

IME 672A: Data Mining & Knowledge Discovery

Assignment 2: Telecom Churn: Data Preprocessing

We have provided the data of a Telecom Company and we have to analyze on which factors the customers are churning their service. So, first we Cleanse the Data Set and after that we Analyze the data on various attributes and find out in which case, they most likely to churn from the service.

1. Data Cleansing:

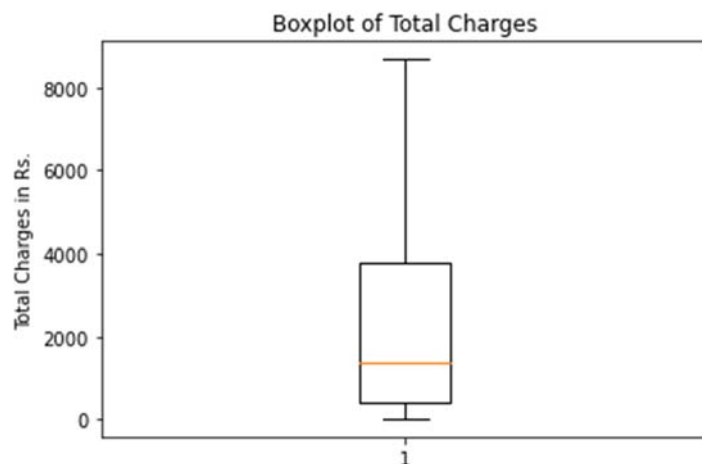
Missing Values:

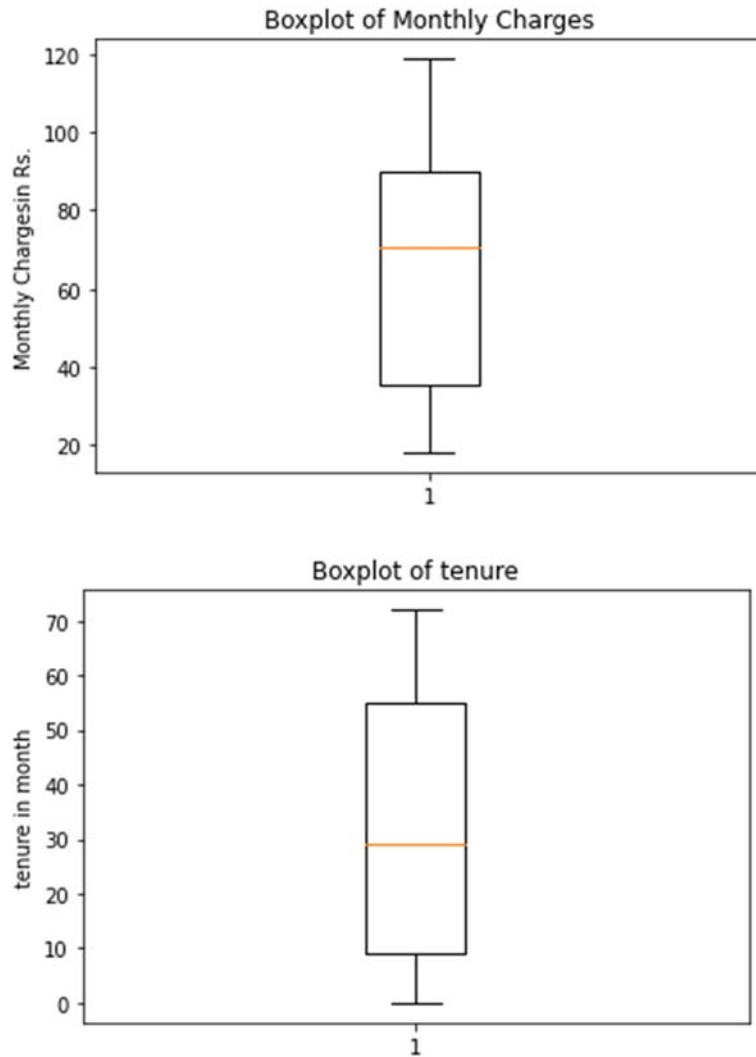
In the dataset we have provided the data of 7,043 customers with 21 attributes. Out of which 3 are numeric and 18 are categorical. After data cleansing, we convert 'TotalCharges' to numeric and got 11 null entries. After that we find out that the null values, we got have the customers with tenure of 0 month. So 'TotalCharges' of null entries is equal to 'MonthlyCharges' for these customers.

