

Pratique R

Psychologie statistique avec R

Yvonnick Noël

Psychologie statistique
avec R

Vj ku' r ci g' k p v g p v k q p c m { ' i g h v ' d r e p m

Yvonnick Noël

Psychologie statistique avec R

 edp sciences

ISBN : 978-2-7598-1736-8

© **2015, EDP Sciences**, 17, avenue du Hoggar, BP 112, Parc d'activités de Courtabœuf,
91944 Les Ulis Cedex A

Imprimé en France

Tous droits de traduction, d'adaptation et de reproduction par tous procédés réservés pour tous pays. Toute reproduction ou représentation intégrale ou partielle, par quelque procédé que ce soit, des pages publiées dans le présent ouvrage, faite sans l'autorisation de l'éditeur est illicite et constitue une contrefaçon. Seules sont autorisées, d'une part, les reproductions strictement réservées à l'usage privé du copiste et non destinées à une utilisation collective, et d'autre part, les courtes citations justifiées par le caractère scientifique ou d'information de l'œuvre dans laquelle elles sont incorporées (art. L. 122-4, L. 122-5 et L. 335-2 du Code de la propriété intellectuelle). Des photocopies payantes peuvent être réalisées avec l'accord de l'éditeur. S'adresser au : Centre français d'exploitation du droit de copie, 3, rue Hautefeuille, 75006 Paris. Tél. : 01 43 26 95 35.

Collection Pratique R
dirigée par **Pierre-André Cornillon**
et **Eric Matzner-Løber**

Département MASS
Université Rennes-2-Haute-Bretagne
France

Comité éditorial

Eva Cantoni

Institut de recherche en statistique
& Département d'économétrie
Université de Genève
Suisse

Pierre Lafaye de Micheaux

Département de Mathématiques
et Statistique
Université de Montréal
Canada

François Husson

Département Sciences de l'ingénieur
Agrocampus Ouest
France

Sébastien Marque

Directeur Département Biométrie
Danone Research, Palaiseau
France

Déjà paru dans la même collection :

Réseaux bayésiens avec R

Jean-Baptiste Denis, Marco Scutati, 2014
ISBN : 978-2-7598-1198-4 – EDP Sciences

Analyse factorielle multiple avec R

Jérôme Pagès, 2013
ISBN : 978-2-7598-0963-9 – EDP Sciences

Psychologie statistique avec R

Yvonnick Noël, 2015
ISBN : 978-2-7598-1736-8 – EDP Sciences

Séries temporelles avec R

Yves Aragon, 2011
ISBN : 978-2-8178-0208-4 – Springer

Régression avec R

Pierre-André Cornillon, Eric Matzner-Løber, 2011
ISBN : 978-2-8178-0184-1 – Springer

Méthodes de Monte-Carlo avec R

Christian P. Robert, George Casella, 2011
ISBN : 978-2-8178-0181-0 – Springer

Vj ku' r ci g' k p v g p v k q p c m { ' i g h v ' d r e p m

REMERCIEMENTS

Cet ouvrage est le fruit d'un lent processus d'évolution s'étalant maintenant sur 15 ans d'enseignement des statistiques et de la psychométrie en cursus de psychologie. Mais sa forme et son contenu doivent beaucoup à un certain nombre de personnes et collègues qui ont parfois initié, souvent stimulé, toujours nourri mes réflexions sur ces questions.

Je pense en particulier à Jean-Michel Petot (Université de Paris X-Nanterre), qui a, sans doute au-delà de ce qu'il souhaitait, encouragé le thésard que j'étais à développer et affiner son goût immodéré des mathématiques et des statistiques, pour traiter de problématiques psychologiques. C'est sans doute là le premier élément de grave dérive mathématique dans mon parcours de psychologue... Je l'ai toujours soupçonné d'être de mèche avec Robert Molimard (Faculté de Médecine de l'Université Paris-Descartes), qui en encadrant mes recherches sur les modèles quantitatifs de l'arrêt du tabac au sein de la Société de tabacologie, m'a poussé sur la même voie. Jacques Juhel (Université de Rennes 2), Paul Dickes et Jean-Luc Kop (Université de Nancy 2), n'ont guère arrangé les choses, en m'invitant à partager avec eux ces moments extrêmement riches et stimulants pour l'esprit, qu'étaient les séminaires fermés de psychologie différentielle à la fin des années 1990, et que sont aujourd'hui les journées MODEVAIIA (*Modélisation de la Variabilité Inter et Intra Individuelle*).

