# OPIM 5604 B15 – Predictive Modeling Assignment     Meghana Kasula (Net ID=mek15120)
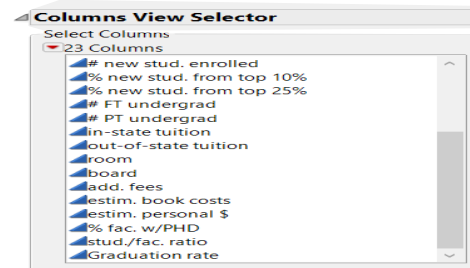
*"The work contained and presented here is my work and my work alone."*

14.4 University Rankings.

a. Use the *Columns Viewer* to produce numeric summaries of all of the variables. Note that many observations are missing some measurements. Clustering methods in JMP omit records with missing values. For the purposes of this exercise, we will assume that the values are missing purely at random (not that this may in fact not be the case). So our first goal is to estimate (impute) these missing values.

Here, we can see the number of missing values in each. Most are continuous values.



b. Select all of the continuous columns, and go to *Cols > Modeling Utilities >Explore Missing Values*. Select the option *Multivariate Normal Imputation,* and click *Yes Shrinkage*. Multivariate normal imputation uses least squares regression to predict the missing values from the nonmissing variables in each row. The imputed values will be highlighted in the data table. Save the data table with imputed values under a new name.

Using the imputation



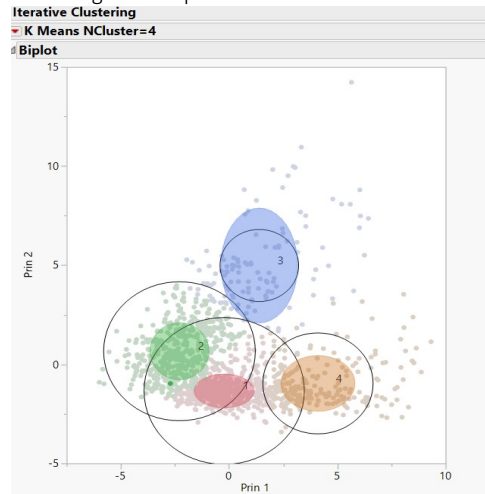method, we successfully filled the missing values with normalizations.



c. Use *k*-means clustering, with only the continuous variables. Use College Name as the label. Form clusters starting with 3 and ending with 8. What is the optimal number of clusters based on the Cubic Clustering Criterion (under *Cluster Comparison*)?

**Cluster Comparison**

| Method | NCluster | CCC | Best |
|---|---|---|---|
| Means Clustering | 3 | -5.5157 | |
| Means Clustering | 4 | -2.7521 | Optimal CCC |
| Means Clustering | 5 | -6.6236 | |
| Means Clustering | 6 | -3.0575 | |
| Means Clustering | 7 | -6.2707 | |
| Means Clustering | 8 | -4.6187 | |

umns Scaled Individually

**ontrol Panel**

K Means NCluster=3
umns Scaled Individually

**Cluster Summary**

| Cluster | Count | Step | Criterion |
|---|---|---|---|
| 1 | 166 | 24 | 0 |
| 2 | 346 | | |
| 3 | 790 | | |

**Cluster Means**

| Cluster | Math SAT | Verbal SAT | ACT | # appli. rec'd | # appl. accepted | # new stud. enrolled | % new stud. from top 10% | % new stud. from top 25% | # FT undergrad | # PT undergrad | in-state tuition | out-of-state tuition | room | board | add. fees | estim. book costs | estim. personal $ | % fac. w/PHD | stud./ fac. ratio | Gradu |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 521.784873 | 457.582968 | 22.3309692 | 8637.29434 | 6061.1874 | 2543.53569 | 28.0439495 | 57.7745846 | 13360.3023 | 3466.96985 | 2835.67794 | 7550.83013 | 2525.97713 | 1968.96324 | 664.779163 | 586.950104 | 1843.74973 | 81.6390453 | 17.1216867 | 54.61 |
| 2 | 572.207666 | 521.323788 | 24.9621106 | 2860.93931 | 1663.07514 | 574.690751 | 41.5429852 | 70.9379022 | 2309.77168 | 365.514715 | 13967.6271 | 14174.2244 | 2967.28965 | 2497.68066 | 363.732525 | 573.436468 | 1130.18133 | 81.6720014 | 11.884104 | 77.05 |
| 3 | 466.560804 | 428.484201 | 20.7924993 | 1462.9746 | 1080.3772 | 498.403142 | 15.2203872 | 40.1894928 | 2276.79124 | 870.264192 | 6236.87755 | 7478.97234 | 2310.41536 | 1878.48044 | 341.61183 | 531.338725 | 1433.16187 | 60.3706927 | 15.6942647 | 53.22 |

