

# Master de Psychologie

## PSY73B : Informatique : traitement des données - TD N°4

### Corrélation et régression

## 16. Corrélation linéaire

### 16.1. Coefficient de corrélation

Dans une expérience de perception, on étudie l'évaluation des longueurs de figures géométriques. Le sujet est invité à évaluer les longueurs des figures, en s'aidant d'une figure de référence dont il connaît la longueur (9 cm).

Dans la condition 1, les figures sont 11 bâtonnets. Les données recueillies pour un sujet sont les suivantes :

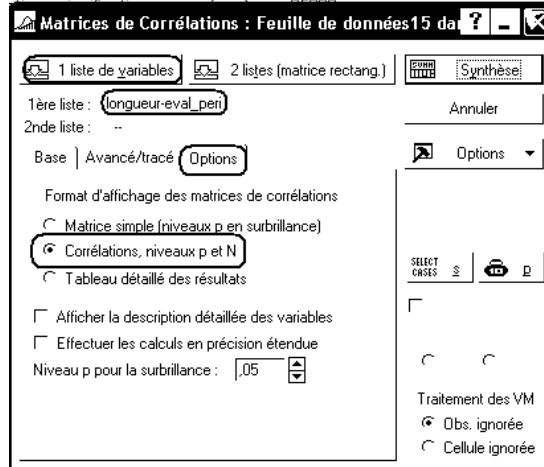
Longueur	2.5	4.6	6.3	7.6	8.5	9.0	9.5	10.4	11.7	13.4	15.5
Eval. long.	2.8	4.4	6.2	7.8	8.2	9.0	9.6	10.6	12.0	13.6	15.2

Dans la condition 2, les figures sont des cercles de périmètres égaux aux longueurs des bâtonnets de la condition 1. L'évaluation du périmètre par un sujet est alors la suivante :

Périmètre	2.5	4.6	6.3	7.6	8.5	9.0	9.5	10.4	11.7	13.4	15.5
Eval. périm.	1.8	3.6	5.8	7.2	8.4	9.0	9.8	11.0	13.2	16.1	21

Ouvrir un nouveau classeur Statistica, y insérer une feuille de données et saisir les données dans quatre variables : Longueur, Eval. Longueur, Périmètre, Eval. Périm. La troisième pourra évidemment être recopiée à partir de la première.

Pour obtenir les coefficients de corrélation entre les différentes variables, on pourra utiliser le menu Statistiques - Statistiques Élémentaires - Matrices de corrélation. Sous l'onglet "Options", on peut par exemple sélectionner "Corrélations, niveaux p et N" :



Sur l'exemple proposé on obtient les résultats suivants :

Corrélations (Feuille de données15 dans Classeur6)				
Corrélations significatives marquées à $p < ,05000$				
N=11 (Observations à VM ignorées)				
Variable	longueur	Eval_long	perimetre	eval_peri
longueur	1,0000	,9982	1,0000	,9867
	p= ---	p=,000	p= ---	p=,000
Eval_long	,9982	1,0000	,9982	,9829
	p=,000	p= ---	p=,000	p=,000
perimetre	1,0000	,9982	1,0000	,9867
	p= ---	p=,000	p= ---	p=,000
eval_peri	,9867	,9829	,9867	1,0000
	p=,000	p=,000	p=,000	p= ---

Les corrélations sont ici très fortes, d'où les niveaux de significativité proches de 0.

Enregistrez le document sous le nom Perception-longueurs.stw. Il sera repris dans le paragraphe sur la régression linéaire.

## 16.2. Alpha de Cronbach

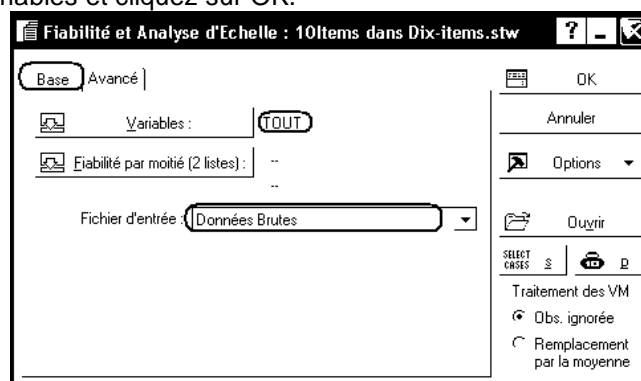
On doit valider un questionnaire pour mesurer les préjugés des individus vis-à-vis des voitures étrangères. On a rédigé un questionnaire posant des questions du type : "les voitures étrangères manquent de personnalité", "Les voitures étrangères se ressemblent toutes", etc... Puis on a soumis ce questionnaire à un groupe de 100 sujets. Les sujets indiquent leur degré d'accord mesuré sur des échelles de 9 points, allant de 1=pas du tout d'accord à 9=tout à fait d'accord.

Ouvrez le classeur Dix-items.stw

On souhaite mesurer la cohérence de cet ensemble de questions à l'aide du coefficient Alpha de Cronbach.

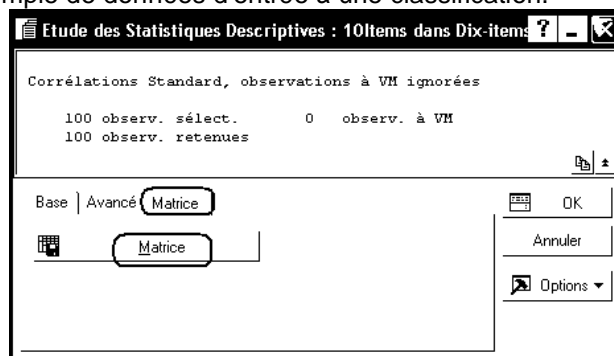
Utilisez le menu Statistiques - Techniques exploratoires multivariées - Fiabilité et analyse d'échelle.

Sélectionnez toutes les variables et cliquez sur OK.



On peut alors afficher les corrélations entre les variables à l'aide du bouton "corrélations". Toutefois, le menu Statistiques - Statistiques Élémentaires - Matrices de corrélation permet également de visualiser quels sont les coefficients de corrélation qui sont significatifs d'un lien entre les variables.

L'onglet "Matrice" permet d'afficher les données dans une feuille de données d'un type particulier, une *matrice*, pour servir par exemple de données d'entrée à une classification.