Dans la foulée, je pense que mes collègues lillois ont cru bien faire en me recrutant sur mon premier poste de psychologie et statistiques à l'Université de Lille 3. J'ai notamment appris beaucoup, en termes de rigueur et de rapport à l'écriture mathématique, des collègues statisticiens Françoise Lefèvre (Université de Lille 3), Catherine Trottier (Université de Montpellier 2) et Sébastien Faure (Université de Lille 3). Le programme de statistiques lillois a, depuis, été une source d'inspiration majeure pour la construction de celui de Rennes. La constitution avec ces collègues d'un groupe de travail *Psychométrie et Econométrie*, en coordination avec l'économiste Frédéric Jouneau (Université de Lille 3), rendait déjà difficile tout retour en arrière. Mais était-il alors vraiment nécessaire que l'équipe de Jean-Claude Darcheville (Université de Lille 3) vienne en remettre une couche, en soumettant à mon esprit déjà torturé par ces questions les derniers modèles dynamiques nonlinéaires de l'apprentissage par renforcement ? Aujourd'hui encore, je me pose la question. Que dire de la diligence et de la gentillesse avec laquelle Wijbrand van Schuur (Université de Groningen), David Andrich (University of Western Australia) et James Roberts (Georgia Institute of Technology) ont accepté d'engager des échanges épistolaires fructueux et réguliers sur les modèles du dépliage ? Je reconnais avoir jeté un peu d'huile sur ce feu-là, en distribuant sans modération, à l'occasion de leurs visites en Europe, ce renforçateur primaire qu'est le vin de Corbières...

A mon arrivée à Rennes en 2003, je ne me suis pas méfié des deux attachés d'enseignement qu'on m'avait adjoints pour les enseignements de statistiques : Olivier Le Bohec (Université de Rennes 2) et Bruno Dauvier (Université d'Aix-Marseille). J'ai cru naïvement que leurs formations d'expérimentaliste, pour l'un, de différen-

tialiste, pour l'autre, m'aideraient à retrouver une voie de recherche plus orthodoxe pour un psychologue. Au gré des discussions, je crois que ma passion pour les modèles quantitatifs du comportement n'a en réalité jamais cessé de s'accroître à leur contact, sans que leurs directeurs de recherche respectifs (Jacques Juhel et Eric Jamet, Université de Rennes 2) ne jugent utile de me mettre en garde. On ne me fera pas croire que leurs efforts répétés, appuyés par Fanny De La Haye et Christophe Quaireau (Université de Rennes 2), au sein du groupe TACIT (*Testing Adaptatif de la Compréhension Implicite de Texte*), pour mettre en œuvre des modèles de réponse à l'item sous n'importe quel prétexte, ne font pas partie d'un vaste plan organisé.

Je ne sais combien leur a coûté (financièrement) la participation à ce mouvement d'autres collègues de Rennes, psychologues (Thierry Marivain, Géraldine Rouxel, Frédéric Devinck, Alessandro Guida, Audrey Noël), ou statisticiens (Laurent Rouvière, Mathieu Emily), des ATER en statistiques (qui feraient mieux de boucler leur thèse), et même des étudiants qui, à travers leurs remarques, questions et jeux de données, n'ont pas cessé de faire évoluer ma pratique de la modélisation, et par conséquent le contenu de cet ouvrage, par ricochet.

J'ai pu, longtemps, maladivement en retenir la parution. C'était compter sans l'énergie et l'autorité d'Eric Matzner-Løber (Université de Rennes 2), qui avec Pierre-André Cornillon, a joué un grand rôle dans le simple fait que ce livre vienne à l'existence.

Comme je ne suis pas rancunier, je dirais simplement à tous... *un grand merci*.

Que celle qui m'a accepté dans sa vie me pardonne d'avoir consacré autant de temps à la rédaction de cet ouvrage. En lisant ces lignes, elle pourra se demander si c'est au final tellement de ma faute...

Rennes, octobre 2012

Y.N.

