

NOTION DE CORRELATION	
Introduction	<i>Coefficient de corrélation de Bravais-Pearson, ou encore coefficient de corrélation linéaire.</i>
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	Le coefficient de corrélation permet de mesurer la liaison ou le lien entre 2 ensembles de données.
Interprétation de la covariance	
Significativité du coefficient de corrélation	Par exemple :
Notion de ddl	-Existe-t-il un lien entre le fait que les enfants mangent des sucreries et leur fréquentation des dentistes ?
Notion d'Hypothèses	-La satisfaction des clients est-elle liée à la température qui règne dans les magasins ?
Notion de seuil de décision et de type d'erreur	-Est-ce que le niveau d'études atteint dépend du milieu social ?
Exemple 02	-Est-ce que la mémorisation des mots d'un texte dépend de la longueur des mots ?
Représentation graphique	-Est-ce que l'impact d'une campagne publicitaire dépend du média choisi ?
Remarques	-Est-ce que le cours du pétrole dépend de celui de l'euro ?
	-Est-ce que le cours de l'euro dépend de celui du pétrole ?
	-Est-ce que le loisir préféré des étudiants dépend de leur sexe ?
	Toutes ces questions mettent en jeu <u>deux</u> variables. Ces deux variables sont observées sur la <u>même</u> population.

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

Ce coefficient peut donc nous aider à émettre des pronostics mais **attention**, ce n'est pas parce qu'une corrélation existe entre 2 séries statistiques qu'il y a un lien de cause à effet entre les deux !

Exemple : une corrélation (probable) entre la taille des enfants entre 5 et 10 et leur score à un test en rapport avec le QI (non normalisé par tranche d'âge) ne signifie pas que « plus on est grand, plus on est intelligent ».

Autre exemple : s'il y a plus de naissances au printemps et en automne, on risque de trouver une corrélation entre le nombre de passage de cigognes et le nombre de naissance.

Dans ces 2 exemples, le lien provient d'une 3^e variable (âge dans le 1^{er} cas et saison dans le 2^e cas).

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

Exemple : Considérons une population de couples (femme , mari) et associons à chaque couple un couple d'observations (âge de la femme ; âge du mari) l'année de leur mariage.

Sur cette (même) population on observe deux variables :

Variable X : âge de la femme l'année de son mariage.

Variable Y : âge du mari l'année de son mariage

On obtient une série double (une série de 10 couples)

X	18	21	21	19	22	20	19	18	22	20
Y	20	24	26	20	24	26	24	20	26	24

NOTION DE CORRELATION

Introduction	<p>On peut évidemment étudier chaque variable <u>indépendamment</u> de l'autre. On se limitera au calcul de la moyenne et de la variance.</p> <p>moyenne=$mX=(18+21+\dots)/10=20$ variance=$\sigma^2X=((18-21)^2+\dots)/10=2$ écart-type=$\text{racine carré}(\sigma^2X)=\sigma X=1,42$</p> <p>moyenne=$mY=(20+24+\dots)/10=23,4$ variance=$\sigma^2Y=((20-25)^2+\dots)/10=5,64$ écart-type =$\text{racine carré}(\sigma^2Y)=\sigma Y=2,37$</p>
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	
Exemple 02	
Représentation graphique	
Remarques	

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

A chaque couple d'observations (x_i, y_i) on associe, dans un repère cartésien, un point (géométrique) d'abscisse x_i et d'ordonnée y_i . L'ensemble des points ainsi obtenu est appelé nuage de points (ou nuage des individus).



NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

A ce nuage de points on ajoute le point $(mX, mY)=(20 ; 23,4)$ appelé centre de gravité du nuage ou plus simplement centre du nuage (on dit aussi point moyen).



NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

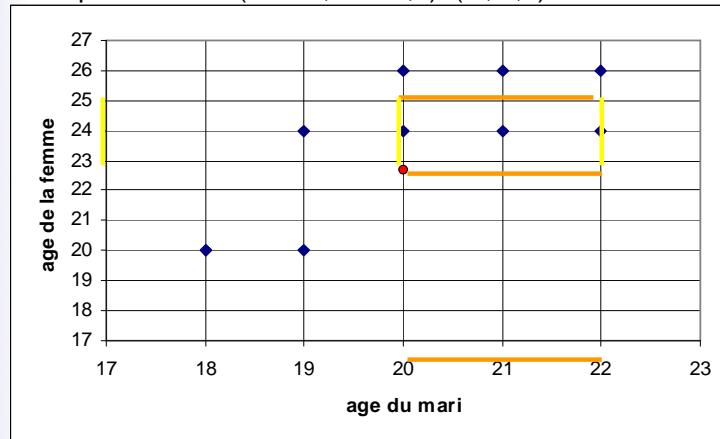
Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

A chaque point (x_i, y_i) on associe son écart par rapport au point moyen $(m_x, m_y)=(20 ; 23,4)$. On obtient un couple d'écarts : $(x_i - m_x, y_i - m_y)$. Exemple au point $(22,26)$ on associe le couple d'écarts : $(22-20, 26-23,4)=(2, 2,6)$



Le couple d'écart exprime les deux variances d'un point.

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

CALCUL DU COEFFICIENT DE CORRÉLATION

Le coefficient de corrélation, noté r , se calcule à partir de 2 recueils de données différents.

$$r = \frac{\text{COV}(x, y)}{\sigma_x \times \sigma_y}$$

Avec σ_x, σ_y = écart type de X et de Y

$$\text{Et } \text{cov}(x, y) = \frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y}) = \frac{\sum x_i y_i}{N} - \bar{x} \bar{y}$$

= covariance de X et Y

= variance (variation) commune aux deux variables/ produit des variances

= moyenne du produit des écarts de deux variables à leur moyenne respective

NOTION DE CORRELATION

	moyenne										
x_i	18	21	21	19	22	20	19	18	22	20	20
y_i	20	24	26	20	24	26	24	20	26	24	23,4
$x_i -$ moyenne X	-2	1	1	-1	2	0	-1	-2	2	0	
$y_i -$ moyenne Y	-3,4	0,6	2,6	-3,4	0,6	2,6	0,6	-3,4	2,6	0,6	
$(x_i -$ moyenne X) * ($y_i -$ moyenne Y)	6,8	0,6	2,6	3,4	1,2	0	-0,6	6,8	5,2	0	covarian ce=SOM ME/10 = 2,6

NOTION DE CORRELATION

Introduction	La covariance permet d'estimer le sens de la variation entre deux variables numériques :
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	► Les variables varient dans le même sens (ou covariant) : les sujets qui ont des valeurs fortes (au dessus de la moyenne) sur une variable, présentent également des valeurs élevées sur l'autre variable. Autrement dit, les sujets les plus grands sont les sujets les plus lourds et, inversement, les sujets plus petits sont les plus légers : dans ce cas, la covariance est de signe positif (+66,831)
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	► Les variables varient en sens inverse : les sujets qui ont des valeurs fortes sur une des deux variables auront tendance à avoir des notes faibles sur l'autre variable. La valeur de la covariance sera alors de signe négatif
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	
Exemple 02	► Les variables ne covariant pas : Parmi les sujets présentant des valeurs fortes sur une variable, on peut observer que, sur l'autre variable, ces mêmes sujets obtiennent des notes fortes, faibles ou moyennes. La valeur de la covariance est proche de la valeur 0.
Représentation graphique	
Remarques	

NOTION DE CORRELATION

Introduction	
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	Significativité du coefficient de corrélation
Notion de ddl	Si le coefficient de corrélation est proche de 1 ou de -1, cela signifie que les variables x et y sont très fortement liées (par une relation de la forme $y = ax + b$). S'il est proche de 0, cela signifie que le lien est peu probable. Un signe positif indique que x et y varient dans le même sens.
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	Mais on ne peut toujours se contenter d'une telle interprétation approximative. Il est nécessaire de savoir si ce lien est réel ou s'il est dû au simple hasard de nos mesures. Autrement dit peut-on le considérer comme reflétant un lien entre les 2 variables sur l'ensemble de la population et non pas seulement sur notre échantillon.
Exemple 02	
Représentation graphique	
Remarques	

