

Project 1 Report

Chaitanya Chembolu – 50248987

Kautuk R Desai - 50247648

Problem Statement

The project aims to evaluate the basic statistics related to probability theory that help us to understand basic machine learning concepts.

These statistics include finding the mean, standard deviation, variance of multiple variables given in the dataset.

Further, we find covariance, and correlation of the variables in order to find the log likelihood for two different cases:

- i. Considering the variables as dependent and
- ii. Independent of each other.

The given dataset is about the computer science program rankings of select universities from sources such as US News, The Chronicle of Higher Education. The dataset consists of the following variables, CS Score, Research Overhead, Admin Base Pay and Tuition with 49 samples each.

Approach

In order to find mean, variance, and standard deviation we use different functions from numpy package. We then calculate the covariance (considering the variables to be dependent, independent) and correlation between the variables.

While finding the variance and covariance, we use numpy functions (`.var()`, `.cov()`) with `ddof` (Delta Degrees of Freedom) as 1. The default implementation of variance (and covariance) uses

$$\sigma^2 = \frac{\sum (X - \mu)^2}{N}$$
`ddof=0` which gives, . Where σ^2 is the variance, and μ is the mean of the variable.

But with ddof = 1, we get the denominator of variance expression as N-1. The expression is shown below,

$$\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N [x(i) - \mu]^2$$

To calculate the log likelihood of the data considering the variables as independent, we followed two approaches

- Using the formula for Gaussian distribution of independent variables and implementing it in python.
- Using the `multivariate_normal.logpdf()` function from `scipy` package. To compute, we use the covariance matrix for independent variables.

Results & Observations

As per the approach mentioned above, the following results were obtained

Statistical Parameter	CS Score	Research Overhead	Admin Base Pay	Tuition
Mean: μ	3.214	53.386	469178.816	29711.959
Variance: σ^2	0.475	12.85	1.4190e+10	3.1368e+07
Standard Deviation: σ	0.676	3.585	119120.615	5600.687

The covariance matrix (dependent variables):

	CS Score	Research Overhead	Admin Base pay	Tuition
CS Score	4.5800e-01	1.1060e+00	3.8798e+03	1.0585e+03
Research Overhead	1.1060e+00	1.2850e+01	7.0279e+04	2.8058e+03
Admin Base pay	3.8798e+03	7.0279e+04	1.4190e+10	-1.6369e+08
Tuition	1.0585e+03	2.8058e+03	-1.6369e+08	3.1368e+07

The corresponding correlation matrix:

	CS Score	Research Overhead	Admin Base pay	Tuition
CS Score	1	0.456	0.048	0.279
Research Overhead	0.456	1	0.165	0.14
Admin Base pay	0.048	0.165	1	-0.245
Tuition	0.279	0.14	-0.245	1