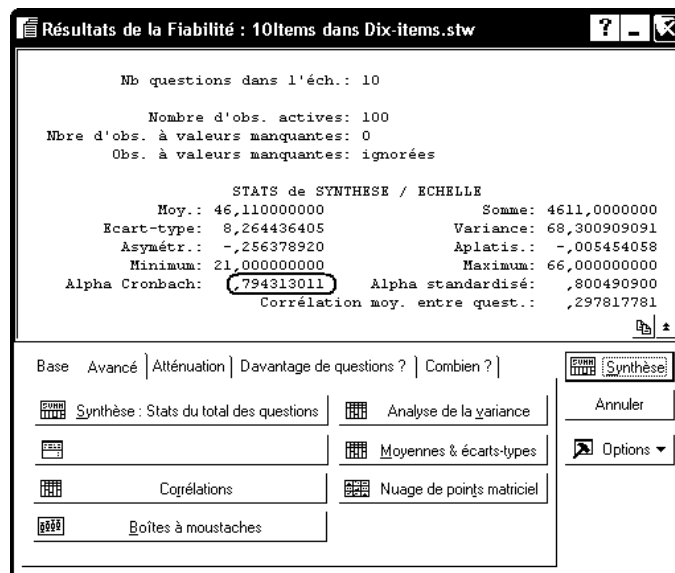


Une telle feuille est caractérisée par la présence d'observations supplémentaires dans le bas du tableau :

	10Items dans Dix-items.stw		
	1 QUEST_1	2 QUEST_2	3 QUEST_3
QUEST_10	0,50041	0,46781	0,29262
Moyennes	4,50000	4,74000	4,70000
Ec-Types	1,44600	1,26027	1,35214
Nb Obs.	100,00000		
Matrice	1,00000		

et s'enregistre dans un format particulier (fichiers d'extension .smx).

Cliquez ensuite sur le bouton OK. On affiche ainsi la fenêtre de dialogue suivante :



La valeur du coefficient Alpha de Cronbach pour l'ensemble des 10 items est 0,79. Le coefficient standardisé est celui que l'on obtiendrait en effectuant une transformation par centrage et réduction sur chaque variable avant de faire la somme.

Le bouton "Synthèse" permet d'avoir des résultats plus détaillés :

Synthèse échelle : Moy.=46,1100 Ec-T.=8,26444 N actif:100 (10Items dans Dix-items.stw)					
Alpha Cronbach :,794313 Alpha Standardisé :,800491					
Corrél. moy. inter-quest.:',297818					
variable	Moy. si supprimé	Var. si supprimé	Ec-T. si supprimé	Corrél. Qst. Tot	Alpha si supprimé
QUEST_1	41,6100	51,9379	7,2068	0,6563	0,7522
QUEST_2	41,3700	53,7931	7,3344	0,6661	0,7547
QUEST_3	41,4100	54,8619	7,4069	0,5492	0,7668
QUEST_4	41,6300	56,5731	7,5215	0,4709	0,7760
QUEST_5	41,5200	64,1696	8,0106	0,0546	0,8249
QUEST_6	41,5600	62,6864	7,9175	0,1186	0,8179
QUEST_7	41,4600	54,0284	7,3504	0,5876	0,7620
QUEST_8	41,3300	53,3211	7,3021	0,6092	0,7590
QUEST_9	41,4400	55,0664	7,4207	0,5025	0,7720
QUEST_10	41,6600	53,7844	7,3338	0,5729	0,7633

On voit, par exemple, que l'on pourrait améliorer le coefficient Alpha en retirant la question 5 ou la question 6.

### 16.3. Corrélation des rangs

Douze étudiants ont complété deux questionnaires conçus pour mesurer (1) l'autoritarisme et (2) la lutte pour le statut social. L'autoritarisme (Adorno et al., 1950) est un concept psychologique ; en résumé, les personnes très autoritaires tendent à être stricts et croient en l'autorité ("droit et ordre").

Les données sont rassemblées dans le fichier Striving.stw.

Utilisez le menu Statistiques - Tests non paramétriques - Corrélations (Spearman, tau de Kendall, Gamma).

Vous obtenez pour le R de Spearman :

Coeffs de Corrélations de Rangs de Spearman		
Cellules à VM ignorées		
Corrélations significatives marquées à p <,05000		
Variable	AUTHORIT	SOCIAL
AUTHORIT	1,0000	0,8182
SOCIAL	0,8182	1,0000

et, pour le tau de Kendall :

	Corrélations du Tau de Kendall Cellules à VM ignorées Corrélations significatives marquées à $p < ,05000$	
Variable	AUTHORIT	SOCIAL
AUTHORIT	1,0000	0,6667
SOCIAL	0,6667	1,0000

Quant à la statistique Gamma, l'aide de Statistica 7 indique :

*Gamma. La statistique Gamma (Siegel & Castellan, 1988) est préférable au R de Spearman ou au Tau de Kendall lorsque les données contiennent de nombreux ex-aequo. En termes d'hypothèses sous-jacentes, Gamma est équivalent au R de Spearman ou au Tau de Kendall ; en termes d'interprétation et de calculs, il est plus proche du Tau de Kendall que du R de Spearman. En résumé, Gamma est également une probabilité ; plus précisément, il se calcule comme la différence entre la probabilité que le rang de deux variables soit identique, moins la probabilité qu'il soit différent, divisé par 1 moins la probabilité d'ex-aequo. C'est pourquoi, Gamma est en fait équivalent au Tau de Kendall, à la différence que les ex-aequo sont ici, explicitement pris en compte.*

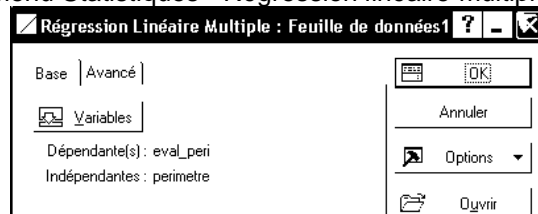
## 17. Régression linéaire à deux ou plusieurs variables

### 17.1. Régression linéaire à deux variables

Ouvrez le classeur Perception-longueurs.stw, créé au paragraphe 15.1. On veut maintenant déterminer la droite de régression de Eval-périmètre par rapport à Périmètre.