NOTION DE CORRELATION

Introduction	
Exemple 01	<u>1°étape</u> : poser les hypothèses
Calcul de la covariance et du coefficient de corrélation	Hypothèse nulle (H_0) : il n'y a pas de lien statistique entre les 2 séries
Interprétation de la covariance	Hypothèse alternative (H_1) : il existe un lien statistique entre les 2 séries
Significativité du coefficient de corrélation	
Notion de ddl	<u>2°étape</u> : calcul de r (appelé r calculé).
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	
Exemple 02	
Représentation graphique	
Remarques	

NOTION DE CORRELATION

Introduction	
Exemple 01	<u>Table du coefficient de corrélation.</u>
Calcul de la covariance et du coefficient de corrélation	Les lignes correspondent au nombre de degrés de liberté (ddl) :
Interprétation de la covariance	avec 2 variables, $ddl = N - 2$ (2 = nombre de variable)
Significativité du coefficient de corrélation	Lecture de r (appelé r_{lu}) correspondant à un risque d'erreur de 5% (ou probabilité de 0.05) de rejeter à tort H_0 .
Notion de ddl	Si $r_{calculé} \geq r_{lu}$ alors on rejette l'hypothèse nulle et accepte H_1 : les deux distributions sont statistiquement liées. On peut l'affirmer avec un risque d'erreur $< 5\%$.
Notion d'Hypothèses	Si $r_{calculé} < r_{lu}$ alors on ne rejette pas l'hypothèse nulle : on ne peut pas dire que les deux distributions sont statistiquement liées. Mais attention : cela ne veut pas dire que l'on accepte H_0
Notion de seuil de décision et de type d'erreur	
Exemple 02	
Représentation graphique	
Remarques	

NOTION DE CORRELATION

Introduction	
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	Notion du degré de liberté (ddl):
Notion d'Hypothèses	Le ddl reflète le nombre d'éléments indépendants entrant dans l'estimation d'une variance. C'est le nombre d'éléments pouvant varier librement, de là le terme de liberté.
Notion de seuil de décision et de type d'erreur	Par exemple si nous voulons estimer la variance de 5 observations, nous devons dans un premier temps estimer la moyenne puisque la variance s'obtient à partir de l'écart à la moyenne. Or si nous fixons la somme des éléments, seuls quatre peuvent varier librement, la valeur du cinquième étant imposée par la valeur de la somme.
Exemple 02	Par exemple la somme est 30 et les quatre premières valeurs sont 5, 6, 2 et 10: alors la valeur du cinquième est forcément 7 pour obtenir la somme de 30 ($30 - (5+6+2+10)$).
Représentation graphique	Donc pour estimer la moyenne d'une population à partir de la moyenne d'un échantillon, nous perdons un degré de liberté.
Remarques	$DDL = \text{nombre d'observations} - \text{nombre de variables}$

NOTION DE CORRELATION

Introduction	
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	
Exemple 02	
Représentation graphique	
Remarques	

Notion d'hypothèses:
Nous souhaitons comparer des résultats obtenus à partir d'un échantillon (i.e. statistiques) à ceux d'une population (i.e. paramètres). Nous voulons faire des inférences d'un échantillon sur une population. En gros ceci revient à dire que nous voulons savoir si notre échantillon est représentatif de la population.
Dans le cas de corrélation, nous voulons savoir si les deux séries sont linéairement liées ou pas. Nous avons deux possibilités, deux alternatives:
OUI: les deux séries sont liées. (H1)
NON: les deux séries ne sont pas liées. (H0)
On pose alors deux hypothèses et on va décider laquelle on accepte:
H1: Hypothèse alternative. C'est une hypothèse statistique inexacte (si l'effet n'est pas nul, alors son intensité est différente de zéro, et il existe une infinité de façon d'être différent de zéro)
H0: Hypothèse nulle. C'est une hypothèse statistique exacte (si l'effet est nul alors l'intensité de l'effet est de zéro, une seule possibilité)

NOTION DE CORRELATION

Introduction	
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	
Notion d'Hypothèses	Notion d'hypothèses: Les hypothèses concernent toujours la population, pas l'échantillon. Comme il n'y a que H1 qui est une hypothèse inexacte alors : On ne peut jamais rejeter H1. Donc: On ne peut jamais accepter H0. ON NE PEUT QUE REJETER ou NON H0.
Notion de seuil de décision et de type d'erreur	Mais à partir de quel moment acceptons H1 ou rejetons H0. Il nous faut bien un critère pour prendre la décision.
Exemple 02	
Représentation graphique	Attention: ne pas rejeter H0 n'est pas équivalent à l'accepter.
Remarques	

