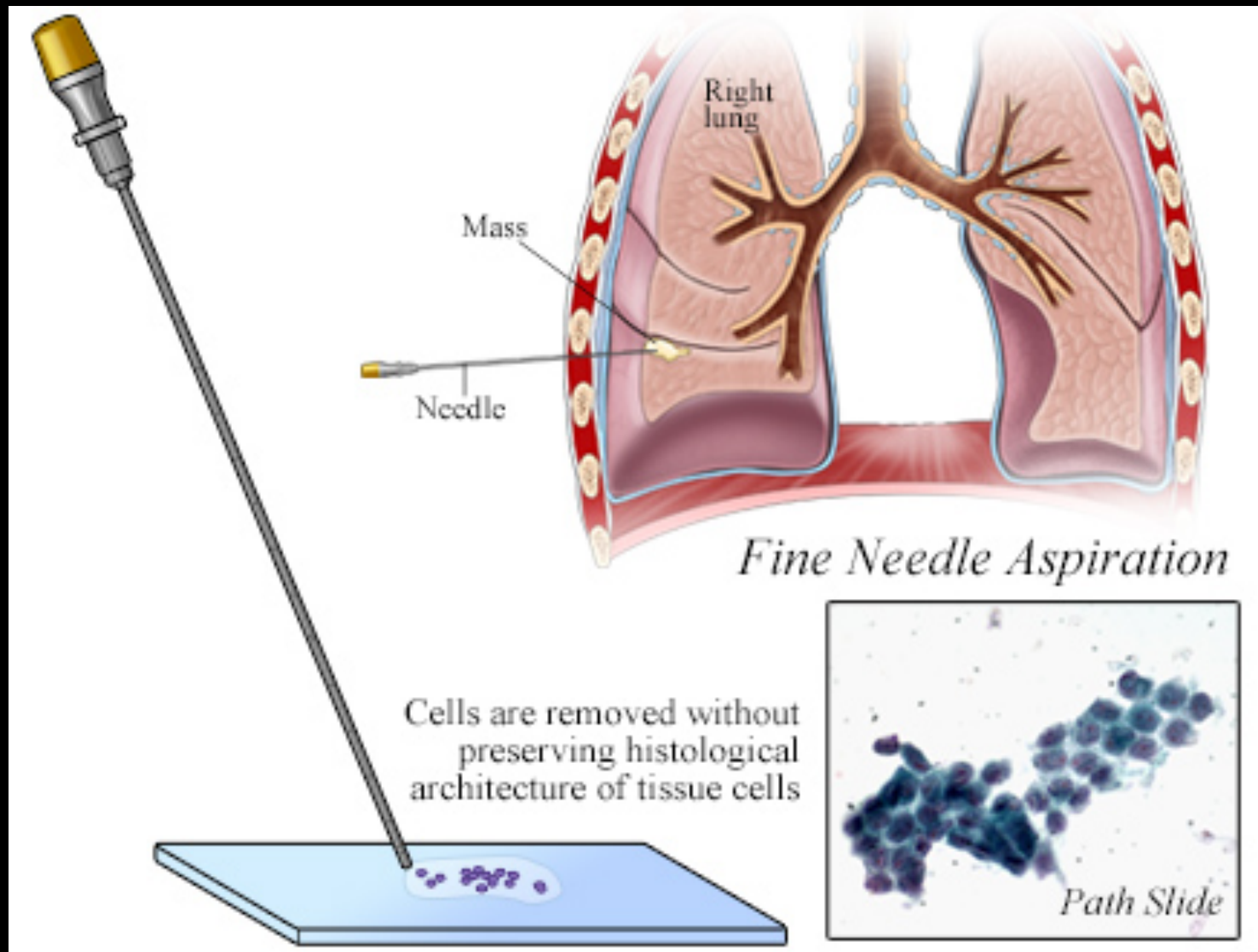


Model to Predict tumor Cell Classification based on Wisconsin Breast Cancer Data Set