#### 17.1.1 Equation de la droite de régression

On peut, pour cela, utiliser le menu Statistiques - Régression linéaire multiple :



On indique Eval\_peri comme variable dépendante, perimetre comme variable indépendante et on clique sur OK.

Le bouton "Synthèse : résultats de la régression" du dialogue suivant permet d'obtenir l'équation de la droite de régression :

Synthèse de la Régression; Variable Dép. : eval_peri R= ,98674804 R?= ,97367170 R? Ajusté = ,97074633 F(1,9)=332,84 p<,00000 Err-Type de l'Estim.: ,94591						
N=11	Bêta	Err-Type de Bêta	B	Err-Type de B	t(9)	niveau p
OrdOrig.			-3,3053	0,7687	-4,2997	0,0020
perimetre	0,9867	0,0541	1,4471	0,0793	18,2438	0,0000

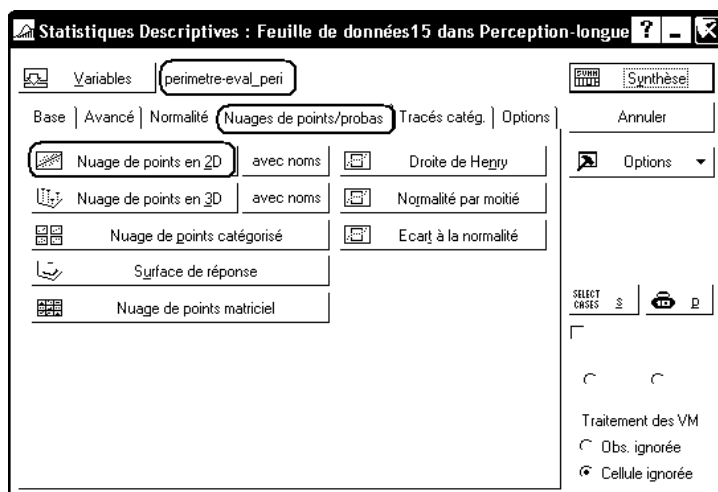
On obtient ainsi comme équation pour la régression :

$$\text{eval\_peri} = -3.3053 + 1.4471 * \text{perimetre}.$$

N.B. Les autres résultats fournis par ce traitement sont décrits infra, dans le paragraphe 16.2.2.

#### 17.1.2 Nuage de points et droite de régression

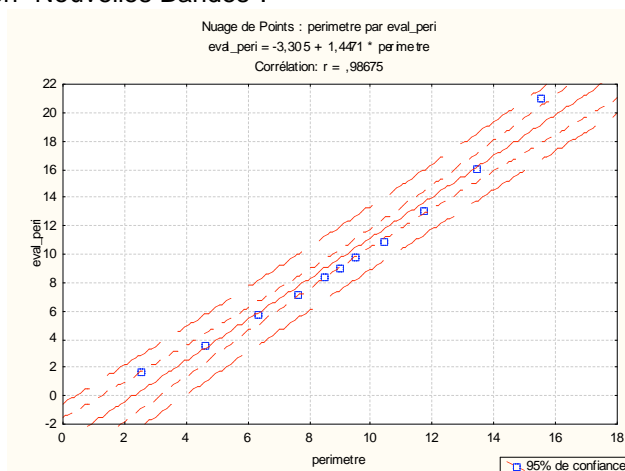
Le plus simple est d'utiliser ici le menu Statistiques - Statistiques Élémentaires - Statistiques Descriptives et l'onglet "Nuages de points/probas" :



Statistica nous affiche le nuage de points, la droite de régression, et les "bandes" donnant l'intervalle de confiance pour la droite de régression, au degré de confiance de 95%. Cet intervalle de confiance correspond aux différentes positions que la droite serait susceptible d'occuper si on recommençait les calculs à partir d'un autre échantillon.

En cliquant sur le graphique à l'aide du bouton droit de la souris, on a accès au menu Propriétés du Graphique (Toutes les Options). L'onglet "Bandes de Régr" permet alors de supprimer les bandes donnant l'intervalle de confiance, ou de leur substituer les représentations graphiques de l'intervalle de détermination, c'est-à-dire la bande du plan qui devrait rassembler 95% des couples (x, y) observés sur la population.

On peut aussi (comme ci-dessous), représenter les deux types de bandes en introduisant un deuxième jeu de bandes à l'aide du bouton "Nouvelles Bandes".



## 17.2. Régression linéaire à plusieurs variables : recherche d'un modèle explicatif

### 17.2.1 Présentation de l'exemple

Exercice adapté à partir de "Les disparités géographiques des dépenses de santé: deux modèles explicatifs pour le secteur libéral", de Roquefeuil, L., Solidarité Santé, N° 4, 1996.

Des variations dans le niveau des dépenses de santé allant du simple au double ont été observées entre les départements. Plusieurs variables peuvent expliquer ce phénomène : la densité des médecins libéraux et la densité de leur clientèle, la morbidité de la population, la proportion de personnes âgées ou l'influence du tiers-payant sur la dépense. Sont étudiées ici :

- l'IDRS ou indicateur des dépenses de remboursement de soins du secteur libéral
- la densité de médecins libéraux dans l'unité géographique concernée
- la mobilité de la clientèle des médecins libéraux : un indicateur de mobilité positif signifie que la valeur des soins "produits" par les médecins de l'unité géographique est supérieure à la valeur des

soins "consommés" par la population de l'unité ; un indicateur négatif au contraire, signifie qu'une partie de la population de l'unité va se faire soigner à l'extérieur de celle-ci.

- la mobilité de la clientèle des médecins spécialistes
- le taux de mortalité, corrigé de la structure par âge de la population totale
- la proportion de personnes âgées de 70 ans et plus
- la part (en %) de dépenses de santé réglées en tiers payant.

Deux niveaux d'unités géographiques sont considérés : les données sont fournies par département et par région.

N.B. Les données figurant dans le fichier sont celles indiquées par l'auteur en annexe de son article, et non des données recréées artificiellement.

### 17.2.2 Etude au niveau départemental

Ouvrez le classeur IDRS.stw et activez la feuille IDRS-Dept.

