



Big Data

Analyse d'un dataset en R

Axel BAYLE, Pierre Prié 5 ISS

I. INTRODUCTION

Pour mettre en pratique les concepts et notions que nous avons vu en cours de Big Data il nous a été proposé d'analyser un dataset d'au moins 1000 entrées. Cette analyse a été réalisée via le langage de programmation R. Nous avons utilisé un dataset qui répertorie les applications sur le Playstore. Pour chaque application nous avons accès à plusieurs informations : le nom, la catégorie, la note, le nombre de revues, la taille, le nombre d'installations, le type (payant ou gratuit), le prix, tranche d'âge visée, le(s) genre(s), date de mise à jour, version actuelle et version d'android. Nous allons essayer grâce à ce dataset de déterminer comment faire une application populaire en termes d'installations et de notes.

II. INFLUENCE DE LA CATÉGORIE

Notre première approche a été de s'intéresser aux catégories d'applications les plus populaires. C'est pourquoi nous commençons par évaluer, pour chaque catégorie, le nombre d'applications sur le Playstore. En effet, cela semble un bon indicateur car il semblerait que plus il y a d'applications d'un type plus c'est une catégorie populaire. On note avec ce premier graphique 1 que la majorité des applications du Playstore sont gratuites et que les catégories qui possèdent le plus d'applications sont Family, Game et Tools. Nous allons maintenant voir quelles catégories sont les plus populaires en termes de note et d'installations.

Pour ce faire nous choisissons de porter notre attention sur le nombre d'installation d'application par catégorie, ainsi que sur la médiane, le premier et le troisième quartile en termes de note. Grâce à cela on a une idée plus claire de la répartition des notes qu'avec la moyenne.

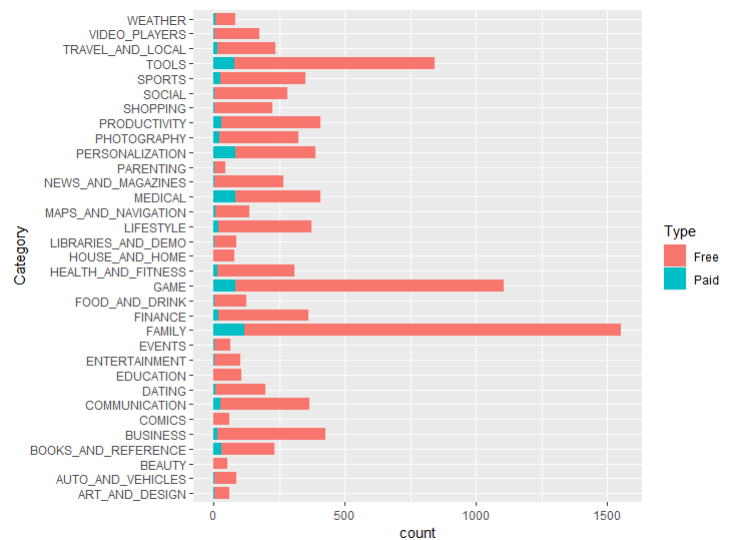


Fig. 1. Nombre d'application par catégorie

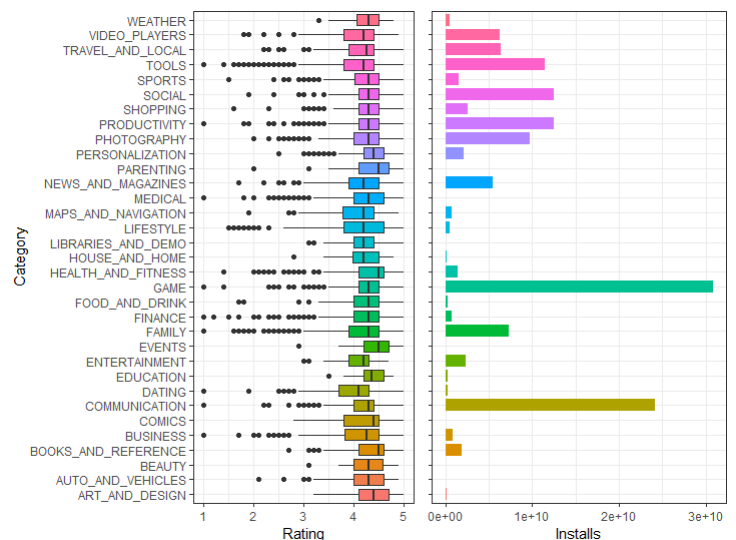


Fig. 2. Note et prix des catégories du playstore

Sur le graphique 2 nous pouvons voir que deux catégories se démarquent des autres quant à leurs nombres