Murali Satuluri

Problem



- **Fine Needle Aspiration is a less invasive alternative to Biopsy.**
- **Cells collected from this test are studied and their features are recorded.**
- **The features of the cell are to be used to predict if the parent tissue is malignant or benign.**

Variables

Input Variables

Output Variable

1.Clump Thickness

2.Uniformity of Cell Sizes

3.Uniformity of Cell Shape

4.Marginal Adhesion

5.Single Epithelial Cell Size

6.Bare Nuclei

7.Bland Chromatin

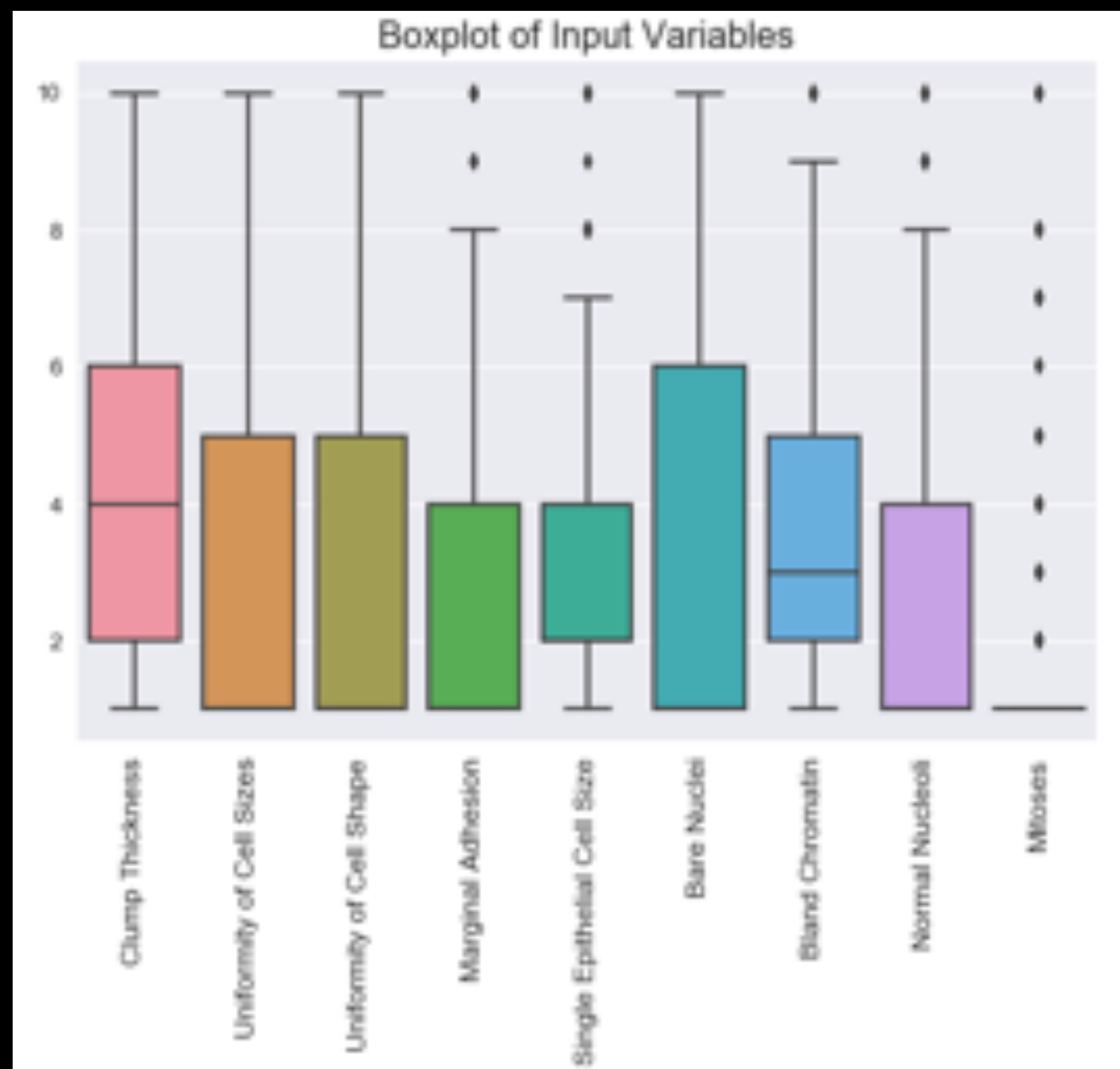
8.Normal Nucleoli

9.Mitoses

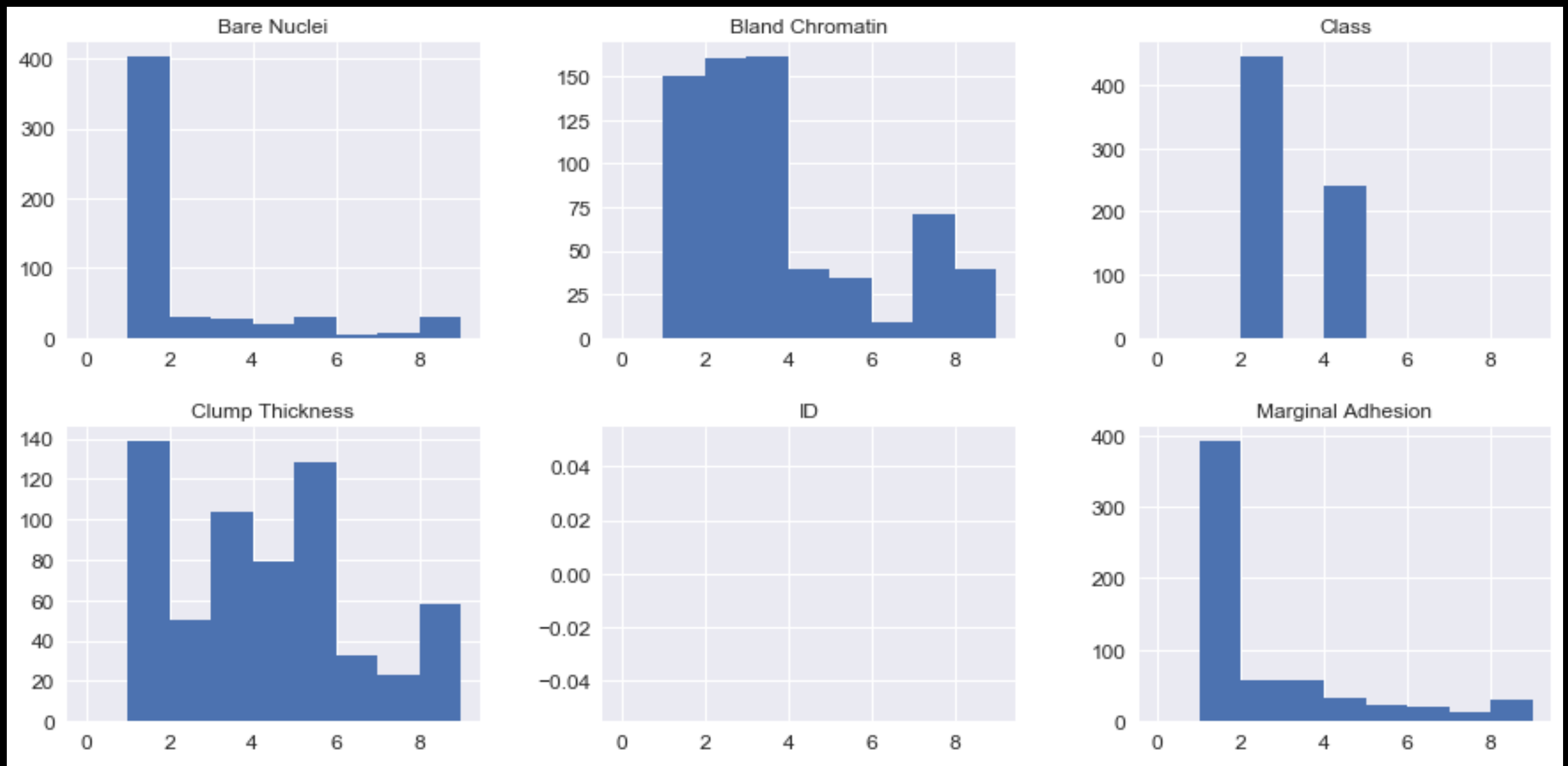
**Cell Classification
(i.e. Benign or Malignant)**