AVANT-PROPOS

A bien des égards, le psychologue apparaît (y compris souvent à ses propres yeux) comme celui qui donne du sens (à une situation, à une conduite). Par on ne sait quel mystère, sa quête de sens semble cependant s'arrêter là où commence son besoin de statistiques. Tout se passe en effet comme si l'enseignement des statistiques en psychologie devait se réduire à la présentation de quelques procédés magiques (un T de Student, un F de Fisher...), dans un ensemble fini de situations problèmes (comparer deux moyennes, deux variances, etc.). Cet enseignement catalogue est aussi ennuyeux pour l'enseignant qu'il l'est pour l'étudiant, autant qu'il est inefficace, car à ne pas comprendre les *modèles* sous-jacents aux statistiques de décision, on court un sérieux risque de les appliquer dans des situations où les attendus du modèle sont simplement absents. Ce risque n'est pas nul, même avec des statistiques très élémentaires et usuelles en psychologie : je ne compte plus le nombre d'articles où l'on utilise des T de Student sur données ordinales, des χ^2 de Pearson sur comptages d'événements non indépendants, etc.

Ces procédés magiques que sont les « statistiques de décision » sont parfois rassemblées dans des opuscules qui ont toujours, peu ou prou, le même titre : *Statistiques appliquées à la psychologie*. Mais comment penser que l'on va aller chercher dans une boîte à outils toute faite les modèles qui donnent du sens aux données particulières d'une étude psychologique ? Cela ne se peut que si l'on a acquis, souvent sans le comprendre, qu'un certain formatage des données permettait de le faire (résumer des données en moyennes et écarts-types de groupes par exemple). Même dans ce cas, c'est une manière de restreindre fortement notre capacité théorique que de préformater les données pour qu'un ensemble réduit d'outils statistiques puissent s'appliquer. C'est bien le modèle qui doit s'adapter à la situation et non l'inverse. On voit ainsi souvent des variables numériques transformées en variables catégorisées... pour pouvoir utiliser l'analyse de la variance !

Cet ouvrage ne parle pas de statistiques appliquées à la psychologie. Il parle de psychologie et de modèles de probabilités que l'on peut extraire de la situation elle-même pour donner du sens aux données. Le point de départ est d'abord et avant tout la situation psychologique (étude de terrain ou expérience de laboratoire), dont on élabore le modèle, en même temps qu'on élabore la théorie psychologique de ce qu'on observe. Il s'agit donc de *psychologie statistique*, à proprement parler, et il n'y a pas dans cette approche de distinction entre théorisation et modélisation, entre psychologie et statistiques, car l'un et l'autre se mènent conjointement. On parlera aussi de *psychométrie* pour désigner cette approche spécifique de la psychologie par modèles de probabilité, sur les mécanismes de réponse, de perception ou de conduite. Cette approche est très ancienne en psychologie, et trouve ses racines notamment dans la psychophysique et dans la modélisation factorielle des compétences cognitives. On verra qu'à renverser ainsi la perspective, partant du psychologique pour en extraire le modèle, il arrivera qu'incidemment on recouvre des modèles connus (modèles de groupe tels qu'ANOVA ou T de Student par exemple), mais pas nécessairement. Cette approche par construction de modèle

permettra aussi d'apercevoir plus rapidement que telle ou telle statistique connue est en réalité inapplicable pour la situation qu'on étudie.

L'approche catalogue stigmatisée ci-dessus est tentante pour l'enseignant(e) de statistiques en psychologie, car les pressions sur lui ou sur elle sont nombreuses pour que son souci pédagogique de développement de la compréhension des modèles sous-jacents passe simplement à la trappe. Et les pressions les plus fortes ne viennent pas nécessairement des étudiants. Il est vrai que les tuteurs de recherche pourront quelque temps cultiver l'illusion qu'en procédant ainsi, leurs étudiants deviennent plus rapidement « autonomes » dans le traitement de leurs données. Le fait qu'en procédant ainsi on incite à une forme de pensée magique (les étudiants ne comprennent guère ce qu'ils font) qui serait jugée inacceptable dans n'importe quelle autre discipline, n'est pas encore le plus lourd tribut à payer : en faisant cela, on ferme la porte à un enseignement de Master digne de ce nom. Si des notions théoriques importantes ne sont pas posées en Licence (combinatoire, algèbre des événements, probabilités élémentaires), la formation aux modèles avancés en Master n'est plus possible, sans empiéter sur l'espace pédagogique déjà restreint pour introduire à la hâte toutes ces notions. Comment initier l'étudiant de psychologie à l'analyse fine des tables de contingence à nombre quelconque d'entrées, ou à l'analyse de la variance en plans factoriels complexes, sans avoir au préalable présenté les lois de probabilité qui en sont le cœur (loi binomiale, loi multinomiale, loi normale) ? Comment former aux modèles d'équations structurales, devenus incontournables y compris dans des études appliquées, sans avoir appris à manier l'algèbre des variances et covariances ? En croyant gagner du temps, on perd en efficacité pédagogique.

