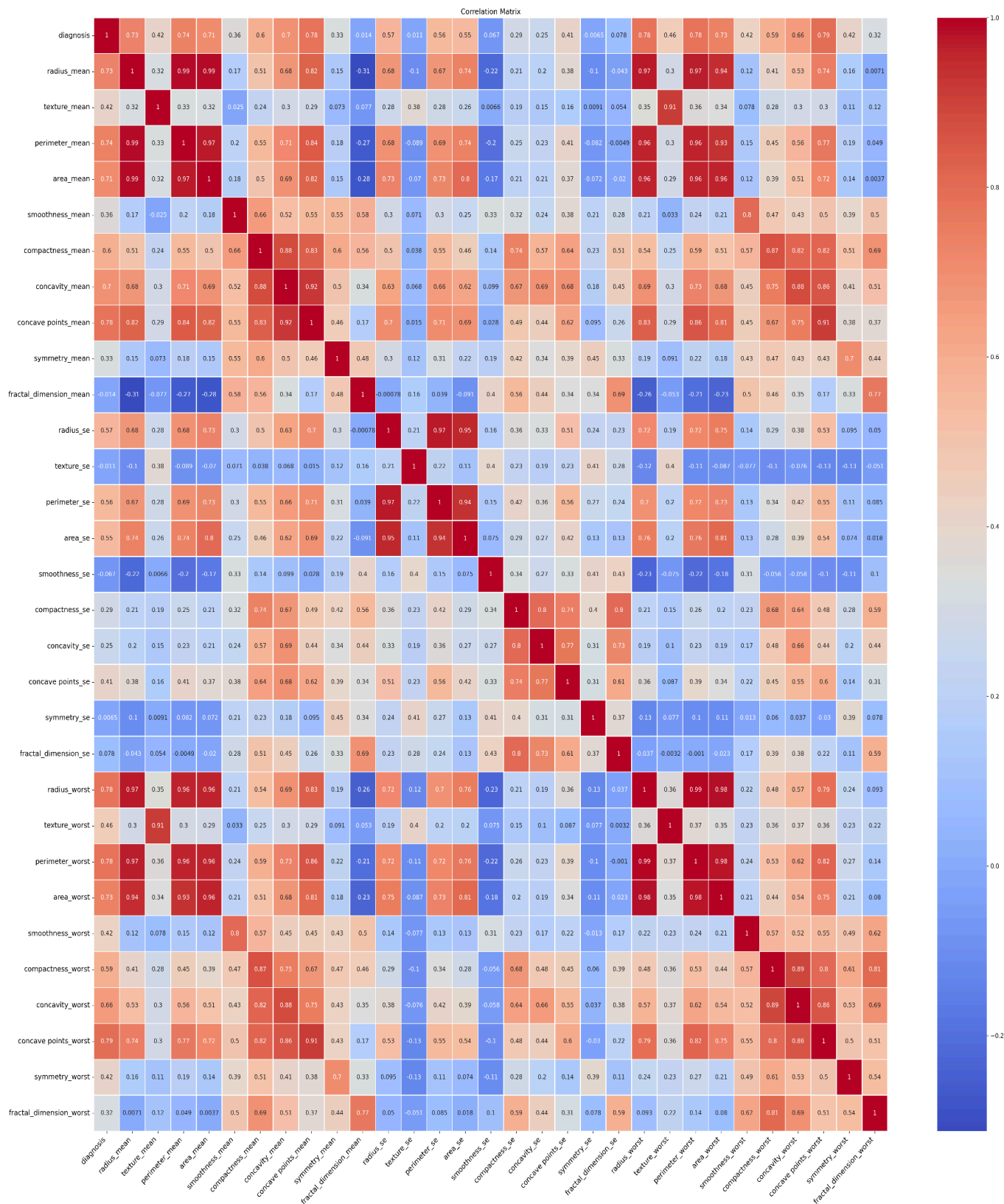## Project Report -1
## Adult Census Data

**1. Provide brief details about the nature of your dataset. What is it about? What type of data are we encountering? How many entries and variables does the dataset comprise?**

1. The WDBC dataset represent the characteristics of the cell nuclei present in the image of brest tissue . The end goal is to find whether it is benign or malignant

2. We are encountering with 31 numerical features like size , shape and texture of the cell , smoothness and 1 categorical features with B or M .

3. The dataset contains 569 entries and 32 columns .

**2. Provide the main statistics about the dataset entries (mean, std, number of missing values, etc.)**

```
data.describe().T
```

