

ÉCHANTILLONS TÉLÉPHONIQUE COMMERCIAUX ET DISTRIBUTION DES CODES D'ÉTAT FINAL PARTIELS

Claude Comeau, Peter Mariolis, Ph.D.¹

RÉSUMÉ

Dans les enquêtes téléphoniques, la distribution des codes d'état final repose habituellement sur un échantillon complet. Toutefois, les échantillons téléphoniques provenant de fournisseurs d'échantillons commerciaux renferment des sous-ensembles identifiables d'enregistrements qui font varier énormément la probabilité d'obtenir un code particulier. En pareil cas, les distributions partielles des codes d'état final pourraient varier indépendamment de celle fondée sur l'échantillon complet. À partir de l'enquête menée en 2000 par le Behavioral Risk Factor Surveillance System (BRFSS), nous examinons dans quelle mesure la distribution des codes d'état final, choisis par territoire pour différents sous-ensembles d'enregistrements, est en corrélation avec celle fondée sur l'échantillon complet. Les corrélations ont tendance à être importantes mais, dans certains cas, elles ne le sont pas, d'où l'opportunité d'examiner les distributions partielles des codes d'état final.

MOTS-CLÉS : Enquêtes téléphoniques; Distribution des codes d'état final; Qualité des données.

1. INTRODUCTION

Dans les enquêtes téléphoniques, la distribution des codes d'état final repose habituellement sur un échantillon complet. Toutefois, les échantillons téléphoniques provenant de fournisseurs d'échantillons commerciaux renferment généralement des sous-ensembles identifiables d'enregistrements qui font varier énormément la probabilité d'obtenir un code particulier. En pareil cas, les distributions partielles (basées sur des sous-échantillons) des codes d'état final pourraient varier indépendamment de celle fondée sur l'échantillon complet. Il existe à cet égard deux distinctions importantes : l'une entre les numéros de téléphone inscrits et non inscrits, l'autre entre les numéros de téléphone d'une banque d'un et plus et ceux d'une banque de zéro.

Les fournisseurs d'échantillons commerciaux ont accès à des bases de données permettant d'identifier les numéros de téléphone de ménages inscrits. Les autres numéros de téléphone ne sont pas inscrits. Une banque de cent est un ensemble de cent numéros de téléphone dont l'indicatif régional, le préfixe et les deux premiers chiffres du suffixe sont les mêmes. D'après les bases de données renfermant les numéros inscrits, les fournisseurs peuvent déterminer le nombre de numéros de téléphone de ménages inscrits dans n'importe quelle banque de cent. Les numéros de téléphone d'une banque d'un et plus sont des numéros (inscrits ou non) d'une banque de cent renfermant au moins un numéro de téléphone de ménage inscrit, alors que ceux d'une banque de zéro sont des numéros (dont aucun n'est inscrit) d'une banque de cent qui ne renferme aucun numéro de téléphone de ménage inscrit. On peut donc établir une distinction entre les numéros de téléphone de ménages inscrits, les numéros non inscrits d'une banque d'un et plus et les numéros de téléphone d'une banque de zéro.

Nous examinons trois taux par sous-échantillon : le taux de réponse CASRO (Council of American Survey Research Organizations), le taux d'identification de ménages et le taux d'achèvement par ménage (CASRO, 1982). Le taux de réponse CASRO est le nombre d'interviews achevées divisé par une estimation

¹ Claude Comeau, Comeau Associates, 177 Ball Hill Road, Milford, New Hampshire, USA 03055
Peter Mariolis, Centers for Disease Control and Prevention (MS K66), 4770 Buford Hwy, NE, Atlanta,
Georgia, USA 30341-3717

du nombre de ménages admissibles compris dans l'échantillon². Le taux d'identification de ménages est le nombre de ménages identifiés divisé par le nombre total d'enregistrements compris dans l'échantillon. Le taux d'achèvement par ménage est le nombre d'interviews achevées divisé par le nombre de ménages identifiés compris dans l'échantillon. Les taux d'identification de ménages et d'achèvement par ménage correspondent aux deux aspects essentiels de la participation à une enquête selon la distinction établie par Groves et Couper (1998) : le contact et la collaboration. L'utilité des mesures fondées sur les codes d'état final partiels ne se limite pas à ces trois taux; nous mentionnons également des utilisations réelles ou possibles d'autres mesures fondées sur les codes d'état final.

