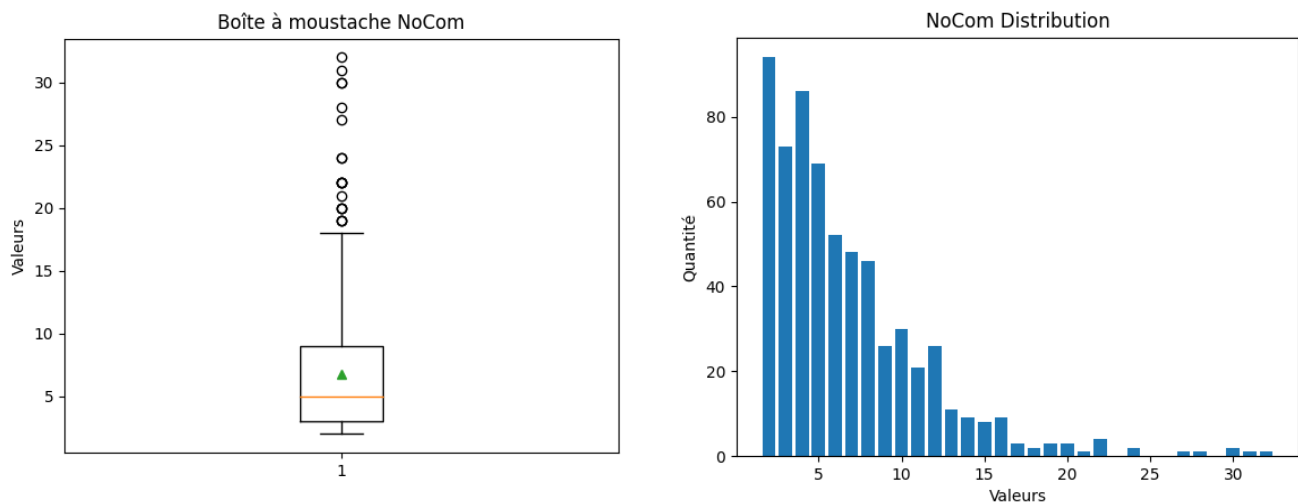


Auteurs : Anthony Grange (20160453) & Luchino Allix-Lastrego (20222844)

TÂCHE 1 :

Voici la visualisation (distribution et boîte à moustaches) de chaque métrique ainsi que ses données pertinentes:
(A noter que les valeurs extrêmes ne sont pas comptabilisées dans le minimum et maximum)

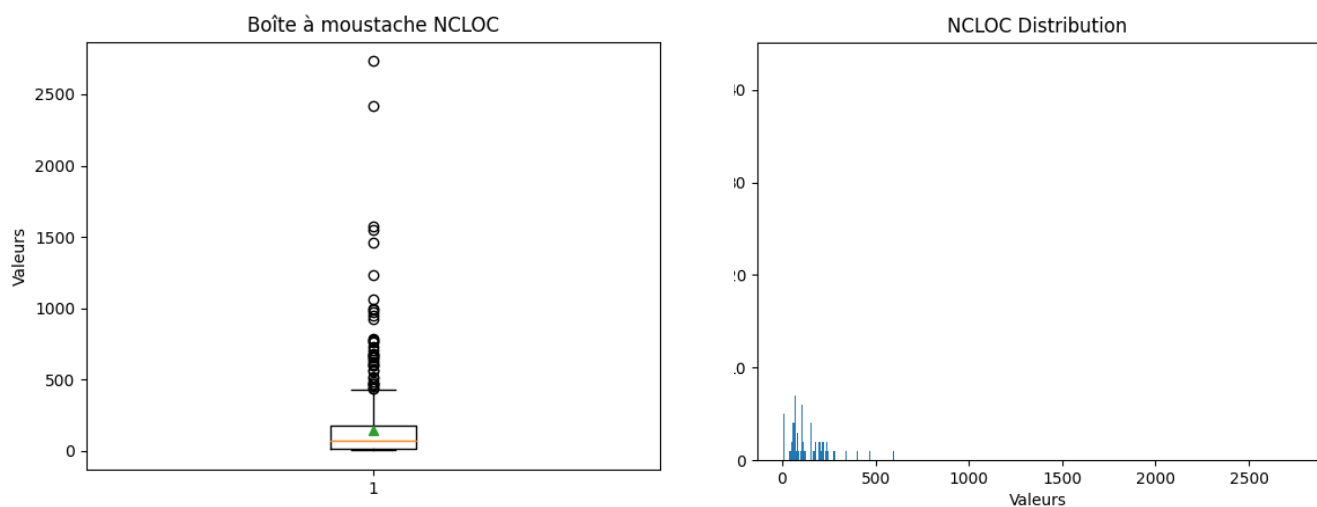
NoCom :