d'installations : Game et Communication, ces catégories ont donc été sélectionnées d'office. Pour la suite de notre étude, nous avons également pré-sélectionné les catégories Photography, Productivity, Social et Tools. Finalement, bien que la catégorie Tools comptabilise plus d'applications, la moyenne de ses notes ainsi que son troisième quartile sont plus faibles que celles des autres applications pré-sélectionnées, cette catégorie n'a donc pas été retenue. On peut noter que contrairement à la catégorie Game qui comptabilise de nombreuses applications (1) mais également de nombreuses installations, les applications de la catégorie Family bien que nombreuses sont elles peu installées. Le nombre d'applications d'une catégorie ne semble donc pas être un critère de succès.

En recoupant les informations que nous apportent les graphiques 1 & 2, nous avons choisi de limiter notre analyse aux catégories qui nous semblent les plus populaires. Ce sont les catégories Game, Communication, Photography, Productivity et Social. Car ce sont celles qui sont les plus installées et celles qui ont une répartition des notes les plus intéressantes.

Maintenant que nous avons choisi les catégories les plus pertinentes pour nous, nous allons étudier les tranches d'âge visées par ces 5 catégories.

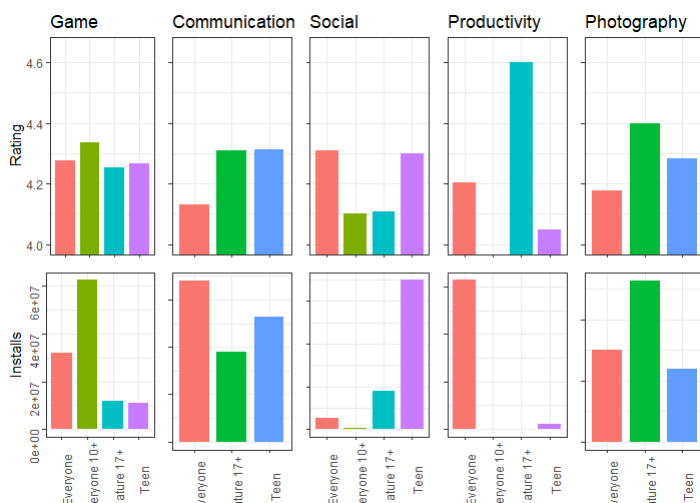


Fig. 3. Populations cibles des catégories

Sur le graphique ci dessus on peut observer plusieurs informations instructives. Tout d'abord, avec les graphiques du haut, on peut remarquer que la catégorie Game est notée de manière quasi-similaire selon que les applications soient destinées à telle ou telle population. Alors que la catégorie Productivity, quant à elle, est notée très différemment selon que l'application cible une population de plus de 17ans ou une population jeune, ou encore tout le monde.

Les catégories Communication, Social et Photographie sont entre ces deux comportements et ont des notes assez similaires malgré que certaines différences existent.

Sur les graphiques du bas, on montre le nombre d'installation selon ces mêmes populations. Cette fois, on se rend compte que chaque catégorie a une population qui se démarque des autres en nombres de téléchargements. Les applications de la catégorie Game sont plus installées lorsqu'elles visent une population 10ans et plus, alors que c'est les applications visant tout le monde qui sont les plus installées dans les catégories Communication et Productivity. Dans la catégorie Social ce sont les applications visant les jeunes qui sont le plus téléchargées. Et finalement, les applications de la catégorie Photographie sont plus installées lorsqu'elles visent un public majeur.

On peut alors se demander si ces deux types de graphiques sont liés. c'est à dire, si les applications d'une catégorie qui visent la population qui installe le plus vont avoir une meilleure note ou pas.