Affichez les statistiques descriptives relatives aux données présentées. Vous devriez obtenir :

Variable	Statistiques Descriptives (IDRS-Dept dans IDRS.stw)				
	N Actifs	Moyenne	Minimum	Maximum	Ecart-type
IDRS	89	3262,76	2237,40	4959,20	495,39
Densité Médecins	89	178,45	123,50	309,10	37,54
Mobilité omnipraticiens	89	4,49	-10,80	31,20	6,18
Mobilité spécialistes	89	-8,52	-61,60	31,00	19,16
Taux de mortalité	89	9,16	8,10	11,00	0,63
Plus de 70 ans	89	11,01	7,00	17,60	2,30
Tiers payant	89	25,58	11,70	56,80	6,60

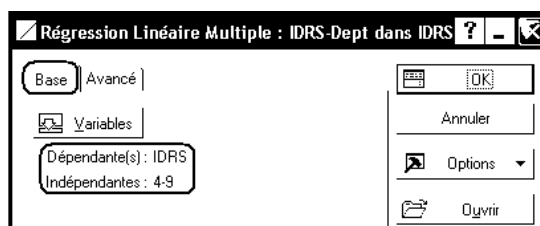
Affichez la matrice des corrélations entre les variables :

Variable	Corrélations (IDRS-Dept dans IDRS.stw) Corrélations significatives marquées à $p < ,05000$ N=89 (Observations à VM ignorées)						
	IDRS	Densité Médecins	Mobilité omnipraticien	Mobilité spécialistes	Taux de mortalité	Plus de 70 ans	Tiers payant
IDRS	1,00	0,67	-0,32	-0,05	-0,13	0,51	0,69
Densité Médecins	0,67	1,00	0,21	0,48	-0,31	0,14	0,36
Mobilité omnipraticiens	-0,32	0,21	1,00	0,39	-0,34	-0,13	-0,28
Mobilité spécialistes	-0,05	0,48	0,39	1,00	-0,11	-0,33	-0,05
Taux de mortalité	-0,13	-0,31	-0,34	-0,11	1,00	-0,38	0,04
Plus de 70 ans	0,51	0,14	-0,13	-0,33	-0,38	1,00	0,15
Tiers payant	0,69	0,36	-0,28	-0,05	0,04	0,15	1,00

Effectuez ensuite une régression linéaire multiple de la variable IDRS sur les autres variables numériques.

Utilisez ensuite le menu Statistiques - Régression Multiple

Sous l'onglet "Base", spécifiez IDRS comme variable dépendante, les 6 autres variables numériques comme variables indépendantes.



Statistica nous affiche alors l'essentiel des résultats de la régression. On peut notamment afficher les résultats de l'ANOVA (bouton ANOVA) montrant qu'ici, le coefficient de régression multiple est

significativement différent de 0, ou encore qu'il existe un lien linéaire significatif entre la variable IDRS et les autres variables :

Effet	Analyse de Variance (IDRS-Dept dans IDRS.stw)				
	Sommes Carrés	dl	Moyennes Carrés	F	niveau p
Régress.	19632968	6	3272161	136,6935	0,0000
Résidus	1962912	82	23938		
Total	21595880				

On peut cliquer sur le bouton OK pour avoir accès à d'autres résultats.

Le bouton "Synthèse de la régression" (onglet "Avancé") affiche les résultats suivants :

Synthèse de la Régression; Variable Dép. : IDRS (IDRS-Dept dans						
R= ,95347109 R²= ,90910713 R² Ajusté = ,90245643						
F(6,82)=136,69 p<0,0000 Err-Type de l'Estim.: 154,72						
N=89	Bêta	Err-Type de Bêta	B	Err-Type de B	t(82)	niveau p
OrdOrig.			-302,094	367,7148	-0,8215	0,4137
Densité Médecins	0,6518	0,0460	8,601	0,6072	14,1659	0,0000
Mobilité omnipraticiens	-0,2355	0,0404	-18,895	3,2444	-5,8239	0,0000
Mobilité spécialistes	-0,1381	0,0450	-3,570	1,1625	-3,0712	0,0029
Taux de mortalité	0,0921	0,0407	72,162	31,9124	2,2613	0,0264
Plus de 70 ans	0,3358	0,0412	72,172	8,8481	8,1567	0,0000
Tiers payant	0,3276	0,0394	24,586	2,9580	8,3117	0,0000

La colonne "B" donne les coefficients de l'équation de régression linéaire. Le modèle fourni par la régression linéaire est le suivant :

$$IDRS = -302 + 8,6 * \text{Dens. Méd} - 18,9 * \text{Mobi Gén} - 3,57 * \text{Mobi Spéc} + 72,2 * \text{Mort.} + 72,2 * \text{Part Agées} + 24,6 * \text{Tiers-P}$$

La valeur de R<sup>2</sup> est de 0,91 : 91% de la variance de la variable IDRS est expliquée par le modèle.

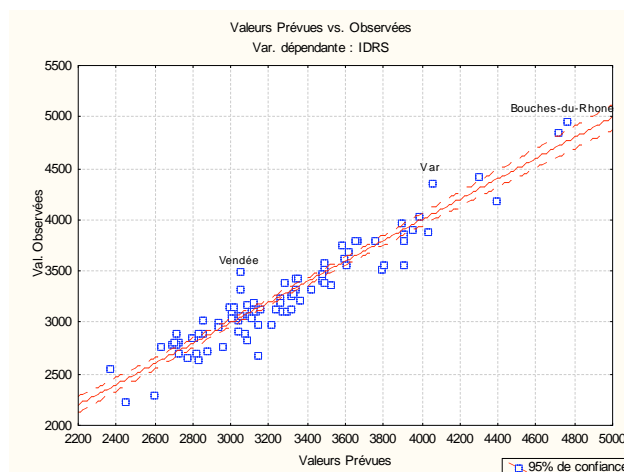
Les coefficients de la colonne "Bêta" sont les coefficients standardisés, c'est-à-dire les coefficients que l'on observerait si on utilisait des variables centrées réduites au lieu des variables observées. On peut également les interpréter comme suit : lorsque "Densité Médecins" augmente d'un écart type, la variable "IDRS" estimée augmente de 0,65 écart type, lorsque la variable "Mobilité omnipraticiens" augmente d'un écart type, "IDRS" diminue de 0,23 écart type.

