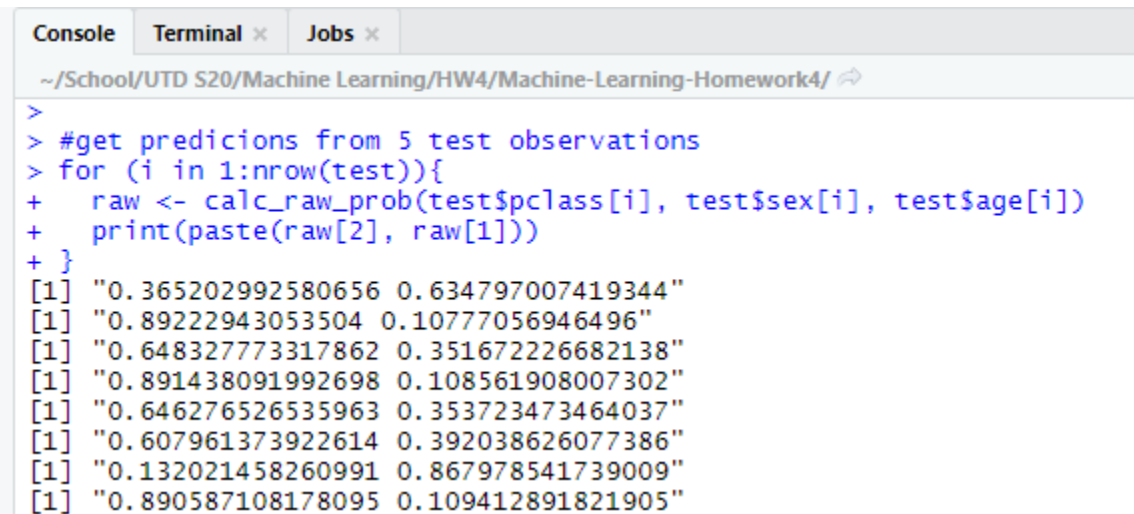


## Naïve Bayes

### In R:

We began by loading the titanic data and splitting into test and train sets. The train set being the first 900 instances and test being the remaining. We ran Naïve Bayes on the data using the e1071 package first to get base results. Then, following the book on Naïve Bayes and the code on git-hub we adapted the Naïve Bayes from scratch to work on our data. The way the algorithm works is explained in the C++ section of Naïve Bayes.



```
~/School/UTD S20/Machine Learning/HW4/Machine-Learning-Homework4/ ↵
>
> #get predictions from 5 test observations
> for (i in 1:nrow(test)){
+   raw <- calc_raw_prob(test$class[i], test$sex[i], test$age[i])
+   print(paste(raw[2], raw[1]))
+ }
[1] "0.365202992580656 0.634797007419344"
[1] "0.89222943053504 0.10777056946496"
[1] "0.648327773317862 0.351672226682138"
[1] "0.891438091992698 0.108561908007302"
[1] "0.646276526535963 0.353723473464037"
[1] "0.607961373922614 0.392038626077386"
[1] "0.132021458260991 0.867978541739009"
[1] "0.890587108178095 0.109412891821905"
```

*Results from the Naïve Bayes from scratch in R.*

### In C++:

The C++ code for Naïve Bayes uses the Armadillo library as well in order to facilitate the programming. All the metrics such as likelihood matrices and raw probabilities are identical in the R script as the C++ code. The program begins by reading the data from the CSV and splitting it into test and train.

```
170 int main(){
171
172     vector<Passenger> passengers = readCSVData();
173     vector<Passenger> test;
174     vector<Passenger> train;
175     //get the first 900 instances as train subjects
176     for(int i=0; i<900; i++){
177         train.push_back(passengers[i]);
178     }
179     //train
180     for(int i=900; i<passengers.size(); i++){
181         test.push_back(passengers[i]);
182     }
183 }
```

The program then gets the counts and priors for the people who survived and perished. After that we're able to compute the likelihoods of both class and sex. This is done by counting survived and perished for each class and dividing by the total count in the train data.

```

89  mat get_likelihood_pclass(vector<Passenger>& pass, mat count){
90      mat likelihood(2,3); likelihood.zeros();
91      for(int sv =0; sv<2; sv++){
92          for(int pc=0; pc<3; pc++){
93              int numSurvived = 0;
94              for(auto p: pass){
95                  if(p.survived==sv and p.pclass==(pc+1)){
96                      numSurvived++;
97                  }
98              }
99              likelihood(sv, pc) = numSurvived/count(0, sv);
100          }
101      }
102      return likelihood;
103  }
104
105  mat get_likelihood_sex(vector<Passenger>& pass, mat count){
106      mat likelihood(2,2); likelihood.zeros();
107      for(int sv =0; sv<2; sv++){
108          for(int sx=0; sx<2; sx++){
109              int numSurvived = 0;
110              for(auto p: pass){
111                  if(p.survived==sv and p.sex==sx){
112                      numSurvived++;
113                  }
114              }
115              likelihood(sv, sx) = numSurvived/count(0, sv);
116          }
117      }
118      return likelihood;
119  }
120

