



Projet Obe Maghreb

Ecole thématique gestion et analyse de données  
20 au 29 avril 2010

## Qualité des données dans les enquêtes épidémiologiques



Pierre Traissac  
UMR 204 « Prévention des malnutritions et pathologies associées »  
IRD, Montpellier, France



1

## Objet

- Assurance qualité des données pour une enquête épidémiologique
- Concept « Erreur d'enquête totale »
- Approche non complètement formalisée / normalisée (e.g. HACCP)
- Présentation : principes généraux
- Détails techniques : documents de référence / mise en pratique

2

## Erreur totale d'une enquête (« total survey error »)

**Question**  
(e.g. de nutrition)

Monde réel  
(e.g. adultes Tunisiens)

3

## Erreur totale d'une enquête (« total survey error »)

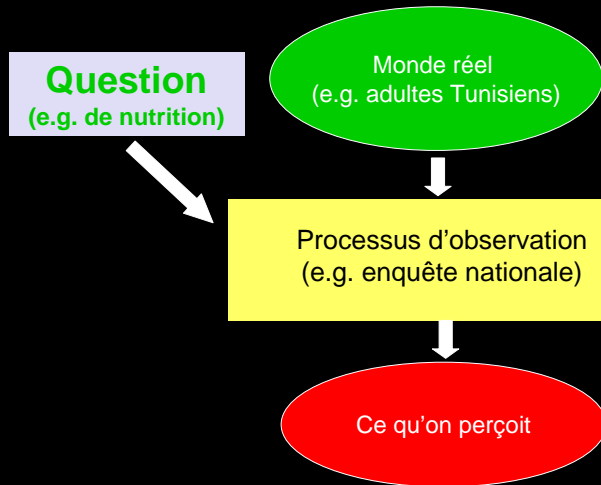
**Question**  
(e.g. de nutrition)

Monde réel  
(e.g. adultes Tunisiens)

Processus d'observation  
(e.g. enquête nationale)

