

VARIABLE ALÉATOIRE (C1)

(18 / 10 / 2021, © Monfort, Dicostat2005, 2005-2021)

[**Note.** On donne souvent à une **variable aléatoire** le même qualificatif que le nom de sa **loi de probabilité** : variable gaussienne et **loi gaussienne**, variable de POISSON - ou variable poissonnienne - et **loi de POISSON**, etc]

Notion fondamentale du **calcul des probabilités**, la notion de **variable aléatoire** correspond à la notion mathématique (théorie de la mesure) d'**application mesurable**.

En **Statistique**, on utilise plutôt le concept « technique » de **statistique**, en association avec les notions de **modèle statistique** et de **problème statistique**.

(i) Soit (Ω, \mathcal{F}) et $(\mathcal{X}, \mathcal{B})$ deux **espaces mesurables** (ou **espaces probabilisables**) et $\xi : \Omega \mapsto \mathcal{X}$ une **application** donnée.

On dit que ξ est une **variable aléatoire**, ou un **aléas**, ou un **élément aléatoire**, défini sur Ω et à valeurs dans \mathcal{X} , ssi ξ est une application $(\mathcal{F}, \mathcal{B})$ -mesurable, ie ssi ξ est compatible avec les **structures** mesurables (Ω, \mathcal{F}) et $(\mathcal{X}, \mathcal{B})$ resp associées à Ω et à \mathcal{X} . Autrement dit, par définition :

$$(1) \quad \xi^{-1}(B) \in \mathcal{F}, \quad \forall B \in \mathcal{B}, \quad \text{ou encore} \quad \xi^{-1}(\mathcal{B}) \subset \mathcal{F},$$

ie l'image inverse de la **tribu** \mathcal{B} par ξ est contenue dans la tribu \mathcal{F} .

Une variable aléatoire (ou **va**) se note à l'aide de symboles variés : eg des lettres grecques ξ, η, ζ ou ε (etc), ou des lettres latines majuscules X, Y, Z ou U, etc :

(a) les lettres grecques peuvent servir à représenter des va « théoriques » (eg des **inobservables**) ;

(b) les lettres latines représentent des va « empiriques », ie des variables constituant des **observations** des va précédentes.

Ainsi, ξ représentera une **vars** et $X = (X_1, \dots, X_N)$ un échantillon constitué de **copies** (ou « répliques ») analogues à la « **variable parente** » ξ , eg des copies indépendantes entre elles et équadistribuées comme ξ (cf **échantillon iid**).

Plus généralement :

(a) un **échantillon aléatoire** quelconque $X = (X_1, \dots, X_N)$ peut aussi être considéré comme une va à valeurs dans un espace produit $\mathcal{X} = \prod_{n=1}^N \mathcal{X}_n$, ou dans un espace puissance $\mathcal{X} = \mathcal{X}_0^N$, ou encore dans un **espace vectoriel** (réel) \mathcal{X} ;

(b) de même, un **processus stochastique** peut être considéré comme une variable aléatoire $X = (X_t)_{t \in T}$ à valeurs dans un espace produit (en général infini) $\mathcal{X} = \prod_{t \in T} \mathcal{X}_t$, ou dans un espace puissance $\mathcal{X} = \mathcal{X}_0^T$.

Si ξ est une va, on note souvent $x = \xi(\omega)$ l'image d'un élément $\omega \in \Omega$ par ξ . On note encore $x = X(\omega)$ l'image d'un élément $\omega \in \Omega$ par X : cette image (ou « valeur ») est appelée **observation** lorsqu'elle peut être « appréhendable » sur des **unités statistiques** (cf **observabilité**, **variable observable**). On l'appelle encore « mesure », quelle que soit la nature de la variable considérée.

(ii) Une va est généralement qualifiée en fonction de la nature (mathématique) de l'**ensemble** de ses valeurs. Ainsi, on distingue trois principaux types de variables :

(a) $\xi : \Omega \mapsto \mathcal{X}$ est une **variable aléatoire quantitative**, ou **variable aléatoire numérique**, (scalaire) ssi l'ensemble \mathcal{X} de ses valeurs est un ensemble numérique (cf **variable quantitative**). Parmi les variables numériques :

(a)₁ ξ est appelée **variable discrète** ssi \mathcal{X} est un **ensemble discret**, ie un ensemble ayant au plus la **puissance du dénombrable** (ie soit fini, soit en bijection avec l'ensemble \mathbf{N} ou l'une de ses parties (cf **loi discrète**). Une variable de cette nature peut encore être qualifiée de **variable « simple »**, ou de **variable « scalaire »** ou de **variable « à une dimension »**.