La présente communication utilise des données d'une enquête menée en 2000 par le Behavioral Risk Factor Surveillance System (BRFSS) pour étudier les liens entre les taux global et partiel dans les trois cas : taux CASRO, taux d'identification de ménages et taux d'achèvement par ménage. À partir des données de la présente communication, nous tentons de répondre à une seule question de base : la distribution des codes d'état final partiels présente-t-elle une variation (linéaire) indépendante de celle des codes d'état final globaux?

2. DONNÉES ET MÉTHODES

L'enquête du BRFSS est menée conjointement par le Centers for Disease Control and Prevention (CDC) et par les départements de la Santé des 50 États américains, du district fédéral de Columbia et de Porto Rico³. Il s'agit d'une enquête téléphonique (habituellement) mensuelle servant à cerner, dans chaque territoire, la prévalence de comportements liés aux maladies chroniques et les pratiques de santé préventive parmi la population civile hors établissement âgée de 18 ans et plus. Les sujets liés à la santé comprennent l'alimentation, le tabagisme, l'activité physique et la consommation d'alcool. Le CDC coordonne l'élaboration d'un ensemble de questions de base qui sont posées par chaque territoire et des ensembles uniformisés de questions portant sur des sujets précis que les territoires peuvent choisir de poser; de plus, chaque territoire est libre de poser d'autres questions de son choix. Le CDC coordonne également l'élaboration de normes concernant les plans d'échantillonnage et les méthodes de collecte des données et offre une assistance technique aux territoires. La base de sondage du BRFSS comprend tous les numéros de téléphone de types NXX 00, 50, 51, 52, et 54, y compris ceux d'une banque de zéro. Dans les cas non résolus, les lignes directrices du BRFSS prescrivent au plus 15 rappels répartis entre les jours de semaine, les soirs de semaine et le week-end. On choisit au hasard un seul adulte par ménage admissible; les interviews par personne interposée ne sont pas autorisées. Les territoires se chargent de la collecte des données. En 2000, 36 territoires ont confié la collecte des données en sous-traitance à des organismes de recherche commerciaux ou universitaires; dans les 16 autres, c'est un service du département de la santé qui a procédé à la collecte. Une fois les données recueillies et soumises à une première vérification, on les envoie au CDC. Ce dernier poursuit la vérification et, à la fin de chaque année, pondère les données et les retourne aux territoires, accompagnées de divers rapports. Le CDC met ensuite l'ensemble de données agrégées à la disposition du public. On trouvera plus de renseignements sur le BRFSS à l'adresse <http://www.cdc.gov/nccdphp/brfss>.

Les données du BRFSS sont publiées dans un ensemble de données annuelles, mais chaque territoire procède chaque mois à la collecte des données. (Il y a des exceptions : le Michigan recueille des données chaque trimestre et certains territoires le font aux quatre mois. De plus, la plupart des enquêtes sont menées en un seul mois; à l'occasion, cependant, la période d'enquête peut s'étendre au mois suivant ou, rarement, au-delà.) On peut donc considérer que les données recueillies chaque mois dans chaque territoire sont le fruit d'une enquête distincte. Telle est la démarche adoptée dans la présente communication.

² On estime le nombre de ménages admissibles en supposant que la proportion de ménages admissibles parmi les enregistrements dont on ignore le statut (pas de réponse ou ligne occupée) égale la proportion de ménages admissibles parmi les enregistrements dont on connaît le statut. Le nombre obtenu à la suite de ce calcul est ajouté au nombre de ménages admissibles identifiés compris dans l'échantillon.

³ Dans la suite du texte, on entend par « territoire » les 50 États américains, le district fédéral de Columbia et Porto Rico.