**Cluster Standard Deviations**

K Means NCluster=4
umns Scaled Individually

**Cluster Summary**

| Cluster | Count | Step | Criterion |
|---|---|---|---|
| 1 | 499 | 43 | 0 |
| 2 | 447 | | |
| 3 | 120 | | |
| 4 | 236 | | |

**Cluster Means**

| Cluster | Math SAT | Verbal SAT | ACT | # appli. rec'd | # appl. accepted | # new stud. enrolled | % new stud. from top 10% | % new stud. from top 25% | # FT undergrad | # PT undergrad | in-state tuition | out-of-state tuition | room | board | add. fees | estim. book costs | estim. personal $ | % fac. w/PHD | stud./ fac. ratio | Gradu |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 488.24207 | 451.083058 | 21.8107265 | 1013.04396 | 769.611935 | 294.872853 | 20.2449832 | 46.779018 | 1245.67535 | 451.389841 | 9836.48775 | 10018.1065 | 2715.32576 | 2160.7133 | 289.800227 | 539.695867 | 1281.82016 | 61.9198127 | 13.3653307 | 63.85 |
| 2 | 458.113179 | 416.898083 | 20.312279 | 2275.316 | 1668.78083 | 831.256941 | 12.7743012 | 37.6496257 | 3997.68474 | 1543.21983 | 3097.40421 | 5623.14823 | 2020.873 | 1677.18582 | 397.782422 | 530.279212 | 1585.65639 | 63.8730084 | 17.8464634 | 44.45 |
| 3 | 536.453461 | 467.627006 | 22.8251692 | 10235.5967 | 7092.91425 | 2846.58333 | 32.6720624 | 62.8512478 | 14923.4932 | 3295.68329 | 3373.77048 | 8133.04167 | 2623.67549 | 2048.07879 | 780.070638 | 599.39493 | 1839.70528 | 83.4256794 | 16.6775 | 59.06 |
| 4 | 594.912615 | 539.322847 | 25.7107798 | 3510.92797 | 1923.50847 | 654.682203 | 47.9672361 | 76.9937802 | 2611.35593 | 334.5297 | 14969.7257 | 15158.6957 | 2958.06658 | 2548.20975 | 381.570371 | 581.906316 | 1102.21091 | 84.9285525 | 11.4601695 | 80.29 |

As we can see, the optimal number of clusters are 4.

**d. Use the biplot, the parallel plot, and other built-in graphical and numeric summaries to explore the clusters. Save the clusters to the data table and use other graphical tools to compare and characterize the clusters. Summarize the key characteristics of each of the four clusters, and try label or name the clusters.**

Following is the Biplot:

Biplot





K Means NCluster=5
umns Scaled Individually

**Cluster Summary**

| Cluster | Count | Step | Criterion |
|---|---|---|---|
| 1 | 106 | 48 | 0 |
| 2 | 270 | | |
| 3 | 499 | | |
| 4 | 418 | | |
| 5 | 9 | | |

**Parallel Coordinate Plot**



The four key characteristics of clusters are:

1.Cluster 4 is having highest SAT and ACT scores, also high Graduation rate.
2.cluster 2 is relatively low or average in every field starting from Sat score to graduation rate.
3.Cluster 3 has high application accepted and enrolled.
4.Cluster 1has lowest application received, accepted and enrolled.