Ainsi, on dit que ξ est :

* une **variable entière non négative** lorsque $\mathcal{X} = \mathbf{N}$; une variable entière non négative (finie ou bornée) lorsque $\mathcal{X} = \{0, 1, \dots, n\} = \mathbf{N}_n$;

* une **variable entière positive** (ou **variable « naturelle »**) lorsque $\mathcal{X} = \mathbf{N}^* = \mathbf{N} \setminus \{0\}$; une variable entière positive (finie ou bornée) lorsque $\mathcal{X} = \{1, \dots, n\} = \mathbf{N}_n^*$;

* une **variable entière signée** lorsque $\mathcal{X} = \mathbf{Z}$; une variable entière signée (finie ou bornée) lorsque $\mathcal{X} = \{-m, \dots, 0, \dots, +n\} = \mathbf{Z}_{mn}$;

* une **variable rationnelle signée** lorsque $\mathcal{X} = \mathbf{Q} = \mathbf{Z} / \mathbf{N}^*$.

Par extension, on peut aussi considérer comme discrète une variable « continue » dont la **loi de probabilité** ne « charge » qu'un ensemble discret \mathcal{X} du type ci-dessus : cette loi est parfois appelée « **distribution de masses** ».

(a)₂ ξ est appelée **variable continue** lorsque l'ensemble \mathcal{X} de ses valeurs possède la puissance du continu. Une variable de cette nature est parfois appelée **variable « scalaire »** ou **variable « à une dimension »**.

Ainsi, on dit que ξ est :

* une **variable réelle signée** lorsque $\mathcal{X} = \mathbf{R}$;

* une **variable réelle non négative** lorsque $\mathcal{X} = \mathbf{R}_+ = \{x \in \mathbf{R} : x \geq 0\}$; une **variable réelle positive** lorsque $\mathcal{X} = \mathbf{R}_+^* = \{x \in \mathbf{R} : x > 0\} = \{x \in \mathbf{R}_+ \setminus \{0\}\}$; il en va de façon analogue pour les valeurs dans \mathbf{R}_- ou \mathbf{R}_-^* ;

* une **variable complexe** lorsque $\mathcal{X} = \mathbf{C}$.

Lorsque $\mathcal{X} = \mathbf{R}$, on parle parfois d'**aléa numérique**. En pratique, le plus souvent (cas scalaire), $\mathcal{X} = \mathbf{D}$ (ensemble des **nombre décimaux**) : en effet, la « représentation machine » d'un nombre mémorisé dans un ordinateur (mémoire vive ou mémoire morte) ne peut, pour des raisons matérielles, qu'être tronquée, même si sa précision est très élevée ;

(b) $\kappa : \Omega \mapsto \mathcal{K}$ est une **variable aléatoire qualitative**, ou une **variable aléatoire attribut**, ssi l'ensemble \mathcal{K} de ses valeurs n'est doté d'aucune propriété en termes de calcul algébrique (cf **variable qualitative**). Autrement dit, \mathcal{K} est un **ensemble non numérique**, ie un ensemble dont les éléments ne peuvent faire l'objet de calculs algébriques. Dans certains cas, les éléments de \mathcal{K} peuvent avoir une « apparence numérique » (eg identifiant, numéro d'ordre, etc), mais les calculs entre ces éléments sont insensés.

Une variable de cette nature peut aussi être qualifiée de « **variable simple** », ou de « **variable scalaire** », ou encore de « **variable à une dimension** » .

On distingue usuellement entre :

(b)₁ **variable nominale**, pour laquelle \mathcal{K} ne possède pas de structure particulière ;

(b)₂ **variable ordinale**, pour laquelle \mathcal{K} est un ensemble ordonné (cf **relation d'ordre**), qui peut lui-même résulter d'un ensemble numérique ;

(c) **variable morphologique** ou **variable objectale**, $\xi : \Omega \mapsto \mathcal{K}$ tq l'ensemble \mathcal{K} de ses valeurs est constitué de « formes » ou d'« objets » particuliers (cf **forme**, **reconnaissance des formes**), eg :

(c)₁ \mathcal{K} est un ensemble de figures ou de formes géométriques : coniques, quadriques, ou, plus généralement, **variétés différentielles**, **graphes**

(arborescence ou réseau), etc. Cet ensemble correspond aux valeurs prises par une **variable aléatoire géométrique** ;

(c)₂ \mathcal{K} est un ensemble (ou « espace ») de **fonctions numériques**. Cet espace correspond à une **variable aléatoire fonctionnelle**.

Dans cette dernière situation, les formes $k \in \mathcal{K}$ (ensemble considéré dépendent souvent d'une variable (vectorielle réelle) θ , ie $\theta \mapsto k = \phi(\theta)$, avec $\theta \in \Theta \subset \mathbf{R}^Q$ (où $Q \in \mathbf{N}^*$). La **loi multidimensionnelle (loi jointe)** de θ détermine alors le caractère aléatoire de la forme ϕ : l'analyse de cette situation se ramène à celle de la loi du **vecteur aléatoire** θ , donc de la loi d'une variable quantitative vectorielle (cf infra).

