

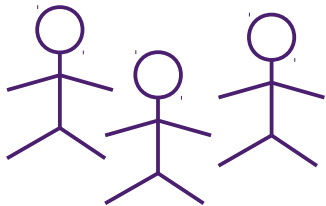
# Alimentation en ligne d'une base de données à partir d'un fichier Excel

Marie-Claude QUIDOZ  
CEFE - Montpellier

Juliette FABRE  
OSU OREME - Montpellier

# Contexte

Producteurs données  
& utilisateurs de R



Gestion des données  
(saisie, traitement, visualisation ..)



Aucune contrainte  
→ pas de typage  
→ pas de structure  
→ pas de contrôle



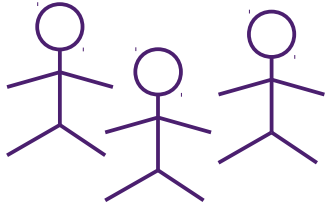
- Incohérences  
- Fautes de frappe



Pérenniser les données  
Offrir des outils de diffusion

# Projet

Producteurs données  
& utilisateurs de R



Gestion des données  
(saisie, traitement, visualisation ..)



Aucune contrainte  
→ pas de typage  
→ pas de structure  
→ pas de contrôle



- Incohérences  
- Fautes de frappe



Insertion automatisée  
Vérification des données

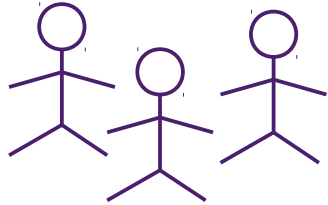


Pérenniser les données  
Offrir des outils de diffusion

# Ajout de deux contraintes

Producteurs données  
& utilisateurs de R

Gestion des données  
(saisie, traitement, visualisation ..)



Pas de manipulation

1

Aucune contrainte  
→ pas de typage  
→ pas de structure  
→ pas de contrôle



- Incohérences  
- Fautes de frappe



Insertion automatisée en ligne  
Vérification des données



2



Pérenniser les données  
Offrir des outils de diffusion

# Fonctions génériques R

Script R de fonctions **utilitaires** et **génériques**

- Interrogation de **bases de données** (RPostgreSQL)
  - Extraction de valeurs : comparaison données / BDD
  - Ordination de jeux de données selon une table : formatage
  - Alternative : RODBC

# Fonctions génériques R

## Script R de fonctions **utilitaires** et **génériques**

- Interrogation de **bases de données** (RPostgreSQL)
  - Extraction de valeurs : comparaison données / BDD
  - Ordination de jeux de données selon une table : formatage
  - Alternative : RODBC
- **Vérification** de données
  - Unicité de valeurs : clés primaires
  - Valeurs manquantes : NOT NULL, clés étrangères
  - Valeurs autorisées : liste, intervalle
  - Format : numérique, date, datetime, nb de caractères, ...

# Fonctions génériques R

## Exemples :

- Les noms de stations existent-ils dans la base ?

```
db_station <- get_db_values(fields = 'station_name',  
                             table = 'station',  
                             schema = 'my_schema')
```

```
check_belong(col = 'station', data = data, value_set = db_station)
```

- Les sample\_id sont-ils fournis, uniques et de 30 caractères max ?

```
check_missing_values(col = 'sample_id', data = data)
```

```
check_unicity(col = 'sample_id', data = data)
```

```
check_nb_character(col = 'sample_id', data = data, nb_char = 30)
```

# Format et règles de validation

- Définition d'un format **standard** pour les fichiers Excel
  - Pour chaque type de données
  - Définition des noms de colonnes, de feuilles
  - flexibilité possible : lignes / colonnes supplémentaires non prises en compte





# Scripts R de traitement

Scripts R de **traitement** (pour chaque type de données)

- **Import** des fichiers avec XLConnect

- multi-plateformes, gestion des caractères spéciaux, des noms de feuilles

- `loadWorkbook()`, `getSheets`, `readWorksheet()`

- **Vérification** des données (fonctions génériques)



- **Formatage** (ajout d'identifiants uniques, traitements, séparation en plusieurs tables, ...) et **insertion**



- **Ecriture des erreurs** dans un fichier

# Interface web

## Interface **web** de soumission de fichier (PHP)

- Upload
- Exécution du script R
- Récupération et téléchargement des fichiers d'erreurs

Bienvenue, admin | Déconnexion | Ma page

OSU OREME DATA  
Site expérimental de publication des données de l'OSU OREME

Ouvrir le menu | Accueil > Truites > Truites > Soumission D'un Jeu De Données

### Soumission d'un jeu de données

Télécharger le [fichier type](#).

Lire le [document](#) décrivant le format approprié des données.

Fichier :  Aucun fichier sélectionné.

© 2011 data.oreme.org | Vos droits et devoirs sur ce site

Envoi du fichier

juliette | Bureau | SO | gek | instrument\_type\_files

Raccourcis	Nom	Taille	Modifié
Rechercher	CTD.pdf	470,2 Ko	07/04/2014
Récemment util...	dcx22f.pdf	632,7 Ko	19/02/2013

juliette

- Bureau
- Système de fich...
- Système de fich...
- Réservé au syst...

Tous les fichiers

Annuler | Ouvrir

Bienvenue, admin | Déconnexion | Ma page

OSU OREME DATA  
Site expérimental de publication des données de l'OSU OREME

Ouvrir le menu | Accueil > Truites > Truites > Upload Trout File

### Soumission d'un jeu de données

Traitement des données : **SERIE\_T\_130102test.xls**

Les données n'ont pas été validées, lisez le [fichier erreur](#)

Attention, certaines données ne respectent pas les règles (mais sont tout de même validées), lisez le [fichier warning](#)

 Détails de la mise à jour :

- Lecture du fichier ...
- Vérification des noms de colonnes ...
- Vérification des données de la feuille Reports ...
- Vérification des données de la feuille Listing ...
- Vérification des données de la feuille Génotypes ...

!! Attention, certaines données ne respectent pas les règles (mais sont tout de même validées), lisez le fichier warning.  
Les données n'ont pas été validées. lisez le fichier erreur.

[Retour](#)

© 2011 data.oreme.org | Vos droits et devoirs sur ce site

Valid CSS & XHTML | Credits | L'équipe

# Conclusions (1)

- Import des fichier Excel : XLConnect
  - Pb avec les fichiers **volumineux** (ex : 13Mo, 150.000 lignes)  
options(`java.parameters = "-Xmx4096m"`) ne suffit pas toujours
  - Bug dans la lecture de champs **heure**
  - Alternative : xlsx
- Fichiers Libre Office (.ods)
  - librairies gnumeric, ROpenOffice, speedR
  - gnumeric, speedR : pas réussi à utiliser
  - ROpenOffice : moins performante ? Difficultés d'import

# Conclusions (2)

- Réalisation dans différents environnements

- Debian GNU/Linux 7 (wheezy)
- R 3.1
- PostgreSQL 9.2.8
- Apache 2.2.22
- PHP 5.3.1011

- Windows Server 2008R2
- R 3.1
- PostgreSQL 9.3.4
- Apache 2.4.9
- PHP 5.5.10

- Application web → gestion des droits nécessaire

→ Upload du fichier, lecture du fichier Excel, exécution du script R

- Généricité du script de vérification

→ difficultés à décrire formellement l'intégralité des règles de validation