

STID 1 - Programmation Statistique

TP2

Manipuler un dataframe

Anthony SARDELLITTI

2023-01-01

Contents

1	Ressources documentaires	1
2	Exercices	1
2.1	Exercice 1 : Importer les données	2
2.2	Exercice 2 : Statistiques descriptives	3
2.3	Exercice 3 : Tris et Selections	3
2.4	Exercice 4 : Tris et Filtres	3
2.5	Exercice 5 : Agrégations	4

1 Ressources documentaires

Pour réaliser ce TP, vous aurez besoin des ressources suivantes :

- Importer un fichier excel
- Fonctions de tests et comparaisons
- Indexation
- Filtres et sélection
- Les fonctions de tests et opérateurs de comparaison
- Trier
- Agréger

2 Exercices

Pour rendre ce TP, voici les étapes à suivre :

1. Créer une branche **tp2** sur votre repository github que vous avez déjà ouvert. Le tp sera pousser dans cette branche
2. Dans cette branche, créer un dossier **tp2** dans lequel vous pourrez pousser le script avec votre code

3. Si vous ne l'avez pas fait pour le TP1, partagez votre repository. Si vous êtes en public, envoyer le lien de votre repository par mail à *anthony.sardellitti@hotmail.fr* **OU** si vous êtes en privé, ajouter cette adresse mail en tant que collaborateur de votre repository.

Pensez à commenter votre code.

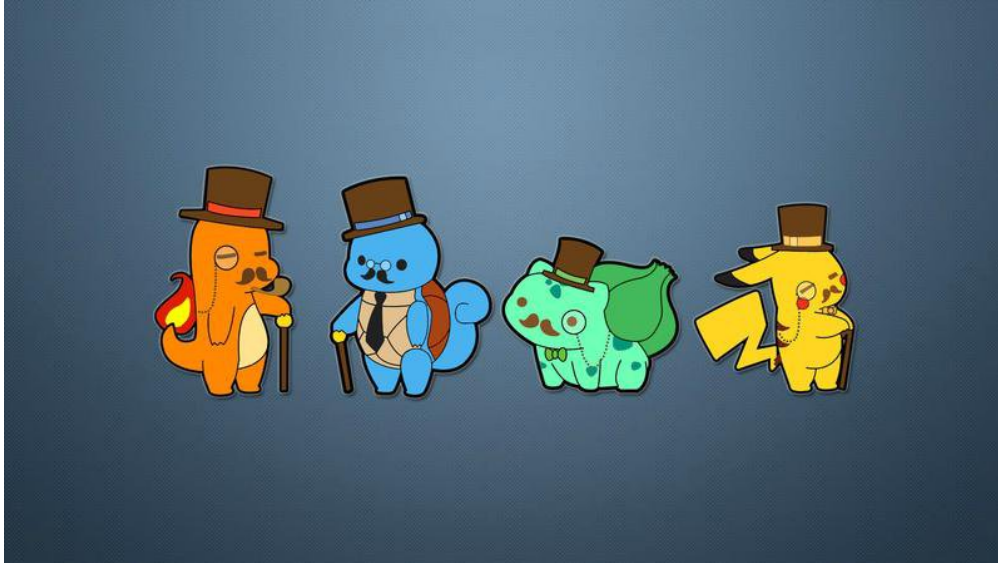


Figure 1: Pokemon

On utilise le fichier `pokemon.xlsx` qui décrit les statistiques des pokemon des deux premières générations. Le fichier est issu du site Kaggle. Il a été adapté pour ce TP. Pour réaliser ce TP, télécharger le fichier en cliquant [ici](#). Voici une description des données :

- `pokedex_number` : numéro du pokemon
- `nom` : nom du pokemon
- `generation` : le numéro de génération dont est issu le pokemon
- `is_legendary` : Oui / Non si le pokemon est légendaire
- `type` : le type du pokemon
- `weight_kg` : le poids du pokemon en kg
- `height_m` : la taille du pokemon en mètre
- `attack` : la puissance d'attaque du pokemon
- `defense` : la puissance de défense du pokemon
- `speed` : la vitesse du pokemon

2.1 Exercice 1 : Importer les données

- a. Importez le jeu de données `pokemon.xlsx` à l'aide du package `readxl`.
- b. Combien de lignes, colonnes sont présentes dans ce dataset (utilisez les fonctions adaptées) ?
- c. Affichez le nom des colonnes.
- d. Affichez le type des colonnes avec la fonction adaptée.
- e. On souhaite analyser les variables `generation`, `is_legendary`, et `type` en tant que variables qualitatives. Modifier le type de ces variables pour les transformer en type factor.
- f. Combien de niveaux (*levels*) sont présents dans ces variables ?
- g. Affichez un résumé des données avec la fonction adaptée.