Si les modalités constituant \mathcal{K} sont déterminées par une « **morphologie** », alors \mathcal{K} est un ensemble de formes que l'on peut distinguer des deux types précédents (sauf eg lorsque \mathcal{K} est lui-même une classification déduite d'un ensemble de formes).

On peut ainsi considérer que la notion de variable morphologique est une notion intermédiaire entre celle de variable qualitative (puisqu'on ne peut effectuer de calculs entre formes) et celle de variable numérique (puisque les formes sont souvent définies ou calculables à partir d'équations numériques tq $k = \phi(\theta)$ ci-dessus).

(iii) On appelle **suite aléatoire**, ou parfois « **liste** » **aléatoire**, (finie) une va ζ dont l'ensemble \mathcal{Z} des « valeurs » est un **ensemble produit (fini)** de la forme $\mathcal{Z} = \prod_{i=1}^n \mathcal{Z}_i$. Chaque composante \mathcal{Z}_i de ce produit est muni d'une **structure** mesurable, éventuellement adaptée à sa nature, ou à sa structure, particulière : \mathcal{Z} peut ainsi constituer un « mélange » d'ensembles numériques ou non, géométriques ou non, unidimensionnels ou non, etc.

On appelle **vecteur aléatoire**, ou **variable aléatoire vectorielle**, une va $\xi : \Omega \mapsto \mathcal{X}$ tq \mathcal{X} est un **espace vectoriel** (réel) muni d'une structure mesurable (resp **espace vectoriel mesurable**). On dit notamment que ξ est un **élément aléatoire** ssi \mathcal{X} est un **espace de BANACH** (généralement muni de sa **tribu borélienne**).

Ainsi, lorsque \mathcal{X} est un espace vectoriel sur un corps \mathbf{K} (en pratique $\mathbf{K} = \mathbf{R}$ ou $\mathbf{K} = \mathbf{C}$) et qu'il est muni d'une base canonique $(e_i)_{i=1, \dots, n}$, la variable ξ se décompose selon $\xi = \sum_{i=1}^n \xi_i e_i$, où les ξ_i sont des va à valeurs dans \mathbf{K} .

(a) si $\mathcal{X} = \mathbf{R}^n$, on dit que ξ est un **vecteur aléatoire réel**, ou une **variable aléatoire vectorielle réelle** ;

(b) si $\mathcal{X} = \mathbf{C}^n$, on dit que ξ est un **vecteur aléatoire complexe**, ou une **variable aléatoire vectorielle complexe** ;

En pratique, on identifie \mathbf{C} à \mathbf{R}^2 , muni de sa tribu borélienne $\mathcal{B}(\mathbf{R}^2)$.

En particulier :

(a) si $\mathcal{X} = \mathbf{K}$, on dit que ξ est une **variable aléatoire réelle scalaire** (vars) si $\mathbf{K} = \mathbf{R}$ ou une **variable aléatoire complexe scalaire** (vacs) si $\mathbf{K} = \mathbf{C}$;

(b) si ξ est une **variable aléatoire matricielle**, ou **matrice aléatoire**, ssi l'ensemble de ses valeurs est $\mathcal{X} = M_{mn}(\mathbf{K})$ (espace vectoriel des matrices de format (m,n) , où \mathbf{K} est un corps donné). Cette notion est à distinguer de celle de **matrice stochastique**.

(iv) Le concept de va est donc très générale, et il en est de même de celui de statistique (qui en résulte).

Ainsi, une va peut prendre ses valeurs dans un ensemble \mathcal{X} dont les éléments sont de nature très variée, eg :

(a) **histogrammes** ;

(b) **matrices** ou **tableaux statistiques** ;

(c) figures géométriques : graphiques, diagrammes, cercles, triangles, etc (cf **probabilité géométrique**). On peut aussi considérer :

(c)₁ des **graphes** : on parle alors de **graphes aléatoires** ;

(c)₂ des **réseaux** : on parle alors de **réseaux aléatoires** ;

(v) Lorsque (Ω, \mathcal{F}) est muni d'une **mesure de probabilité** définie sur \mathcal{F} , l'**espace probabilisé** (Ω, \mathcal{F}, P) devient, par ξ , un espace probabilisé $(\mathcal{X}, \mathcal{B}, P^\xi)$, où P^ξ est, par définition, la **loi de probabilité** de ξ (ie l'image de P par ξ) (cf aussi **mesure image**). Dans ce cas, les qualificatifs précédents sont utilisés en considérant le **support** $\text{Supp } P^\xi$ de la loi de ξ , et non pas l'ensemble $\xi(\Omega) \subset \mathcal{X}$ des valeurs possibles de ξ , ie son image $\text{Im } \xi$ (image de Ω par ξ).

