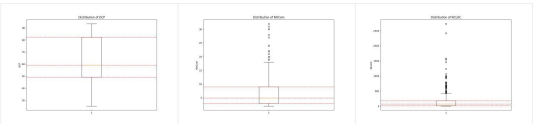


Question N.01

```
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd

df = pd.read_csv('/Users/ethanmai/Downloads/freeschart-stats.csv')
data = df[['NOCm']]
#creating plot
plt.boxplot(data)
#show the median of the data
plt.axhline(data.median(), color='r', linestyle='dashed', linewidth=3)
#show the first quartile of the data
plt.axhline(data.quantile(0.25), color='r', linestyle='dashed', linewidth=1)
#show the third quartile of the data
plt.axhline(data.quantile(0.75), color='r', linestyle='dashed', linewidth=1)
```



Question 2

```
NOCm = df[['NOCm']]
NCLOC = df[['NCLOC']]

#remove outliers with IQR
q1 = NOCom.quantile(0.25)
q3 = NOCom.quantile(0.75)
Iqr = q3 - q1
NOCm = NOCom[~((NOCm < (q1 - 1.5 * Iqr)) | (NOCm > (q3 + 1.5 * Iqr)))]

q1 = NCLOC.quantile(0.25)
q3 = NCLOC.quantile(0.75)
Iqr = q3 - q1
NCLOC = NCLOC[~((NCLOC < (q1 - 1.5 * Iqr)) | (NCLOC > (q3 + 1.5 * Iqr)))]

#remove the missing values
if len(NCLOC) > len(NOCm):
    NCLOC = NCLOC[:len(NOCm)]
else:
    NOCom = NOCom[:len(NCLOC)]

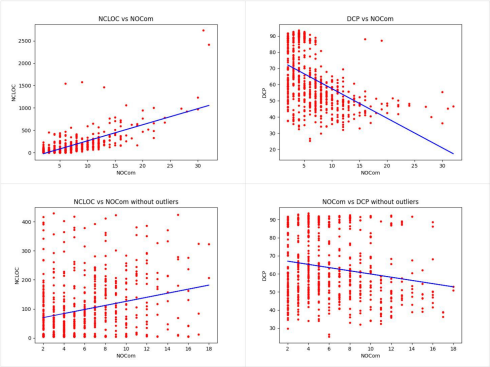
diff = len(NCLOC) - len(NOCm)

#plot the data
plt.scatter(NOCm, NCLOC, color='red', markers='s', s=10)
#calculate the correlation coefficient
corr = np.corrcoef(NOCm, NCLOC)
print(corr)

#plot the regression line
a, b = np.polyfit(NOCm, NCLOC, 1)
plt.plot(NOCm, m=NOCm, b=NCLOC, color='blue')

#add title
plt.title("NCLOC vs DCP without outliers")
#add label
plt.xlabel("NOCm")
#add label
plt.ylabel("NCLOC")
plt.show()
```

Dataset	correlation	m	b
NOCm vs DCP	-0.487753	-1.8300216935322569	75.80952917110602
NOCm vs DCP without outliers	-0.12691989	-0.6289842618132071	67.53760367576922
NOCm vs NCLOC	0.71457166	-36.946998232489396	-87.59249915923612
NOCm vs NCLOC without outliers	0.24695678	7.000198678430157	55.9549436446627



Question 3

L'hypothèse est la suivante: « les classes qui ont été modifiées plus de 10 fois sont mieux commentées que celles qui ont été modifiées moins de 10 fois »

Afin de valider ou réfuter cette hypothèse nous allons regarder les variables suivantes:

- NOCm : le nombre de commit
- DCP : la densité des commentaires

Ici on regarde DCP et pas NCLOC car seulement la densité des commentaires dans le code est pertinente (car un petit fichier peut avoir un plus petit nombre de commentaire qu'un gros fichier tout en ayant une plus grosse densité de commentaire)

Si l'hypothèse est valide, on peut s'attendre à une corrélation positive entre NOCom et DCP.

Or selon les calculs, on observe une corrélation de -0.487753 et -0.12691989 si on retire les valeurs aberrantes. Aussi, le DCP moyen pour les fichier dont le nombre de commit est inférieur à 10 est de 63.629675572519076 contre 62.70324074074074 pour les fichiers modifiés plus de dix fois ce qui réfute l'hypothèse.

Les menaces à la validité:

Dans cette analyse le manque de contexte peut nuire à la validité de notre conclusion: Ici, nous avons choisi de comparer le nombre de modifications d'un fichier à sa densité de commentaires mais nous n'avons pas pris en compte d'autres facteurs qui peuvent influencer la documentation d'un programme. Par exemple, la complexité d'un programme:

Si un programme est très complexe à développer, alors grandes sont les chances qu'il a été modifié plusieurs fois, sans pour autant avoir été commenté proportionnellement. Ou encore, la taille d'un fichier: Si un fichier est relativement petit, alors chaque ligne de commentaire augmente significativement la densité de commentaires par ailleurs ce n'est pas imaginable d'un petit fichier exige moins de modifications lors de son implémentation. Tout ces facteurs ignorés lors de cette analyse peut donc nuire à sa validité. Pour assurer une conclusion plus valide il faudrait donc ajouter ces facteurs dans notre analyse afin d'avoir un meilleur contexte.