Les données utilisées pour notre étude comprennent 1 674 110 enregistrements provenant de 49 territoires (le district fédéral de Columbia et les 50 États américains, sauf le Minnesota et le Wisconsin) qui ont utilisé en 2000 un plan d'échantillonnage à partir de listes. Les enregistrements individuels ont été agrégés en 564 enquêtes mensuelles, ce qui a donné un maximum de 12 enregistrements par territoire. Le plan d'échantillonnage stratifie les numéros de téléphone selon qu'ils sont compris dans une banque d'un et plus ou dans une banque de zéro. En général, les numéros de téléphone d'une banque d'un et plus sont échantillonnés quatre fois plus que ceux d'une banque de zéro. Les enregistrements échantillonnés mensuellement sont produits chaque trimestre. On mesure la composition d'un échantillon selon les pourcentages de numéros inscrits, de numéros non inscrits d'une banque d'un et plus et de ceux d'une banque de zéro. La figure 1 présente les formules de calcul du taux CASRO, du taux d'identification de ménages et du taux d'achèvement par ménage en fonction des codes d'état final utilisés par le BRFSS.

Notre stratégie analytique vise à déterminer de façon distincte, pour chaque territoire, les coefficients de corrélation de Pearson entre chaque mesure globale et ses différentes mesures partielles. Cette stratégie neutralise les écarts entre les équipes de collecte, les caractéristiques de la population et d'autres facteurs dans la mesure où ils sont liés aux écarts entre territoires. Chaque coefficient de corrélation est fondé sur 4 à 12 enregistrements. Les principales données présentées ici sont les distributions de ces coefficients de corrélation pour chaque mesure des codes d'état final.

Figure 1. Formules de calcul des mesures des codes d'état final selon les codes de règlement définitif du BRFSS pour 2000	
Taux de réponse CASRO	
$\frac{01}{\left[(01 + 02 + 07 + 09) + \frac{(01 + 02 + 07 + 09)}{(01 + 02 + 07 + 09) + (03 + 05 + 06 + 08 + 11)} \times (04 + 10) \right]}$	
Taux d'identification de ménages	
$\frac{(01 + 02 + 06 + 07 + 08 + 09 + 11)}{(01 + 02 + 03 + 04 + 05 + 06 + 07 + 08 + 09 + 10 + 11)}$	
Taux d'achèvement par ménage	
$\frac{01}{(01 + 02 + 06 + 07 + 08 + 09 + 11)}$	
Codes de règlement définitif du BRFSS	
01 Interview achevée	07 Répondant choisi absent pendant la période d'interview
02 Interview refusée	08 Barrière linguistique
03 Numéro hors service	09 Interview achevée dans le cadre du questionnaire
04 Pas de réponse	10 Ligne occupée
05 Pas une résidence privée	11 Répondant incapable de communiquer à cause d'une incapacité physique ou mentale
06 Pas de répondant admissible à ce numéro	

3. RÉSULTATS

Pour l'ensemble des territoires et des mois, une moyenne de 21,0 % des enregistrements échantillonnés consiste en numéros de téléphone inscrits (étendue: 8,7 % à 32,9 %) (figure 2). La moyenne par territoire va de 10,4 % à 30,3 %, la moyenne globale étant de 21,0 %. L' étendue moyenne du pourcentage de numéros inscrits dans un territoire est de 3,5 points de pourcentage (et varie : 1,4 à 5,8 points). Comme l'échantillon a été stratifié par banque de numéros de téléphone, la distribution des numéros de téléphone d'une banque d'un et plus et d'une banque de zéro par territoire et par mois varie peu d'un mois à l'autre. Ainsi, la variabilité de la distribution des numéros de téléphone non inscrits d'une banque d'un et plus par

territoire et par mois est presque identique à celle de la distribution des numéros de téléphone inscrits par territoire et par mois.