Médiane : 5.0 – Moyenne : 6.7 – Minimum 2.0 – Maximum 18.0 – Q1 : 3.0 – Q3 : 9.0

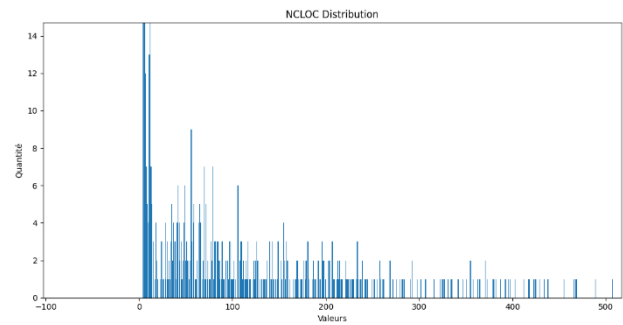
Analyse de la distribution : Clairement décroissante, la plupart des fichiers ont peu de commits, il y a très peu de fichiers avec beaucoup de commits. Elle n'est donc pas normale.

NCLOC :

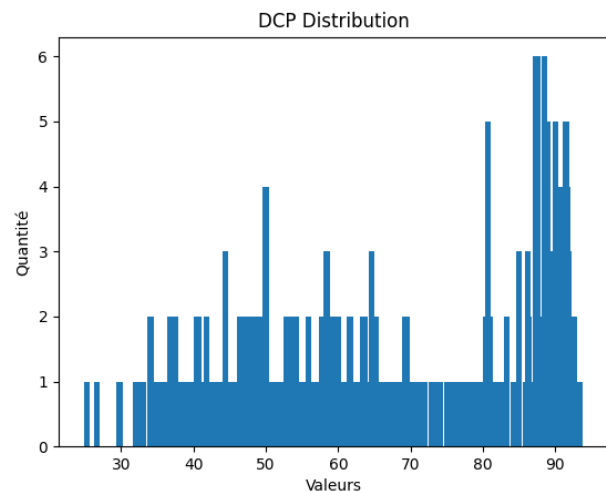
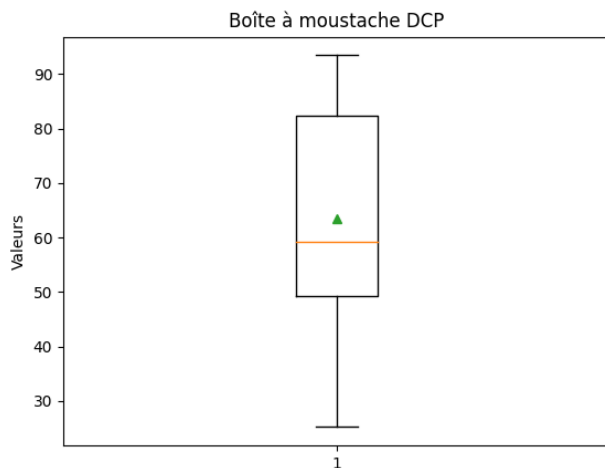


Médiane : 71.50 – Moyenne : 145.7 – Minimum 4.0 – Maximum 429.0 – Q1 : 12.0 – Q3 : 180.0

Analyse de la distribution : Également décroissante, si on zoom sur l'image on s'aperçoit que le nombre de ligne de code est très bas pour une très grosse majorité des fichiers, puis une légère augmentation dans le nombre de fichiers possédant aux alentours de 75 lignes de codes pour ensuite baisser. Si on faisait abstraction du premier pic, on aurait une distribution normale légèrement aplatie à droite.