Data Cleaning

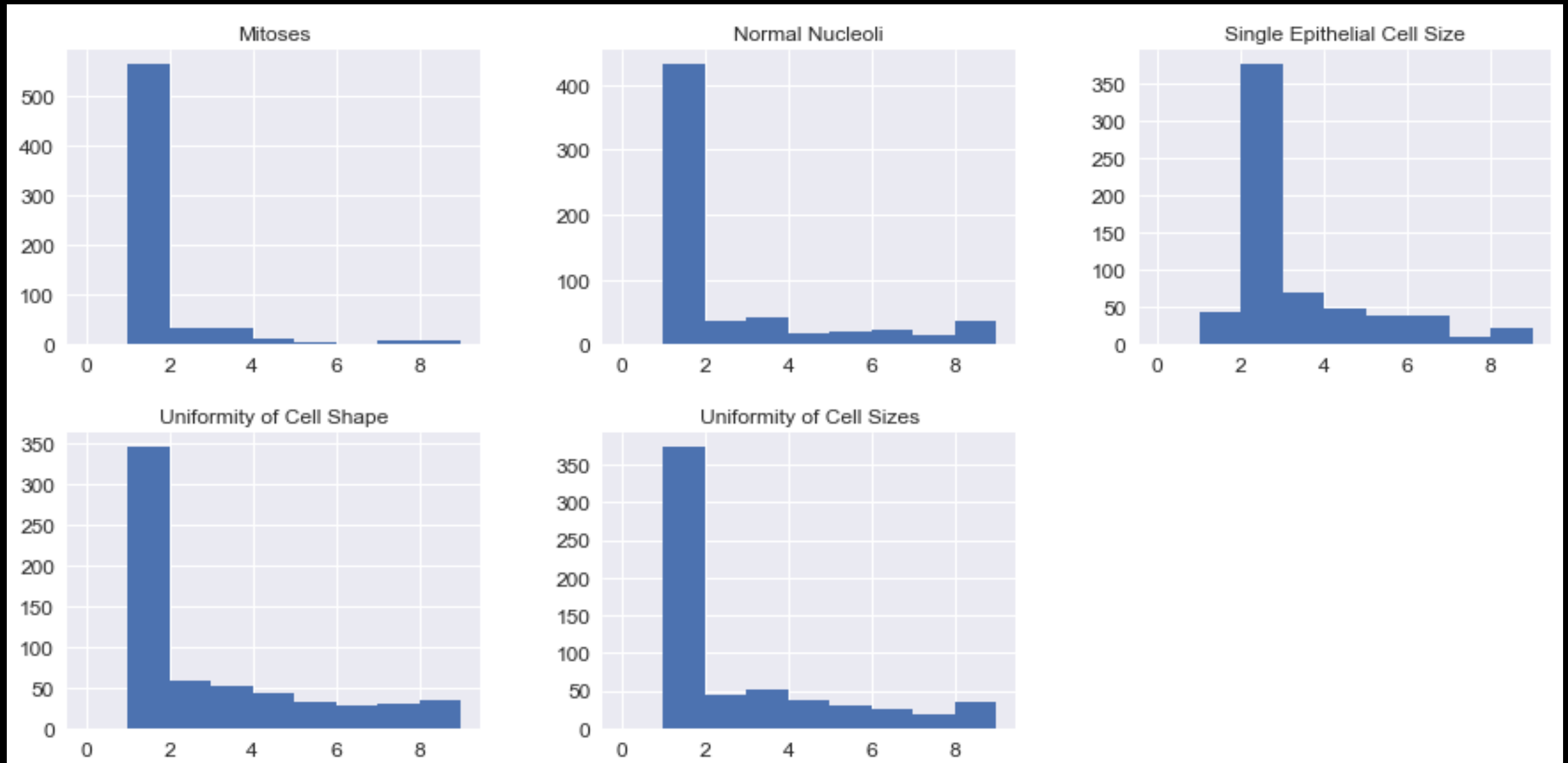
- There are 699 samples. 16 of them containing NA values were removed.
- Boxplots of the remaining data are shown below.
- Variables are not normally distributed. Since there is no clear basis on which the extreme data can be called outliers, no further data were removed.