Les corrélations entre les taux de réponse CASRO global et partiel de numéros inscrits vont de 0,14 à 0,99 (médiane : 0,89); les corrélations pour tous les territoires, sauf deux, sont de 0,59 et plus (tableau 1). Les corrélations entre les taux CASRO global et partiel de numéros non inscrits d'une banque d'un et plus sont plus faibles et plus variables : elles vont de -0,14 à 0,98 (médiane : 0,78); quatre territoires présentent des corrélations comprises entre -0,14 et 0,24. Les corrélations entre les taux CASRO global et partiel de numéros d'une banque de zéro sont encore plus faibles et plus variables : elles vont de -0,52 à 0,76 (médiane : 0,12) et, contrairement aux autres distributions, celle des corrélations d'une banque de zéro n'est pas très asymétrique.

Les corrélations entre les taux global et partiel d'identification de ménages avec numéros inscrits vont de -0,05 à 0,999 (médiane : 0,74) (tableau 2). Les corrélations entre les taux global et partiel d'identification de ménages avec numéros non inscrits d'une banque d'un et plus sont plus importantes et moins variables : si elles vont de -0,63 à 0,97 (médiane : 0,78), l'intervalle interquartile est de 0,16 (0,24 pour les enregistrements inscrits) et tous les territoires, sauf deux, présentent des corrélations de 0,35 et plus. Les corrélations entre les taux global et partiel d'identification de ménages d'une banque de zéro sont encore plus faibles et plus variables : elles vont de -0,43 à 0,87 (médiane : 0,24) et, contrairement aux autres distributions, celle des corrélations d'une banque de zéro n'est pas très asymétrique.

Les corrélations entre les taux global et partiel d'achèvement par ménage avec numéros inscrits vont de 0,25 à 0,99 (médiane : 0,91); dans cinq territoires, les corrélations sont d'au plus 0,61 (tableau 3). Les corrélations entre les taux global et partiel d'achèvement par ménage non inscrits d'une banque d'un et plus sont plus faibles et plus variables : elles vont de -0,06 à 0,98 (médiane : 0,76); quatre territoires présentent des corrélations comprises entre -0,06 et 0,28. Les corrélations entre les taux global et partiel d'achèvement par ménage d'une banque de zéro sont encore plus faibles et plus variables : elles vont de -0,53 à 0,81 (médiane : 0,20) et, contrairement aux autres distributions, celle des corrélations d'une banque de zéro n'est pas très asymétrique.

4. OBSERVATIONS

Les résultats de notre étude montrent qu'en général, la distribution des codes d'état final pour les numéros inscrits et non inscrits d'une banque d'un et plus est grandement corrélée avec celle des codes d'état final globaux correspondants mais que, dans certains cas, elle ne l'est pas. On peut en déduire qu'il convient sans doute de considérer séparément la distribution des codes d'état final pour les numéros inscrits et celle des numéros non inscrits d'une banque d'un et plus. La situation est plus claire dans le cas des numéros d'une banque de zéro : la plupart des territoires présentent des associations faibles ou modérées entre la distribution des codes d'état final d'une banque de zéro et le taux global, de sorte que le calcul de distributions distinctes des codes d'état final pour les numéros d'une banque de zéro donne généralement des distributions des codes d'état final pour les numéros d'une banque de zéro qui fournissent des données sans lien statistique avec les distributions correspondantes des codes d'état final pour l'ensemble des enregistrements.

Une façon d'utiliser les taux partiels consiste à les comparer entre eux. Le taux de réponse CASRO et le taux d'achèvement par ménage sont généralement le plus élevés pour les numéros inscrits, un peu moins élevés pour les numéros non inscrits d'une banque d'un et plus et plutôt faibles pour les numéros d'une banque de zéro (données non présentées). Les écarts par rapport à cette tendance pourraient justifier un examen approfondi de l'exactitude des données. Il faut poursuivre les travaux pour examiner les rapports entre les mesures partielles ainsi que les facteurs à l'origine de tendances différentes.