|  | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| diagnosis | 569.0 | 0.372583 | 0.483918 | 0.000000 | 0.000000 | 0.000000 | 1.000000 | 1.00000 |
| radius_mean | 569.0 | 14.127292 | 3.524049 | 6.981000 | 11.700000 | 13.370000 | 15.780000 | 28.11000 |
| texture_mean | 569.0 | 19.289649 | 4.301036 | 9.710000 | 16.170000 | 18.840000 | 21.800000 | 39.28000 |
| perimeter_mean | 569.0 | 91.845255 | 24.091002 | 43.790000 | 75.210000 | 86.180000 | 103.800000 | 188.50000 |
| area_mean | 569.0 | 654.889104 | 351.914129 | 143.500000 | 420.300000 | 551.100000 | 782.700000 | 2501.00000 |
| smoothness_mean | 569.0 | 0.096404 | 0.014038 | 0.052630 | 0.086410 | 0.095920 | 0.105300 | 0.16340 |
| compactness_mean | 569.0 | 0.104341 | 0.052813 | 0.019380 | 0.064920 | 0.092630 | 0.130400 | 0.34540 |
| concavity_mean | 569.0 | 0.088799 | 0.079720 | 0.000000 | 0.029560 | 0.061540 | 0.130700 | 0.42680 |
| concave points_mean | 569.0 | 0.048919 | 0.038803 | 0.000000 | 0.020310 | 0.033500 | 0.074000 | 0.20120 |
| symmetry_mean | 569.0 | 0.181131 | 0.027428 | 0.106000 | 0.161900 | 0.179200 | 0.195700 | 0.30400 |
| fractal_dimension_mean | 569.0 | 0.062790 | 0.007065 | 0.049960 | 0.057690 | 0.061540 | 0.066120 | 0.09744 |
| radius_se | 569.0 | 0.405172 | 0.277313 | 0.111500 | 0.232400 | 0.324200 | 0.478900 | 2.87300 |
| texture_se | 569.0 | 1.215555 | 0.551626 | 0.360200 | 0.833900 | 1.108000 | 1.473000 | 4.88500 |
| perimeter_se | 569.0 | 2.866059 | 2.021855 | 0.757000 | 1.606000 | 2.287000 | 3.357000 | 21.98000 |
| area_se | 569.0 | 40.337079 | 45.491006 | 6.802000 | 17.850000 | 24.530000 | 45.190000 | 542.20000 |
| smoothness_se | 569.0 | 0.007041 | 0.003003 | 0.001713 | 0.005169 | 0.006380 | 0.008146 | 0.03113 |
| compactness_se | 569.0 | 0.025449 | 0.017918 | 0.002252 | 0.012950 | 0.020420 | 0.032450 | 0.13540 |
| concavity_se | 569.0 | 0.031841 | 0.030226 | 0.000000 | 0.014980 | 0.025890 | 0.042050 | 0.39600 |
| concave points_se | 569.0 | 0.011796 | 0.006170 | 0.000000 | 0.007638 | 0.010930 | 0.014710 | 0.05279 |
| symmetry_se | 569.0 | 0.020542 | 0.008266 | 0.007882 | 0.015160 | 0.018730 | 0.023480 | 0.07895 |
| fractal_dimension_se | 569.0 | 0.003795 | 0.002646 | 0.000895 | 0.002248 | 0.003187 | 0.004558 | 0.02984 |
| radius_worst | 569.0 | 16.269190 | 4.833242 | 7.930000 | 13.010000 | 14.970000 | 18.790000 | 36.04000 |
| texture_worst | 569.0 | 25.677223 | 6.146258 | 12.020000 | 21.080000 | 25.410000 | 29.720000 | 49.54000 |
| perimeter_worst | 569.0 | 107.261213 | 33.602542 | 50.410000 | 84.110000 | 97.660000 | 125.400000 | 251.20000 |
| area_worst | 569.0 | 880.583128 | 569.356993 | 185.200000 | 515.300000 | 686.500000 | 1084.000000 | 4254.00000 |
| smoothness_worst | 569.0 | 0.132369 | 0.022832 | 0.071170 | 0.116600 | 0.131300 | 0.146000 | 0.22260 |
| compactness_worst | 569.0 | 0.254265 | 0.157336 | 0.027290 | 0.147200 | 0.211900 | 0.339100 | 1.05800 |
| concavity_worst | 569.0 | 0.272188 | 0.208624 | 0.000000 | 0.114500 | 0.226700 | 0.382900 | 1.25200 |
| concave points_worst | 569.0 | 0.114606 | 0.065732 | 0.000000 | 0.064930 | 0.099930 | 0.161400 | 0.29100 |
| symmetry_worst | 569.0 | 0.290076 | 0.061867 | 0.156500 | 0.250400 | 0.282200 | 0.317900 | 0.66380 |
| fractal_dimension_worst | 569.0 | 0.083946 | 0.018061 | 0.055040 | 0.071460 | 0.080040 | 0.092080 | 0.20750 |

## 3. Provide the correlation matrix and write your inference about it.



Correlation Matrix

1.Features like radius_mean, texture_mean, perimeter_mean, and area_mean show higher mean values compared to other features.

2.In the correlation matrix concave_point_worst , concavity_point_mean are highly correlated with diagnosis . That is it strongly predict the malignancy .

3.Other than concave_point_worst, perimeter_mean , radius_mean, area_mean ,area_worst, perimeter_worst,radius worst are strongly correlated with malignancy .

4. From the row radius_mean the perimeter_mean , area_mean are highly correlated to each other . This indicates that it contains the redundant values .

To reduce the redundancy we need to take PCA to combine into one feature .

5. These redundant data contains in perimeter_mean , area_mean,radius_worst,perimeter_worst columns respectively .

6.concavity_se , compactness_se show weak correlations with diagnosis.

7. fractional_dimension_mean  is likely associated with benign tumors .