Par exemple, on pourra vérifier que

$$Beta(Densité) = \frac{Ecart\ type(Densité)}{Ecart\ type(IDRS)} \times B(Densité) = \frac{37,54}{495,39} \times 8,601 = 0,6518$$

Les valeurs de t sont obtenues en divisant la valeur correspondante de B par son erreur type. Autrement dit, on teste si le coefficient B est significativement différent de 0.

Sous l'onglet "Nuage", on pourra obtenir différentes représentations graphiques dont, par exemple, le graphique illustrant l'adéquation entre les valeurs observées et les valeurs théoriques :



Remarque : dans l'article cité supra on étudiait ensuite, pour chaque département, la part de la dépense expliquée par chaque variable. L'étude concluait d'ailleurs sur une variabilité importante selon le département. Ainsi, dans l'Aisne, la densité de médecins explique 37% de la dépense, alors qu'elle explique 64% de la dépense dans les Hautes-Alpes. Par exemple, pour les premiers départements, on obtient :

	Densité Médecins	Mobilité omnipraticiens	Mobilité spécialistes	Taux de mortalité	Plus de 70 ans	Tiers payant	Résidu
Ain	46,1%	-5,0%	6,7%	27,9%	27,0%	23,1%	-12,7%
Aisne	36,8%	-1,0%	2,2%	23,6%	19,8%	28,2%	-0,3%
Allier	44,8%	-2,0%	-1,2%	20,0%	28,8%	19,7%	-1,2%
Alpes Hte Prov	48,0%	-5,2%	2,1%	16,5%	23,3%	19,3%	4,0%
Hautes-Alpes	63,7%	-13,7%	2,5%	20,8%	27,3%	16,7%	-7,1%
Alpes Maritimes	59,9%	-7,2%	-1,1%	13,7%	24,2%	14,0%	3,3%
Ardèche	41,7%	-3,1%	4,8%	21,2%	29,1%	16,5%	-0,3%

Ce genre de tableau ne paraît pas simple à produire à l'aide de Statistica, et il est sans doute préférable de recopier les données et certains résultats dans Excel pour réaliser ce travail.

### 17.3. Coefficients de corrélation partielle

Ouvrir le fichier Eval-Cours.stw. Il rassemble les données figurant dans un exercice des fiches de TD de Statistiques. L'énoncé accompagnant ces données est rappelé dans le rapport contenu dans le classeur.

- Calculer les paramètres descriptifs des six variables.
  - Calculer les coefficients de corrélation des variables prises deux à deux.
- Quels sont les couples de variables pour lesquels la corrélation est significative, au seuil de 5% ?
- On veut estimer la variable Qual-Glob en utilisant comme prédicteurs les 5 autres variables.
- Déterminer l'équation de régression et le coefficient de corrélation. Vous devriez obtenir :

$$\text{Qual-Glob} = - 1.19 + 0.763 \text{ Péda} + 0.132 \text{ Exam} + 0.489 \text{ Connai} - 0.184 \text{ Rés} + 0.000525 \text{ Inscr}$$

Les coefficients de corrélation partielle entre la variable Qual-Glob et chacun des prédicteurs, lorsque les autres variables sont contrôlées, peuvent être obtenus à partir du bouton "Corrélations partielles" du dialogue "Résultats". On obtient les résultats suivants :

Variable	Variables dans l'équation ; VD: Qual-Glob (Données dans Eval-Cours.stw)						
	Bêta ds	Corrél. Partiel.	Corrél. Semipart	Tolérance	R?	t(44)	Niveau p
Pédagogie	0,6620	0,6545	0,4281	0,4182	0,5818	5,7421	0,0000
Examen	0,1061	0,1213	0,0604	0,3246	0,6754	0,8107	0,4219
Connaissance	0,3251	0,4751	0,2670	0,6746	0,3254	3,5813	0,0008
Résultat	-0,1055	-0,1656	-0,0830	0,6197	0,3803	-1,1137	0,2715
Inscription	0,1242	0,1990	0,1004	0,6534	0,3466	1,3472	0,1848



### 17.3.1 Calcul "à la main" des coefficients de corrélation partielle

On veut calculer le coefficient de corrélation partielle entre la variable Qual-Glob et la variable Pédagogie. Nous allons procéder en trois étapes :

- Déterminez les résidus de la régression de la variable Qual-Glob par rapport aux 4 autres variables (Examen, Connaissance, Résultat, Inscription).
- Déterminez de même les résidus de la régression de la variable Pédagogie par rapport aux 4 autres variables.
- Créez une nouvelle feuille de données et collez dans les deux premières colonnes de cette feuille les colonnes "Résidus" des feuilles de résultats précédentes.
- Supprimez les 4 dernières observations qui viennent d'être collées (il s'agit de paramètres descriptifs des résidus, sans intérêt ici).
- Calculez enfin le coefficient de corrélation entre les deux variables ainsi définies. Vous devriez retrouver le résultat, à savoir :  $r=0,65$ .

### 17.3.2 Corrélations partielles et neutralisation de l'effet d'un facteur

Ouvrez le fichier Coping.stw du répertoire Corrélations-Partielles.

On a relevé les valeurs de deux variables numériques, RSS et DI, sur 12 sujets, 6 hommes et 6 femmes :

	Sujet	Sexe	RSS	DI	RSS centrée par sexe	DI centrée par sexe
1	s1	F	5	0	-4	-1,33
2	s2	F	8	0	-1	-1,33
3	s3	F	9	2	1	-0,33
4	s4	F	10	1	0	0,67
5	s5	F	10	3	3	0,67
6	s6	F	12	2	1	1,67
7	s7	H	2	0	-3,33	-0,17
8	s8	H	4	0	-1,33	-0,17
9	s9	H	6	0	0,67	-0,17
10	s10	H	6	0	0,67	-0,17
11	s11	H	6	0	0,67	-0,17
12	s12	H	8	1	2,67	0,83

On constate un effet important du facteur sexe : pour les deux variables, les scores des hommes et ceux des femmes sont notablement différents.

Calculer le coefficient de corrélation des variables RSS et DI. On obtient :  $r = 0,77$ . Ce coefficient est difficile à interpréter, car il est dû à la fois au lien éventuel entre RSS et DI et à l'effet du facteur Sexe.