Si l'on peut en outre présenter sur des problèmes simples (comparaisons de deux proportions par exemple) une approche de la modélisation qui restera valable sur les problèmes de dimension supérieure (comparaison de plus de deux proportions ou comparaisons de distributions catégorisées), l'effort d'apprentissage initial devient rapidement rentable pour l'étudiant. On verra que dans cet ouvrage, tout en présentant les outils classiques de sélection de modèles (valeur p), on mettra l'accent sur une unique statistique de décision, applicable dans tous les cas (le *facteur de Bayes*). Un effort particulier a été fait pour fournir des expressions du facteur de Bayes dans tous les types de problèmes, qui restent calculables avec une machine de poche.

Ce livre donne ainsi une place non négligeable à l'approche dite « bayésienne ». On cherche dans cette approche à calculer la probabilité qu'un modèle soit vrai pour des données particulières, par opposition avec l'outil classique pearsonien qu'est la valeur p , qui calcule la probabilité des données, *d'après un modèle choisi comme arbitrairement vrai* (« l'hypothèse nulle »). La signification de la probabilité n'est évidemment pas du tout la même dans les deux cas. Elle est particulièrement simple à interpréter dans le cas bayésien (on garde le modèle le plus probablement vrai) et délicate à manipuler dans le cas de la valeur p (que faire du modèle quand la probabilité des données est assez élevée *quand on le suppose vrai* ?). Les outils et modèles bayésiens ont littéralement explosé dans la littérature statistique interna-

tionale, tant les avantages de l'approche sont séduisants. Les voix sont nombreuses aussi en psychologie, qui insistent sur ces avantages, et il n'est plus possible de passer cela sous silence dans un manuel actuel de psychologie statistique. Il est très surprenant que ces outils ne soient pas plus utilisés en France, où assez tôt des chercheurs ont déployé de gros efforts pour divulguer les principes bayésiens, développer des logiciels spécialisés, et mettre en avant tous les bénéfices scientifiques et pédagogiques de la méthode. L'apport d'Henri Rouanet et de Dominique Lépine, dès les années 1970 (Rouanet & Lépine, 1976), et avec eux de Bruno Lecoutre (Lecoutre, 1984, 1996) est à cet égard tout à fait remarquable, précédant de 40 ans les appels récents au passage au bayésien (Wagenmakers *et al.*, 2010, 2011 ; Kruschke, 2010). Nous mettons à disposition du lecteur, sous forme de bibliothèques R avec interface graphique, certains des outils mis en avant par ces auteurs.

Le besoin de statistiques et de modèles de probabilités est croissant dans la psychologie contemporaine, et ce besoin n'est plus cantonné aux disciplines expérimentales. L'auteur de ces lignes, clinicien quantitativiste de formation, se souvient encore de ses premières expériences de psychologue en centre de soins pour toxicomanes, en 1997. La méconnaissance de la méthodologie élémentaire et des statistiques rendait les praticiens (tant médecins que psychologues), simplement incapables de tirer profit de l'abondante littérature américaine sur la délivrance de méthadone, qui ne faisait que démarrer en France. Il a été assez difficile de convaincre, données chiffrées à l'appui, que la politique de délivrance de méthadone qui venait juste d'être mise en place était sous-dosante et qu'en agissant ainsi on programmait littéralement la rechute (constat rapidement fait chaque lundi matin). De la même façon, devant la nécessité de sélectionner les patients les plus motivés (le nombre de place étant limité), il n'a guère été possible de convaincre que la motivation au changement et sa maturation dans les conduites addictives sont mesurables et permettent de faire de véritables prédictions comportementales (Noël, 1999, 2009). Il ne serait pas difficile de trouver des exemples équivalents dans le domaine du recrutement, de la santé, de l'éducation ou de la criminologie. Il est grand temps que les psychologues de terrain prennent conscience de ce que le déficit de formation statistique et méthodologique a un véritable impact sur la qualité de la pratique, ne serait-ce que pour se tenir informé des avancées dans la littérature. L'effort de formation pour le psychologue de terrain ou le chercheur n'est certes pas négligeable, mais tout à fait réalisable (en session de deux ou trois jours par an par exemple).