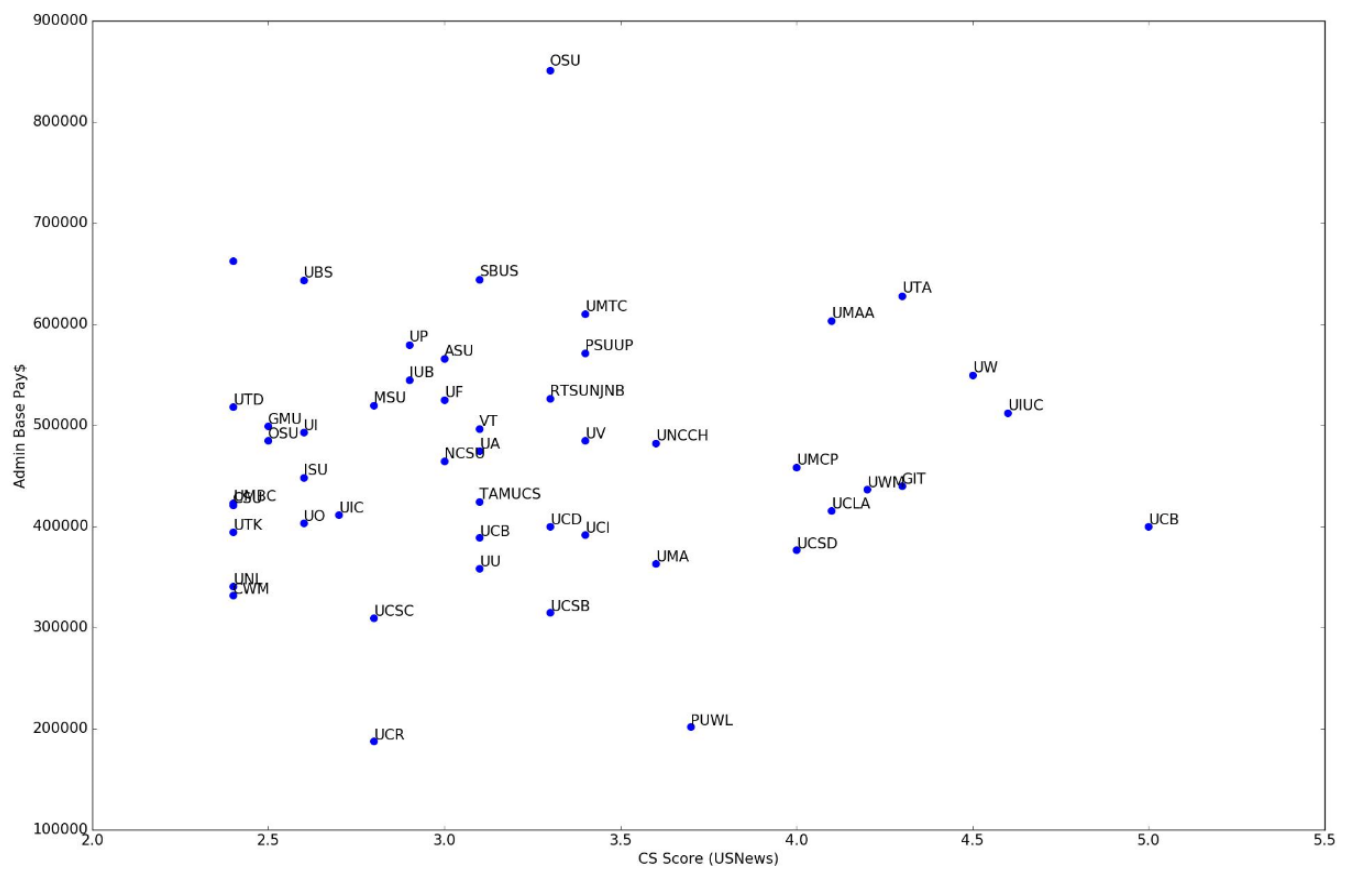
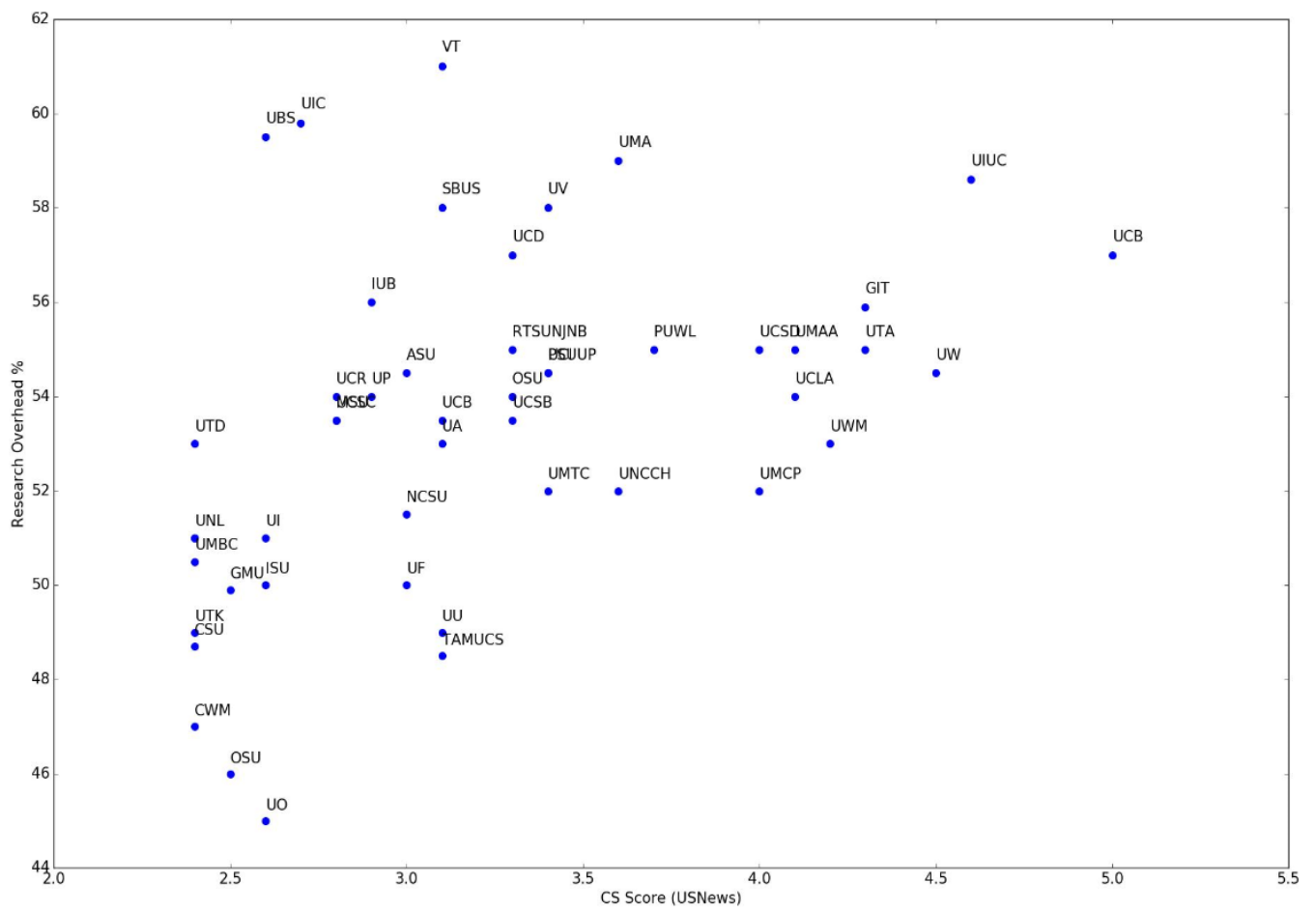
The Covariance matrix (independent variables):