DCP :



Médiane : 59.2 – Moyenne : 63.5 – Minimum 25.2 – Maximum 93.4 – Q1 : 49.3 – Q3 : 82.3

Analyse distribution : Une majorité des fichiers a une densité de commentaires autour de 90%, le restant des fichiers se situent plus à l'entour de 50%. Si on fait abstraction du pic de fin, on pourrait voir une très nette distribution normale.

TÂCHE 2 :

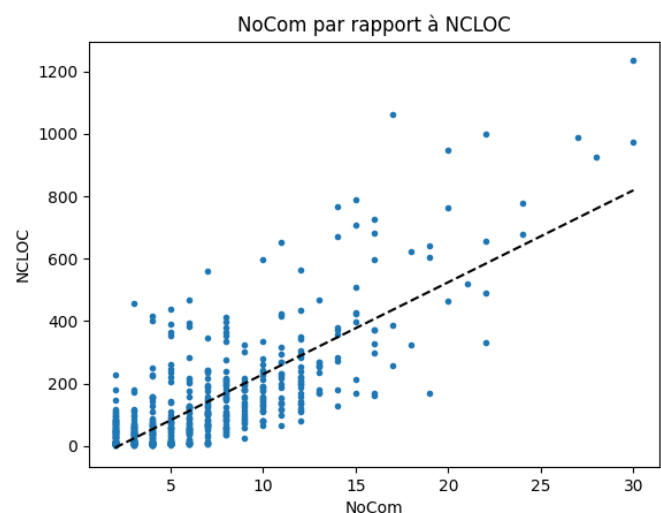
(Nous avons choisi arbitrairement 5 chiffres significatifs pour les valeurs affichée ci-bas.)

Corrélation entre NoCom & NCLOC

On voit sur le graphique que les points semblent corrélés, la droite de régression (sans les points considérés comme extrêmes, c-à-d NCLOC > 1400) a comme équation :

$$y = 36.047x - 97.992$$

Elle semble bien donner une direction générale aux points, donc une corrélation n'est pas à rejeter.



Vu qu'on ne peut pas affirmer que les métriques suivent une distribution normale, nous devons passer par le coefficient de corrélation de Spearman. Il vaut : $\rho = 0.68803$

Encore une fois, les métriques semblent corrélées mais le lien n'est pas très robuste. Donc on peut en déduire que le nombre de commit par fichier semble corrélé avec le nombre de lignes de codes (non vides ni commentaire) par fichier.

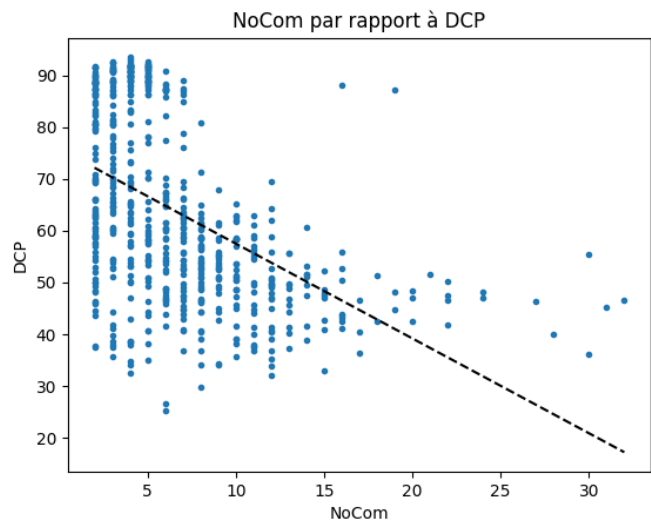
Cette corrélation ne serait pas absurde car avec un commit en général on rajoute du code à un fichier, mais il faut prendre en compte le refactoring, ou encore les commits concentré sur les commentaires, sans oublier les commits qui retirent du code qui n'est plus ou pas utilisé (cette réflexion ne montre pas qu'il y a une corrélation, juste que si elle existait, ceci pourrait la justifier).