Les écarts entre les taux partiels soulèvent également la question des liens entre chaque taux et l'exactitude des données. La distribution des codes d'état final peut être influencée par des facteurs étrangers à l'exactitude des données, par exemple la mesure dans laquelle les numéros hors service émettent un triton. Il est possible que la distribution des codes d'état final pour les numéros inscrits soit moins influencée par

des facteurs étrangers que celle des codes d'état final pour d'autres mesures partielles. Les numéros inscrits ont déjà été ceux de ménages. Si les taux de changement des numéros qui ne sont plus ceux de ménages sont relativement constants d'un territoire à l'autre, alors les écarts entre les taux de réponse parmi les numéros inscrits correspondraient plus directement aux écarts entre les taux de réponse réels que ne le feraient les écarts entre les taux globaux ou d'autres taux partiels. Encore une fois, il faut poursuivre les travaux pour déterminer dans quelle mesure cette hypothèse se vérifie.

Nous avons abordé le calcul de taux partiels en fonction des taux pour une étude complète. Le calcul de taux partiels peut être encore plus utile lorsqu'on examine les résultats par intervieweur. L'utilisation de taux partiels permettrait de neutraliser les écarts dans l'échantillon utilisé par les intervieweurs individuels. Par exemple, une tendance anormale de la distribution des codes d'état final pour certains intervieweurs permettrait de repérer les intervieweurs qui ont besoin d'une amélioration.

Enfin, on peut souligner une application des taux partiels dans la production. Comme les taux d'efficacité varient d'un sous-échantillon à l'autre, on peut estimer avec plus de précision le nombre d'interviews achevées à partir d'un ensemble donné d'enregistrements échantillonnés en calculant et en appliquant séparément les taux d'efficacité pour chaque sous-échantillon plutôt qu'en appliquant des taux d'efficacité globaux. Les fournisseurs d'échantillons commerciaux peuvent indiquer sur demande à quel sous-échantillon appartient chaque enregistrement.

BIBLIOGRAPHIE

Council of American Survey Research Organizations (1982), *Report of the CASRO Completion Rates Task Force*, New York: Audits and Surveys Company.

Groves, R. M., et M. P. Couper (1998), *Nonresponse in Household Interview Surveys*, New York: Wiley.

Figure 2. Pourcentage d'enregistrements échantillonnés qui consistent en numéros de téléphone inscrits, par territoire et par mois de présentation