	SeniorCitizen	tenure	MonthlyCharges	TotalCharges
count	7043.000000	7043.000000	7043.000000	7043.000000
mean	0.162147	32.371149	64.761692	2279.798992
std	0.368612	24.559481	30.090047	2266.730170
min	0.000000	0.000000	18.250000	18.800000
25%	0.000000	9.000000	35.500000	398.550000
50%	0.000000	29.000000	70.350000	1394.550000
75%	0.000000	55.000000	89.850000	3786.600000
max	1.000000	72.000000	118.750000	8684.800000

Noise/Outliers:

Change the "No phone service" and "No internet service" with "No" in different attributes.





There is no Outlier in the attributes 'tenure', 'MonthlyCharges', 'TotalCharges' with IQR method.

2. Correlation:

Chi-Square test:

For Chi Square test if $p \leq 0.05$ then the attributes are Dependent or reject null hypothesis(H_0)

Chi Square test results for different attributes w.r.t Churn:

Chi-Square Test for gender

The Churns are independent on gender : Independent (fail to reject H_0)

Chi-Square Test for SeniorCitizen

The Churns are dependent on SeniorCitizen : Dependent (reject H0)

Chi-Square Test for Partner

The Churns are dependent on Partner : Dependent (reject H0)

Chi-Square Test for Dependents

The Churns are dependent on Dependents : Dependent (reject H0)

Chi-Square Test for PhoneService

The Churns are independent on PhoneService : Independent (fail to reject H0)

Chi-Square Test for MultipleLines

The Churns are dependent on MultipleLines : Dependent (reject H0)

Chi-Square Test for InternetService

The Churns are dependent on InternetService : Dependent (reject H0)

Chi-Square Test for OnlineSecurity

The Churns are dependent on OnlineSecurity : Dependent (reject H0)

Chi-Square Test for OnlineBackup

The Churns are dependent on OnlineBackup : Dependent (reject H0)

Chi-Square Test for DeviceProtection

The Churns are dependent on DeviceProtection : Dependent (reject H0)

Chi-Square Test for TechSupport

The Churns are dependent on TechSupport : Dependent (reject H0)

Chi-Square Test for StreamingTV

The Churns are dependent on StreamingTV : Dependent (reject H0)

Chi-Square Test for StreamingMovies

The Churns are dependent on StreamingMovies : Dependent (reject H0)

Chi-Square Test for Contract

The Churns are dependent on Contract : Dependent (reject H0)

Chi-Square Test for PaperlessBilling

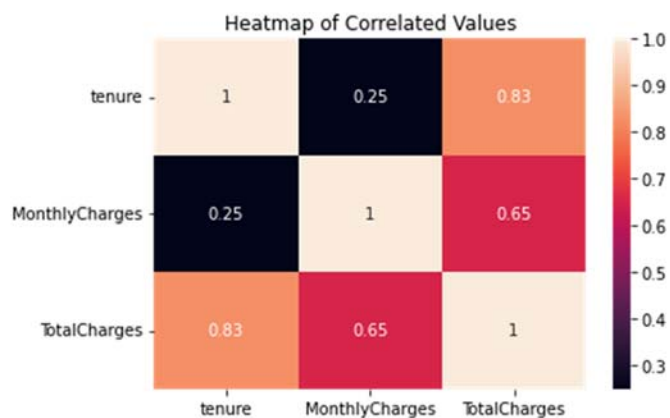
The Churns are dependent on PaperlessBilling : Dependent (reject H0)

