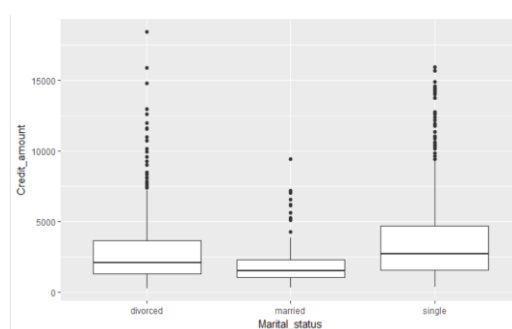
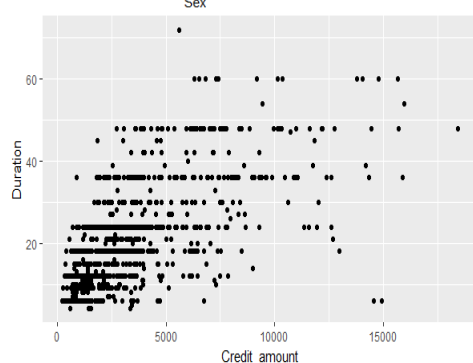
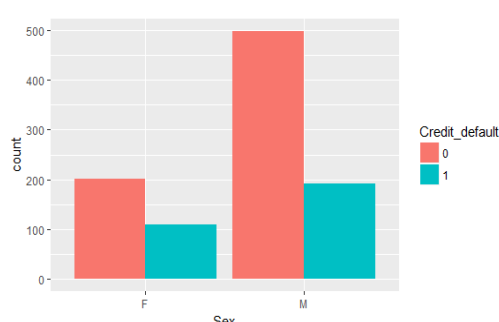
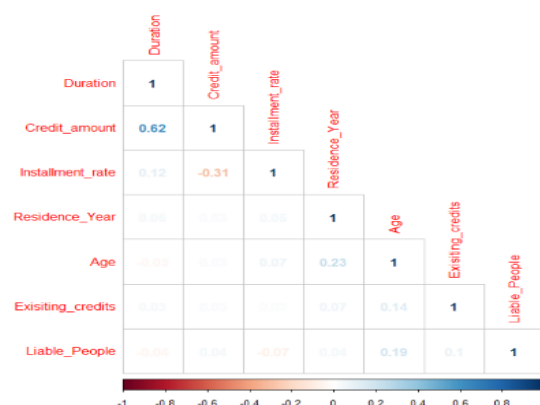


