

## M2 Génie Logiciel

### Rapport TD3 : MapReduce

#### Exercice 1 : Filtrage

Nous avons vu que le fichier est structuré par ligne ainsi

*Country, City, AccentCity, Region, Population, Latitude, Longitude*

mais que toutes les villes n'indiquent pas forcément toutes les informations.

Le compte de nombre de lignes sur le fichier initial est de **3173959** et nous obtenons **47980** villes indiquant leur population.

#### Exercice 2 : Compteurs

Nous créons nos compteurs dans un enum WCP :

- **NB\_CITIES** qui compte le nombre de villes valide ; une ville valide est une ville qui possède un nom dans le fichier initial.
- **NB\_POP** qui compte le nombre de villes avec leur population renseignée ;
- **TOTAL\_POP** qui compte la population totale du fichier.

On incrémente **NB\_CITIES** dans la fonction map car on doit considérer toutes les villes et pas seulement celles avec une population renseignée.

Nous avons choisi d'y placer aussi **NB\_POP**, bien qu'elle aurait pu être placée au début du reducer.

On compte ensuite, dans la fonction reduce, la population totale. une fois les villes filtrées.

On utilise donc la méthode **getCounter(Enum)** dont l'argument est de la forme **WCP.type** pour accéder à un compteur spécifique.

Résultats :

```
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
TP3$WCP
  NB_CITIES=3173959
  NB_POP=47980
  TOTAL_POP=2289584999
File Input Format Counters
  Bytes Read=151149418
File Output Format Counters
  Bytes Written=1635596
```

Nous savons que d'autres groupes ont utilisés la méthode **getCounter**(String,String) pour parvenir à ce résultat, mais nous ne parvenons pas à comprendre où instancier les compteurs avec cette méthode, et donc utiliser les compteurs à la fois dans le Mapper et le Reducer.