

7CCSMDM1 Data Mining Coursework 1

1.1 -

(i) number of instances - 48,842

(ii) number of missing values - 6,465

(iii) fraction of missing values over all attribute values - $6,465/633,108 = 1.02\%$

(iv) number of instances with missing values - 3,620

(v) fraction of instances with missing values over all instances - $3,620/48,842 = 7.41\%$

1.4 -

The results from when the missing results are simply deleted tend to have an increased error rate of around 3-5% with the data used in the coursework when compared to the other methods used. The method of replacing every Nan value for just the string 'missing' and the method of taking the mode for each attribute value and replacing the Nan's with that value produce very close error rates of around 16% when tested multiple times. The two methods are so close that depending on the run, one can slightly edge the other with a margin of 0.1%. This shows that both methods are almost equally as effective as one another when being used in decision trees.

2.1 -

	Fresh	Milk	Grocery	Frozen	Detergents Paper	Delicatessen
Mean	12000.297 727272728	5796.2659 09090909	7951.2772 72727273	3071.9318 18181818	2881.4931 818181817	1524.8704 545454545
Min	3	55	3	25	3	3
Max	112151	73498	92780	60869	40827	47943

2.1 -

Below is all the scatter plots for all the different cluster pairs for the attributes in the wholesale customers data set.

For the fresh attribute, there's a correlation all attributes except the following pairs:

Fresh/Frozen

Fresh/Delicatessen

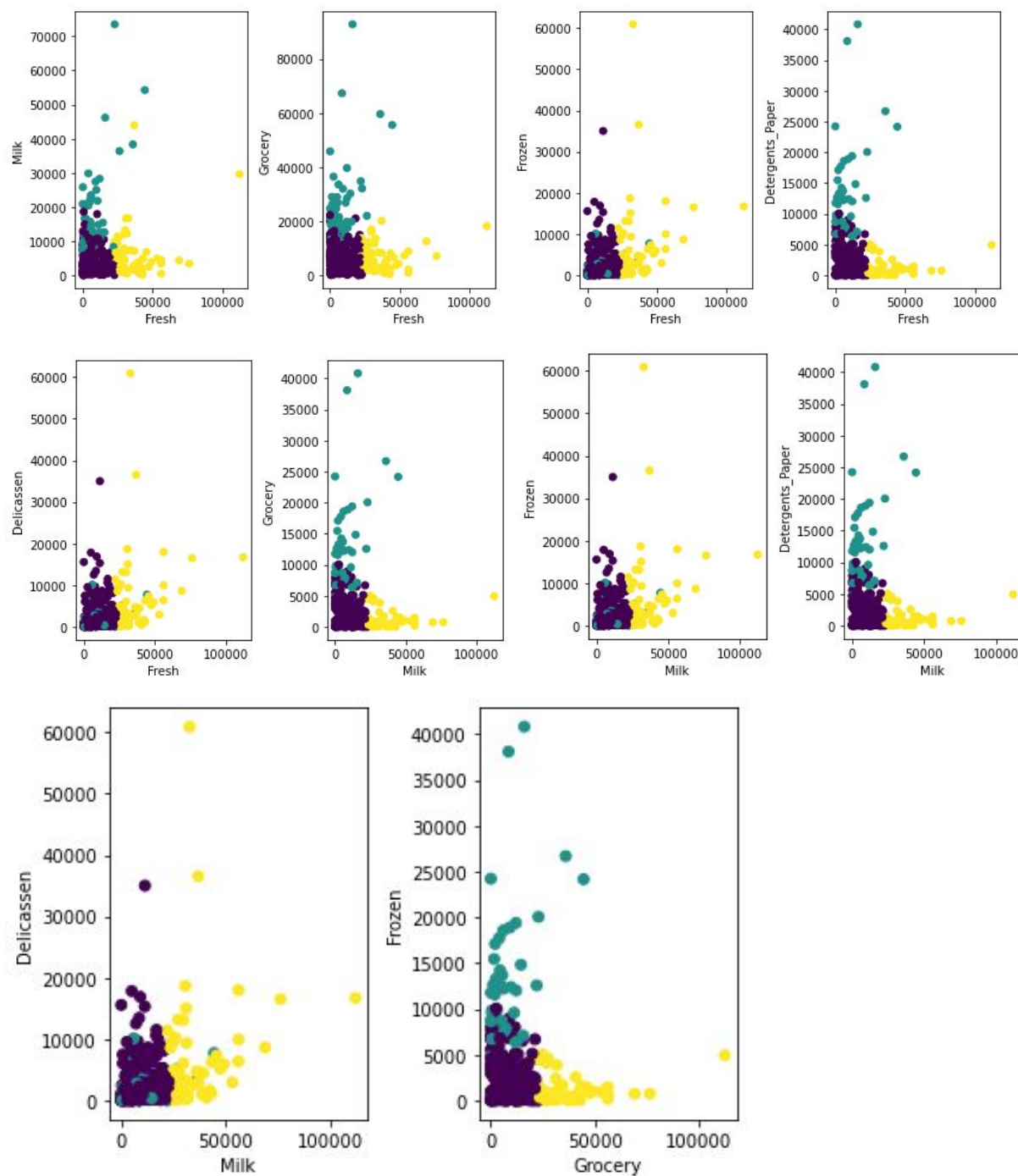
Milk/Frozen

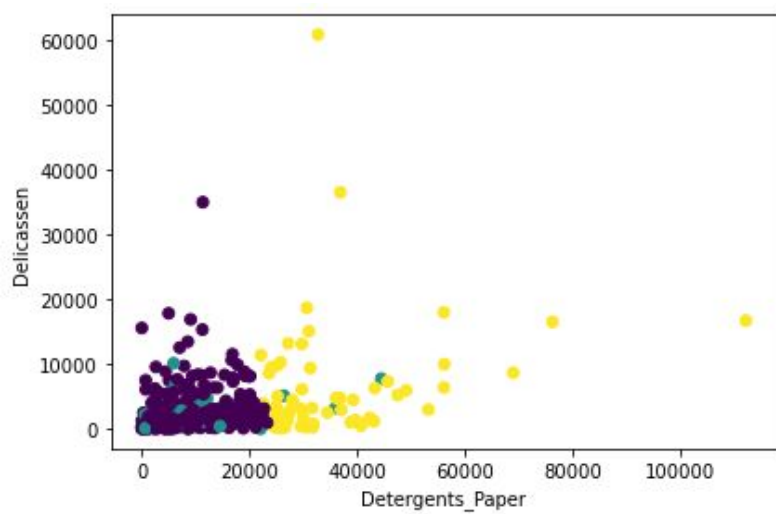
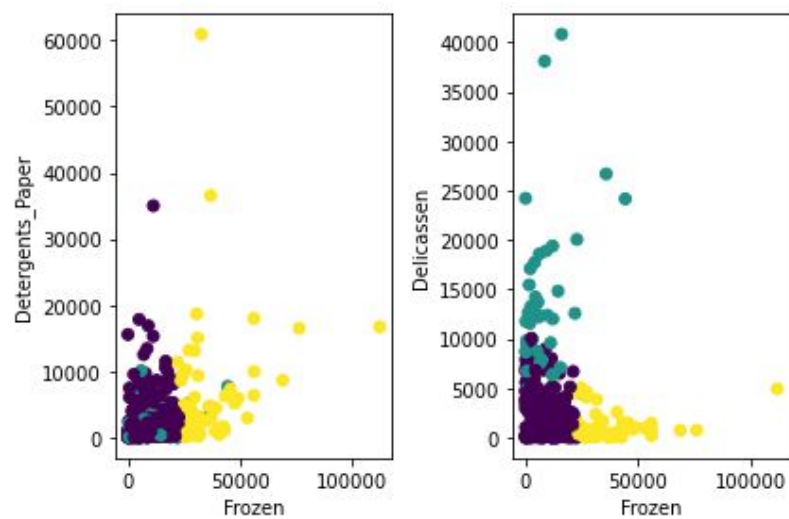
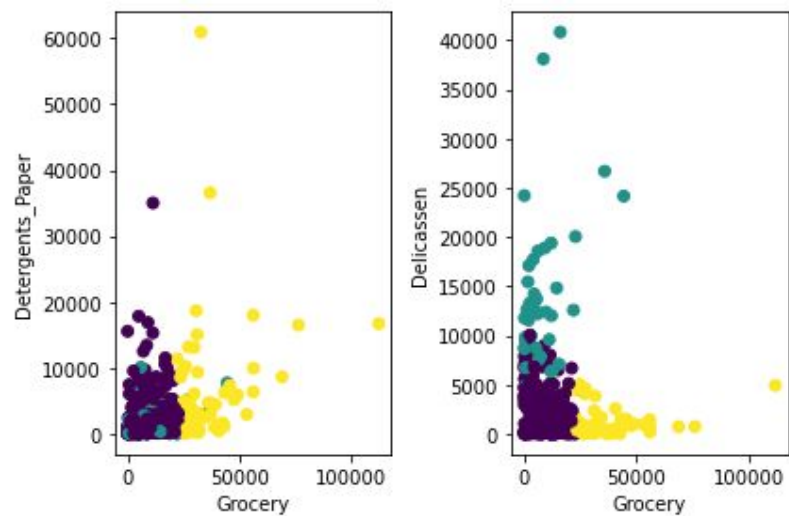
Milk/Delicatessen

Grocery/Detergents

Frozen/Detergents

Detergents/Delicatessen





2.3 -

As the clusters increased the within cluster scores decrease and the inverse happens with the between cluster scores, this would mean logically that, the clusters are more intense and the items are closer together within each cluster, and as the clusters increase, the distance between the clusters increase so the clusters are further apart. The ratio also increases as the number of clusters increases.

Below is the table of scores from the difference clusters.

	K = 3	K = 5	K = 10
BC	45239982498.87720 5	29454313223.72842	13756106907.27934 8
WC	1924917595.390871 8	12710533365.63535 7	81065549205.12784
BC/WC	0.042549034925878 8	0.431533856148298 87	5.893058970211275