NOTION DE CORRELATION	
Introduction	<p><u>Notion de risque d'erreur:</u> Le seuil d'improbabilité est le seuil de significativité = α Généralement on utilise comme seuil .05 ou .01. $\alpha = .05$ identifie la notion d'évènement improbable avec événements ayant moins de 5 chances sur 100 de se produire. α est de l'ordre du subjectif, c'est notre seuil de décision, je suis sûr de ce que je dis à 95% de chance.</p> <p><u>Deux types d'erreurs:</u> -<u>erreur de type I</u>: accepter H1 alors qu'elle n'est pas vraie. Dire qu'il y a un effet alors que cet effet n'existe pas dans la réalité. La probabilité de faire une erreur de type I correspond à α,</p> <p style="text-align: center;">⇒ Pour minimiser Erreur de type I: il faut un α petit.</p> <p>-<u>erreur de type II</u>: rejeter H1 alors qu'elle est vraie. Ne pas accepter l'existence d'un effet alors que cet effet existe. La probabilité de faire une erreur de type II correspond à β, qui ne peut être connu car H1 est inexacte.</p> <p style="text-align: center;">⇒ Pour minimiser Erreur de type II: il faut un α grand.</p>
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	
Exemple 02	
Représentation graphique	
Remarques	

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

Les deux erreurs varient en sens inverse en fonction du seuil α

décision expérimentale	Etat de la nature (inconnu)	
	H0 vraie (H1 fausse)	H0 fausse (H1 vraie)
Non rejet de H0	Non rejet correct probabilité = $1-\alpha$	Erreur de type II probabilité = β
Rejet de H0	Erreur de type I probabilité = α	Détection correct probabilité = $1-\beta$

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques



Covariance = 2,6

r calculé = $2,6 / (1,42 * 2,37) = 0,77$

r lu (avec $\alpha = 0.05$ et $ddl = 8$) = 0,6319

r calculé > r lu donc rejet de H_0 , accepte H_1

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques

Ex : Sur une population de 7 individus, on mesure 2 caractères X et Y qui prennent les valeurs suivantes :

X	5	-2	0	-1	1	2	-3
Y	9	4	0	1	1	4	9

x_i	y_i	$x_i y_i$	x_i^2	y_i^2
5	9	45	25	81
-2	4	-8	4	16
0	0	0	0	0
-1	1	-1	1	1
1	1	1	1	1
2	4	8	4	16
-3	9	-27	9	81
$\Sigma=2$	$\Sigma=28$	$\Sigma=18$	$\Sigma=44$	$\Sigma=196$

NOTION DE CORRELATION

Introduction	x_i	y_i	$x_i y_i$	x_i^2	y_i^2	<p>Calcul de cov (x,y) : $x = \sum x_i/N = 2/7 = 0,2857 = 0,29$ $y = \sum y_i/N = 28/7 = 4$ $cov(x,y) = 18/7 - 0,29 \cdot 4 = 2,57 - 1,16 = 1,41$</p> <p>Calcul de σ_x et σ_y : $\sigma_x^2 = \sum x_i^2/N - x^2 = 44/7 - 0,29^2 = 6,29 - 0,0841 = 6,21 \rightarrow \sigma_x = \sqrt{6,21} = 2,49$ $\sigma_y^2 = \sum y_i^2/N - y^2 = 196/7 - 4^2 = 28 - 16 = 12 \rightarrow \sigma_y = \sqrt{12} = 3,46$</p> <p>Calcul de r : $r \text{ calculé} = 1,41 / (2,49 \cdot 3,46) = 0,16$ Le degré de corrélation entre ces deux variables avoisine le 0, on peut donc penser qu'elles ne sont pas très liées mais on vérifie $ddl : N-2 = 7-2 = 5$ $r_{lu} = 0,7545$ donc $r \text{ calculé} = 0,16 < r_{lu} : 0,7545$ \Rightarrow on ne rejette pas H_0 : on ne peut conclure à un lien statistique entre X et Y</p>
Exemple 01	5	9	45	25	81	
Calcul de la covariance et du coefficient de corrélation	-2	4	-8	4	16	
	0	0	0	0	0	
Interprétation de la covariance	-1	1	-1	1	1	
	1	1	1	1	1	
Significativité du coefficient de corrélation	2	4	8	4	16	
	-3	9	-27	9	81	
Notion de ddl	$\Sigma=2$	$\Sigma=28$	$\Sigma=18$	$\Sigma=44$	$\Sigma=196$	
Notion d'Hypothèses						
Notion de seuil de décision et de type d'erreur						
Exemple 02						
Représentation graphique						
Remarques						