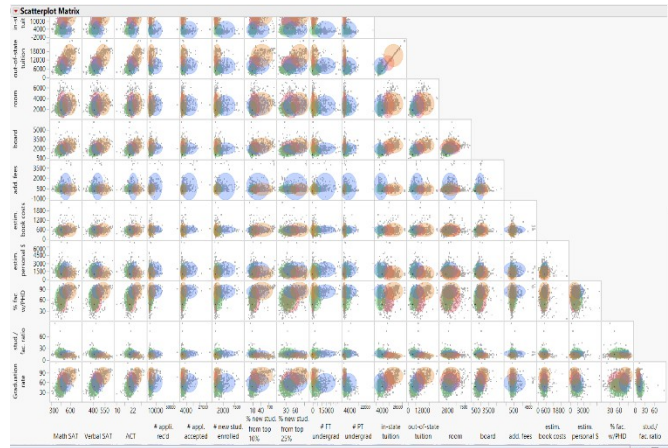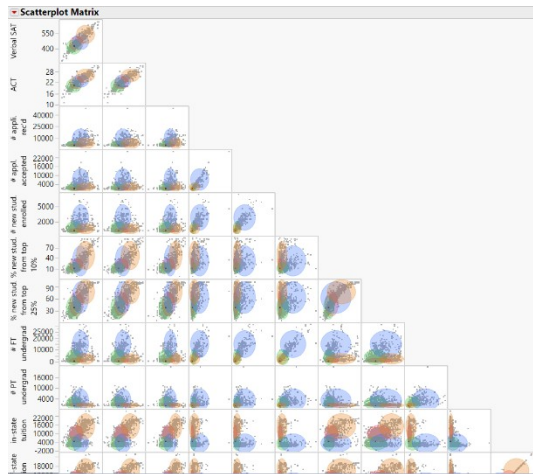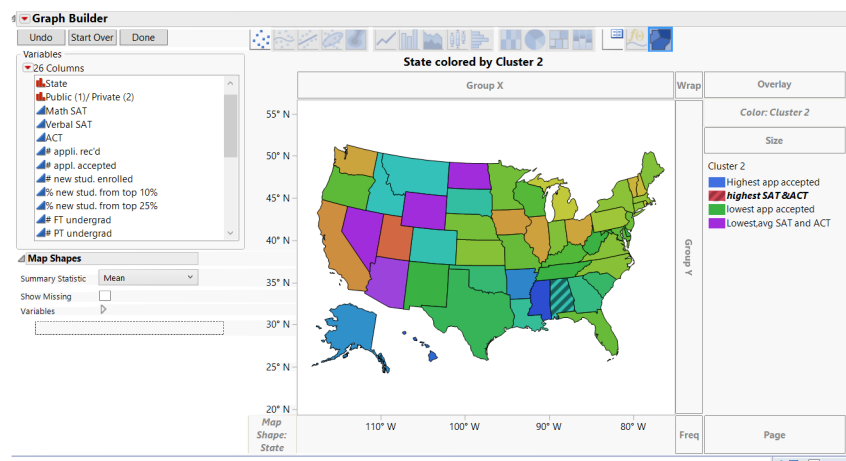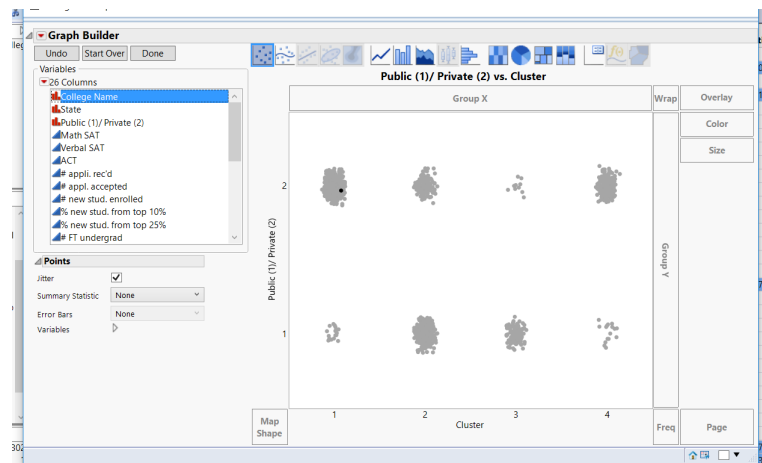| | ut-of-state tuition | room | board | add. fees | estim. book costs | estim. personal $ | % fac. w/PHD | stud./fac. ratio | Graduation rate | Cluster | Distance | Cluster 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 7560 | 1620 | 2500 | 130 | 800 | 1500 | 76 | 11.9 | 15 | 1 | 13.741212143 | highest SAT &ACT |
| 2 | 5226 | 1800 | 1790 | 155 | 650 | 2304 | 67 | 10 | 43.305222237 | 2 | 9.428937207 | Lowest,avg SAT and ACT |
| 3 | 5226 | 2514 | 2250 | 34 | 500 | 1162 | 39 | 9.5 | 39 | 2 | 10.873059871 | Lowest,avg SAT and ACT |
| 4 | 5226 | 2600 | 2520 | 114 | 580 | 1260 | 48 | 13.7 | 32.613560022 | 2 | 35.141165434 | Lowest,avg SAT and ACT |
| 5 | 3400 | 1108 | 1442 | 155 | 500 | 850 | 53 | 14.3 | 40 | 2 | 8.2442844558 | Lowest,avg SAT and ACT |
| 6 | 5600 | 1550 | 1700 | 300 | 350 | 1095.7842316 | 52 | 32.8 | 55 | 2 | 13.908998951 | Lowest,avg SAT and ACT |
| 7 | 4440 | 1755.2252521 | 1516.0911954 | 124 | 300 | 600 | 72 | 18.9 | 51 | 2 | 11.025005794 | Lowest,avg SAT and ACT |
| 8 | 3000 | 1960 | 1331.2033605 | 84 | 500 | 1666.0243071 | 48 | 18.7 | 15 | 2 | 7.8049221279 | Lowest,avg SAT and ACT |
| 9 | 6300 | 1782.5573008 | 1625.3496868 | 349.24133736 | 600 | 1908 | 85 | 16.7 | 69 | 3 | 5.6476419422 | Highest app accepted |
| 10 | 11660 | 2050 | 2430 | 120 | 400 | 900 | 74 | 14 | 72 | 4 | 7.2662320678 | lowest app accepted |
| 11 | 2970 | 1336 | 1200 | 20 | 500 | 1522.5754269 | 62 | 19.4 | 76 | 2 | 6.3498439966 | Lowest,avg SAT and ACT |
| 12 | 8080 | 1380 | 2540 | 100 | 500 | 1100 | 63 | 11.4 | 44 | 1 | 4.7034967116 | highest SAT &ACT |
| 13 | 2610 | 1030 | 1570 | 85 | 570 | 1500 | 66 | 20.1 | 33 | 2 | 2.884049357 | Lowest,avg SAT and ACT |
| 14 | 5780 | 2083.1191021 | 1670.3373083 | 357.66702339 | 512.55736262 | 1350.5342725 | 70 | 17.9 | 43 | 2 | 4.0156745543 | Lowest,avg SAT and ACT |
| 15 | 1740 | 1162 | 1287 | 243 | 570 | 2100 | 58 | 18.8 | 36 | 2 | 5.3601753888 | Lowest,avg SAT and ACT |
