

CSE 574: Project 1

Project Report

Submission Date: Sep 28 2015

Suramrit Singh/5016 9918/suramrit@buffalo.edu

1. Overview

Given a data set in an excel spreadsheet (UniversityData.xls), consisting of several variables, the project required the application of MATLAB to evaluate statistics, namely: mean and variance of univariate distributions and covariance and correlation coefficient of pairs of variables.

The project also involved application of the statistics to construct compact representations of Bayesian Networks and checking the goodness of these representations using the concept of likelihood

1.1 Background

Terms and Concepts:

Mean: In probability, mean gives us a central tendency of a probability distribution or a random variable.

It is represented as

$$\mu = \sum xP(x)$$

where x is the value of a random variable and $P(x)$ is the probability

For a data set it is defined as the sum of the values divided by the total no of values.

$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$

Variance: For a data set sample variance is computed as:

$$\sigma^2 = \left(\frac{1}{N-1} \right) * \sum [x(i) - \mu]^2$$

Variance gives us a measure of dispersion of the data from the mean. A higher value of variance implies that the data set is scattered away from the mean.

Covariance: For a pair of variables X_1, X_2 with samples $x_1(i), x_2(i)$ the covariance is given by

$$\sigma_{12} = \frac{1}{N-1} * \sum [x_1(i) - \mu_1][x_2(i) - \mu_2]$$

Covariance gives a measure of how 2 variables change together. Since there is no square term in covariance, it can also be a negative value.

For a set of n variables $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$, a covariance matrix is given as:

$$\Sigma = \begin{bmatrix} \sigma_{11} & \dots & \dots & \sigma_{1d} \\ \vdots & \ddots & \vdots & \vdots \\ \sigma_{1d} & \dots & \dots & \sigma_{dd} \end{bmatrix}$$

Likelihood and Log likelihood:

The likelihood of a set of parameter values or events given certain outcomes/samples, is equal to the probability of those observed outcomes given those events.

Generally natural log of the likelihood is taken for easier estimation and analysis. For n independent samples of a probability density, the log likelihood is given by

$$L(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \log \sum p(\mathbf{x}_i)$$

Data set:

The given data set for this project is multivariate with 4 variables, CS Score, Research Overhead, Base Pay, and Tuition. They have been represented by X_1, X_2, X_3, X_4 throughout the document.

2. Tools Used

Matlab and related libraries were used for the computation and analysis of the statistical data. The data used in the project was imported from an excel spreadsheet (UniversityData.xls). Wolfram alpha was also intermittently used for computations and plotting.

3. Tasks

3.1 Mean, variance and standard deviation calculation

The values for each variable were imported from the excel sheet and a resultant 49×4 matrix was made with each column containing the values of the random variables.

The mean of this matrix was calculated using: `$\mu = \text{mean}(rvMatrix)$` ;

Here ***rvMatrix*** is a **49x4** matrix containing the variable values.

The resultant matrix, **mu**, was a vector of dimension 1x4 with each column containing the mean for each variable.

After the calculation of mean, variance was calculated using: ***variance = var(rvMat)***

The resultant matrix, **variance**, was a vector of dimension 1x4 with each column containing the variance for each variable.

The standard deviation was calculated taking the square root of each column of vector variance.

3.1.2 Results

The results from section 3.1.1 are tabulated below. Mu1 corresponds to mean of variable X1 and so on. Mean:

Mu1	3.214
Mu2	53.3650612
Mu3	469178.816326531
Mu4	29711.9591836735

Tab 1

Variance:

Var1	0.4575
Var2	12.6160629
Var3	14183920820.9031
Var4	31367695.789966

Tab 2

Standard deviation:

Sigma1	0.676387462923434
Sigma2	3.5519097574643
Sigma3	119120.614592534
Sigma4	5600.68708195396

Tab 3

3.2 Computing covariance and correlation (Matlab Implementation)

In order to calculate the covariance and correlation matrices, the following implementations were used in matlab:

For covariance , ***covarianceMat = cov(rvMat);***

For correlation, ***correlationMat = corrcoef(rvMat);***

The covariance and correlation matrices obtained are:

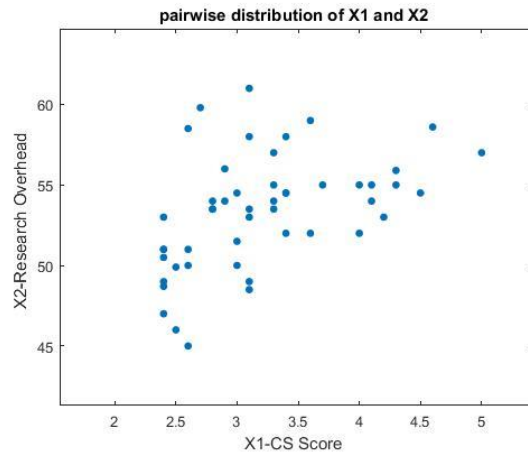
	X1	X2	X3	X4
X1	0.4575	1.11842261904762	3879.78184523809	1058.47976190476
X2	1.11842261904762	12.6160629251701	66651.6643282313	2975.82980442177
X3	3879.78184523809	66651.6643282313	14189720820.9031	-163685641.257653
X4	1058.47976190476	2975.82980442177	-163685641.257653	31367695.7899660

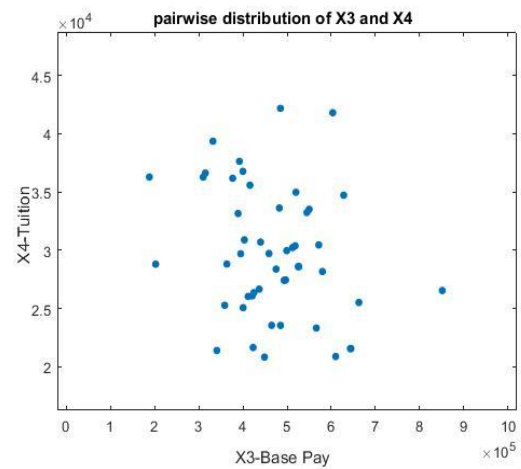
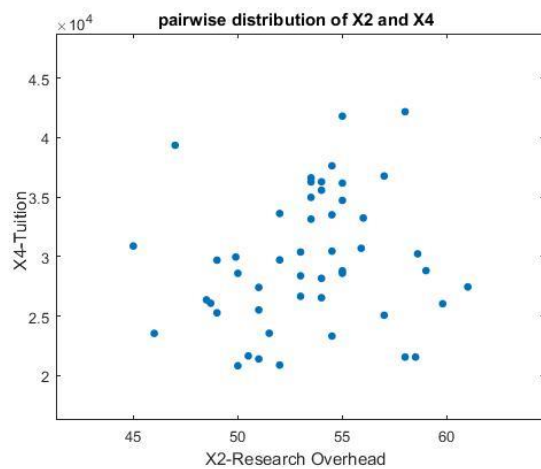
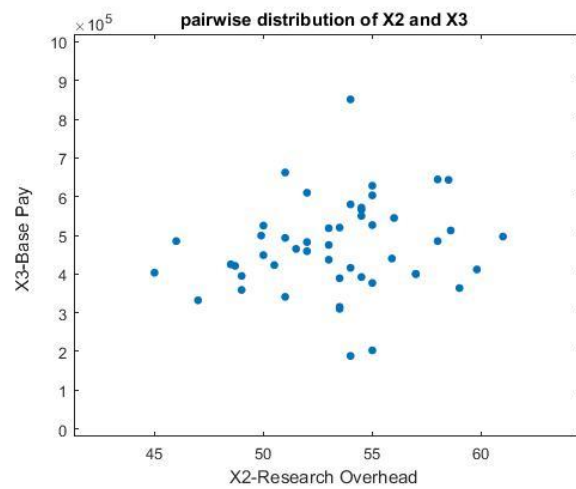
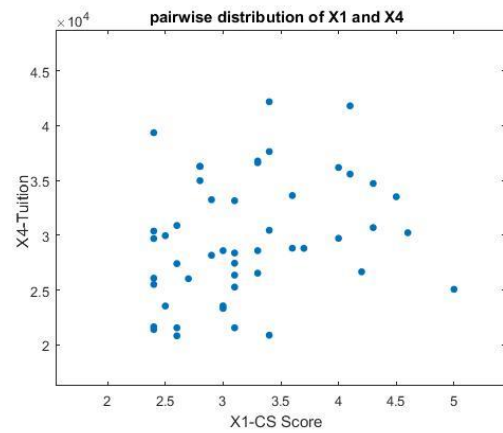
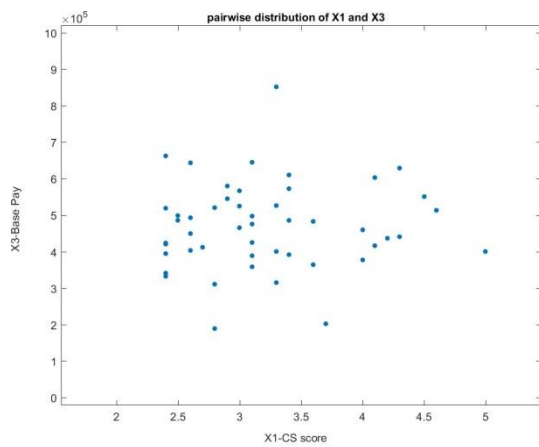
	X1	X2	X3	X4
X1	1	0.465530853186398	0.0481531643615803	0.279412416164465
X2	0.465530853186398	1	0.157529592573249	0.149590790502660
X3	0.0481531643615803	0.157529592573249	1	-0.245347902755565
X4	0.279412416164465	0.149590790502660	-0.245347902755565	1