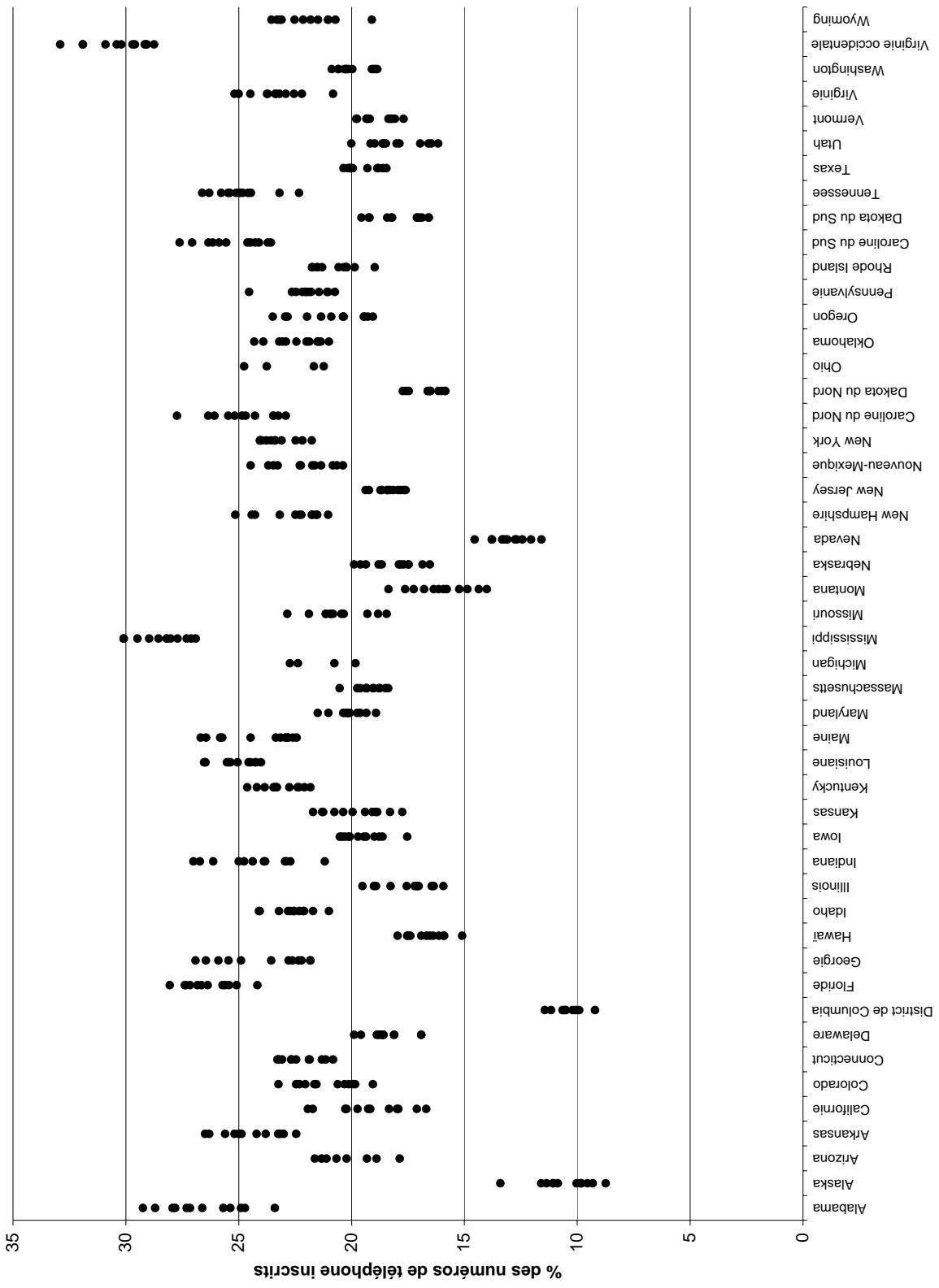


Tableau 2. Distributions des corrélations entre les taux global et partiel d'identification de ménages			Enregistrements d'une banque de zéro		
Enregistrements inscrits			Enregistrements non inscrits d'une banque d'un et plus		
Quantile	Estimation	Diagramme à deux dimensions	Quantile	Estimation	Diagramme à deux dimensions
100 % (maximum)	0,9992285		100 % (maximum)	0,969071	
99 %	0,9992285		99 %	0,969071	
95 %	0,9756642		95 %	0,964890	
90 %	0,9351920		90 %	0,947907	
75 % Q3	0,8488764		75 % Q3	0,894336	
50 % (médiane)	0,7425905		50 % (médiane)	0,783756	
25 % Q1	0,6096933		25 % Q1	0,633399	
10 %	0,3287779		10 %	0,449564	
5 %	0,2038621		5 %	0,354852	
1 %	-0,0469052		1 %	-0,628246	
0 % (minimum)	-0,0469052		0 % (minimum)	-0,628246	
Diagramme arborescent	n°	Diagramme à deux dimensions	Diagramme arborescent	n°	Diagramme à deux dimensions
10 0	1		9 00002355677	12	
9 689	3		8 345556788899	12	+-----+
9 0014	4		7 044566778	9	*-----*
8 5689	5	+-----+	6 0033479	7	+-----+
8 0113	5		5 337	3	
7 56677	6		4 256	3	
7 00224	5	*-----*	3 5	1	
6 589	3	+	2		
6 11244	5	+-----+	1		
5 9	1		0		
5 3	1		-0 3	1	0
4 5	1		-1		0
4 334	3		-2		
3 5	1		-3		
3 3	1		-4		
2			-5		
2 00	2	0	-6 3	1	*
1 5	1	0	-----+-----+		
1			Multiplier diagramme arborescent par 10**-.1		
0					
0					
-0					
-0 5	1	0			
-----+-----+					
Multiplier diagramme arborescent par 10**-.1					