Cet ouvrage couvre un programme de Licence (années 1, 2 et 3). Typiquement, les chapitres 1 et 2 constituent le programme de L1, sur un semestre. Les contenus de ces deux chapitres correspondent assez bien avec un programme de statistiques dites « descriptives » dans une approche traditionnelle, mais introduisent des distinctions sur les niveaux de mesure qui étendent la théorie classique de la mesure de Stevens, pour permettre l'introduction des modèles de distributions à partir du chapitre 7. Ils sont une préparation déguisée à la distinction des familles de distribution de probabilité, qui sont introduites dans les chapitres suivants. Les chapitres 3 à 8 (jusqu'aux modèles binomiaux à un paramètre et en introduisant

la méthode bayésienne de façon élémentaire) correspondent à un programme de L2, sur un semestre. Les chapitres 8 (à partir des modèles binomiaux à deux paramètres) à 15 correspondent à un programme de L3, sur deux semestres. Compte tenu des contraintes de temps, variables selon les universités, il n'est certainement pas possible d'amener en cours tous les contenus de cet ouvrage, qui fournit des explications et des approfondissements qui pourront être négligés à un premier niveau de lecture. La perspective bayésienne peut être présentée de façon simplifiée, par exemple en n'utilisant que l'approximation *BIC* du facteur de Bayes pour la sélection de modèles.

L'ensemble des procédures peut être mis en œuvre avec les bibliothèques **AtelierR** et **R2STATS** pour **R**. Le logiciel **R** est devenu en quelques années un outil incontournable, adopté par de nombreuses écoles et universités dans le monde. La gratuité du logiciel n'est pas la seule raison de son succès : il a aussi une dynamique de développement excellente, avec un accès aisé à de nombreux documents d'aide en ligne et une liste de diffusion extrêmement active. Pour accélérer l'apprentissage des étudiants, nous avons doté ces deux bibliothèques d'une interface graphique en **GTK** (*Gnome Toolkit*), permettant l'installation sur systèmes Linux, Mac ou Microsoft Windows¹. On trouvera également sur la page web de **R2STATS** des exercices corrigés de façon détaillée, qui mettent en œuvre toutes les procédures décrites dans l'ouvrage sur des données réelles de psychologie expérimentale et appliquée.

1. Voir les instructions d'installation sur la page web : <http://yvonnick.noel.free.fr/r2stats>

Table des matières

Remerciements	vii
Avant-propos	ix
1 Description sur une variable	1
1.1 Processus de mesure	1
1.2 Structure de la mesure	4
1.2.1 Variable qualitative nominale	5
1.2.2 Variable qualitative ordinale	7
1.2.3 Variable quantitative discrète	10
1.2.4 Variable quantitative continue	16
1.2.5 Notion de rapport et d'intervalle	21
1.3 Synthèse	22
2 Description de liaison	23
2.1 Lien entre une variables numérique et une variable catégorisée	23
2.1.1 Plans d'analyse	23
2.1.2 Comparaison des indices de centralité	24
2.1.3 Comparaison des dispersions	27
2.1.4 Comparaison de distributions en blocs	32
2.1.5 Situer un individu	35
2.2 Lien entre deux variables numériques	38
2.2.1 Covariance empirique	38
2.2.2 Coefficient de corrélation de Bravais-Pearson	40
2.2.3 Corrélation et causalité	45
2.3 Lien de deux variables catégorisées	46
2.3.1 Cotes, rapport de cotes et lograpports de cotes	46
2.3.2 Rapport de vraisemblance	48
3 Algèbre des événements	51
3.1 Notion d'ensemble	52
3.2 Intersection et union	54