In the above tables, $(i, j)^{th}$ cell show for each variable pair X_i and X_j s the covariance and correlation values for variables respectively.

3.2.1 Pairwise data plot

In order to understand the correlation between the variables better, pair wise data was plotted and the following plots were observed.





3.2.1.1 Inference

The scatter plot for X1 and X2 can be observed to be more closely packed than other variable pairs. Also there is a positive correlation between the variables.

For other variables it cannot be directly inferred as to what their correlation is since the scatter plot is more dispersed.

3.3 Computing Log Likelihood when variables are independent

The combined probability of the variables will be represented as: $P(X1, X2, X3, X4)$

Assuming that the variables are independent of each other, the combined probability can be represented as:

$$P(X1, X2, X3, X4) = P(X1) * P(X2) * P(X3) * P(X4) \quad (1)$$

Therefore the Log Likelihood for the data will be:

$$\text{Log Likelihood} = \text{Loglikelihood}(X1) + \text{Loglikelihood}(X2) + \text{Loglikelihood}(X3) + \text{Loglikelihood}(X4) \quad (2)$$

The Log Likelihood of a variable is calculated by first finding the normal probability density for the variables, in matlab this was implemented by fusing the function ($\text{normpdf}(\text{rvMat}(:, i), \text{mu}(i), \text{sigma}(i))$)) [4]

This function computes the pdf at each of the values in a column vector for a random variable , using the normal distribution with mean μ and standard deviation σ

Then we take the log of the values obtained and sum them to get a scalar value. The summed values of each variable are then summed to obtain the log likelihood for the distribution.

This procedure can be implemented in matlab as:

$$\text{sum}(\text{sum}(\log(\text{normpdf}(\text{rvMat}(:, i), \text{mu}(i), \text{sigma}(i))), 1), 1) \quad \text{for } i = [1, 4]$$

3.3.1 Result

The value of Log likelihood is computed to be ≈ -13146

3.4 Finding a Bayesian Network and its Log likelihood

3.4.1 Finding valid Bayesian Network

Since there are 4 variables, the Bayesian network for the distribution will be represented by a 4x4 binary matrix. From a total of 2^{16} possible combinations of 4X4 binary matrices, we need to find the matrices which are representations of a directed acyclic graphs (with no self loops). In order to list all the possible 4X4 binary matrices, we take all the integers from 1 to 2^{16} , convert them to a binary string of 16 digits and then reshape the string to a 4X4 matrix. This can be stored in a 3d matrix array where each index of the array stores a possible combination for a 4X4 binary matrix. The matlab implementation for the procedure will be:

$$\mathbf{matrices} = \text{reshape}(\text{dec2bin}((0:2^{16}-1), 16).', '0', 4, 4, [1]))$$

dec2bin() converts a decimal number to binary, *reshape()* converts the 16 digit binary number into a 4x4 matrix. The result stored in '*matrices*' will be a 3-d array of matrices of size 4X4 [1][2].