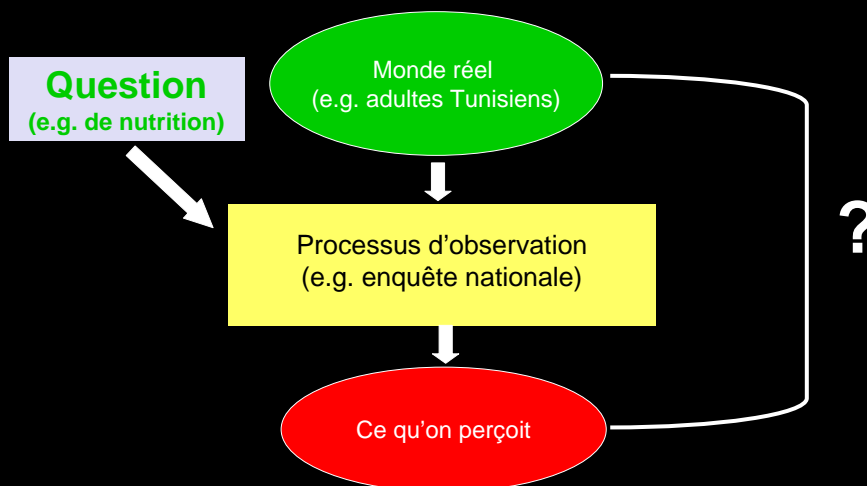
4

## Erreur totale d'une enquête (« total survey error »)



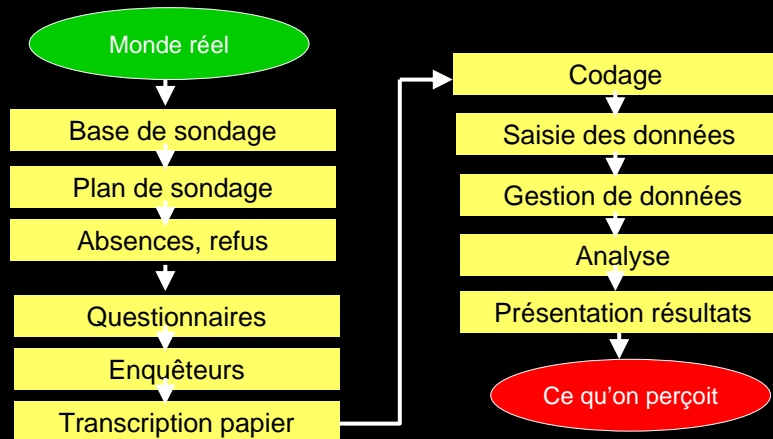
5

## Erreur totale d'une enquête (« total survey error »)



6

## Erreur totale d'une enquête (« total survey error »)



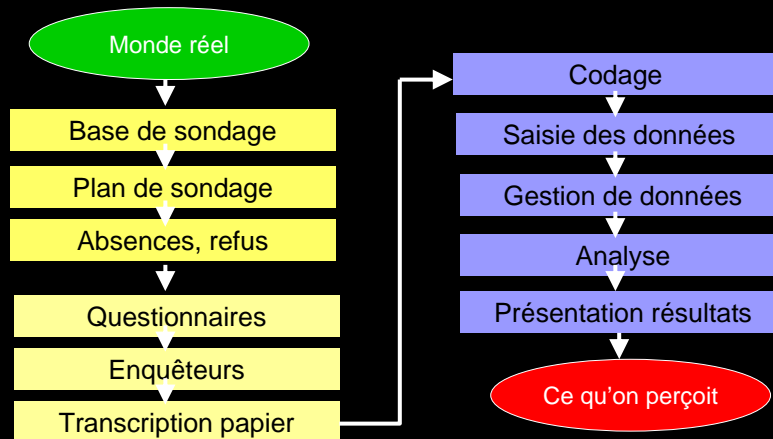
7

## Erreur totale d'une enquête (« total survey error »)

- Erreurs lors de la collecte de l'information sur le terrain  
(« response errors »)
- Erreurs dans les traitements ultérieurs  
(« processing errors »)

8

## Erreur totale d'une enquête (« total survey error »)



9

## Erreurs lors de la collecte de l'information (« response errors »)

### ■ Erreurs liées à la sélection des sujets

- population cible v.s. population source
- variance d'échantillonnage (vs. recensement)
- plan de sondage théorique (base de sondage et plan de sondage)
- mise en œuvre terrain (sélection des individus, absences, refus)

### ■ Erreurs de mesure

- qualités intrinsèques des instruments de mesure et questionnaires (précision, validité, reproductibilité)
- mise en œuvre pratique (administration du questionnaire, prise des mesures, transcription)

10

## Erreurs dans la chaîne de traitement (« processing errors »)

- Transcription des mesures/observations/réponses
- Codage (le cas échéant)
- Saisie (lecture, frappe), validation, apurement
- Gestion des données
  - Sélections, fusion de tableaux de données
  - Calcul de nouvelles variables, indices, recodages
- Analyses (choix des méthodes, mise en oeuvre)
- Préparation de tableaux / graphiques
- Documents écrits / supports visuels

11

## Erreurs dans la chaîne de traitement (« processing errors »)

- Transcription des mesures/observations/réponses
- Codage (le cas échéant)
- Saisie (lecture, frappe), validation, apurement
- Gestion des données
  - Sélections, fusion de tableaux de données
  - Calcul de nouvelles variables, indices, recodages
- Analyses (choix des méthodes, mise en oeuvre)
- Préparation de tableaux / graphiques
- Documents écrits / supports visuels

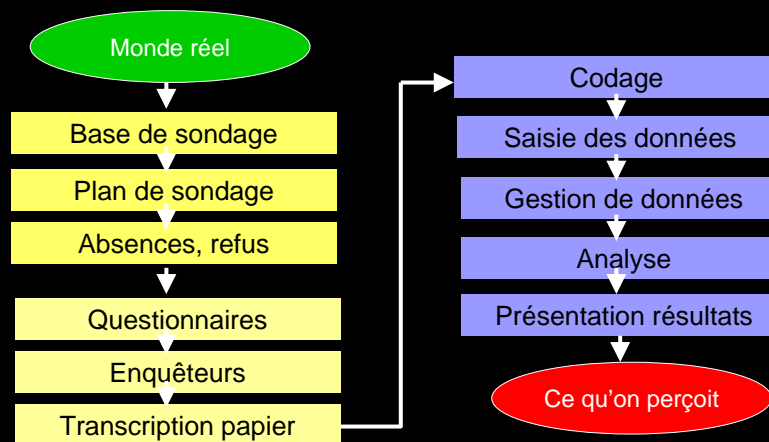
12

## Principe de base n°1

« Une chaîne est aussi solide que le plus faible de ses maillons »