3.3	Algèbre sur les ensembles	55
3.3.1	Ordre des opérateurs	55
3.3.2	Distributivité	56
3.3.3	Lois de De Morgan	56
3.3.4	Tableau de synthèse	57
3.4	Application : le jeu de la sélection de cartes	57
4	Calcul des probabilités	59
4.1	Notion intuitive	59
4.1.1	Probabilité connue	60
4.1.2	Probabilité inconnue	61
4.2	Probabilité conjointe, conditionnelle et marginale	61
4.3	Règles de calcul	64
4.3.1	Formules de Bayes	64
4.3.2	Loi du produit	65
4.3.3	Loi de l'addition	66
4.3.4	Théorème des probabilités totales	68
4.3.5	Tableau de synthèse	69
4.4	Dénombrements	69
4.4.1	Permutations	70
4.4.2	Arrangements	70
4.4.3	Combinaisons	71
4.4.4	Répartition en classes identifiées	72
4.4.5	Tableau de synthèse	73
4.5	Probabilités sur un ensemble non dénombrable	73
4.5.1	Simulation d'un processus uniforme	74
4.5.2	Probabilité ponctuelle dans une loi continue	76
4.5.3	Construction de la densité uniforme	77
4.5.4	Notion d'intégrale	78
4.6	Applications	79
4.6.1	Sally Clark est-elle coupable ?	79
4.6.2	Sensibilité et spécificité des tests psychologiques	81
5	Espérances et moments	87
5.1	Espérance mathématique et théorie des jeux	88
5.2	Variance et gestion des risques	93
5.3	Algèbre des covariances	97
5.4	Application : l'analyse factorielle	99
5.5	Tableaux de synthèse	104

6	Notion de modèle	105
6.1	La statistique inférentielle	105
6.2	Démarche d'hypothèse	106
6.3	Un exemple neuropsychologique	108
6.3.1	Modélisation	108
6.3.2	La valeur p	110
6.3.3	Seuil de décision et erreur de type I	110
6.3.4	Erreur de type II	111
6.4	Probabilité des données ou probabilité du modèle	112
6.4.1	Le facteur de Bayes	112
6.4.2	Probabilités a posteriori des modèles	115
7	Modèles binomiaux	117
7.1	Modèles à un paramètre	117
7.1.1	Test d'hypothèse	117
7.1.2	Estimation d'une probabilité inconnue	125
7.2	L'approche bayésienne	128
7.2.1	Historique critique sur la démarche par valeur p	129
7.2.2	Facteur de Bayes pour l'inférence sur une probabilité	132
7.2.3	Estimation bayésienne d'une probabilité	138
7.2.4	Le critère d'information bayésien (<i>BIC</i>)	145
7.3	Modèles à deux paramètres	149
7.3.1	La théorie de la dissonance cognitive	149
7.3.2	Facteur de Bayes pour la comparaison de deux probabilités	151
7.3.3	Approximation par la différence des <i>BIC</i>	152
7.4	Modèles à trois paramètres	156
7.4.1	Antécédents d'abus sexuels et délinquance	156
7.4.2	Facteur de Bayes pour la comparaison de trois probabilités	158
7.4.3	Approximation par la différence des <i>BIC</i>	161
7.5	Modèles généraux et factoriels	164
7.5.1	Etiquetage social négatif et « compliance »	164
7.5.2	Modélisation	165
7.5.3	Comparaison de tous les modèles possibles	167
8	Modèles multinomiaux	169
8.1	Construction de la loi multinomiale	169
8.2	Modèles sur une variable catégorisée	170
8.2.1	Comparaison à une distribution multinomiale fixée	170
8.2.2	Comparaison à une alternative multinomiale structurée	172
8.3	Modèles sur deux variables catégorisées	172
8.3.1	Comparaison de distributions multinomiales	173
8.3.2	Examen des liaisons locales	175
8.3.3	Modèle de l'indépendance	176
8.3.4	Contrastes dans une table de contingence	179