Il apparaît que la corrélation n'est pas systématique, en effet, même si c'est le cas pour la catégorie Photography et Game, on s'aperçoit que la catégorie Productivity déroge à la règle et ne suit pas du tout ce schéma. Les applications de Productivity sont les plus installées lorsqu'elles visent tout le monde mais elles ont tendance à être mieux notées lorsqu'elles visent un public majeur.

Ainsi alors que pour les catégories Game et Photography les critères de notes et de nombres d'installations ont le même effet sur la population à viser. Dans les autres catégories il faudra choisir si nous privilégions le critère de la note ou du nombre d'installation dans la caractérisation d'une application populaire, pour choisir ensuite la population à viser avec l'application.

III. INFLUENCE DU PRIX

Nous avons maintenant voulu observer l'influence du prix des applications sur leurs notes moyennes et leurs nombres d'installations. Pour cela, pour chacune des 5 catégories étudiées, nous avons calculé le prix moyen des applications :

- sur l'ensemble d'entre elles pour la catégorie
- en les filtrant et en ne gardant que celles dont la note est supérieure au troisième quartile des notes de la catégorie
- en les filtrant et en ne gardant que celles dont le nombre d'installation est supérieur au troisième quartile des installations de la catégorie

Nous avons d'abord réalisé ce graphique sur l'ensemble des populations d'application, mais nous avons observé une trop forte influence des applications gratuites sur le résultat du fait de leur surpopulation par rapport aux applications payantes, comme nous avons pu le voir avec le graphique 1. Nous avons donc réalisé un pré-filtrage des applications en ne gardant que les payantes pour réaliser notre étude

sur l'influence du prix. Nous avons obtenue le graphique suivant :

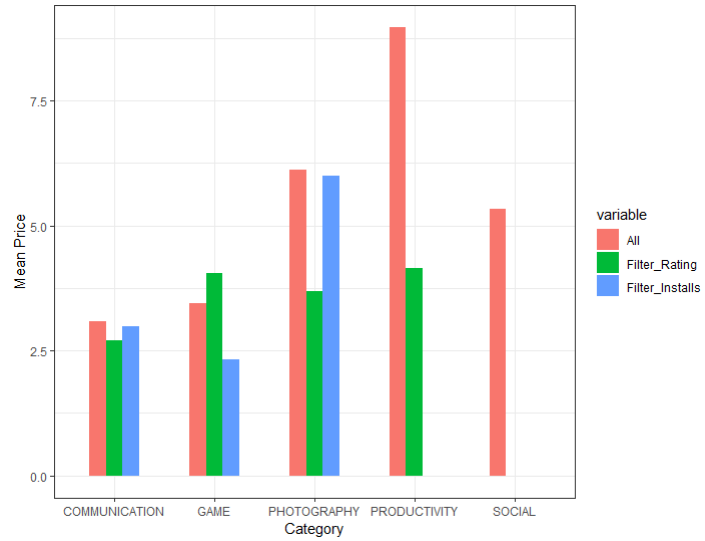


Fig. 4. Prix moyen des applications par catégorie : sans filtre, filtré par le 3ème quartile des notes, filtré par le 3ème quartile des installations

Pour la catégorie Communication, nous pouvons voir que le prix moyen des applications est de 3 dollars, ce prix reste quasi-identique chez les applications qui ont le plus de succès (en note et nombre d'installation). Pour la catégorie Game, le prix moyen des applications est de 3.5 dollars, ce prix est identique pour les applications les mieux notées, mais plus faible (2.1 dollars) pour les plus installées. Le résultat est identique pour la catégorie Photography mais cette fois-ci, c'est les applications les mieux notées qui ont un prix plus faible par rapport à la moyenne qui est de 6 dollars. Nous obtenons des résultats très différents pour les catégories Productivity et Social. Premièrement, pour Productivity, le prix moyen des applications est beaucoup plus important que pour les autres catégories et est de 9 dollars. On voit que le prix moyen des applications les mieux notées est moitié plus faible (4.5 dollars), pour les plus installées le prix moyen est là de 0 euros. Ce résultat se retrouve avec la catégorie social mais cette fois-ci, toutes les applications à succès sont gratuites.