*Comment neutraliser l'effet du sexe dans le calcul de l'intensité du lien entre RSS et DI ?*

1) Faites une régression multiple de DI sur les variables Sexe et RSS, afin d'évaluer les coefficients de corrélation partielle. Vous devriez obtenir :

Variable	Variables dans l'équation ; VD: DI (Donnees dans Classeur1)						
	Bêta ds	Corrél. Partiel.	Corrél. Semipart	Tolérance	R?	t(9)	Niveau p
Sexe	0,111458	0,129428	0,082673	0,550186	0,449814	0,391578	0,704480
RSS	0,694655	0,631059	0,515257	0,550186	0,449814	2,440493	0,037334

Ainsi, après neutralisation de l'effet du sexe, la corrélation entre RSS et DI n'est que de  $r'=0,63$ . Elle reste cependant significative.

2) De manière équivalente, on peut remplacer chaque valeur observée de RSS et DI par son écart algébrique à la moyenne par sexe correspondante. C'est ce qui a été fait dans les colonnes "RSS centrée par sexe" et "DI centrée par sexe".

Par exemple, RSS vaut 5 sur le sujet féminin s1, et vaut 9 pour les femmes. La valeur de "RSS centrée par sexe" sur le sujet s1 est donc :  $5 - 9 = -4$ .

Le coefficient de corrélation des deux variables centrées par sexe est alors exactement le coefficient de corrélation partielle précédent :

Corrélations (Données dans Classeur1) Corrélations significatives marquées à $p < ,05000$ N=12 (Observations à VM ignorées)		
Variable	RSS centrée par sexe	DI centrée par sexe
RSS centrée par sexe	1,000000	0,631050
DI centrée par sexe	0,631050	1,000000

## 18. Régression linéaire pas à pas

### 18.1. Principe de la méthode

Les données sont formées par une VD Y et plusieurs variables explicatives  $X_1, X_2, \dots, X_p$ .

On choisit, parmi les variables explicatives, celle qui est le mieux corrélée à Y. Pour simplifier les notations, nous supposons qu'il s'agit de la variable  $X_1$ .

On calcule l'équation de régression linéaire de Y sur  $X_1$  :  $Y = b_1 X_1 + b_0$ .

On calcule alors les résidus :  $R_1 = Y - b_1 X_1 - b_0$

On choisit, parmi les variables explicatives restantes, celle qui est le mieux corrélée à  $R_1$ . Nous supposons ici qu'il s'agit de la variable  $X_2$ .

On calcule l'équation de régression linéaire de Y sur  $X_1$  et  $X_2$  :  $Y = b'_1 X_1 + b_2 X_2 + b'_0$ .

On calcule les nouveaux résidus :  $R_2 = Y - (b'_1 X_1 + b_2 X_2 + b'_0)$  et on poursuit la méthode jusqu'à ce que les variables explicatives restantes ne soient plus significativement corrélées aux résidus.

### 18.2. Présentation de l'exemple

Exercice adapté à partir de "Intelligence pratique ou traditionnelle : Que mesure l'entrevue structurée situationnelle ?", Durivage A., St-Martin J., Barette J., Revue européenne de Psychologie Appliquée, 1995, vol. 45 n° 3, pp. 171-178.

L'objectif de l'étude consiste à explorer le construit sous-jacent à l'entrevue structurée situationnelle lors de la sélection du personnel. Constitue-t-elle une mesure de l'intelligence traditionnelle (QI) ou de connaissances tacites associées théoriquement à de l'intelligence pratique.

Méthodologie : l'entrevue structurée situationnelle et les tests ont été administrés à 48 candidats potentiels à un poste de responsable des bénévoles dans un centre hospitalier psychiatrique du Québec. Les variables suivantes ont été recueillies :

- Score à l'entrevue structurée : échelle de 0 à 40
- Score au BGTA : batterie générale de tests d'aptitude mesurant l'intelligence traditionnelle. Ce score est ici donné sous la forme d'une variable centrée et réduite
- Scores sur les dimensions "Organisation", "Impulsivité", "Compréhension", "Altruisme" des tests de personnalité de Jackson (échelles de 0 à 20)
- Age (de 20 à 41 ans dans l'expérience originale)
- Le nombre d'années d'expérience de travail à temps plein (de 0 à 21 ans dans l'expérience originale).

### 18.3. Régression linéaire pas à pas de Entrevue sur les autres variables

Chargez le classeur Entrevue-structuree.stw, qui contient des données générées artificiellement, conformes aux résultats indiqués par les auteurs.

Affichez les statistiques descriptives relatives aux données présentées. Vous devriez obtenir :

Variable	Statistiques Descriptives (Entrevue dans Entrevue-structuree.stw)				
	N Actifs	Moyenne	Minimum	Maximum	Ecart-type
Entrevue	48	30,09	20,77	38,58	4,43
Comprehension	48	12,18	5,09	19,02	3,31
Organisation	48	13,06	6,05	20,70	3,69
Altruisme	48	12,00	2,49	19,58	4,25
Impulsivite	48	9,51	2,66	16,48	3,41
BGTA	48	-0,00	-2,08	2,31	1,01
Age	48	25,00	16,11	38,06	4,75
Anciennete	48	8,60	0,00	18,62	4,02

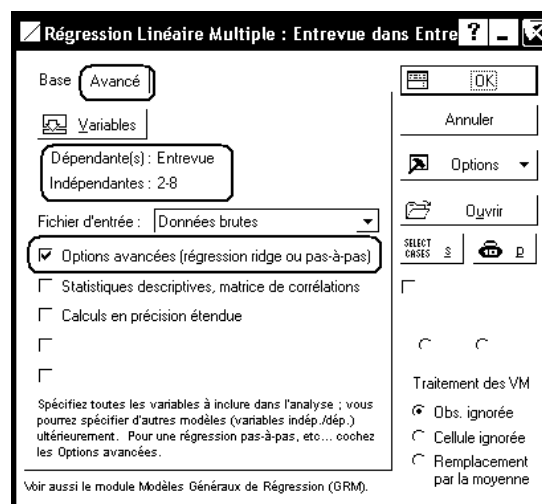
Affichez la matrice des corrélations entre les différentes variables. Quelles sont les corrélations qui apparaissent significatives ?