8.4	Modèles sur trois variables catégorisées	182
9	Modèles gaussiens	187
9.1	Construction de la loi normale	187
9.1.1	La loi normale	187
9.1.2	Changement d'échelle et d'origine d'une loi normale	193
9.2	Inférence sur une moyenne : variance connue	196
9.2.1	Distribution d'une moyenne d'échantillon	196
9.2.2	Comparaison à une valeur théorique	199
9.2.3	Approche bayésienne	201
9.3	Inférence sur une variance	206
9.3.1	Construction d'un estimateur de variance sans biais	207
9.3.2	La loi de χ^2	209
9.3.3	Comparaison à une valeur théorique	212
9.3.4	Approche bayésienne	213
9.4	Inférence sur une moyenne : variance inconnue	222
9.4.1	La loi de Student	223
9.4.2	Test de comparaison à une norme	224
9.4.3	Analyse d'une différence test-retest	225
9.4.4	Inférence sur la taille de l'effet	227
9.4.5	Approche bayésienne	234
9.5	Inférence sur deux moyennes d'échantillons indépendants	248
9.5.1	Statistique de Student	250
9.5.2	Mesure de la taille d'effet (g de Hedges)	255
9.5.3	Approche bayésienne	255
9.6	Inférence sur des variances d'échantillons indépendants	259
9.6.1	Construction d'un modèle de groupe (ANOVA)	263
9.6.2	Calcul pratique sous R/R2STATS	273
9.6.3	Test des hypothèses de l'ANOVA	277
9.6.4	Comparaisons spécifiques (contrastes)	282
9.6.5	Approche bayésienne	288
A	Compléments techniques	299
A.1	Les fonctions exponentielle et logarithme	299
A.2	Maximisation d'une vraisemblance binomiale	302
A.3	La loi Beta-binomiale	303
A.4	Formules exactes du facteur de Bayes	305
A.5	Maximisation d'une vraisemblance gaussienne	311
A.6	Lois <i>a posteriori</i> sur les paramètres d'une loi normale	312
	Bibliographie	314
	Index	321

Chapitre 1

Description sur une variable

1.1 Processus de mesure



S.S. Stevens (1906-1973)

En 1946, le psychologue américain Stanley Smith Stevens publia un article (Stevens, 1946) dans la revue *Science*, où il proposait une théorie synthétique des niveaux de mesure en science. Toutes les mesures n'ont pas les mêmes propriétés ni la même richesse, et les distinctions qu'il a posées entre quatre niveaux de mesure vont nous servir de base dans ce chapitre. Nous irons cependant un peu au-delà, pour établir une typologie étendue permettant dans les parties ultérieures de cet ouvrage de lier directement les types de mesures et les types de modèles statistiques. A chaque type et niveau de mesure est associé un ensemble d'opérations statistiques appropriées.

Un exemple

Un psychologue clinicien travaillant auprès de patients toxicomanes recueille systématiquement lors du premier entretien les 6 informations suivantes :

1. le sexe ;
2. la motivation au travail de sevrage, rapportée par le patient sur une échelle à quatre réponses possibles : (a) pas du tout motivé, (b) plutôt pas motivé, (c) plutôt motivé, (d) motivé ;
3. les scores du patient sur plusieurs échelles de personnalité ;
4. la température corporelle (la fièvre est l'un des symptômes du manque) ;
5. un score d'anxiété fondé sur le comptage des manifestations somatiques de l'anxiété ;
6. son âge.

Au fil des années, il a enregistré ces données pour 750 patients. Il présente ces informations dans un tableau qui a la forme suivante :

N° Patient	Sexe	Motivation	Anxiété	Age	...
1	H	a	2	23	...
2	F	c	10	27	...
3	H	d	12	30	...
...

Chaque patient est ainsi caractérisé par un profil d'attributs qui lui est propre. Les patients apparaissent en lignes et les classes d'information qui les décrivent apparaissent en colonnes. Dans une même colonne, on peut voir apparaître des valeurs différentes selon les sujets pour une même classe d'information. Cette mise en relation d'un ensemble de sujets avec un ensemble de descripteurs (appelés aussi *variables*) est l'opération de base du travail statistique.

A partir de cette procédure d'enregistrement d'informations, le psychologue va pouvoir se poser plusieurs types de questions :

1. sur chaque variable prise isolément : l'échantillon est-il composé plutôt de femmes ou d'hommes ? Leur degré de motivation est-il plutôt élevé ou faible ? etc. Il est intéressant de pouvoir résumer les caractéristiques saillantes de son échantillon,
2. sur les liens qui unissent certaines variables entre elles : la motivation est-elle différente chez les hommes et les femmes ? Ce qui revient à dire : y a-t-il un lien entre sexe et motivation au sevrage ? Les patients plus âgés sont-ils plus motivés au sevrage ? etc. Il est intéressant de connaître les associations entre caractéristiques, pour pouvoir ultérieurement leur donner du sens dans une *théorie psychologique*.

Ces deux types de questions renvoient à deux grands axes du travail statistique du psychologue : décrire, résumer et modéliser ce que des individus *ont en commun*, décrire, résumer et modéliser ce en quoi des individus *diffèrent*.

