

# Projet : Statistiques et ACP

## Consignes :

- **Résoudre** les différents exercices à l'aide du **logiciel R sous R-studio**
- Écrire les **scripts** pour chaque traitement
- Présenter les résultats dans un **pdf** avec tous les tracés et tableaux de calcul
- **Commenter** et **interpréter** les résultats lors de la soutenance
- **Barème** : exos 1,2,3 : 10 pts, exo 4 : 10pts

## Exercice 1

On donne la série unidimensionnelle suivante, correspondant à la répartition des entreprises du secteur automobile en fonction de leur chiffre d'affaire en millions d'euros.

Chiffres d'affaires	[0, 0.25[	[0.25, 0.5[	[0.5, 1[	[1, 2.5[	[2.5, 5[	[5, 10[
Nombres d'entreprises	137	106	112	154	100	33

1. Calculer le chiffre d'affaire moyen et l'écart-type de la série.
2. Construire l'histogramme des fréquences.
3. Construire les deux polygones (croissant et décroissant) des fréquences cumulées
4. Calculer la médiane et la proportion d'entreprises dont le chiffre d'affaire est supérieur à 3 millions d'euros.

## Exercice 2

Récupérer le fichier data3.txt et l'ouvrir. Ce fichier représente 11 individus (numérotés de 1 à 11) sur lesquels sont mesurés 8 variables  $x_1, x_2, x_3, x_4, y_1, y_2, y_3, y_4$ . On notera que les variables  $x_1, x_2, x_3$  sont identiques.

1. Récupérer les données à l'aide d'un tableau excel afin de mieux les visualiser.
2. Sous R-studio, à l'aide de la commande `read.table`, récupérer le fichier de données.
3. Calculer la moyenne et la variance des 8 variables  $x_1, \dots, x_4$  et  $y_1, \dots, y_4$
4. Calculer les covariances et les coefficients de corrélation des couples  $(x_1, y_1), \dots, (x_4, y_4)$ . Que constate-t-on ?
5. Tracer la représentation des couples  $(x_1, y_1), \dots, (x_4, y_4)$ . Commenter.
6. Centrer et réduire  $x_1$  et  $y_1$ . Reprendre les questions 3, 4 et 5 pour  $(x_1, y_1)$

## Exercice 3

L'indice moyen d'un salaire a évolué de la façon suivante :

année	1	2	3	4	5	6	7
indice	165	176	193	202	222	245	253

1. Représenter cette série statistique par un nuage de points.
2. En utilisant la méthode des moindres carrées, calculer l'équation de la droite représentant l'indice en fonction de l'année.
3. Comment pourrait-on prévoir l'indice à l'année 9?

## Exercice 4 : Salles de cinéma

On mesure pour 94 départements de France métropolitaine les dix variables décrites ci-dessous. De plus, il est fourni la matrice de corrélation des variables ainsi qu'un extrait du tableau de données concernant les 20 départements les plus peuplés.

**popu** population du département (en millions d'habitants)  
**entr** nombres d'entrées réalisées (en millions)  
**rece** recettes (en millions d'euros)  
**sean** nombre de séances (en milliers)  
**comm** nombre de communes équipées de salles de cinéma  
**etab** nombre de cinémas en activité  
**salle** nombre de salles en activité  
**faut** nombre de fauteuils disponibles  
**artes** nombre de salles d'art et essai (passant des films indépendants)  
**multi** nombre de multiplexes (au moins 8 salles)

### Corrélations

	popu	entr	rece	sean	comm	etab	salle	faut	artes	multi
popu	1.00	0.71	0.68	0.77	0.61	0.77	0.85	0.87	0.72	0.80
entr	0.71	1.00	1.00	0.99	0.19	0.76	0.93	0.91	0.70	0.63
rece	0.68	1.00	1.00	0.98	0.15	0.74	0.91	0.89	0.67	0.59
sean	0.77	0.99	0.98	1.00	0.27	0.80	0.96	0.94	0.72	0.70
comm	0.61	0.19	0.15	0.27	1.00	0.75	0.49	0.53	0.64	0.52
etab	0.77	0.76	0.74	0.80	0.75	1.00	0.91	0.91	0.85	0.67
salle	0.85	0.93	0.91	0.96	0.49	0.91	1.00	0.99	0.79	0.79
faut	0.87	0.91	0.89	0.94	0.53	0.91	0.99	1.00	0.80	0.81
artes	0.72	0.70	0.67	0.72	0.64	0.85	0.79	0.80	1.00	0.55
multi	0.80	0.63	0.59	0.70	0.52	0.67	0.79	0.81	0.55	1.00