Variable	Corrélations (Entrevue dans Entrevue-structuree.stw)							
	Entrevue	Comprehension	Organisation	Altruisme	Impulsivite	BGTA	Age	Anciennete
Entrevue	1,00	0,48	-0,14	0,02	0,03	-0,06	0,40	0,33
Comprehension	0,48	1,00	-0,18	0,00	0,12	0,07	0,01	0,04
Organisation	-0,14	-0,18	1,00	0,18	-0,51	-0,12	0,02	0,06
Altruisme	0,02	0,00	0,18	1,00	-0,19	-0,22	-0,05	0,00
Impulsivite	0,03	0,12	-0,51	-0,19	1,00	-0,22	-0,17	-0,28
BGTA	-0,06	0,07	-0,12	-0,22	-0,22	1,00	0,03	0,01
Age	0,40	0,01	0,02	-0,05	-0,17	0,03	1,00	0,83
Anciennete	0,33	0,04	0,06	0,00	-0,28	0,01	0,83	1,00

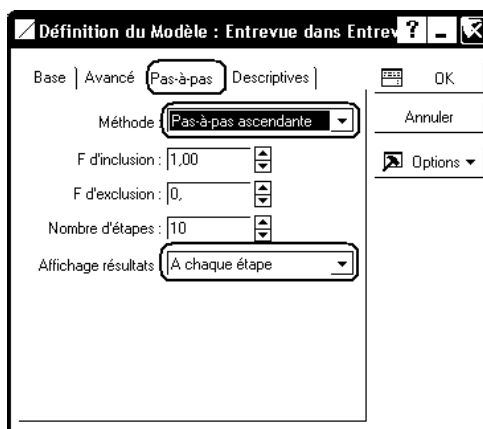
### 18.3.1 Exécution de la procédure

Utilisez ensuite le menu Statistiques - Régression Multiple

Sous l'onglet "Avancé", spécifiez Entrevue comme variable dépendante, les 7 autres variables comme variables indépendantes. Cochez l'option "régression ridge ou pas-à-pas".



Dans le dialogue suivant, activez l'onglet "pas-à-pas" et sélectionnez la méthode "pas à pas ascendante", et l'affichage des résultats à chaque étape :



A la première étape, Statistica affiche les résultats suivants :

Résultats Régress. Multiple (Etape 0)			
Var dép. : Entrevue	R Multiple = 0,0000000	F = 0,000000	
	R <sup>2</sup> = 0,0000000	dl = 0,47	
Nb d'obs. : 48	R <sup>2</sup> ajusté = 0,0000000	p = -0,00000	
	Erreur-type de l'estim. : 4,426350465		
Etape 0 : Aucune variable dans l'équation			
(bêta significatifs en surbrillance)			

Cliquez sur "suivant". On obtient :

Résultats Régress. Multiple (Etape 1)			
Var dép. : Entrevue	R Multiple = ,48000000	F = 13,77131	
	R <sup>2</sup> = ,23040000	dl = 1,46	
Nb d'obs. : 48	R <sup>2</sup> ajusté = ,21366957	p = ,000556	
	Erreur-type de l'estim. : 3,925078427		
Ord.Orig : 22,282917179	Err.-Type: 2,178731	t( 46) = 10,227	p = ,0000
Comprehension bêta=,480			
(bêta significatifs en surbrillance)			

puis :

Résultats Régress. Multiple (étape 2 , sol. finale)			
pas d'autre F d'inclusion au seuil spécifié			
Var dép. : Entrevue	R Multiple = ,62177055	F = 14,18071	
	R <sup>2</sup> = ,38659861	dl = 2,45	
Nb d'obs. : 48	R <sup>2</sup> ajusté = ,35933633	p = ,000017	
	Erreur-type de l'estim. : 3,542915919		
Ord.Orig : 13,138962343	Err.-Type: 3,341282	t( 45) = 3,9323	p = ,0003
Comprehension bêta=,476		Age bêta=,395	
(bêta significatifs en surbrillance)			

Il ne reste plus de variable significativement corrélée aux résidus, et Statistica substitue le bouton "OK" au bouton "Suivant". Cliquez sur ce bouton.

### 18.3.2 Analyse des résultats

Sous l'onglet "Avancé", le bouton "Synthèse de la régression" permet d'obtenir les résultats suivants :

Synthèse de la Régression; Variable Dép. : Entrevue (Entrevue d						
R= ,62177055 R²= ,38659861 R² Ajusté = ,35933633						
F(2,45)=14,181 p<,00002 Err-Type de l'Estim.: 3,5429						
N=48	Bêta	Err-Type de Bêta	B	Err-Type de B	t(45)	niveau p
OrdOrig.			13,14	3,34	3,93	0,0003
Comprehension	0,48	0,12	0,64	0,16	4,08	0,0002
Age	0,40	0,12	0,37	0,11	3,39	0,0015

Ainsi, l'équation de régression obtenue par ce modèle est :

$$\text{Entrevue} = 0,64 * \text{Comprehension} + 0,37 * \text{Age} + 13,14$$

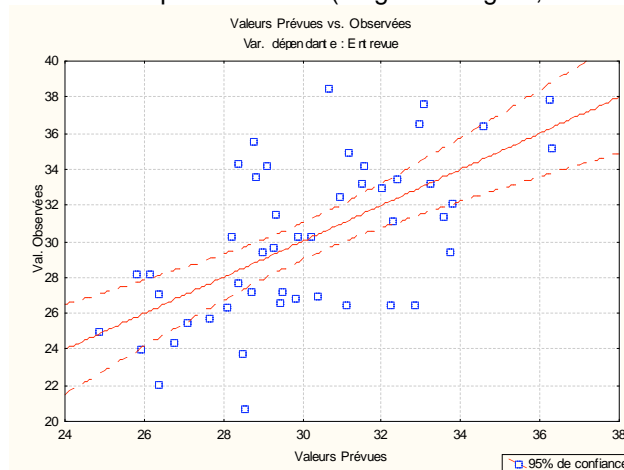
Ce modèle explique 38% de la variance de la variable Entrevue. Les coefficients de la colonne "Bêta" sont les coefficients standardisés, c'est-à-dire les coefficients que l'on observerait si on utilisait des variables centrées réduites au lieu des variables observées. On peut également les interpréter comme suit : lorsque "Comprehension" augmente d'un écart type, la variable "Entrevue" estimée augmente de 0,48 écart type, lorsque la variable "Age" augmente d'un écart type, "Entrevue" augmente de 0,4 écart type.