Input Variables

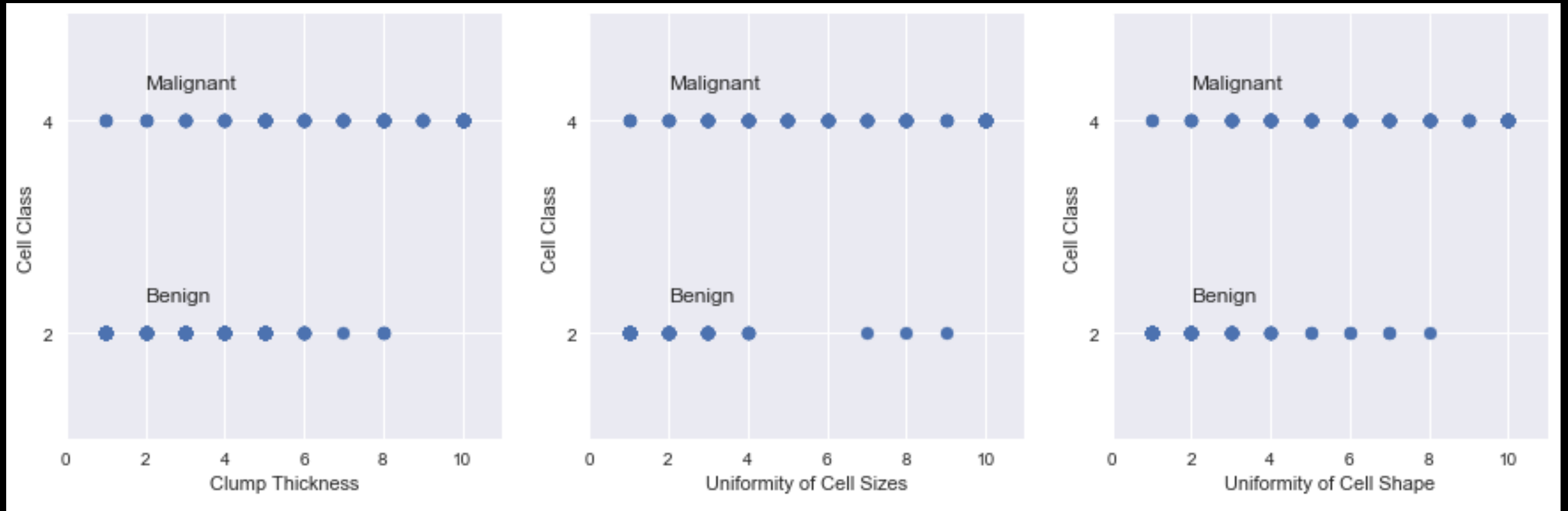


Input Variables



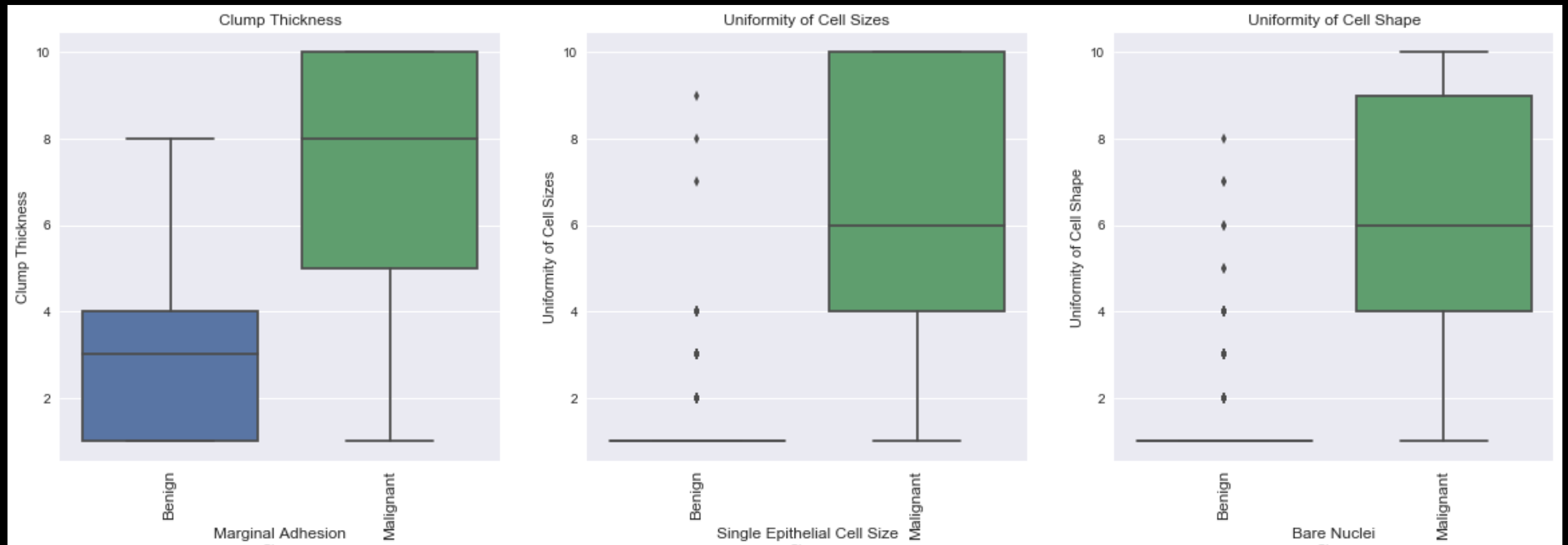
Variables are not normally distributed. They are skewed to the left. These variables clearly seem correlated.

