



```
[12]: # replace NaN with median value
b = pd.DataFrame(mpg_df.horsepower, dtype="category")
imp_median = SimpleImputer(strategy="median")
f_median = imp_median.fit_transform(b)
print(f_median)

[130.]
[165.]
[150.]
[150.]
[145.]
[120.]
[170.]
[225.]
[190.]
[170.]
[160.]
[150.]
[225.]
[95.]
[95.]
[97.]
[78.]
[88.]
[46.]
[87.]
[90.]
[95.]
[113.]
[90.]
[215.]
[200.]
[140.]
[193.]
[88.]
[90.]
[95.]
[93.5]
[100.]
[90.]
[100.]
[88.]
[100.]
[165.]
[175.]
[153.]
[150.]
[180.]
[95.]
[110.]
[72.]
[150.]
[88.]
[90.]
[70.]
[76.]
[65.]
[130.]
[60.]
[70.]
[95.]
[80.]
[54.]
[90.]
[165.]
[175.]
[150.]
[153.]
[150.]
[208.]
[155.]
[160.]
[190.]
[95.]
[150.]
[130.]
[140.]
[150.]
[112.]
[76.]
[97.]
[86.]
[69.]
[86.]
[92.]
[97.]
[80.]
[89.]
[175.]
[150.]
[145.]
[170.]
[150.]
[198.]
[150.]
[158.]
[150.]
[175.]
[105.]
[100.]
[88.]
[95.]
[46.]
[150.]
[167.]
[170.]
[180.]
[100.]
[75.]
[72.]
[105.]
[75.]
[75.]
[97.]
[70.]
[88.]
[95.]
[115.]
[53.]
[86.]
[81.]
[92.]
[79.]
[83.]
[140.]
[150.]
[120.]
[152.]
[100.]
[105.]
[90.]
[52.]
[60.]
[70.]
[53.]
[100.]
[78.]
[110.]
[95.]
[71.]
[70.]
[102.]
[150.]
[89.]
[108.]
[120.]
[180.]
[145.]
[130.]
[150.]
[68.]
[80.]
[58.]
[70.]
[145.]
[145.]
[145.]
[130.]
[110.]
[105.]
[100.]
[98.]
[180.]
[170.]
[190.]
[149.]
[78.]
[88.]
[75.]
[89.]
[63.]
[83.]
[67.]
[78.]
[97.]
[110.]
[110.]
[48.]
[66.]
[52.]
[70.]
[60.]
[110.]
[150.]
[139.]
[105.]
[95.]
[85.]
[88.]
[100.]
[90.]
[105.]
[85.]
[120.]
[145.]
[165.]
[139.]
[140.]
[68.]
[95.]
[97.]
[75.]
[85.]
[105.]
[85.]
[97.]
[103.]
[125.]
[115.]
[71.]
[68.]
[115.]
[85.]
[88.]
[90.]
[110.]
[130.]
[129.]
[138.]
[135.]
[155.]
[142.]
[125.]
[150.]
[71.]
[65.]
[80.]
[77.]
[125.]
[71.]
[90.]
[70.]
[70.]
[65.]
[69.]
[90.]
[115.]
[115.]
[90.]
[76.]
[60.]
[65.]
[90.]
[90.]
[88.]
[79.]
[92.]
[75.]
[65.]
[105.]
[65.]
[48.]
[48.]
[67.]
[67.]
[67.]
[93.5]
[62.]
[62.]
[132.]
[150.]
[88.]
[93.5]
[72.]
[84.]
[84.]
[92.]
[110.]
[58.]
[64.]
[60.]
[67.]
[65.]
[62.]
[68.]
[63.]
[65.]
[65.]
[74.]
[93.5]
[75.]
[100.]
[74.]
[80.]
[76.]
[116.]
[120.]
[110.]
[105.]
[110.]
[88.]
[88.]
[85.]
[88.]
[84.]
[90.]
[92.]
[93.5]
[74.]
[68.]
[68.]
[63.]
[70.]
[88.]
[75.]
[67.]
[67.]
[110.]
[85.]
[92.]
[112.]
[96.]
[84.]
[90.]
[86.]
[52.]
[84.]
[79.]
[82.]]

In [13]: # variance after replacing with the median value
var_median = f_median.var()
print(var_median)

1457.2982752960781

In [14]: # replace NaN with mode value
c = pd.DataFrame(mpg_df.horsepower, dtype="category")
imp_mode = SimpleImputer(strategy="most_frequent")
f_mode = imp_mode.fit_transform(c)
print(f_mode)

[130.]
[165.]
[150.]