```

We then wrote a function that returns the mean and variances of the continuous variable age. Those numbers can then be plugged into a probability density function to get the likelihood of a new age value.

```

149  const double pi = 3.14159265358979323846;
150  //function to calculate age likelihood
151  //run like this: calc_age_lh(6, 25.9, 138)
152  double calc_age_lh(double age, double mean, double var){
153      return 1/sqrt(2*pi*var)*exp(-(pow(age-mean, 2)/(2*var)));
154  }
155

```

*Function used to calculate the age likelihood.*

Finally, we write the function that calculates the raw probabilities for survived and perished. The numerator for perished is the product of the likelihoods of perishing for each class and the prior of perishing. The denominator for both is the total probability of perishing and surviving. The numerator for survived is the same as perished save for the likelihoods and prior being for surviving.

```

156 //function used to calculate raw probabilities
157 mat calc_raw_prob(int pclass, int sex, double age, mat& lh_pclass, mat&lh_sex, mat&age_mean_var, mat&apriori){
158     mat raw_prob(1,2); raw_prob.zeros();
159
160     pclass-=1;//for indexing purposes
161     double num_s = lh_pclass(1, pclass)*lh_sex(1, sex)*apriori(0, 1)*calc_age_lh(age, age_mean_var(0,1), age_mean_var(1,1));
162     double num_p = lh_pclass(0, pclass)*lh_sex(0, sex)*apriori(0, 0)*calc_age_lh(age, age_mean_var(0,0), age_mean_var(1,0));
163     double denominator = num_s + num_p;
164     raw_prob(0,1) = num_s / denominator;
165     raw_prob(0,0) = num_p / denominator;
166
167     return raw_prob;
168 }
169

```

We then run that function for each instance in the test set to get their raw probabilities. The values we got are identical to the ones from the R script.

```

PS C:\Users\emanu\Documents\School\U
0.3652 0.6348

0.8922 0.1078

0.6483 0.3517

0.8914 0.1086

0.6463 0.3537

0.6080 0.3920

0.1320 0.8680

0.8906 0.1094

```

*Raw probabilities from C++.*

```

Console Terminal x Jobs x
~/School/UTD S20/Machine Learning/HW4/Machine-Learning-Homework4/ ↗
>
> #get predictions from 5 test observations
> for (i in 1:nrow(test)){
+   raw <- calc_raw_prob(test$class[i], test$sex[i], test$age[i])
+   print(paste(raw[2], raw[1]))
+ }
[1] "0.365202992580656 0.634797007419344"
[1] "0.89222943053504 0.10777056946496"
[1] "0.648327773317862 0.351672226682138"
[1] "0.891438091992698 0.108561908007302"
[1] "0.646276526535963 0.353723473464037"
[1] "0.607961373922614 0.392038626077386"
[1] "0.132021458260991 0.867978541739009"
[1] "0.890587108178095 0.109412891821905"

```

*Raw probabilities from R.*

The R program was timed using `proc.time()` and the C++ program with a PowerShell script.

```

> endTime <- proc.time()
> #elapsed time
> endTime-startTime
      user  system elapsed
      0.05    0.00    0.37
>

```

*Runtime from Naïve Bayes was 50 milliseconds and...*

```

PS C:\Users\emanu\Documents\School\UTD S20\Machi
Days           : 0
Hours          : 0
Minutes        : 0
Seconds        : 0
Milliseconds   : 11
Ticks          : 114683
TotalDays      : 1.32734953703704E-07
TotalHours     : 3.18563888888889E-06
TotalMinutes   : 0.000191138333333333
TotalSeconds   : 0.0114683
TotalMilliseconds : 11.4683

Press Enter to continue...:

```

*11 milliseconds for C++.*

The run time for C++ is a little under 5 times faster. The confusion matrix was also calculated in the C++ code. We got the same accuracy as R.

```

Confusion Matrix:
  69.0000  25.0000
  10.0000  42.0000

Accuracy: 76.0274%
Specificity: 62.6866%
Sensitivity: 87.3418%

```

*C++ confusion matrix.*

#### Confusion Matrix and Statistics

```

              Reference
Prediction  0  1
0  69  25
1  10  42

Accuracy : 0.7603
95% CI : (0.6827, 0.827)
No Information Rate : 0.5411
P-Value [Acc > NIR] : 3.612e-08

Kappa : 0.5089

McNemar's Test P-Value : 0.01796

Sensitivity : 0.8734
Specificity : 0.6269

```

*R confusion matrix.*

<https://www.mathsisfun.com/data/standard-deviation.html> was used as a quick reference because I had forgotten how to calculate variance. <https://en.cppreference.com/w/> was used to program in C++. [http://arma.sourceforge.net/docs.html#example\\_prog](http://arma.sourceforge.net/docs.html#example_prog) was used to program using armadillo matrices.