| 16 | 4000 | 1250 | 1450 | 325.00415671 | 600 | 1000 | 35 | 16.7 | 15 | 2 | 11.213085286 | Lowest,avg SAT and ACT |
| 17 | 6150 | 2095.5742417 | 1724.2056108 | 70 | 550 | 1200 | 59 | 16.6 | 52 | 1 | 3.6632156653 | highest SAT &ACT |
| 18 | 6639 | 2510.948436 | 1877.9294656 | 368 | 600 | 2500 | 52 | 14.5 | 48.170262761 | 2 | 9.4520321104 | Lowest,avg SAT and ACT |
| 19 | 8236 | 3700 | 1992.8614531 | 375.22676839 | 569 | 1650 | 74 | 14.7 | 61 | 1 | 6.8077720701 | highest SAT &ACT |
| 20 | 11478 | 2450 | 2338 | 645 | 500 | 900 | 80 | 14.4 | 70 | 1 | 4.2858433395 | highest SAT &ACT |
| 21 | 4460 | 1683 | 1071 | 100 | 600 | 1000 | 47 | 14.3 | 47 | 2 | 6.3967329946 | Lowest,avg SAT and ACT |
| 22 | 5666 | 1424 | 1540 | 418 | 1000 | 1400 | 56 | 15.5 | 46 | 2 | 15.406536705 | Lowest,avg SAT and ACT |
| 23 | 2883 | 1250 | 1120 | 75 | 300 | 1285.7347103 | 50 | 23 | 48 | 2 | 10.157458707 | Lowest,avg SAT and ACT |
| 24 | 6735 | 3395 | 1929.0460475 | 416.85424464 | 600 | 1425 | 70 | 12.2 | 65 | 1 | 3.1697497268 | highest SAT &ACT |
| 25 | 5424 | 1600 | 1930 | 632.22298916 | 580 | 1654 | 80 | 17.3 | 50 | 3 | 4.120257679 | Highest app accepted |
| 26 | 4440 | 1935 | 3240 | 291 | 750 | 2200 | 96 | 6.7 | 33 | 2 | 23.522743717 | Lowest,avg SAT and ACT |
| 27 | 4960 | 2500 | 1000 | 425.97178278 | 600 | 2100 | 83 | 12.7 | 38 | 2 | 12.582456221 | Lowest,avg SAT and ACT |
| 28 | 3108 | 1530 | 1320 | 183 | 480 | 1470 | 80 | 13.3 | 43 | 2 | 8.4278202872 | Lowest,avg SAT and ACT |
| 29 | 5715 | 1650 | 2780 | 351.47070975 | 600 | 2500 | 74 | 18.1 | 42 | 2 | 6.8894221716 | Lowest,avg SAT and ACT |
| 30 | 3552 | 840 | 1470 | 84 | 500 | 2800 | 54 | 42.6 | 41.571164124 | 2 | 30.513029669 | Lowest,avg SAT and ACT |
| 31 | 3216 | 2328 | 1518.1658046 | 143 | 450 | 1657.5046877 | 47 | 16.2 | 38.032822502 | 2 | 11.85747515 | Lowest,avg SAT and ACT |
| 32 | 8644 | 2382 | 1540 | 120 | 500 | 800 | 79 | 12.6 | 54 | 4 | 12.12055574 | lowest app accepted |
| 33 | 3460 | 1727.9970114 | 1446.9707085 | 60 | 450 | 1000 | 57 | 19.6 | 48 | 2 | 3.4389670594 | Lowest,avg SAT and ACT |
| 34 | 3720 | 1763.8339436 | 1493.5750728 | 130 | 500 | 1674.0079002 | 71 | 20 | 39 | 2 | 2.784528263 | Lowest,avg SAT and ACT |