2.2 Exercice 2 : Statistiques descriptives

- Déterminer la moyenne de la variable `weight_kg`.
- Déterminer la médiane de la variable `weight_kg`.
- Déterminer les quartiles de la variable `height_m`.
- Déterminer les déciles de la variable `height_m`.
- Déterminer la variance et l'écart-type de la variable `weight_kg`.
- Déterminer un tri à plat pour compter les effectifs des modalités de chaque variable *factor* en triant chaque sortie par ordre décroissant.

2.3 Exercice 3 : Tris et Selections

Pour chaque question suivante, affectez le résultat de la requête dans un objet puis calculez sa dimension. Exemple :

```
#Sélectionnez les deux premières colonnes du data frame
requete_0 <- pokemon[,1:2]
dim(requete_0)
```

```
## [1] 251  2
```

- Sélectionnez la colonne `nom` et `is_legendary`.
- Sélectionnez les 50 premières lignes et les deux premières colonnes.
- Sélectionnez les 10 premières lignes et toutes les colonnes.
- Sélectionnez toutes les colonnes sauf la dernière.
- Triez le dataset par ordre alphabétique et afficher le `nom` du pokemon de la première ligne.
- Triez le dataset par `weight_kg` en ordre **décroissant**, et afficher le `nom` du pokemon de la première ligne.
- Triez le dataset par `attack` en ordre **décroissant** puis par `speed` en ordre **croissant**, et afficher le `nom` des pokemons des 10 premières lignes.

2.4 Exercice 4 : Tris et Filtres

Pour chaque question suivante, affectez le résultat de la requête dans un objet puis calculez sa dimension. Pour faciliter la lecture, sélectionnez la colonne `nom` et les colonnes concernées par le filtre. Exemple :

```
#Sélectionnez les pokemons de type feu
requete_0 <- pokemon[ pokemon$type == "fire", c("nom","type")]
dim(requete_0)
```

```
## [1] 20  2
```

- Filtrez sur les pokemons qui ont 150 ou plus d'`attack` puis trier le résultat par ordre décroissant d'`attack`.
- Filtrez sur les pokemons de `type` *dragon,ghost,psychic* et *dark*.
- Filtrez sur les pokemons de `type` *fire* avec plus de 100 d'`attack`, puis trier le résultat par ordre décroissant d'`attack`.
- Filtrez sur les pokemons qui ont entre 100 et 150 de `speed`. Les trier par `speed` décroissant.
- Filtrez sur les pokemons qui ont des valeurs manquantes sur la variable `height_m`.
- Filtrez sur les pokemons qui ont des valeurs renseignées à la fois pour la variable `weight_kg` et la variable `height`.
- Filtrez sur les pokemons pesant plus de 250 kg et affichez le résultat pour vérifier.

2.5 Exercice 5 : Agrégations

Pour chaque question suivante, affectez le résultat de la requête dans un objet puis calculez sa dimension.
Exemple :

```
#Calculez la vitesse moyenne par generation
requete_0 <- aggregate(x = speed ~ generation, data = pokemon , FUN = mean)
dim(requete_0)
```

```
## [1] 2 2
```

- Calculez l'**attack** moyenne en fonction de la variable **type**, puis filtrez sur les 3 types avec les moyennes les plus élevées.
- Calculez le nombre de pokemon par **type** , puis triez par ordre décroissant ces effectifs.
- Calculez la médiane de **weight_kg** par **type**.
- Calculez le nombre de pokemon par **type** et **generation**
- Calculez la moyenne de chaque critère (**weight_kg**, **height_m**, **attack**, **defense** et **speed**) en fonction de chaque **type**.