Les valeurs de t sont obtenues en divisant la valeur correspondante de B par son erreur type. Autrement dit, on teste si le coefficient B est significativement différent de 0.

On peut également obtenir le tableau des valeurs observées et des valeurs estimées de la variable Entrevue (Onglet "Avancé", bouton Synthèse : Résidus et prévisions)

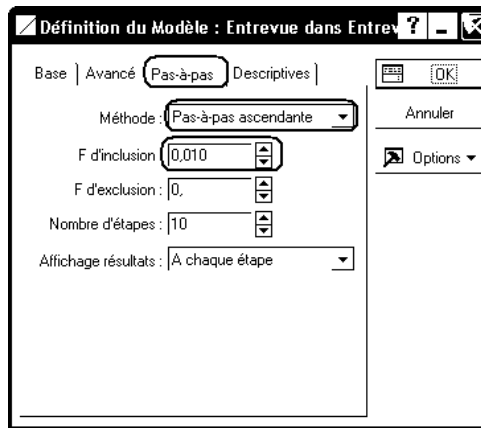
N° d'Obs.	Valeurs Prévues & Résidus (Entre Var. dépendante : Entrevue		
	Valeur Observée	Valeur Prévue	Résidus
s1	32,47	30,92	1,54
s2	28,22	26,14	2,09
s3	35,16	36,28	-1,12

Diverses représentations graphiques peuvent être obtenues. Un nuage de points tri-dimensionnel Entrevue - Comprehension - Age serait peu lisible. En revanche, on pourra construire un graphique comparant les valeurs observées aux valeurs estimées par le modèle (Onglet "Nuages", bouton "Prévues v/s observées") :



### 18.3.3 Variantes

On peut souhaiter recueillir également des informations concernant les variables qui ont été écartées du modèle. Reprenez par exemple l'étude, en indiquant cette fois 0,01 comme valeur limite de F pour inclure une variable :



La régression pas à pas est alors faite sur toutes les variables, avec les résultats suivants. On notera que l'ajout des 5 variables restantes ne permet pas vraiment d'augmenter la part de variance expliquée (40% au lieu de 38%). On notera que, lorsqu'est introduite la variable "Ancienneté", fortement corrélée à l'âge, ni "Ancienneté" ni "Age" ne restent significatifs.

Synthèse de la Régression; Variable Dép. : Entrevue (Entrevue dar R= ,63694958 R? = ,40570477 R? Ajusté = ,30170311 F(7,40)=3,9009 p<,00250 Err-Type de l'Estim.: 3,6988						
N=48	Bêta	Err-Type de Bêta	B	Err-Type de B	t(40)	niveau p
OrdOrig.			14,3833	6,0676	2,3705	0,0227
Comprehension	0,4728	0,1249	0,6314	0,1668	3,7852	0,0005
Age	0,4415	0,2203	0,4115	0,2053	2,0038	0,0519
BGTA	-0,1208	0,1360	-0,5293	0,5957	-0,8886	0,3795
Organisation	-0,1001	0,1496	-0,1201	0,1795	-0,6692	0,5072
Altruisme	0,0259	0,1298	0,0269	0,1351	0,1992	0,8431
Impulsivite	-0,0402	0,1629	-0,0523	0,2117	-0,2468	0,8064
Anciennete	-0,0557	0,2279	-0,0614	0,2509	-0,2446	0,8080

#### 18.4. Exercice à rendre par mail

Source des données : Source : [http://www.amstat.org/publications/jse/jse\\_data\\_archive.html](http://www.amstat.org/publications/jse/jse_data_archive.html)

Dans les débats autour des réformes du système éducatif, il a été fréquemment affirmé que la dépense d'éducation par élève n'était pas un gage de réussite. A l'appui de cette affirmation, certains ont souligné que parmi les dix états américains dont la dépense moyenne par élève était la plus basse en 1994/95, quatre se trouvaient parmi les dix états qui avaient les meilleurs résultats au SAT.

Pour étudier cette question, des statisticiens américains ont extrait du *Digest of Education Statistics* les données suivantes :

Nom de l'état

Dépense par élève moyenne dans les écoles publiques élémentaires et secondaires, en milliers de dollars (1994-95)

Taux d'encadrement : ratio élève par enseignant

Salaires annuels moyens estimés des enseignants des écoles publiques élémentaires et secondaires

Pourcentage d'inscrits au SAT parmi les élèves satisfaisant les conditions d'inscription

Score moyen observé au SAT verbal

Score moyen observé au SAT mathématique

Score moyen observé global au SAT global.

Les données correspondantes se trouvent dans le classeur Statistica SATdata.stw du serveur des salles de TD.

1) Etudier la corrélation entre la dépense moyenne par élève et le score moyen observé au SAT. La corrélation est-elle significative ? Interpréter le signe du coefficient de corrélation.

En effectuant une régression linéaire du score sur la dépense moyenne, retrouver le résultat indiqué dans la source :

*"every \$1,000 increase in spending per student per year is associated with a decline of nearly 21 points in the average statewide SAT score, an estimate that easily reaches conventional levels of statistical significance ( $p < .01$ )."*

Représenter le nuage de points et la droite de régression.

2) Etudier la corrélation entre les variables "Pourcentage Inscrits au SAT" et "Score global au SAT". Représenter le nuage de points correspondant. Que peut-on en conclure ?

3) Etudier la régression linéaire multiple de la variable "Score global au SAT" sur les deux variables "Pourcentage Inscrits au SAT" et "Dépense par élève". Analyser les résultats obtenus. Compléter cette étude par le calcul et l'interprétation des coefficients de corrélation partiels.

Retrouver ainsi la conclusion :

*"With a robust  $R^2$  and slope coefficients that are both highly statistically significant ( $p < .01$ ), it is now clear that while the bulk of variation in statewide SAT scores is attributable to the percentage of students taking the exam, increased spending on public education is in fact associated with better academic performance."*