e. Use the categorical variables that were not used in the analysis (State and Private/Public) to characterize the different clusters. (Hint: Create a geographic map to explore the clusters geographically.) Is there any relationship between the clusters and the categorical information?

The clusters spread can be viewed in the graph beside.
The legend is rightly mentioned and denoted.



In general, the amount of observations in private colleges are high and also, clust 1 and 4 dominantly has private institutes.

### f. Can you think of other external information that might explain the contents of some or all of these clusters?
1. Private institutes are in general have less rigorous criteria in selection than public institute. That is why we can observe more observations in private institutes.
2. In cluster 1, 2 we observed that low scores of ACT and SAT led to low acceptance rate which is a natural phenomenon.
3. Similarly, high score in SAT and ACT led to more acceptance range, observed in cluster 4.
4. Cluster 4, could have the observations of good academicians.

### g. Consider Tufts University. Which cluster does Tufts belong to? Which other universities is Tufts similar to, based on the clustering and the categorical variables?
Tufts is from cluster 4.
Hence, using local data filter, following are from the same cluster as Tufts.

**Local Data Filter**

Clear  Start Over  Favorites ▼

☑ Show ☑ Include
447 matching rows
☐ Inverse

▼ State (51)

HI (5)
IA (29)
ID (6)
IL (49)
IN (42)
KS (20)
KY (24)
LA (20)
MA (56)
MD (23)
ME (14)
MI (36)
MN (25)
MO (35)
MS (15)

▼ Public (1)/ Private (2) (2)
1 | 2

▼ Cluster 2 (4)
Highest app accepted (120)
highest SAT &ACT (499)
lowest app accepted (236)
Lowest,avg SAT and ACT (447)

Add Filter Columns
▼ 26 Columns
# FT undergrad
# PT undergrad
in-state tuition
out-of-state tuition
room
board
add. fees
estim. book costs
estim. personal $