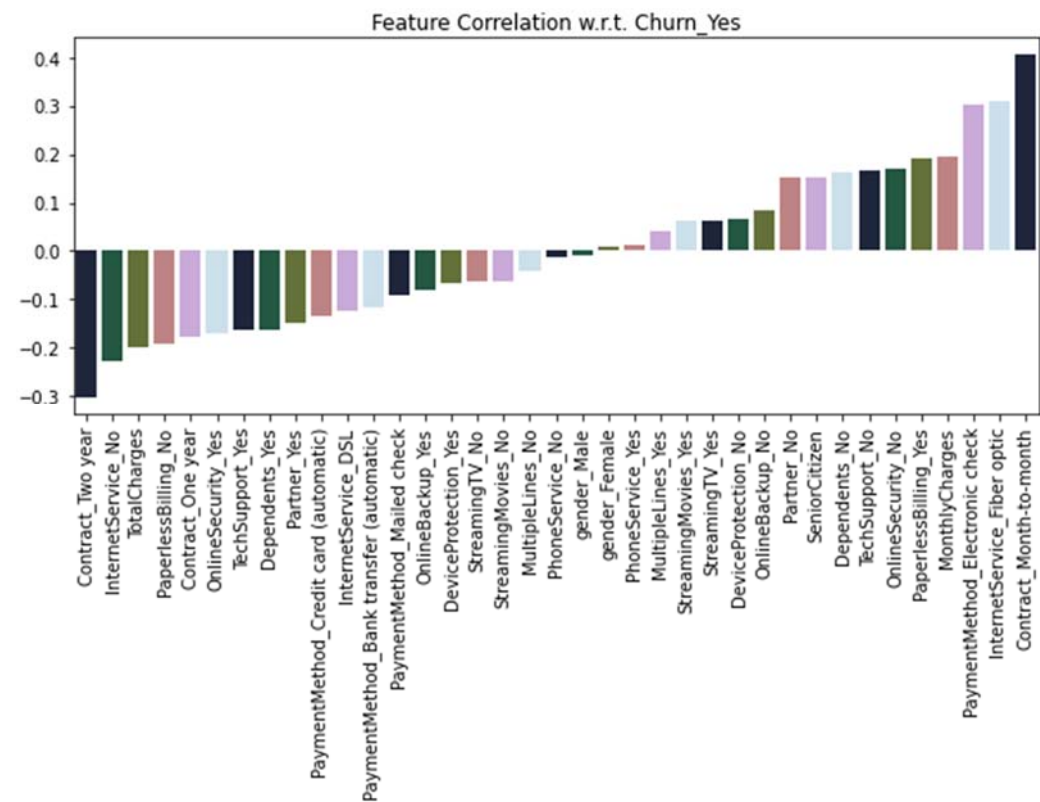
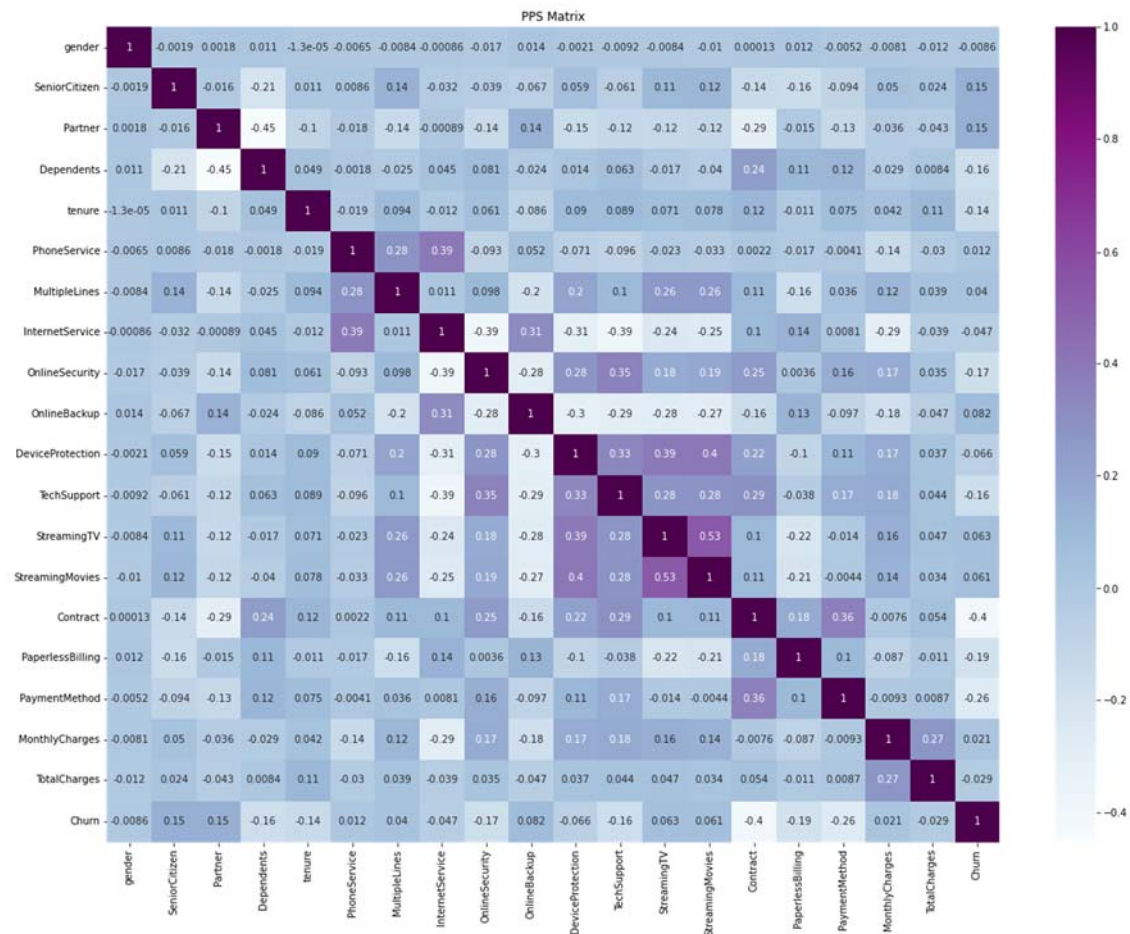
Chi-Square Test for PaymentMethod