Pour résumer ce graphique, le prix n'a que peu d'influence pour les catégories Communication, Game et Photographie sur les succès des applications (variation de 1 à 2 dollars) alors que pour les catégories Productivity et Social, le prix influence le succès des applications car ce sont les applications gratuites qui fonctionnent.

On pourra également noter que de plus en plus d'applications deviennent gratuite, notamment des jeux, mais proposent alors des micros-transactions ingame. Le dataset ne prenant pas en compte ce genre d'informations, l'étude du prix des applications en fonction de leurs succès pour la catégorie Game peut être faussé.

IV. INFLUENCE DE LA TAILLE

La dernière étape est de chercher si il existe un lien entre la taille d'une application et sa popularité. Pour ce faire nous avons choisi de différencier selon les deux critères de popularité.

A. sur le critère du nombre d'installations

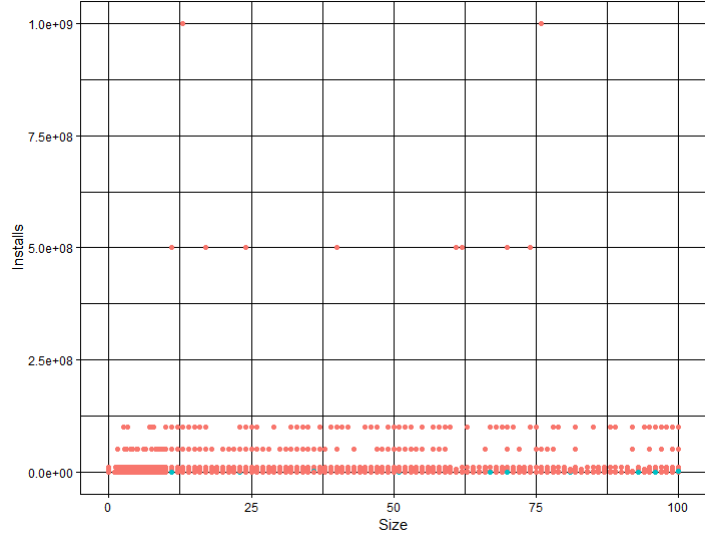


Fig. 5. Nombre d'installations en fonction de la taille

En affichant les applications sur un graphe, le constat est qu'il n'y a pas de lien entre installations et taille. En effet on trouve quasiment le même nombre d'installations quelle que soit la taille de cette application. On remarque aussi que la plupart des applications comptent peu d'installations et que même les plus installées ont des taille variables.

B. sur le critère de la note

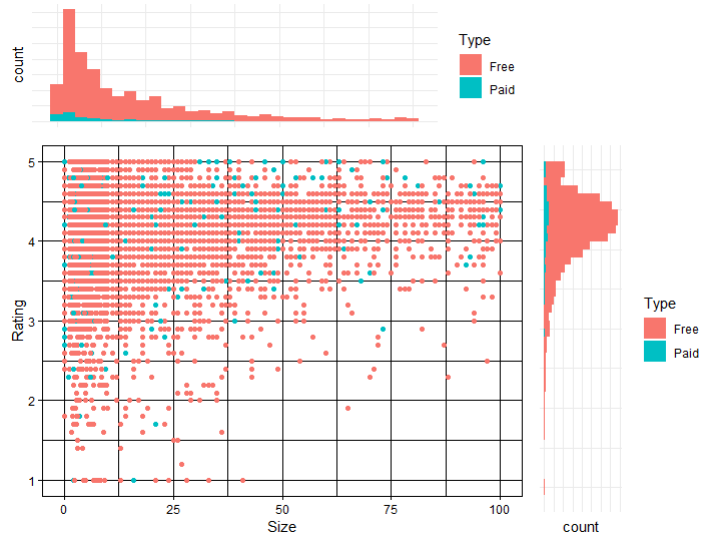


Fig. 6. Note en fonction de la taille

Contrairement au graphique 5, le graphique 6 montre un réel changement selon la taille.

Enfin, n'oublions pas que ce qui rend une application populaire est quand même le contenu qu'elle propose et la publicité qu'elle met en place pour se vendre.