Once all the combinations for the 4X4 binary matrices have been obtained we can check the matrices that satisfy the condition for representing directed acyclic networks. This can be done in matlab by

$$\mathbf{graphisdag}(\text{sparse}(\mathbf{matrices}(:, :, i))) \quad \text{:for } i\text{th matrix in 'matrices' } \quad [3]$$

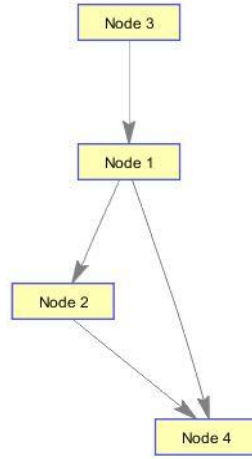
Once a Bayesian network has been obtained, we then try and find the loglikelihood for that network and compare it to the earlier found value in sec 3.3.

3.4.2 Sample Bayesian Network and Log likelihood calculation

A Bayesian network combination found by the method described in sec 3.4.1 was

$$\mathbf{M} = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

The Bayesian network represented by this matrix is



Possible Bayesian Network for given variables

For this particular matrix, the combined probability, determined by the dependencies of the network, will be given by:

$$P(n1, n2, n3, n4) = P(n3) * P(n1|n3) * P(n2|n1) * P(n4|n1, n2, n3) \quad (3)$$

Using Bayes theorem,

$$= P(n3) * \left(\frac{P(n1, n3)}{P(n3)} \right) * \left(\frac{P(n2, n1)}{P(n1)} \right) * \left(\frac{P(n1, n2, n3, n4)}{P(n1, n2, n3)} \right) \quad (4)$$

Here $P(n3)$ is the univariate probability distribution for the variable $X3$ and rest of the terms are combined probability distributions for the variables.

Thus taking the log of each multiplicative term in (4), and taking the sum of the values, we obtain the Log Likelihood for the given Bayesian network (using (2)).

3.4.3 Matlab implementation

Once the dependencies between the nodes are determined, we can calculate $P(X1, X2, X3 \dots Xn)$ for an n variable distribution in matlab by using the *mvnpdf* function. the matlab implementation for this is:

$$\mathbf{mvnpdf}(\mathbf{data\ values\ of\ dependent\ variables}, \mathbf{mu}, \mathbf{covariance})) \quad (5)$$

mvnpdf computes returns the density of the multivariate normal distribution with mean *mu* for the variables and covariance *sigma*[5]

We then calculate the log of density multivariate distribution and then take sum of the values to obtain the log likelihood for the Bayesian network.

Thus the entire procedure can be implemented in matlab as:

$$\text{sum}(\log(\text{mvnpdf}(\text{data values of dependent variables}, \mu, \text{covariance})))) \quad (6)$$

Once ‘sum’ for each term in (4) was calculated using (6) we then add the values obtained for each term to obtain the Log likelihood