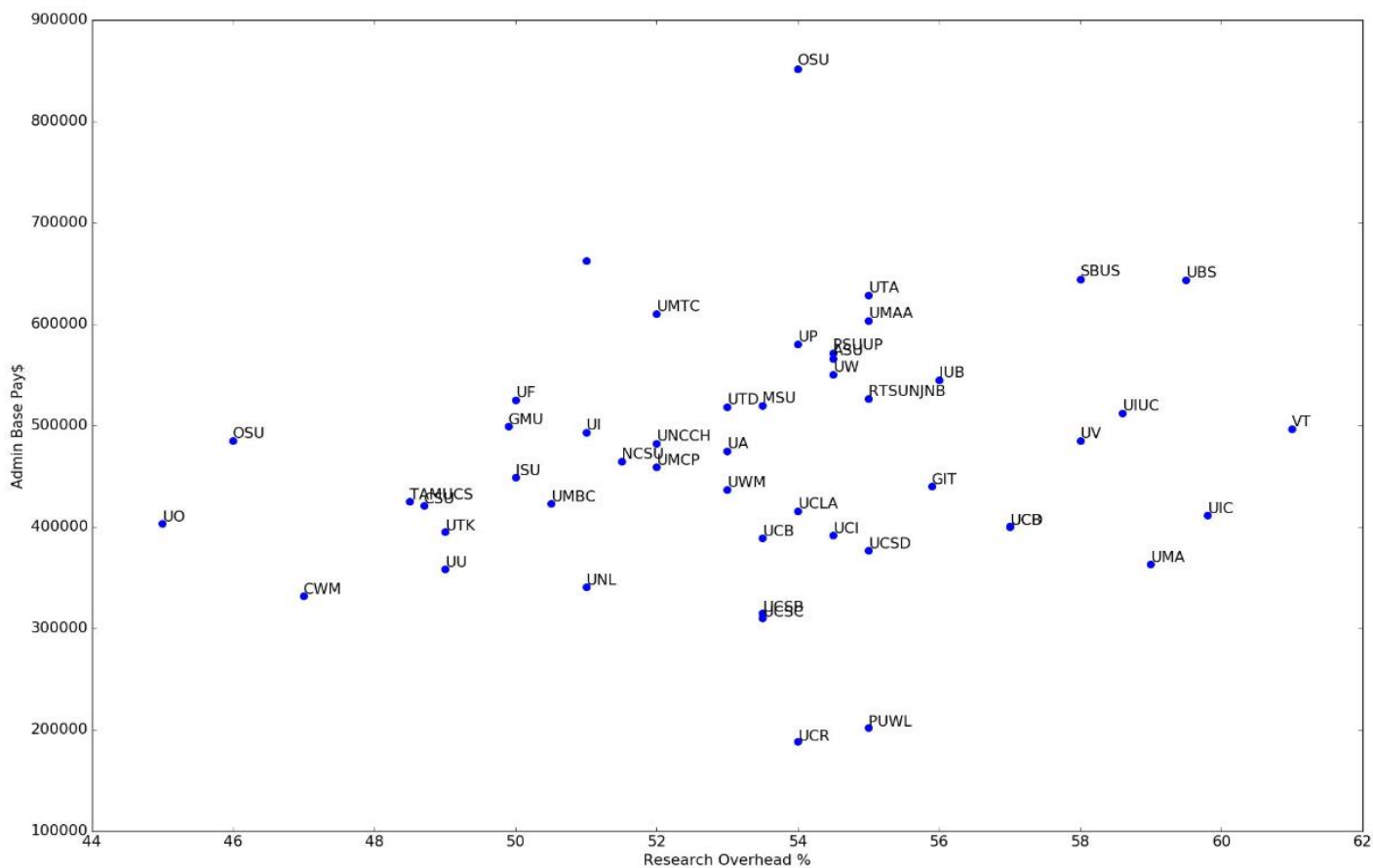
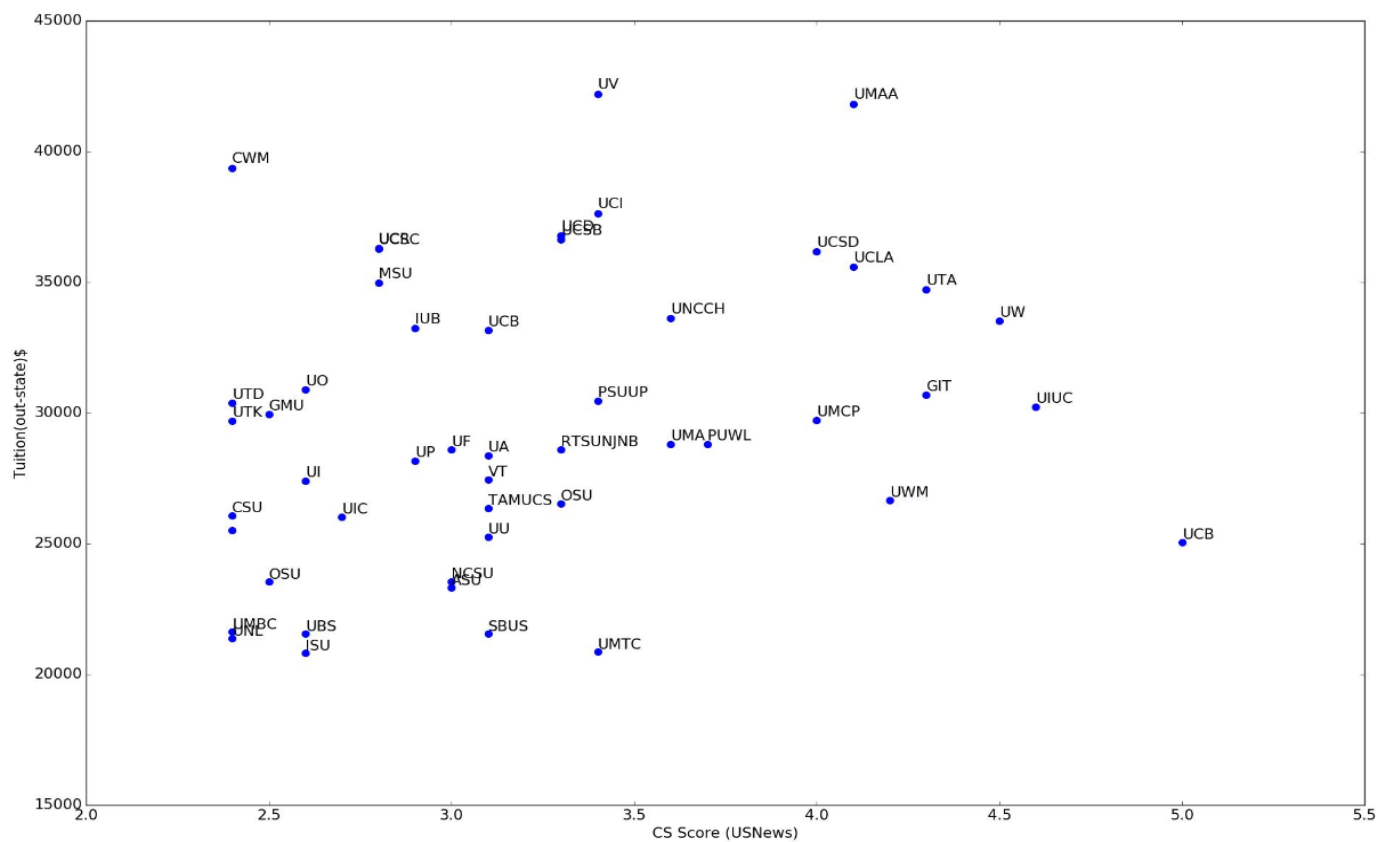
	CS Score	Research Overhead	Admin Base pay	Tuition
CS Score	4.5800e-01	0	0	0
Research Overhead	0	1.2850e+01	0	0
Admin Base pay	0	0	1.4190e+10	0
Tuition	0	0	0	3.1368e+07