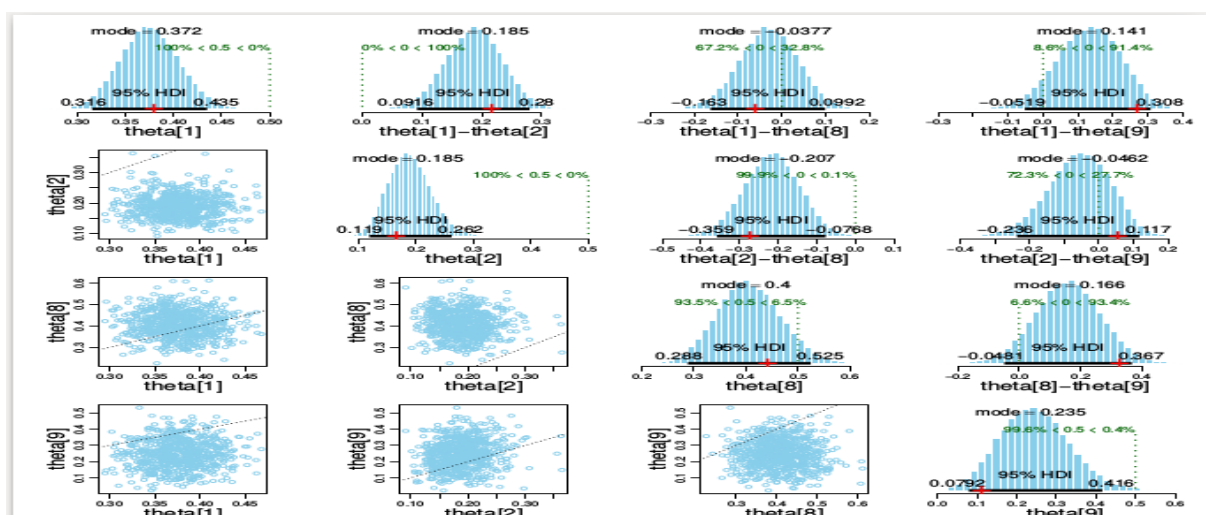
UCI German Dataset

The data that we have chosen to work with is regarding German credit data. This data has almost 1,000 observations and 21 attributes. The main variable of interest is whether or not the client would be classified as good credit or bad credit. A secondary variable of interest is credit amount. The attributes are the status of a checking account, how long the account has been open for, credit history, purpose of the loan, if employed, how long they have held a job, savings amount, marital status, residence, and age, among others. Banks would be interested in what type of applicant will default and applicants would be interested in how much they can ask for.

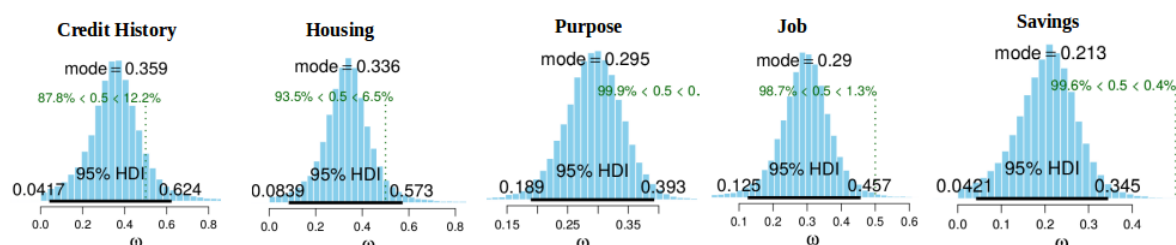


We began by doing initial exploratory data analysis to understand the data. We ran both structure and summary functions in R. Then we used bar plots, correlation matrix, scatter plots and boxplots to understand the different variable types that we had. We have a bar plot showing the gender based on credit default. For credit default 1 means default and 0 means no default. We can see that there are more male applicants and that the majority of the applicants did not default. Next, there is a duration scatter plot by credit amount. This shows that they have a positive correlation. Next, there is a marital status and credit amount box plot. The credit amount is fairly similar between single, divorced and married. However, there does appear to be less variance amongst the married crowd. The correlation matrix confirms this theory, showing a correlation of .62. This is the highest correlation of all of the numerical attributes.





We used MCMC to assess the likelihood of default based on the different categories within categorical variables. Doing that, we can get a sense of how likely is an applicants to default and also the tendency of default for the group as a whole. The figure shows the result for the category Purpose and the mode of four possible loan purposes: 1 - Car (new), 2 - Car (used), 8 - Education, and 9 - Retraining. Due to the classes being unbalanced with 70% non-defaulters and 30% defaulters, all the modes are close to the default rate of the dataset. We can see from Theta 1, that the mode of the probability of someone of this subgroup default is 0.37, and also that no shrinkage is occurring, as it matches the proportion in the data (as indicated by the red cross). On the other hand, Theta 9 presents a noticeable shrinkage away from the raw proportion towards the mode of the group (which is 0.3). Notice also that the widths of Theta 9's 95% HDIs is broad. This occurs because there is only a little amount of data for this specific subgroup leading it to be strongly influenced by the metrics of the group. Finally, Theta 8, Education, is a purpose that the financial institutions should pay attention to. When Theta 8 is compared with Theta 2 and Theta 9 we can see that zero is not contained within the 95% HDI. Therefore Theta 8 significantly differs and banks should realize that applicants whose purpose of a loan is education are more likely to default than those buying cars or paying for retraining.



Five categorical features with qualitative information about bank account holders are organized by their group mode from largest to smallest. We can see that there is much variation

amongst the groups. This analysis will help financial institutions to decipher which applicants would be more likely to default.

