

## Rapport sur la réparation des données (we rate dogs)

Cette préparation s'est déroulée suivant plusieurs phases, que j'expose de manière suivante :

### Phase 1 : Collection des données

Dans cette phase je procède d'abord au téléchargement, par la méthode requests, du fichier 'image-predictions.tsv'

Ensuite j'utilise pandas pour effectuer le chargement des trois datasets 'twitter-archive-enhanced.csv' dans le dataframe tweets, 'image-predictions.tsv' dans le dataframe predict\_images et 'tweet-json.txt' retweet\_and\_favorite.

### Phase 2 : Evaluation des données

Dans cette partie j'ai effectué une lecture visuelle et programmation pour comprendre mes données et détecter des potentielles insuffisances ou défauts de qualité ou de rangement (ordre) dans les jeux de données. Dans cette procédure nous avons énuméré les problèmes de qualité et d'ordre suivants (voir tableau 1 les défauts de qualité et tableau 2 pour les défauts de rangement):

*Tableau 1: Défauts de qualité*

N°	Remarque
<b>tweets</b>	
1	Eliminer les retweets
2	Eliminer les replies
3	La valeur 'a' dans la colonne name n'est précise
4	La colonne timestamp est une string au lieu de datetime
5	Il existe des colonnes non nécessaires pour une analyse
6	None pour définir pour désigner NAN dans les colonnes doggo, floofer, pupper, puppo et name.
<b>La table predict_images</b>	
7	Des doublons sur la colonne jpg_url
8	Les colonnes p1,p2,p3 ne sont intuitivement parlantes
<b>Table retweet_and_favorite</b>	
9	L'identifiant (id) du tweet diffère des deux autres tables (tweet_id)

*Tableau 2: Défauts de rangement*

N°	Remarques
1	La table <b>retweet_and_favorite</b> doit faire partie de la table <b>tweets</b>
2	la table <b>predict_images</b> doit faire partie de la table <b>tweets</b>
3	Quatre valeurs en une colonne (doggo floofer pupper puppo) sur la table tweets.

### Phase 3 : Nettoyage des données

Dans cette phase j'ai essayé de résoudre les différents problèmes des données remarqués plus haut. Mes propositions peuvent être comme suit :

*Tableau 3: Résolution des problèmes de qualité*

N°	Remarque
<b>tweets</b>	
1	Récupérer les retweets et les supprimer
2	Récupérer les tweet réponses et les supprimer
3	Supprimer les enregistrements avec nom='a'
4	Convertir la colonne timestamp au type datetime
5	Éliminer les colonnes qui n'ajoute pas de la valeur à l'analyse
6	Remplacer les colonnes contenant None par du vide ( pas par NaN pour faciliter la concaténation des colonnes doggo, floofer, pupper, puppo)
<b>La table predict_images</b>	
7	Supprimer les doublons sur la colonne jpg_url
8	Renommer les colonnes p1, p2, p3 en prediction (_1, _2, _3) pour plus de clarté.
<b>Table retweet_and_favorite</b>	
9	Renommer l'identifiant, préalablement nommer id, en tweet_id.

*Tableau 4: Résolution des problèmes d'ordre*

N°	Remarques
1	Faire un merge sur le tweet_id pour joindre la table <b>retweet_and_favorite</b> à la table <b>tweets</b>
2	Faire un merge sur le tweet_id pour joindre la table <b>predict_images</b> à la table <b>tweets</b> .
3	Regrouper les colonnes doggo floofer pupper puppo en une colonne etape_chien et supprimer les colonnes originales (doggo floofer pupper puppo). Remplacer les valeurs vides de la colonne etape_chien par des NaN.

### Phase 4 : Sauvegarde

Cette phase consiste à enregistrer dans un fichier csv 'twitter\_archive\_master.csv' le dataframe contenant les données nettoyées.