(vi) Par ailleurs :

(a) on suppose souvent que $(\mathcal{X}, \mathcal{B}) = (\Omega, \mathcal{F})$ et $\xi = \text{id}_\Omega$ (identification de l'**espace fondamental** de départ et de l'**espace d'observation** image) ;

(b) la notion de variable aléatoire ne dépend pas du fait que l'**espace probabilisable** (Ω, \mathcal{F}) est muni d'une (mesure de) probabilité P . Cependant, on considère parfois l'espace probabilisé (Ω, \mathcal{F}, P) et l'ensemble $\mathcal{A}(\Omega, \mathcal{F})$ des

applications mesurables (va au sens précédent) pour définir la **relation d'équivalence** suivante sur $\mathcal{A}(\Omega, \mathcal{F})$:

$$(2) \quad \xi \sim \eta \Leftrightarrow \xi = \eta \text{ (P - p.s.)}$$

On appelle alors **variable aléatoire** tout élément de l'**espace quotient** $A(\Omega, \mathcal{F}) = \mathcal{A}(\Omega, \mathcal{F}) / \sim$, constitué des classes de va (au sens précédent) qui sont P-presque sûrement égales entre elles. Par suite, les concepts usuels (**espérance mathématique**, **moment**, **mode**, **quantile**, etc) sont définis à partir d'éléments de $A(\Omega, \mathcal{F})$;

(c) le concept de variable aléatoire généralise celui d'**événement aléatoire**. En effet, la variable aléatoire $\mathbf{1}_A$ (**variable indicatrice** de l'événement $A \in \Omega$) est une **vars** particulière, puisqu'elle prend ses valeurs dans $[0, 1] \subset \mathbf{R}$ (cf aussi **fonction indicatrice**) ;

(d) une **variable aléatoire** peut être **observable** (ie « visible ») (cf **observabilité**) ou **non observable** (ie **inobservable** ou « invisible »). Ainsi, on peut être conduit à distinguer (eg en psychologie) entre « **variables manifestes** » et « **variables latentes** ».

En effet, il arrive souvent qu'une théorie (ou un modèle), dans sa **forme initiale**, relie, sous des formes tq :

$$(3) \quad \omega_g = \rho_g(t_h) \quad (g = 1, \dots, G) \quad (h = 1, \dots, H)$$

une liste ω_g de variables observables à une liste t_h de variables inobservables, en sorte que (si l'élimination des t_h est possible), dans la **forme finale** de cette théorie (ou de ce modèle) seules les variables observables ω_g apparaissent (ie soient manifestes) :

$$(4) \quad \sigma(\omega_1, \dots, \omega_G) = 0.$$

Ce procédé permet, en théorie, de réaliser l'**inférence statistique** recherchée, mais comporte des risques. Ainsi, en supposant que la « vraie » relation (inconnue) entre observables est $\omega_2 = \delta / \omega_1$ alors que la **spécification** initiale s'écrit $\omega_1 = \alpha \cdot i_1 + \beta$ et $\omega_2 = \gamma \cdot i_1^2$, l'élimination de i_1 conduit à une relation $\omega_2 = \gamma \cdot ((\omega_1 - \beta) / \alpha)^2 + \beta$ différente de la vraie (l'une est hyperbolique, l'autre quadratique, en ω_1).

(vii) On donne souvent à une variable aléatoire le **nom** de sa loi de probabilité : ainsi, on parle de **variable uniforme** (dont la loi est uniforme, qu'elle soit discrète ou continue), de variable binômiale (dont la loi est la **loi binômiale**), de **variable de POISSON** (dont la loi est la **loi de POISSON**), de **variable normale** ou de **variable gaussienne** (dont la loi est la **loi normale** ou **loi gaussienne**), etc.

(viii) Les va utilisées sont souvent des transformées de variables « initiales » : leurs « valeurs » (éléments de \mathcal{X}) sont donc obtenues par **changement de variable**

aléatoire. De même, l'observation des va sur des **unités statistiques** fait souvent, ensuite, l'objet de **transformation des données** initiales (celles de l'**échantillon**) : « tabulation » (ie **agrégation**), groupement de données en classes (cf **groupement de classes**), « représentations » graphiques, etc (cf aussi **disposition**).

(ix) La notion de variable aléatoire est généralement associée à celle d'**espace probabilisable** (ou d'espace probabilisé). Elle se généralise en celle de **statistique**, laquelle est associée à la notion de **modèle statistique**, ou de **représentation statistique**, ie de **famille** d'espaces probabilisés.