	depart	popu	entr	rece	sean	comm	etab	salle	faut	artes	multi
D59	Nord	2.555	6.868	37.459	174	35	48	151	34230	18	5
D75	Paris	2.125	30.439	192.244	698	1	92	368	72752	38	5
D13	Bouches du Rhone	1.836	6.651	39.197	193	28	49	155	27488	20	2
D69	Rhone	1.579	6.992	37.359	193	33	52	141	27023	30	3
D62	Pas de Calais	1.442	2.976	15.903	123	23	28	111	22053	10	5
D92	Hauts de Seine	1.429	3.978	21.701	107	33	39	89	20762	24	2
D93	Seine St Denis	1.383	4.803	25.543	127	26	32	97	21168	16	4
D78	Yvelines	1.354	4.625	26.700	130	29	35	95	19302	14	2
D33	Gironde	1.287	5.057	24.555	171	43	52	154	32714	22	5
D76	Seine Maritime	1.239	3.366	19.172	128	23	30	108	23436	11	3
D94	Val de Marne	1.227	4.052	23.967	122	30	36	89	23448	17	2
D77	Seine et Marne	1.194	3.298	19.964	87	30	33	89	19029	10	2
D44	Loire Atlantique	1.134	4.383	22.283	136	39	52	118	25570	26	4
D91	Essonne	1.134	2.242	11.185	84	30	33	77	14414	16	1
D95	Val d'Oise	1.105	2.160	11.581	66	21	22	56	12939	7	1
D38	Isere	1.094	3.993	22.619	142	36	49	133	25116	18	4
D31	Haute Garonne	1.046	4.674	25.114	106	25	32	82	16797	18	2
D67	Bas Rhin	1.026	3.569	20.338	124	13	18	78	16615	8	3
D57	Moselle	1.023	3.187	18.251	88	19	22	79	17866	9	3
D6	Alpes Maritimes	1.011	3.520	21.731	111	23	42	94	16764	8	1

## Partie 1 : première approche

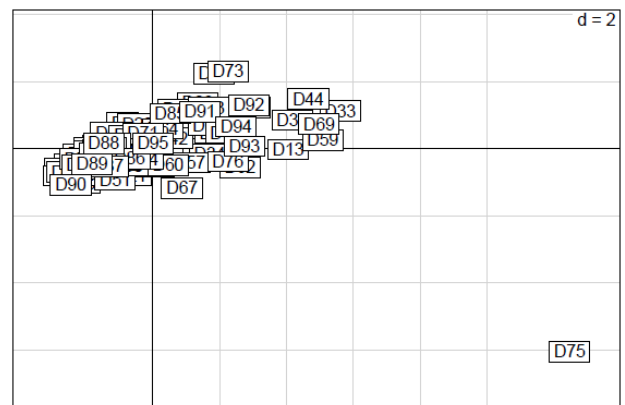
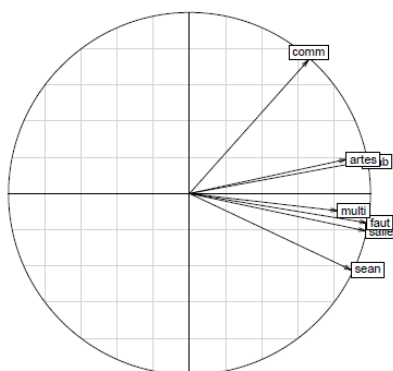
1. A l'aide du tableau des données exo4.txt, calculer la matrice de corrélation des variables et vérifier vos résultats
2. Que pouvez-vous dire à propos des corrélations entre les variables ? Commentez leurs valeurs.
3. En comparant les données brutes de Paris (D75) avec les autres fournies ici, que peut-on dire de ses particularités ?

## Partie 2 : ACP version1

On s'intéresse pour l'analyse aux variables concernant l'offre de cinéma : sean, comm, etab, salle, faut, artes et multi. On effectue une ACP sur les données centrées-réduites, et on obtient les valeurs propres suivantes

5.47	0.87	0.45	0.17	0.02	0.01	0.01
------	------	------	------	------	------	------

On donne le cercle des corrélations des variables ainsi que le plan de composantes principales avec les projections des individus :



4. Retrouver les deux tracés du dessus avec le logiciel R
5. Donner une interprétation rapide de la première composante principale à partir uniquement du cercle des corrélations. Que se passe-t-il sur la seconde composante principale ?
6. En remarquant la projection des individus sur le premier plan principal, que peut-on observer ? En s'appuyant sur la question 2, expliquer pourquoi Paris est particulier à la fois sur le premier et le second axe.
7. On se propose de diviser chaque donnée par la population du département popu (c'est-à-dire que la variable entr sera exprimée par habitant, rece en euros par habitant, sean pour 1000 habitants, etc.). Expliquer en quoi cette approche est intéressante.