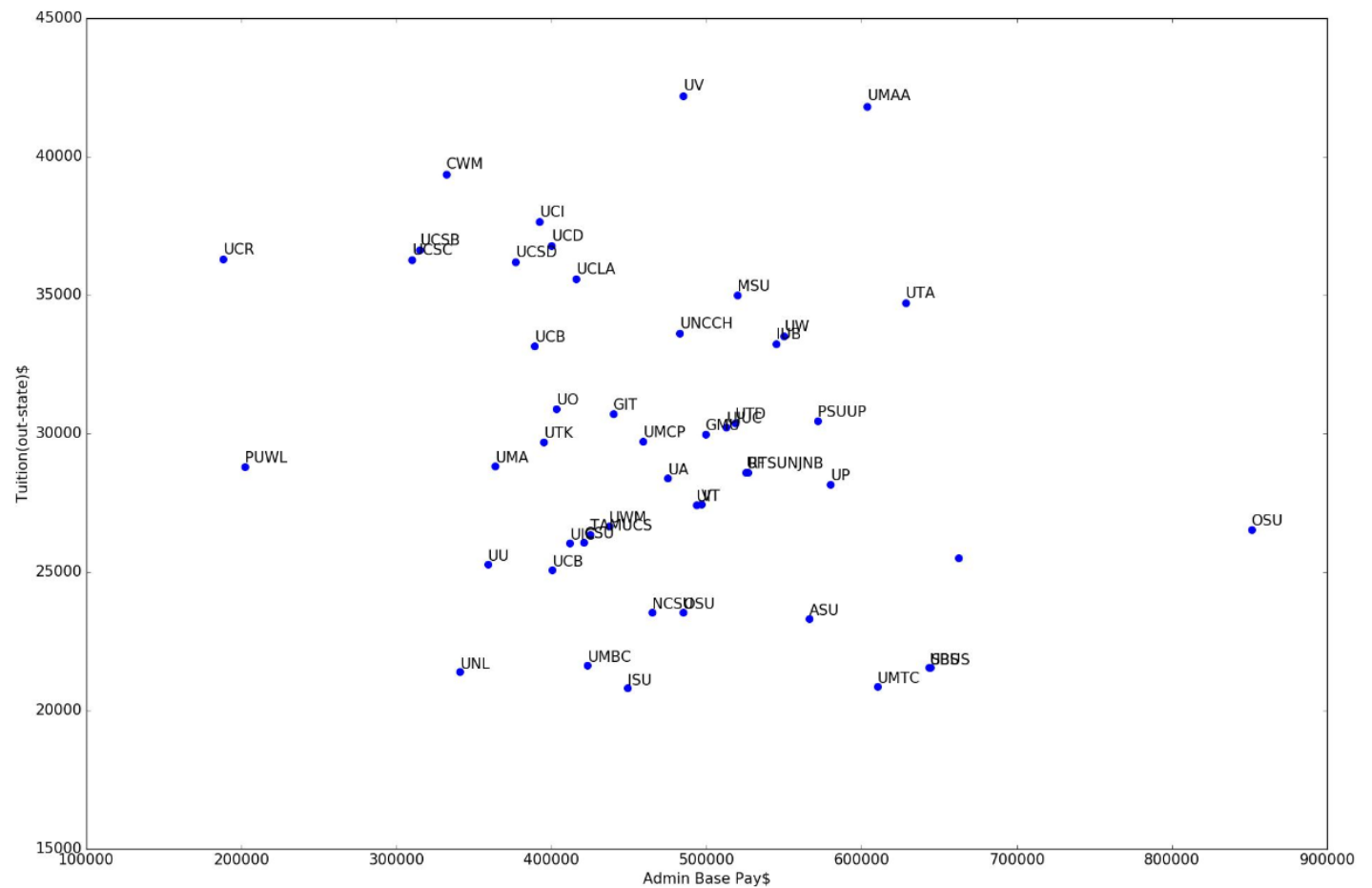
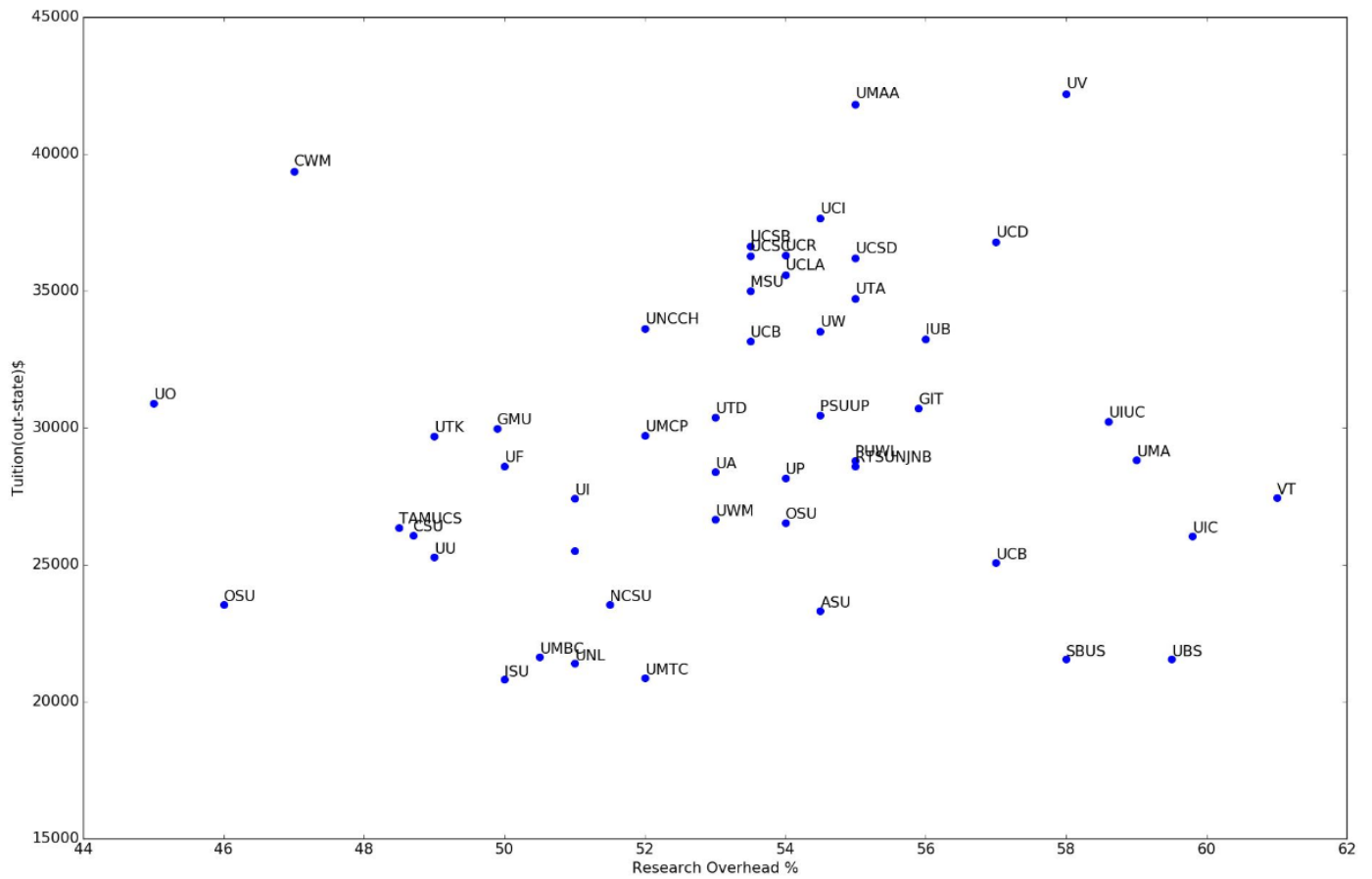
The log likelihood of the given dataset:

Variable Relation	Compute Method	Log Likelihood value
Independent	Formula: $\text{loglikelihood} = \sum_{i=1}^N \log p(\mathbf{x}_i)$ Where, $p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \Sigma ^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right]$	-1315.119
Independent	Function: >> multivariate_normal.logpdf()	-1265.25
Dependent	Function: >> multivariate_normal.logpdf()	-1262.327

The pairwise data showing the label associated with each data point of the variables is shown below,







Observation & Conclusion

From the pairwise data plots and correlation matrix, we can infer that CS Score and Research Overhead are the most correlated pair whereas CS Score and Admin Base Pay are the least correlated pair. Also, the log likelihood of the data is higher when the variables are dependent compared to when variables are independent.