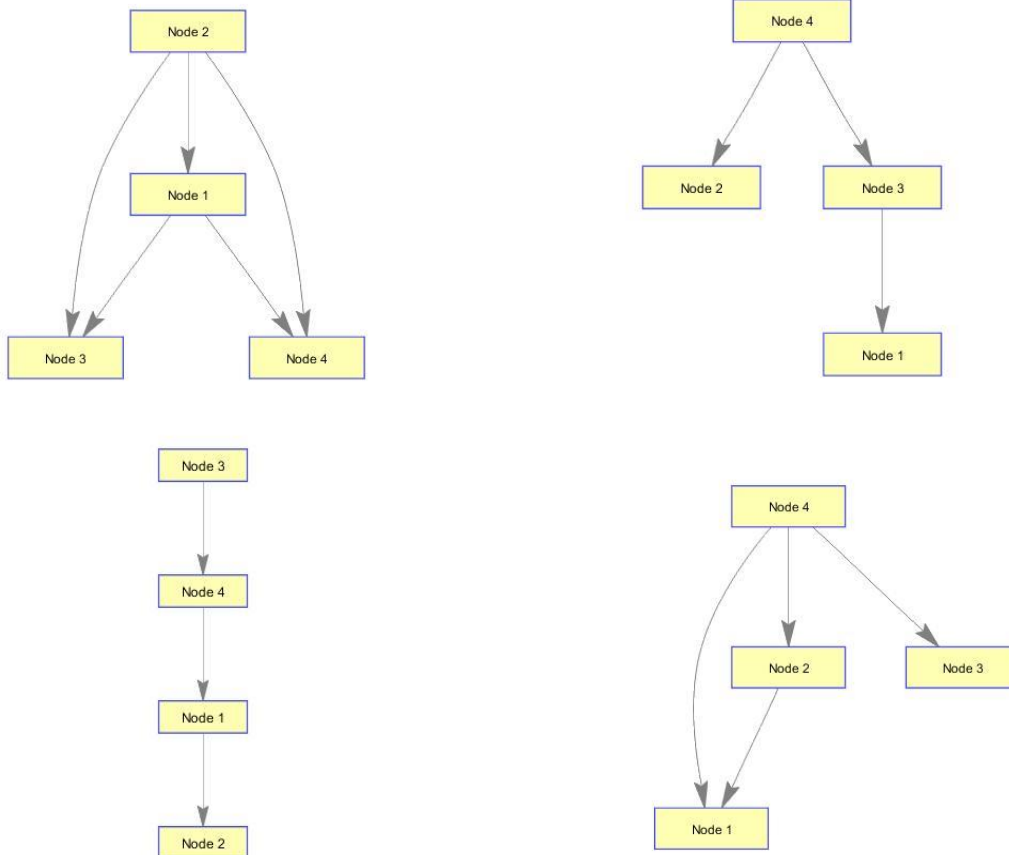
3.4.4 Results and Inference

The log likelihood obtained for the Bayesian network of section 3.4.2 is -1306.6227 which is better than the obtained in section 3.3.

Acc to definition of Likelihood, thus it can be inferred that our variables have a higher likelihood of having dependencies as described by the Bayesian network than being independent of each other.

3.5.5 Further Optimization

Some samples of other possible Bayesian networks are



It was found that there are more than 500 possible matrices(out of the exhaustive 2^{16} combinations) which can describe a directed and acyclic Bayesian network. Thus we can exhaustively search the log likelihood of all such possible matrices, compare the values and evaluate the Bayesian network that results in the maximum log likelihood for our set of variables.

It is possible to calculate log likelihood for all such valid Bayesian networks and then find the optimum value.

4. Conclusion

Using matlab we were able to get different statistical properties (mean, variance, covariance etc) for a multivariate data set of 4 variables. We found the covariance and correlation matrices for the variables and plotted their pair wise distribution to infer about their correlation.

We also calculated the log likelihood value for the distributions by first assuming that the variables were independent. We then found various possible Bayesian networks for the variables, and for a particular Bayesian network calculated the log likelihood value. The log likelihood value for this network was found to be greater than when the variables were considered to be independent thus implying that there was a better likelihood of the variables being dependent as described by the given Bayesian Network than being independent.

This can also be further optimized by further analyzing the log-likelihood value for all possible Bayesian networks for the given 4 variables and finding the maximum value, which would give us the optimal Bayesian network.

5. References

- 1] Reshape, <http://www.mathworks.com/help/matlab/ref/dec2bin.html> Reshape
- 2] Dag test, <http://www.mathworks.com/help/bioinfo/ref/graphisdag.html>
- 3] Mvnpdf, <http://www.mathworks.com/help/stats/mvnpdf>.
- 4] Normpdf <http://www.mathworks.com/help/stats/normpdf.html>
- 5] Aaron Goebel goaaron, Mihir Mongia mmongia Non-linear Reconstruction of Genetic Networks Implicated in AML Pathology.