### Partie 3 : ACP version2

On normalise les données comme indiqué à la question précédente. On effectue l'ACP sur les données centrées-réduites sur les nouvelles variables normalisées, mais en utilisant les parts de population comme poids des individus. On obtient les données suivantes : valeurs propres, corrélations avec les quatre premiers axes et, pour une sélection (arbitraire) de 20 départements parmi les 94, les poids des individus, leurs coordonnées sur les 4 premiers axes, ainsi que la qualité de leur représentation par les 4 premiers sous espaces.

Valeurs propres								Weight	Axis1	Axis2	Axis3	Axis4	Axis1	Axis1:2	Axis1:3	Axis1:4		
[1] 3.71 2.04 0.75 0.37 0.07 0.04 0.02								D4 0.0024	D4	-6.19	3.93	-0.74	0.37	D4	68.6	96.2	97.2	97.4
								D5 0.0021	D5	-13.55	4.11	0.43	-3.56	D5	85.5	93.4	93.5	99.4
								D9 0.0024	D9	-2.47	3.95	-0.36	1.00	D9	26.1	93.2	93.7	98.1
								D15 0.0026	D15	-1.55	3.04	-0.44	0.30	D15	19.9	95.9	97.5	98.3
								D23 0.0021	D23	-1.83	3.24	-0.64	1.01	D23	21.9	90.2	92.9	99.4
Corrélations								D28 0.0070	D28	2.56	1.00	-0.50	-0.78	D28	77.8	89.5	92.5	99.8
								D32 0.0030	D32	-5.65	5.55	-0.69	2.54	D32	45.7	89.7	90.3	99.6
								D38 0.0188	D38	-1.54	-1.16	0.87	0.25	D38	51.7	81.0	97.3	98.6
								D40 0.0056	D40	-7.00	1.56	0.90	0.91	D40	91.4	95.9	97.4	99.0
								D46 0.0027	D46	-3.44	4.45	-0.69	2.02	D46	32.6	87.2	88.5	99.8
sean -0.50 -0.75 -0.41 0.04 comm -0.68 0.64 0.30 -0.12 etab -0.91 0.34 0.08 -0.17 salle -0.93 -0.30 -0.09 -0.09 faut -0.92 -0.33 -0.01 -0.08 artès -0.64 0.54 -0.21 0.50 multi -0.22 -0.68 0.66 0.24								D48 0.0013	D48	-3.01	3.49	0.42	-1.10	D48	39.6	92.6	93.4	98.7
								D53 0.0049	D53	-1.55	0.62	1.22	1.59	D53	33.7	39.1	59.9	95.4
								D67 0.0176	D67	1.28	-1.53	0.31	0.18	D67	39.2	95.7	98.1	98.9
								D73 0.0064	D73	-10.33	1.52	1.79	-2.82	D73	88.7	90.6	93.3	99.9
								D74 0.0108	D74	-4.55	-0.75	1.93	0.12	D74	81.6	83.9	98.6	98.7
								D75 0.0365	D75	-3.39	-4.30	-2.67	0.11	D75	30.9	80.7	99.9	99.9
								D80 0.0095	D80	1.85	1.16	-0.36	-1.12	D80	55.4	77.4	79.4	99.6
								D90 0.0024	D90	-2.41	-1.29	-1.40	-2.20	D90	25.4	32.6	41.2	62.5
								D94 0.0211	D94	0.53	-0.16	-0.13	-0.05	D94	47.5	51.8	54.6	55.0
								D95 0.0190	D95	2.36	0.34	-0.09	-0.57	D95	92.1	93.9	94.1	99.5

- Commentez la nouvelle répartition de l'inertie. Combien d'axes principaux retient-on ? La situation est-elle meilleure qu'avec la première analyse ?
- Quelles sont les variables qui déterminent les axes que l'on retient ? Précisez les critères utilisés. Y a-t-il un effet de taille ?
- Parmi les départements dont les données sont fournies ci-dessus, quels sont ceux qui déterminent les axes que l'on retient ? Précisez les critères utilisés. Y a-t-il des départements sur-représentés ?
- Comment peut-on interpréter les axes à partir des deux questions précédentes ?
- Parmi les départements dont les données sont fournies ci-dessus, quels sont ceux dont la qualité de représentation est mauvaise sur l'espace propre retenu ? Précisez les critères utilisés.