Add
Cancel

**Tabulate**

Cluster 2 = Lowest,avg SAT and ACT

| Public (1)/ Private (2) | College Name | State | Cluster 2 | N |
|---|---|---|---|---|
| 1 | Adams State College | CO | Lowest,avg SAT and ACT | 1 |
| | Alabama Agri. & Mech. Univ. | AL | Lowest,avg SAT and ACT | 1 |
| | Alabama State University | AL | Lowest,avg SAT and ACT | 1 |
| | Albany State College | GA | Lowest,avg SAT and ACT | 1 |
| | Alcorn State University | MS | Lowest,avg SAT and ACT | 1 |
| | Angelo State University | TX | Lowest,avg SAT and ACT | 1 |
| | Arkansas State University | AR | Lowest,avg SAT and ACT | 1 |
| | Arkansas Tech University | AR | Lowest,avg SAT and ACT | 1 |
| | Armstrong State College | GA | Lowest,avg SAT and ACT | 1 |
| | Auburn University at Montgomery | AL | Lowest,avg SAT and ACT | 1 |
| | Augusta College | GA | Lowest,avg SAT and ACT | 1 |
| | Austin Peay State University | TN | Lowest,avg SAT and ACT | 1 |
| | Bemidji State University | MN | Lowest,avg SAT and ACT | 1 |
| | Black Hills State University | SD | Lowest,avg SAT and ACT | 1 |
| | Bloomsburg Univ. of Pennsylvania | PA | Lowest,avg SAT and ACT | 1 |
| | Bluefield State College | WV | Lowest,avg SAT and ACT | 1 |
| | Boise State University | ID | Lowest,avg SAT and ACT | 1 |
| | Bowie State University | MD | Lowest,avg SAT and ACT | 1 |
| | Bridgewater State College | MA | Lowest,avg SAT and ACT | 1 |
| | California State Univ. at Bakersfield | CA | Lowest,avg SAT and ACT | 1 |
| | California State Univ. at Dominguez Hills | CA | Lowest,avg SAT and ACT | 1 |
| | California State Univ. at Hayward | CA | Lowest,avg SAT and ACT | 1 |
| | California State Univ. at Los Angeles | CA | Lowest,avg SAT and ACT | 1 |
| | California State Univ. at San Bernardino | CA | Lowest,avg SAT and ACT | 1 |
| | California State University at Fresno | CA | Lowest,avg SAT and ACT | 1 |
| | California State University at Stanislaus | CA | Lowest,avg SAT and ACT | 1 |
| | California University of Pennsylvania | PA | Lowest,avg SAT and ACT | 1 |
| | Central Connecticut State University | CT | Lowest,avg SAT and ACT | 1 |
| | Central Missouri State University | MO | Lowest,avg SAT and ACT | 1 |
| | Central State University | OH | Lowest,avg SAT and ACT | 1 |
| | Central Washington University | WA | Lowest,avg SAT and ACT | 1 |
| | Chadron State College | NE | Lowest,avg SAT and ACT | 1 |
| | Cheyney University of Penn. | PA | Lowest,avg SAT and ACT | 1 |
| | Chicago State University | IL | Lowest,avg SAT and ACT | 1 |
| | Christopher Newport University | VA | Lowest,avg SAT and ACT | 1 |
| | Clarion University of Pennsylvania | PA | Lowest,avg SAT and ACT | 1 |
| | Cleveland State University | OH | Lowest,avg SAT and ACT | 1 |
| | Clinch Valley Coll. of the Univ. of Virginia | VA | Lowest,avg SAT and ACT | 1 |
| | Coastal Carolina University | SC | Lowest,avg SAT and ACT | 1 |
| | College of Charleston | SC | Lowest,avg SAT and ACT | 1 |
| | Columbus College | GA | Lowest,avg SAT and ACT | 1 |
| | Coppin State College | MD | Lowest,avg SAT and ACT | 1 |
| | CUNY - City College | NY | Lowest,avg SAT and ACT | 1 |

h. Return to the original data table (with the missing values). Run the same *k*means cluster analysis using this data.

i. Compare the results to those achieved after imputing missing values in terms of the number of clusters and the characteristics of the clusters? What are the key differences?



The main difference is that , the number of optimum clusters has drastically increased to 11.
Also, due to a lot of missing values, the clusters has to be more in lesser observations due to the vagueness and obscurity.
ii. In the initial analysis, we assumed that the values were missing at random and imputed the missing values. Describe why this approach was, or was not, a good strategy.
This approach is good, but the rows with missing values will be not considered for analysis. Hence it would omit some useful data.

Compute the metrics for the following rules: (Please utilize Bank-2.jmp)
a. CKING → SVG
Formula for
Support =P[A intersection B]/Total = (4329/7991) =0.543425
Confidence = P(A intersection B)/P (A) = (4329/7991)/(6855/7991)= 0.634562
Expected Confidence =P(B) = 4944/7991=0.6134
Lift =confidence/extreme confidence =0.63/0.61 =1.02001

**b. (CKING,SVG) →CD**
Support =P(A intersection B)/Total = (1139/7991) =0.1467
Confidence = P(A intersection B)/P (A) = (1139/7991)/(4329/7991)= 0.2623
Extreme Confidence =P(B) = 1960/7991=0.24345
Lift =confidence/extreme confidence =0.26/0.24 =1.0700