

Note: every value less than 1000 handled as rubbish value

1-We have 13 feature, 148653 sample of data(rows) There is no missing data in any dependent column, but there is 609 values missing in basepay feature, 4 missing values in overtimePay, 4 missing values in otherPay and 36163 missing values in benefits.

After deleting rubbish values: 145060 sample of data(rows)  
There is no missing data in any dependent column, but there is 215 values missing in basepay feature, 0 missing values in overtimePay, 0 missing values in otherPay and 35363 missing values in benefits feature.

2-mean= 74768.32197169267

median= 71426.60999999999

mode= 0

min= -618.13

max= 567595.43

range= 568213.56

standard deviation= 50517.00527394987

after deleting rubbish values:

mean= 74954.38346771419

median= 71553.0

mode= 18594.0

min= 0.3

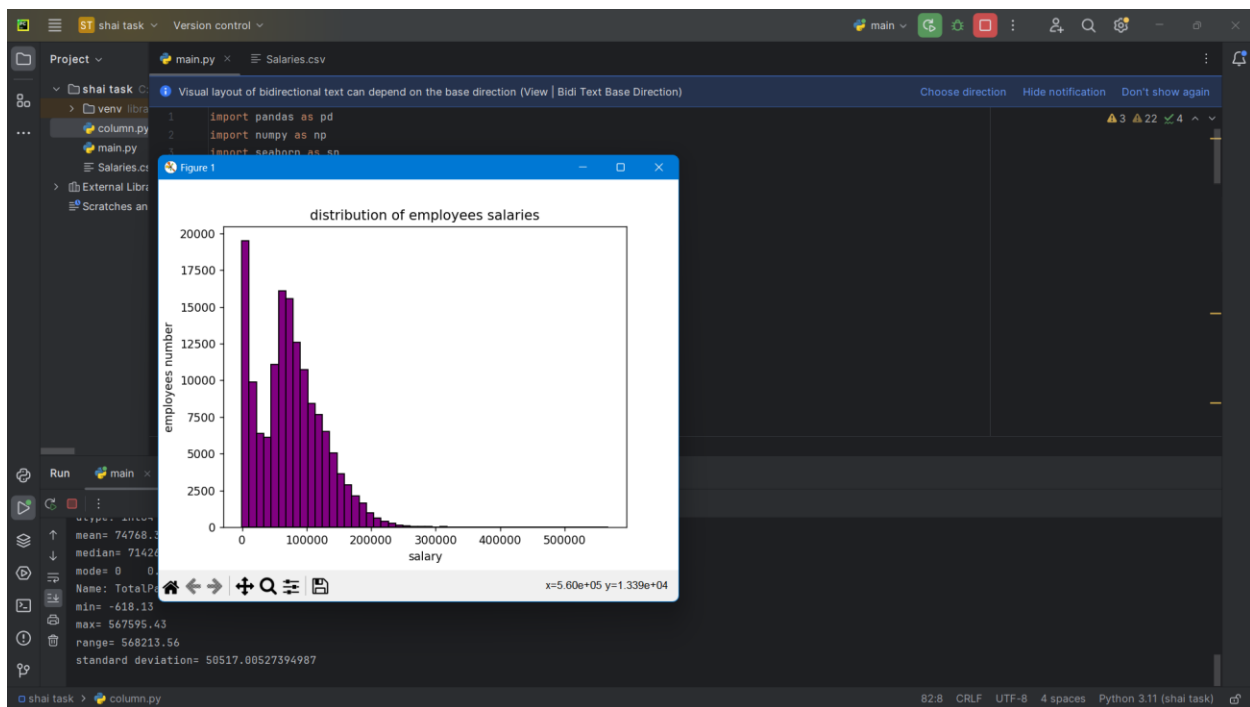
max= 567595.43

range= 567595.13

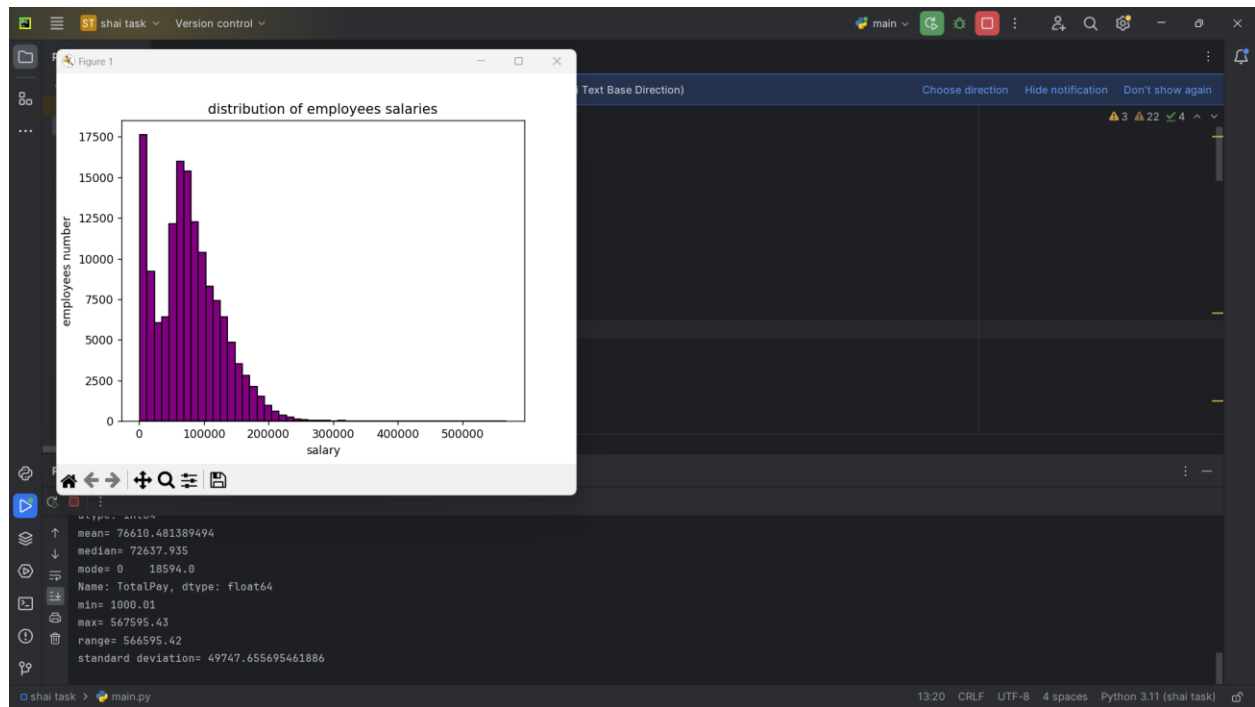
standard deviation= 50441.7662188465

3- I used mean to fill the missing data because using mode is not practical because its value=0, note: we can also use median

4-value 0 is the most frequent value and the period around 100000 containing the largest number of employees and after 220000 there is less than 1% of the total number of employees.

