

# TP Pandas

Pour apprendre à manipuler pandas, nous partirons d'un jeu de données open data hébergé sur le site : [data.gouv.fr](http://data.gouv.fr)

Il s'agit d'une liste de fréquentation au cours du temps de différents musées franciliens. Premières étapes :

- Déposer les fichiers envoyés « freq\_musees.csv » et « communes.csv » dans votre dossier de travail Jupiter
- Charger le fichier freq\_musees.csv dans un dataframe nommé « df »
- Afficher un extrait du dataframe
- Afficher les 13 premières lignes du dataframe
- Afficher seulement les colonnes « etablissements » et « localite » à l'écran
- Compter le nombre de musées disponibles dans ce dataframe
- Le premier row est pas terrible, enlevez-le
- Nous ne nous intéresserons pas à certaines des colonnes du dataframe, autant les enlever. Enlever les colonnes « grat\_06 », « grat\_07 », « grat\_08 », « grat\_09 », « grat\_10 », « commentaires », « wgs84 » en créant un nouveau dataframe « df\_clean »
- Créer un nouveau dataframe « df\_test » d'après « df\_clean » contenant uniquement les colonnes « num\_ref », « etablissements » et « total\_06 ».
- Vérifier le type de la colonne « total\_06 » de « df\_test ». Que constatez-vous ?
- Créer un data frame df\_test2 à partir de df\_test

- Caster la colonne « total\_06 » en int. Que constatez-vous ?
- Le format du fichier csv de départ n'est pas très propre (par exemple, le nombre 288179 est écrit avec un espace dans le fichier : 288 179). Au moment du cast, Pandas considère ce nombre comme un string et le passe à null, ce que nous voulons éviter. Nous devons donc modifier la colonne pour enlever les espaces. Faites-le sur le data frame df\_test.
- Enfin, passez le type de « total\_06 » de « df\_test » à Integer, le problème devrait être réglé.
- Maintenant que vous avez compris comment faire le cast proprement sur une colonne, créez un nouveau dataframe « df\_clean\_and\_cast » en appliquant le cast sur toutes les colonnes le nécessitant à partir de « df\_clean » (qui est complet).  
Liste des colonnes à caster : « total\_06 », « total\_07 », « total\_08 », « total\_09 », « total\_10 », « evolution\_07\_06\_en », « evolution\_08\_07\_en », « evolution\_09\_08\_en », « evolution\_10\_09\_en »
- Trouver combien il y a de musées parisiens
- Afficher par ordre décroissant la liste des musées les plus fréquentés pendant l'année 2010
- Afficher par ordre décroissant la liste des musées parisiens les plus fréquentés pendant l'année 2010
- Afficher la liste des musées ayant une fréquentation supérieure à 1000000 millions de personnes en 2010
- Afficher la liste des musées ayant une fréquentation entre 500000 et 1000000 millions de visiteurs en 2010
- Créez un nouveau dataframe « df\_enrich » à partir de « df\_clean\_and\_cast » contenant une colonne supplémentaire : la somme de toutes les fréquentations de l'année 2006 à 2010 (somme\_freq)
- Créez un nouveau dataframe « df\_enrich2 » à partir de « df\_enrich » écartant tous les musées dont la fréquentation n'a pas été remontée au moins une année (de sorte à écarter tous les « null »). Notez ensuite le nombre de musées suite au filtre appliqué.

- Calculez la moyenne de fréquentation sur les 5 années pour chacun des musées et ajouter une nouvelle colonne « moy\_freq » au dataframe « df\_enrich2 »
- Créez un nouveau dataframe « df\_enrich3 » à partir de « df\_enrich2 » ne gardant que les musées qui ont une fréquentation qui a augmenté d'année en année de 2006 à 2010. Affichez les ainsi que leur nombre.
- Renommer la colonne « localite » du « df\_enrich3 » en « ville »
- A partir du « df\_enrich2 », créer un nouveau dataframe « df\_enrich4 » et calculez l'évolution de fréquentation des musées entre les années 2006 et 2010 dans une nouvelle colonne « evolution\_10\_06\_en »
- Créez un nouveau dataframe « df\_enrich5 » à partir de « df\_enrich4 ». Créez trois nouvelles colonnes à partir de la colonne « num\_ref ». Une colonne « departement\_id » contenant les deux premiers chiffres, une colonne « code\_insee » contenant les 5 premiers et une colonne « num\_musee\_insee » contenant les deux derniers chiffres.
- Le fichier communes.csv contient la liste des communes de France.
- Créez un dataframe « df\_communes » à partir du fichier communes.csv et ne gardez que les champs « code\_insee », « commune », « département », « region », « population »
- Renommer la colonne « Code INSEE » en « code\_insee » pour avoir le même nom que sur df\_enrich5
- Faire un join entre le dataframe « df\_communes » et le « df\_enrich5 » sur la colonne « code\_insee » dans un nouveau dataframe « df\_join ».
- Faire un dataframe « df\_group\_by\_code\_insee » contenant la somme des fréquentations des musées par code\_insee
- Enregistrer le dataframe « df\_group\_by\_code\_insee » dans un fichier csv