Scatter Plots



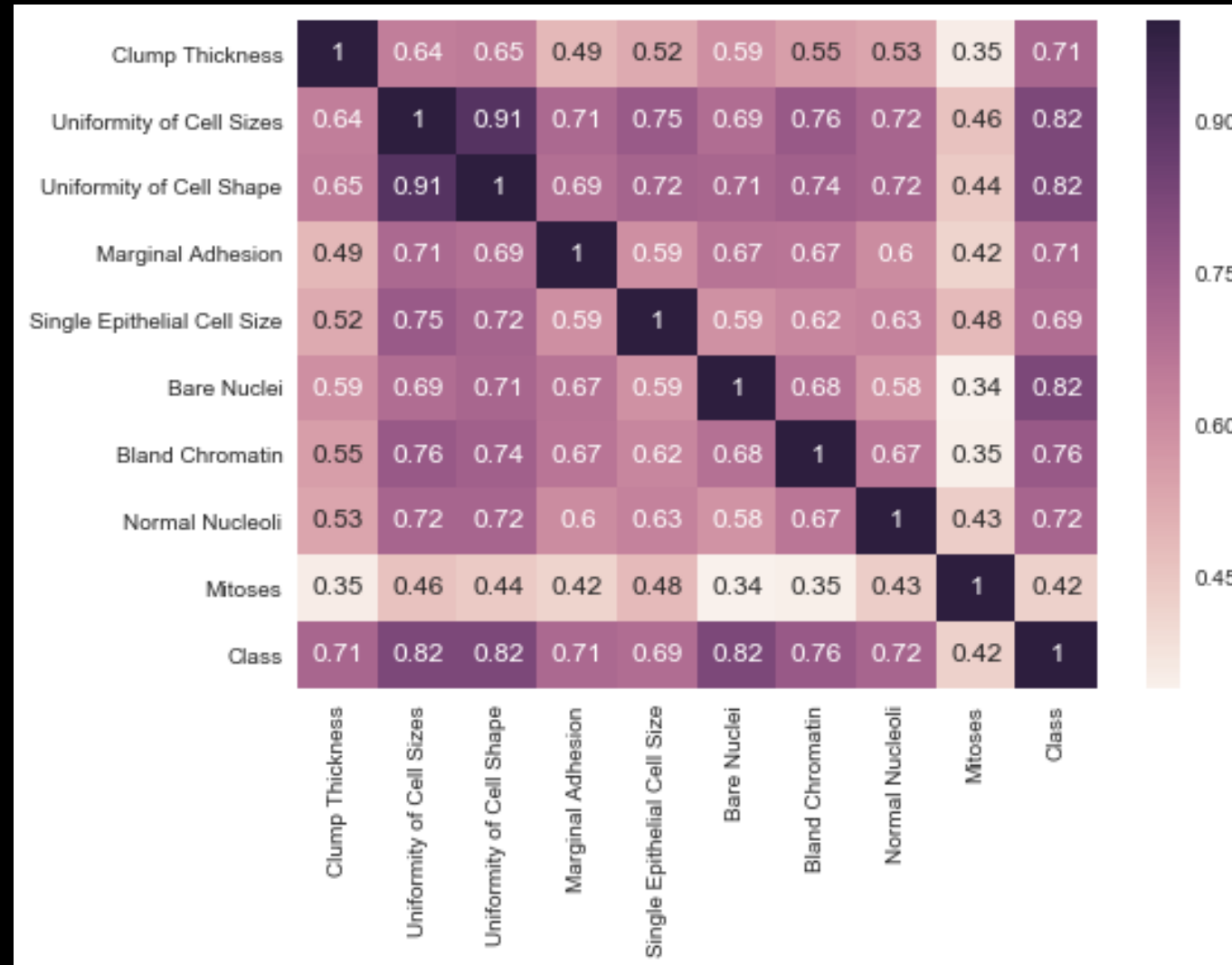
Scatter Plots don't help

Box Plots



Box plots filtered by Cell Type Classification help see the impact of various input variables on the output variable.

Correlation Matrix



- Most of the input variables are correlated to the 'Uniformity of the cell size & shape'.
- There is also good correlation between the cell class & input variables.

Logistic Regression Model

**Linear Logistic Regression Model was built using all variables.
This model has 87% Accuracy.**

Variable Name	Coeff	std err	z	P> z	Conf. Interval	
					[0.025	0.975]
Clump Thickness	-0.3225	0.059	-5.445	0.000	-0.439	-0.206
Uniformity of Cell Sizes	0.9437	0.139	6.809	0.000	0.672	1.215
Uniformity of Cell	0.1804	0.112	1.617	0.106	-0.038	0.399
Marginal Adhesion	0.1780	0.079	2.239	0.025	0.022	0.334
Single Epithelial Cell	-0.7894	0.105	-7.514	0.000	-0.995	-0.583
Bare Nuclei	0.4921	0.065	7.615	0.000	0.365	0.619
Bland Chromatin	-0.5463	0.095	-5.751	0.000	-0.732	-0.360
Normal Nucleoli	0.3595	0.078	4.598	0.000	0.206	0.513
Mitoses	-0.2500	0.089	-2.814	0.005	-0.424	-0.076

- **Cell's tend to be benign if the input variables are lower in value.**
- **Build a model using the most significant variables (shown in blue cells).**

Logistic Regression Model

Variable Name	coef	std err	z	P> z	Conf. Interval	
					2.5	97.5
Uniformity of Cell Sizes	1.0600	0.110	9.604	0.000	0.844	1.276
Single Epithelial Cell	-0.8451	0.098	-8.608	0.000	-1.038	-0.653
Bland Chromatin	-0.4952	0.081	-6.100	0.000	-0.654	-0.336
Bare Nuclei	0.4701	0.060	7.820	0.000	0.352	0.588

Conclusions

- **Cell attributes from Fine Needle Aspiration procedure can be used to predict if the cell was malignant or benign.**
- **Cell attributes - Bare Nuclei, Uniformity of cell size & single epithelial cell & Bland Chromatin were found to be the most important.**
- **Model built using just the four variables listed above was able to predict the cell type with more than 80% accuracy.**