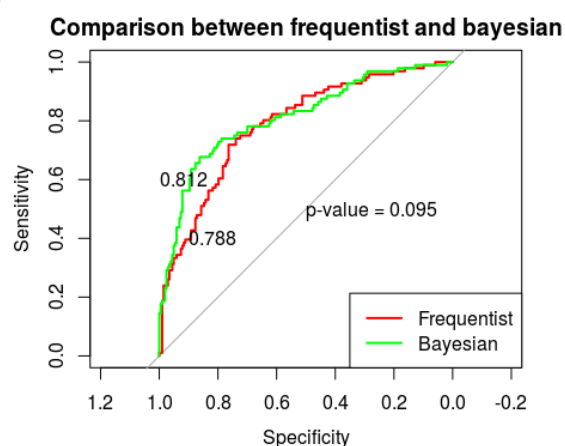
We developed a model to predict the probability of whether or not an applicants will default on their loan based on the available data. To accomplish this, we first fit a logistic regression model using the *glm* function available in R. The overall accuracy for the frequentist model is about 73% and tells how often this classifier is correct. The precision is 0.63 and measures how often the model is correct in predicting a positive outcome. The recall is 0.44 and tells us how often false negatives occur. Our objective for this simple model is to have a frequentist estimation for comparison to a Bayesian approach. It is worth noting that in this problem we are more concerned about the false negative rates (recall) instead of the false positives rates. This is because a bank (or any other financial institution lending money) is interested in minimizing the chances of not being paid back. It is more problematic to state that a person is reliable when it is not (false negative), instead of saying that a person is not reliable and, in fact, it is (false positive). A way to handle this problem using logistic, is to modify the threshold for the decision of classifying as positive. In other words, one can increase the rate of false positives to decrease the rate of false negatives. This is possible by lowering the decision threshold. For example we can move it from 0.5 to 0.3. Hence, any probability outcome from the logistic model greater than 0.3 would be classified as a "default risk". In our following attempt, we evaluate the same problem, predicting the chance of a person default to pay a loan or not, using the Bayesian logistic regression approach. To be able to estimate the parameters of the linear function of the logistic model, i.e., the posterior of our Bayesian inference, we rely on the MCMC

```
model {
  for ( i in 1:Ntotal ) {
    # In JAGS, ilogit is logistic:
    y[i] ~ dbern( mu[i] )
    mu[i] <- ( guess*(1/2) + (1.0-guess)*ilogit(zbeta0+
      sum(zbeta[1:Nx]*zx[i,1:Nx])) )
  }
  # Priors vague on standardized scale:
  zbeta0 ~ dnorm( 0 , 1/2^2 )
  for ( j in 1:Nx ) {
    zbeta[j] ~ dnorm( 0 , 1/2^2 )
  }
  guess ~ dbeta(1,9)
}
```

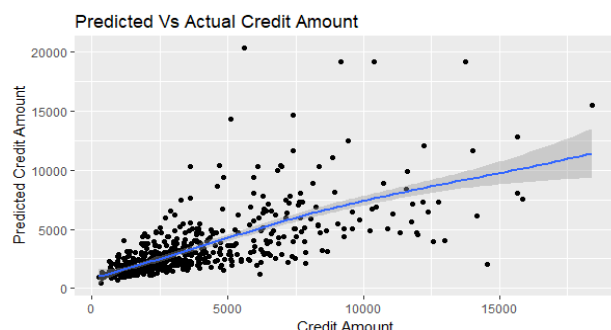
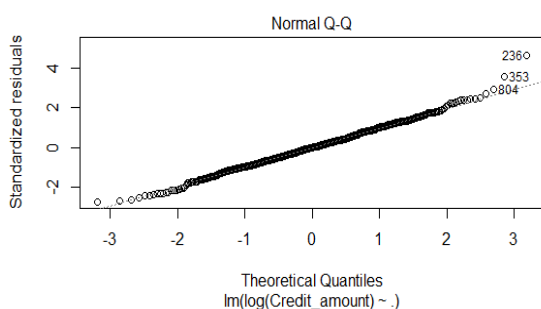
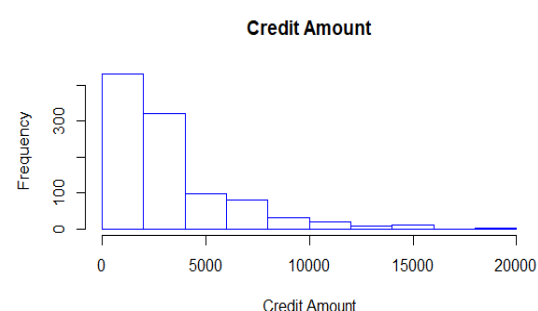
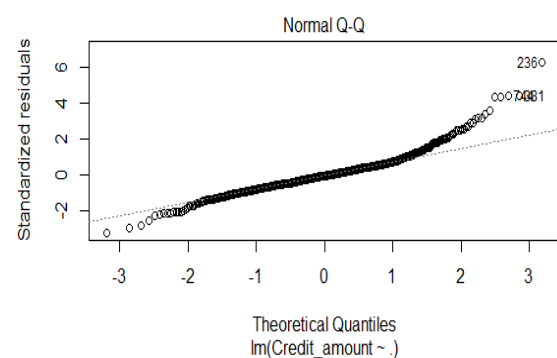
strategy. A snippet of the model written using the software JAGS is also presented. It can be observed that we built a robust model against outliers, once we are employing the “guess” variable. This model generates the estimates for all the coefficients, which we then use to make predictions