13

## Erreur totale d'une enquête (« total survey error »)



14

## Principe de base n°1

### « Une chaîne est aussi solide que le plus faible de ses maillons »

- ne pas négliger certaines étapes a priori moins « nobles » : aspects « techniques » (e.g. mesures, codage, saisie, recodages) vs. « conceptuels » (e.g. échantillonnage, questionnaire, analyse)
- planification, organisation, écriture de guides opératoires, formation des opérateurs, tests pour TOUTES les étapes : sélection des ménages, administration du questionnaire, mesures, codage, saisie, apuration, gestion de données, analyse

15

## Principe de base n°2

### « Confiance n'exclut pas contrôle »

- De son propre travail, de celui des autres
- Procédures de contrôle, détection et correction d'erreurs
- Toutes les étapes de la chaîne de collecte / traitement
- Aussi pour gestion et analyse de données (cf exercices)

16



## Principe de base n°3

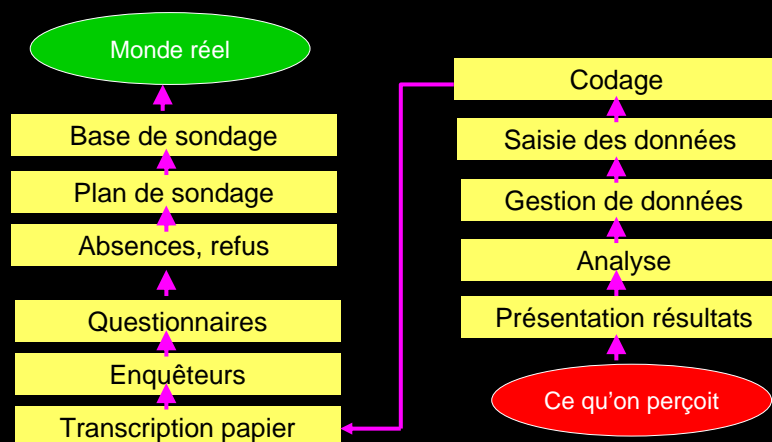
### « Traçabilité »

- But : pouvoir à tout moment remonter la chaîne de collecte/traitement jusqu'à la « donnée de base » (pour vérification, détection et correction erreurs)

17

## Principe de base n°3

### « Traçabilité »



18

## Principe de base n°3

### « Traçabilité »

- But : pouvoir à tout moment remonter la chaîne de collecte/traitement jusqu'à la « donnée de base » (pour vérification, détection et correction erreurs)
- Moyen : **documentation écrite** de tous les processus, étapes de traitement, difficultés éventuelles de mise en œuvre etc.

19

## Traçabilité : documentation écrite

- **Éléments concernant la méthodologie** (base de sondage, échantillonnage, mesures, questionnaire)
- **Modes opératoires** : guides pour enquêteurs (choix des individus, mesures), pour opérateurs de saisie.
- « Journal de bord » des activités de terrain et/ou de laboratoire et de bureau (problèmes rencontrés, modes de résolution)

20

## Traçabilité : documentation écrite

- Description précise des processus d'accès aux données (saisie, vérification, apurement)
- Archivage des questionnaires « papier »
- Description précise des processus de traitement (e.g. QFCA → score de diversité alimentaire DQI, HEI ...)
- Dictionnaires de variables
- Logiciels : programmes et non mode interactif

21

## Traçabilité : documentation écrite

- Description précise des processus d'accès aux données (saisie, vérification, apurement)
- Archivage des questionnaires « papier »
- Dictionnaires de variables
- Logiciels : programmes et non mode interactif
- Description précise des processus de traitement (e.g. QFCA → ingéré en divers nutriments)

22

- Dictionnaire de variables

Ordre dans le fichier	nom	Contenu	format	Unités / codes	n	Pgm de création
1	nusaiad	Numéro séquentiel de saisie adulte	Numérique 4.	Sans objet	250	adultes.qes
2	idadu	Identifiant unique d'adulte	Numérique 4.	Sans objet	250	adultes.qes adultes.chk
3	idmen	Identifiant de ménage	Numérique 4.	Sans objet	250	adultes.qes adultes.chk
4	numad	Numéro adulte dans ménage	Numérique 2.	Sans objet	250	adultes.qes
5	datenq	Date d'enquête	Date (dmy 10)	jj/mm/aaaa	250	adultes.qes
6	sexe	Sexe de la personne	Numérique 1.	1: masculin 2: féminin	250	adultes.qes
7	age	Age de la personne	Numérique 3.	années révolues	245	adultes.qes
8	statmat	Statut matrimonial	Numérique 1.	1: célibataire 2: marié 3: .....	240	adultes.qes
9	gross	Grossesse visible	Numérique 1.	1:oui 2:non	145	adultes.qes
10	tailled	Taille debout	Numérique 5.1	cm	238	adultes.qes