[150.]
[140.]
[198.]
[220.]
[215.]
[225.]
[190.]
[170.]
[160.]
[150.]
[225.]
[95.]
[95.]
[97.]
[78.]
[88.]
[46.]
[87.]
[90.]
[95.]
[113.]
[215.]
[200.]
[193.]
[88.]
[90.]
[95.]
[100.]
[100.]
[165.]
[175.]
[150.]
[153.]
[180.]
[95.]
[110.]
[72.]
[150.]
[88.]
[90.]
[70.]
[76.]
[65.]
[130.]
[60.]
[70.]
[95.]
[80.]
[54.]
[90.]
[165.]
[175.]
[150.]
[153.]
[150.]
[208.]
[155.]
[160.]
[190.]
[95.]
[150.]
[130.]
[140.]
[150.]
[112.]
[76.]
[97.]
[86.]
[69.]
[86.]
[92.]
[97.]
[80.]
[89.]
[175.]
[150.]
[145.]
[170.]
[150.]
[198.]
[150.]
[158.]
[150.]
[175.]
[105.]
[100.]
[88.]
[95.]
[46.]
[150.]
[167.]
[170.]
[180.]
[100.]
[75.]
[72.]
[105.]
[75.]
[75.]
[97.]
[70.]
[88.]
[95.]
[115.]
[53.]
[86.]
[81.]
[92.]
[79.]
[83.]
[140.]
[150.]
[120.]
[152.]
[100.]
[105.]
[90.]
[52.]
[60.]
[70.]
[53.]
[100.]
[78.]
[110.]
[95.]
[71.]
[70.]
[102.]
[150.]
[89.]
[108.]
[120.]
[180.]
[145.]
[130.]
[150.]
[68.]
[80.]
[58.]
[70.]
[145.]
[145.]
[145.]
[130.]
[110.]
[105.]
[100.]
[98.]
[180.]
[170.]
[190.]
[149.]
[78.]
[88.]
[75.]
[89.]
[63.]
[83.]
[67.]
[78.]
[97.]
[110.]
[110.]
[48.]
[66.]
[52.]
[70.]
[60.]
[110.]
[150.]
[139.]
[105.]
[95.]
[85.]
[88.]
[100.]
[90.]
[105.]
[85.]
[120.]
[145.]
[165.]
[139.]
[140.]
[68.]
[95.]
[97.]
[75.]
[85.]
[105.]
[85.]
[97.]
[103.]
[125.]
[115.]
[71.]
[68.]
[115.]
[85.]
[88.]
[90.]
[110.]
[130.]
[129.]
[138.]
[135.]
[155.]
[142.]
[125.]
[150.]
[71.]
[65.]
[80.]
[77.]
[125.]
[71.]
[90.]
[70.]
[70.]
[65.]
[69.]
[90.]
[115.]
[115.]
[90.]
[76.]
[60.]
[65.]
[90.]
[90.]
[88.]
[79.]
[92.]
[75.]
[65.]
[105.]
[65.]
[48.]
[48.]
[67.]
[67.]
[67.]
[93.5]
[62.]
[62.]
[132.]
[150.]
[88.]
[93.5]
[72.]
[84.]
[84.]
[92.]
[110.]
[58.]
[64.]
[60.]
[67.]
[65.]
[62.]
[68.]
[63.]
[65.]
[65.]
[74.]
[93.5]
[75.]
[100.]
[74.]
[80.]
[76.]
[116.]
[120.]
[110.]
[105.]
[110.]
[88.]
[88.]
[85.]
[88.]
[84.]
[90.]
[92.]
[93.5]
[74.]
[68.]
[68.]
[63.]
[70.]
[88.]
[75.]
[67.]
[67.]
[110.]
[85.]
[92.]
[112.]
[96.]
[84.]
[90.]
[86.]
[52.]
[84.]
[79.]
[82.]]

In [15]: # variance after replacing with the mode value
var_mode = f_mode.var()
print(var_mode)

1486.29231585061

In [16]: # find lowest variance from mean, median and mode imputing strategies
m1(var_mean,var_median,var_mode)

Out[16]: 1455.511639831812

Conclusion:
Among mean, median and mode imputation methods mean have a lowest variance which is 1455.511639831812. Whereas in with median and mode method, we get 1457.2982752960781 and 1486.29231585061 variance, respectively.

Variance means how far the given data is spread out. Zero variance means data values in sample are identical. When data is very close to the mean, variance become very low. So, when we replace the NaN with mean value, we get lowest variance.

Problem 3

In [17]: from sklearn import datasets
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
```