Corrélation entre NoCom & DCP

Les points qui représentent le nombre de commits d'un fichier par rapport à la densité de commentaire de ce fichier ne semble pas vraiment corrélé à première vue. La droite de régression calculée est la suivante :

$$y = -1.8300 + 75.809$$

Elle ne semble pas vraiment donner une direction générale, les points sont trop éparpillés. Aucun point n'a été exclu pour améliorer la régression car il n'y a, à première vue, aucun points (ou au moins, un nombre suffisamment petit de point) qui changeraient drastiquement la direction de la droite.



Ceci se confirme avec le calcul du coefficient de Spearman. Celui-ci vaut : $\rho = -0.53352$ Il indique une très légère corrélation linéaire inverse, mais clairement pas assez suffisant pour tirer une conclusion positive sur la corrélation entre le nombre de commits et la densité de commentaire par fichier. Cela peut se comprendre, car ce n'est pas parce qu'on fait un commit qu'on rajoute plus ou moins de commentaires (cette réflexion ne montre pas qu'il ne peut pas y avoir de corrélation, juste qu'elle est peu probable).

TÂCHE 3 :

Conception d'une quasi-expérience :

Nous voulons étudier l'échantillon de données du projet *jfreechart* grâce aux données fournies dans le fichier *jfreechart-stats.csv*.

Le choix d'une quasi-expérience est pertinent car il est impossible de tester physiquement l'hypothèse.

L'hypothèse est donnée : " les classes qui ont été modifiées plus de 10 fois sont mieux commentées que celles qui ont été modifiées moins de 10 fois."

On peut définir *NoCom* en tant que variable indépendante car on peut la manipuler en choisissant les fichiers avec des nombres de commit différents pour évaluer notre hypothèse. *NoCom* est pertinent car chaque commit effectué correspond à un changement dans le code.

On peut définir *DCP* en tant que variable dépendante pour voir la variation suivant *NoCom*. *DCP* est pertinent car il correspond à la densité de commentaire.

La quasi-expérience consiste ensuite à calculer le taux de corrélation entre *NoCom* et *DCP*, si ce taux de corrélation se trouve entre -1 et -0.5 on peut conclure que nos deux variables sont corrélées mais que lorsque *NoCom* augmente, la densité de commentaire diminue. Si le taux de corrélation se trouve entre -0.5 et 0.5 alors nos deux variables ne sont pas corrélées. Dans ces deux cas on peut conclure que l'hypothèse est réfutée. Finalement si le taux de corrélation se trouve entre 0.5 et 1 alors nos deux variables sont corrélées et l'hypothèse est validée pour *jfreechart*. On peut également simplement regarder les fichiers avec moins de 10 *NoCom* et comparer les *DCP* avec les fichiers de 10 *NoCom* ou plus. Cette méthode répondrait plus précisément à la question mais il serait plus difficile de généraliser le résultat.

Il existe des menaces à la validité de construction. Effectivement la densité de commentaire ne prend pas en compte la qualité de ces commentaires. Il ne semble pas avoir de menaces à la validité interne ou externe. Il existe cependant une menace à la validité de conclusion puisque le nombre de commit n'est pas le seul facteur lié à la densité de commentaire.

Nous avons donc effectué la quasi-expérience en calculant le taux de corrélation entre *NoCom* et *DCP* et nous avons trouvé un taux de corrélation de -0.533 , ce qui montre que dans notre cas *NoCom* et *DCP* sont très légèrement corrélées négativement. Ainsi plus le nombre de commit augmente plus la densité de commentaire a une légère tendance à légèrement diminuer. On réfute donc notre hypothèse.