using the *sigmoid* function. To compare the result, we use the same train/test set created before. The Bayesian approach shows a better performance than the frequentist, with an overall accuracy of 80%, a recall of 0.55 and a precision of .77.

Finally, we plot the ROC curve of the two models. This is a commonly used graph that summarizes the performance of a classifier over all possible



thresholds. It is generated by plotting the True Positive Rate (y-axis) against the False Positive Rate (x-axis) as you vary the threshold for assigning observations to a given class. We print the AUC (Area Under the Curve) for both models. An AUC higher than 0.5 means that the classifier is better than chance. Here, both have a good AUC with the Bayesian classifier slightly better. The p-value informs that the difference between them is not significant.



In order to provide advice to the applicants as well as the bank a regression problem exploring how the attributes relate to credit amount was conducted. This analysis began by using a linear model and all of the attributes. The results showed an adjusted R-squared value of .6. This is an okay result but we would like a model that explains more of the variance. The Q-Q plot of the residuals of this model show that the distribution is not normal. An ANOVA test was performed to see which attributes were the most important when it comes to credit amount. A p-value of .001 was used to determine which attributes would be deemed important and kept in the model. This left the model with six attributes: status checking, duration, purpose, installment rate, property, and job. A test of normality was performed on credit amount and it was found to be heavily skewed, as showed in the histogram. This led the next model to utilize the log function to normalize the response variable of Credit

Amount. The final frequentist regression model only used the six main attributes and the log of credit amount. This resulted in a model with an adjusted R-squared of .63, which is slightly better than the initial model.

After making feature selection and adjusting the data, we performed a Bayesian Analysis using the same attributes and with credit amount being

the response variable. BPM, BMA, MPM and HPM were all utilized. It was found that BPM was the best model based on the RMSE values. The next step of the analysis is to compare the Bayesian and Frequentist method. The best model for the frequentist method had an RMSE of 2080.457 while the

Bayesian Method for Data Science

Instructor: Yuxiao Huang

Group 2: Krystin Sinclair, Ye-in Jeon, Christian Cleber Masdeval Braz

best model of the Bayesian Analysis had a RMSE of 2019.237. The predicted versus actual scatterplot with fit line shows the best model available which is the Bayesian BPM.

In conclusion, there is much advice that can be given to both banks and loan applicants. For banks, they should be wary of long term loans because for every month increase in duration of the loan the odds of defaulting increase by a factor of 1.02. As for loan applicants, the applicants that receive the higher credit amount were by those that have management positions or jobs with high qualifications and when the purpose of the loan is to purchase a car. The analysis and use of a variety of methods showed that Bayesian Analysis is a very useful tool. It provided better models than frequentist method and for future research Bayesian Analysis should be utilized.

References

Dataset

<http://home.cse.ust.hk/~qyang/221/Assignments/German/>

Articles

<https://loans.usnews.com/beyond-credit-scores-factors-that-affect-a-loan-application>

<https://studentloanhero.com/featured/personal-loan-purpose-happens-change/>

<https://www.cbsnews.com/news/5-things-that-can-torpedo-your-mortgage-application/>

<http://www.ijbf.uum.edu.my/images/pdf/5no1ijbf/6ijbf51.pdf>

<https://www.sciencedirect.com/science/article/pii/S0883902688900183>

<https://dl.acm.org/citation.cfm?id=131259>

<https://www.emeraldinsight.com/doi/pdfplus/10.1108/eb013696>

<http://www.rcmloan.com/credit-building-solana-beach-ca/>