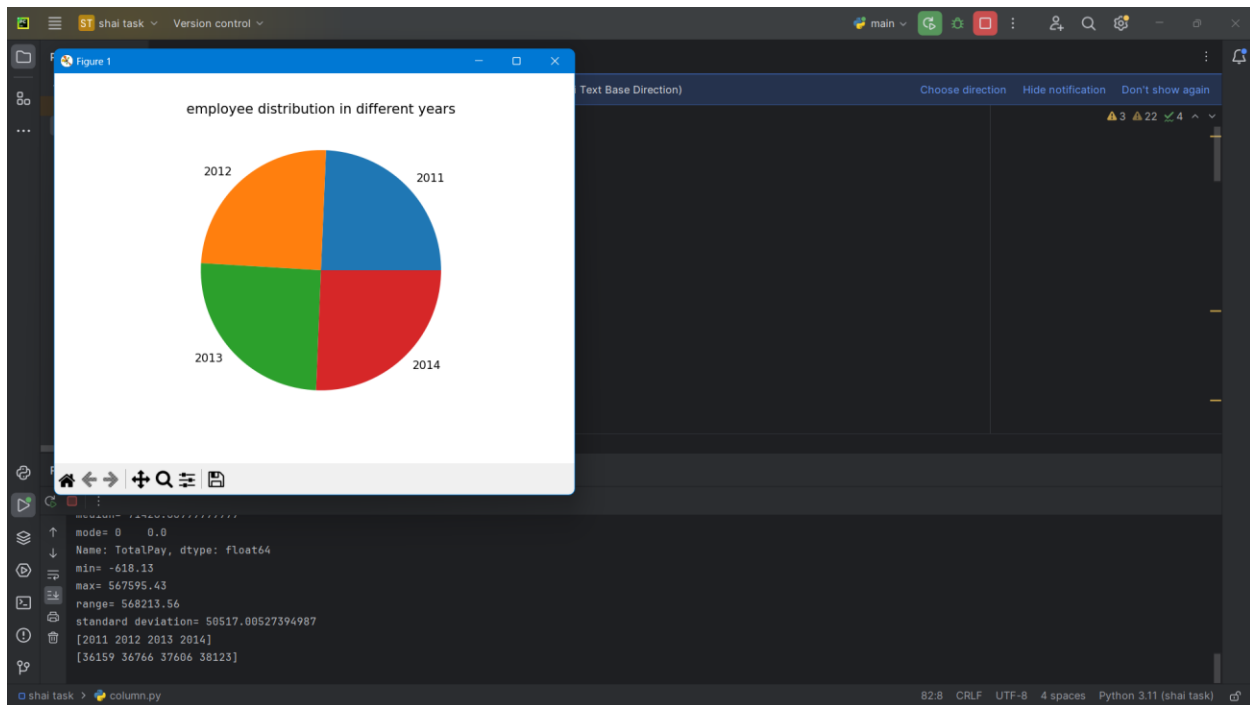


After deleting rubbish data:

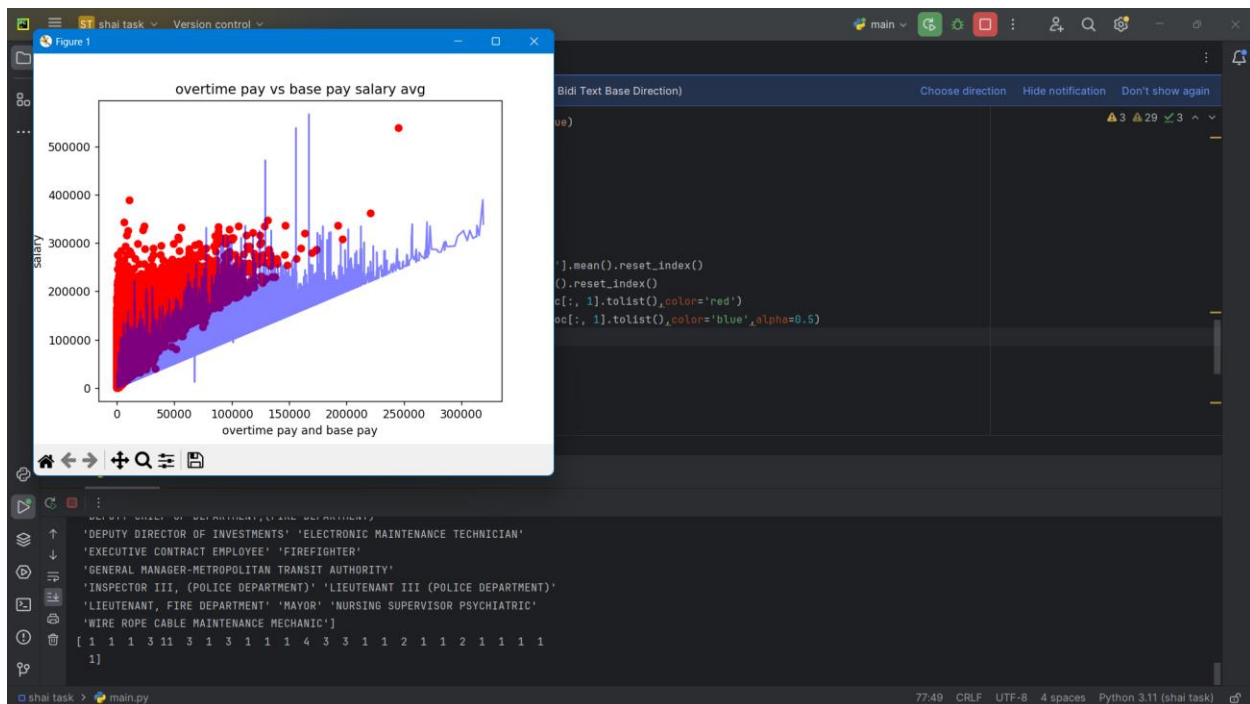


The rest of the charts havent any noticable change.

And the proportion of employees in different years is approximately one quarter per year from 2011 to 2014.



5)



6) total pay have high dependency with base pay but have much lower dependency with other pay.