Ces deux aspects sont toujours manipulés conjointement en psychologie, car il s'agit d'une dialectique irréductible (l'un ne va pas sans l'autre) mais peuvent se trouver différemment pondérés selon les champs. Le psychologue peut selon les cas décider d'insister sur les principes communs du comportement humain, et dégager des *lois*, ou s'attacher à mettre en avant la *différence individuelle*. Les statistiques fournissent des outils de description pour l'un et l'autre mouvement de la pensée, qui ne seront donc pas pour nous dissociés. Sur des données concrètes, nous chercherons à poser un modèle statistique.

Définition 1.1 (Modèle statistique)

On appelle modèle statistique une représentation simplifiée de la réalité qui rend compte à la fois de l'homogénéité et de la variabilité d'un phénomène psychologique au moyen de probabilités.

La notion de probabilité sera présentée plus loin. L'objectif des deux premiers chapitres de ce manuel est d'étudier les procédures de description et de résumé qui permettent de répondre à ce type de question d'un point de vue intuitif. Les procédures naturelles vues dans cette partie serviront d'ancrage pour la compréhension des chapitres suivants, consacrés à l'inférence et à la modélisation probabiliste.

Terminologie et notations

Nous fixons dans cette partie certains points de vocabulaire et conventions de notation. La structure en lignes et colonnes permet de présenter commodément ce qu'on appelle des *variables*. La notion de variable intègre deux aspects :

1. elle est définie en référence à une population d'individus pour laquelle elle est pertinente. Ces « individus statistiques » ne sont pas nécessairement des personnes. On peut étudier des unités statistiques qui sont des groupes de sujets (des familles, des classes d'école, des villes...) ou des objets. Les expressions « individus statistiques », « unités statistiques » et « observations statistiques » sont donc génériques et synonymes ;
2. à chaque individu, elle affecte une *valeur* (ou modalité), dans un ensemble de valeurs possibles qui lui est propre. Elle vient ainsi qualifier ou quantifier des caractères ou attributs propres à l'individu.

Une variable est donc une association qu'on établit entre un ensemble d'individus pour laquelle elle est pertinente et un ensemble de valeurs possibles qui lui est propre. Cette opération d'association élémentaire un-pour-un entre deux ensembles s'appelle *application* en mathématiques.

Définition 1.2 (Application)

On appelle application une relation entre un ensemble de départ (appelé domaine de définition) et un ensemble d'arrivée (appelé image), qui à chaque élément de l'ensemble de départ associe un élément et un seul de l'ensemble d'arrivée.

Prendre une mesure comportementale correspond implicitement à une opération d'association de ce type.

Définition 1.3 (Variable)

Une variable est une application d'un ensemble d'individus statistiques I vers un ensemble de valeurs observables U .

Il est courant de donner un nom symbolique aux variables. Nous garderons cette convention courante en statistique de dénommer les variables par des lettres majuscules (par exemple X) et par une lettre minuscule (par exemple x_i , la valeur prise par l'individu statistique i pour cette variable X , avec une numérotation des sujets par $i = 1, \dots, N$). A chaque individu de I est associé par la variable X une valeur observable de U_X , ce que l'on peut écrire symboliquement :

$$\begin{aligned} X : I &\rightarrow U_X \\ s_i &\rightarrow x_i. \end{aligned}$$

Vj ku'rci g'kpvgpvkqpcmf 'igh'dnc pm

Crédits photographiques

- page 1 : photo de S. Stevens. Harvard University Archives, HUP Stevens, S.S. (3a).
- page 81 : photo de H. Murray. Harvard University Archives, HUP Murray, H. (3).
- page 87 : photo de John von Neumann. This information has been authored by an employee or employees of the Los Alamos National Security, LLC (LANS), operator of the Los Alamos National Laboratory under Contract No. DE-AC52-06NA25396 with the U.S. Department of Energy.
- page 142 : *Portraits de suspects*, Busey & Loftus, 2006. Reprinted from Trends in Cognitive Psychology, 11 (3), Busey, T.A. and Loftus, G.R., Cognitive Science and the Law, 111-117, Copyright (2007), with permission from Elsevier.
- page 145 : photo de G. Schwarz Archives of the Mathematisches Forschungsinstitut Oberwolfach.