NOTION DE CORRELATION

Introduction	
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	
Exemple 02	
Représentation graphique	<u>Représentation graphique.</u>
Remarques	Si l'on représente sur un graphique les scores de chaque sujet dans un repère orthogonal dont les axes correspondent aux 2 variables : Le lien entre les 2 variables est visible dans le fait que le nuage de point se rapproche fortement d'une droite (appelée droite de régression) de type $ax + b$ de pente non nulle. Si la pente est positive alors la corrélation est positive, si elle est négative alors la corrélation est négative. <i>Remarque</i> : une droite horizontale ou verticale signifierait que les données ne dépendent que d'une seule des 2 variables.

NOTION DE CORRELATION

Introduction

Exemple 01

Calcul de la covariance et du coefficient de corrélation

Interprétation de la covariance

Significativité du coefficient de corrélation

Notion de ddl

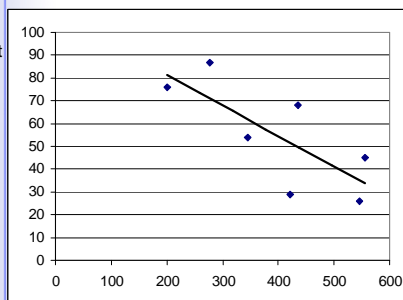
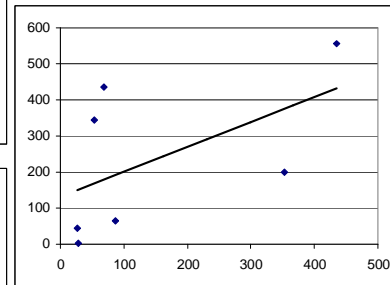
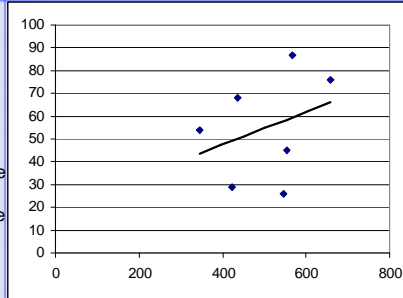
Notion d'Hypothèses

Notion de seuil de décision et de type d'erreur

Exemple 02

Représentation graphique

Remarques



NOTION DE CORRELATION

Introduction	<p>REMARQUE:</p> <p>Propriétés de la covariance :</p> <p>Symétrie : $\text{cov}(X,Y) = \text{cov}(Y,X)$</p> <p>$\text{Cov}(X,X)=\text{var}(X)$</p> <p>Les vecteurs dont on veut connaître la covariance doivent être de taille identique</p> <p>Linéarité:</p> <p>On calcule la régression linéaire par la méthode des moindres carrés, c'est-à-dire qu'on cherche la droite $y=a'x+b'$ qui minimise les distances entre la droite et tous les points. C'est-à-dire qu'on minimise $\sum (y_i - y'_i)^2$.</p> <p>La régression linéaire Dx/y a pour forme: $y=a'x+b'$ Où $a' = \text{cov}(x,y) / \text{var}(y)$ et $b' = \text{moyenne}(X) - (a' * \text{moyenne}(Y))$</p> <p>La régression linéaire Dy/x a pour forme: $y=a'x+b'$ Où $a' = \text{cov}(x,y) / \text{var}(x)$ et $b' = \text{moyenne}(Y) - (a' * \text{moyenne}(X))$</p>
Exemple 01	
Calcul de la covariance et du coefficient de corrélation	
Interprétation de la covariance	
Significativité du coefficient de corrélation	
Notion de ddl	
Notion d'Hypothèses	
Notion de seuil de décision et de type d'erreur	
Exemple 02	
Représentation graphique	
Remarques	