23

## Traçabilité : documentation écrite

- Description précise des processus d'accès aux données (saisie, vérification, apurement)
  - Archivage des questionnaires « papier »
- Dictionnaires de variables
  - Logiciels : programmes et non mode interactif
  - Description précise des processus de traitement (e.g. QFCA → ingéré en divers nutriments)

24

# Traçabilité : gestion et analyse

**Indice de masse corporelle en 4 classes**

	N	%
1	101	5.5
2	774	41.9
3	537	29.0
4	423	22.9
.	14	0.8
<b>Total</b>	<b>1849</b>	<b>100.0</b>

**Mode interactif : non!**

25

# Traçabilité : gestion et analyse

## Programmes : oui !

**Indice de masse corporelle en 4 classes**

	N	%
1	101	5.5
2	774	41.9
3	537	29.0
4	423	22.9
.	14	0.8
<b>Total</b>	<b>1849</b>	<b>100.0</b>

```

1 //*****
2 // PROGRAMME : gen_femmes2.pgm (pgm Epidata Analysis)
3 // DATE : 07/05/2007
4 // AUTEUR : P. Traissac, IRD UR106 Nutrition, Alimentation, Sociétés
5 // LIEU : IRD Tunis
6 // BUT : Création de nouvelles variables, recodage (inc, inc_c3, âge...)
7 // créer ainsi la table femmes2.rec
8 // TD formation Epidata Analysis Tunis 15 au 17 mai 2007
9 // (exercices repris de formation Stata TRAHINA PT et SED 11/2007)
10 //*****
11 // DONNEES EN ENTREE :- femmes.dbf
12 // u.s. = femme (n=1849)
13 //
14 // EN SORTIE :- femmes2.rec données femmes avec nouvelles variables ajoutées
15 // u.s. = femme (n=1849)
16 //*****
17 //*****
18 // lecture tableau de données
19 // read peut lire format Epidata (EpiInfo) .rec ou format Dbase .dbf
20 // au 07/05/2007 problème de lecture variables réelles avec décimales
21 // en format .rec donc on lit un fichier au format .dbf
22 //
23 // D:\MLS_Tunisie\INSP\FormEpidataAna\td_ga_an
24 read "femmes.dbf" *close
25
26 // création nouvelle variable inc et codage en classes
27 gen inc=poids/(taille/100)^2
28 label inc "Indice de masse corporelle en kg/m2"
29 describe inc
30 define inc4 #
31 label inc4 "Indice de masse corporelle en 4 classes"
    
```

26

## Traçabilité : gestion et analyse

### ■ Entête des programmes (e.g. Epidata Analysis)

```
//*****  
// Nom du programme : gen_femmes2.pgm  
// Type de programme : Epidata Analysis  
// Auteur : PT, IRD, UR 106 Nutrition, Alimentation, Sociétés  
// Date : 08/05/2007  
// Lieu : IRD, Montpellier – IRD, Tunis  
// But : Calculs et recodages données femmes en préalable à analyse  
//       facteurs de risque de l'obésité  
//  
// Données en entrée :  
//       femmes.rec (u.s. : femme adulte, n=1849)  
//  
// Données en sortie : femmes2.rec (n=1763)  
//  
// Remarques diverses : pgm repris de exercices Stata SED et PT  
//*****  
À partir d'ici les instructions de programmation proprement dites
```

27

## Traçabilité : gestion et analyse

### ■ Commentaires dans les programmes

```
// tableau d'entrée issu de la double saisie  
read "femmes.rec" /close  
//       création nouvelle variable imc  
gen imc=poids/(taille/100)^2  
label imc "Indice de masse corporelle en kg/m2 "  
//       imc en 4 classes (bornes OMS : maigre, surpoids, obésité)  
//       cf rapport « Physical status : use and interpretation of anthropometry »  
//       OMS 1995 (page 357)  
define imc4 #  
label imc4 "Indice de masse corporelle en 4 classes"  
//       à ce stade elle n'a que des valeurs manquantes  
recode imc to imc4 lo-18.4999=1 18.50-24.9999=2 25.0-29.9999=3 30-hi=4  
//       on peut assigner des labels aux codes ainsi créés  
labelvalue imc4 /1="maigre" /2="normal" /3="surpoids" /4="obèse"  
//       codage obésité (>=30 kg/m2)  
define obese #  
recode imc to obese lo-29.999=0 30-hi=1  
//       sauvegarde nouvelles variables dans table femmes2  
savedata "femmes2.rec" /replace
```

28

## Traçabilité : documentation écrite

- Description précise des processus d'accès aux données (saisie, vérification, apurement)
- Archivage des questionnaires « papier »
- Dictionnaires de variables
- Logiciels : programmes et non mode interactif
- Description précise des processus de traitement (e.g. QFCA → ingéré en divers nutriments)

29

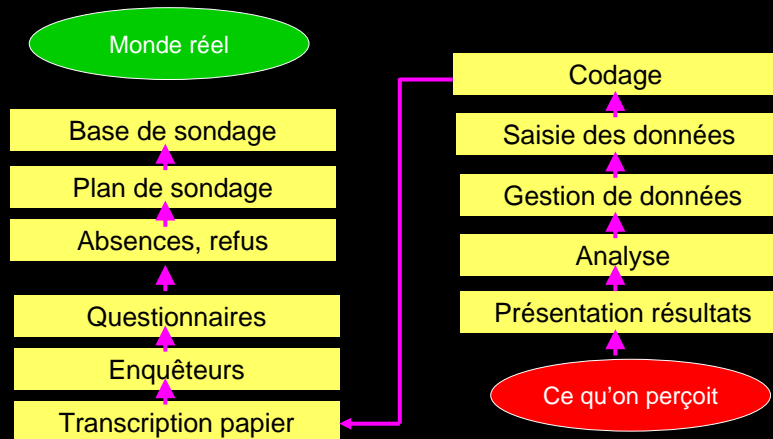
## Traçabilité : gestion et analyse

- Documenter les opérations sur les fichiers de données

Données en entrée	Programme	En sortie	Commentaires
menages.rec n=1763, p= 39	gen_menages2.pgm (pgm Epidata Ana.)	menages2.rec n=1763, p= 45	Calcul score biens Terciles par milieu .....
femmes.rec n=1849, p=36 menages2.rec n=1763, p= 39	gen_femmes2.pgm (pgm Epidata Ana.)	femmes2.rec n=1849, p=43 fem2men.rec n=1849, p=90	Calcul imc, âge, recodages .... Fusion données ménages
fem2men.rec n=1849, p=90	ana1_fem.pgm (pgm Epidata Ana.)	ana1_fem.log (fichier résultats)	Analyse bivariée facteurs de risque obésité

30

## Principe de base n°3 « Traçabilité »



31

## Un mot sur la sécurité des données

- Tous types :
  - protocoles, questionnaires, fichiers de données, programmes, résultats d'analyses, documents, rapports, courriels, images, ...
- Protection accès non autorisés (serrures , mots de passe, cryptage)
- Protections anti-virus (internet +++)
- Protection contre pertes accidentelles de données :
  - sauvegardes régulières (CD-ROM, DVD, bandes, disque amovible)
  - doublage et mise en sécurité des supports de sauvegarde
- Archivage questionnaires papier (forme papier ou informatique)
- Motiver les personnes

32



## Conclusion

- Erreur d'enquête totale
- Une chaîne est aussi solide que le plus faible de ses maillons
- Confiance n'exclut pas contrôle
- **Traçabilité (documentation ++)**
- Formation, motivation des personnes
- Sécurité des données

33

## Quelques références

- Adelf-Aeerna-Aderest-Epiter (2003). *Recommendations. Déontologie et bonnes pratiques en épidémiologie*. ADEL, AEEMA, ADEREST, EPITER: 31.
- Bennett S., Myatt M. et al. (2001). *Data management for surveys and trials. A practical primer using EpiData*. The EpiData Association.
- Biemer P., Lyberg L. (2003). *Introduction to survey quality*. Wiley & Sons.
- Juul S. (2001). *Take good care of your data*. Aarhus, Department of Epidemiology and Social Medicine, University of Aarhus: 56.
- Traissac P. (2003). *Quelques éléments de bonnes pratiques pour la saisie, gestion et analyse des données*. Projet TAHINA. INCO Med – ICA3-CT-2002-10011
- United Nations Statistics Division (2003). *Household Surveys in Developing and Transition Countries: design, implementation and analysis*. <http://unstats.un.org/unsd/HHsurveys/index.htm>

34