The Churns are dependent on PaymentMethod : Dependent (reject H0)

Chi-Square Test for SeniorCitizen

The Churns are dependent on SeniorCitizen : Dependent (reject H0)

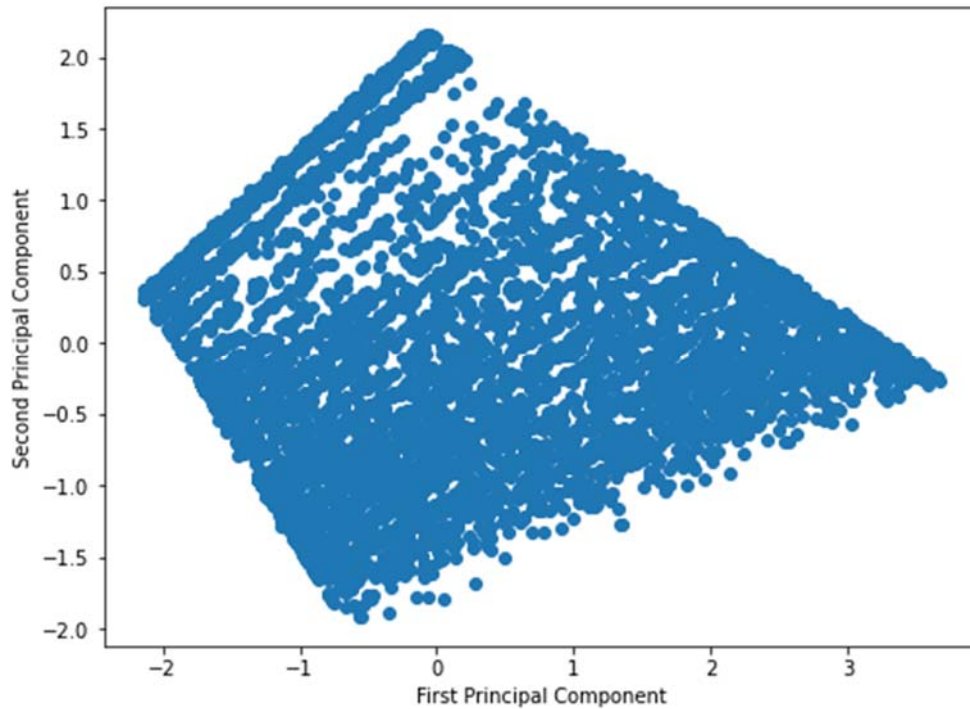




3. Data Reduction:

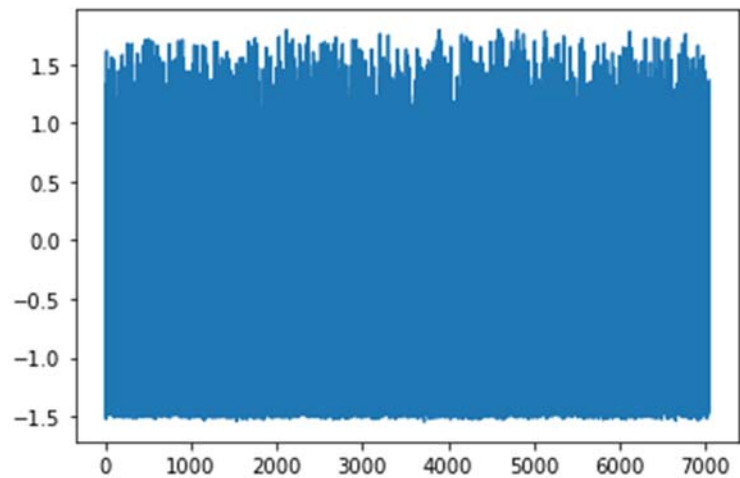
Data Reduction using Principal Component Analysis (PCA):

We reduce our data from 'tenure', 'MonthlyCharges' and 'TotalCharges' to PC1 and PC2.



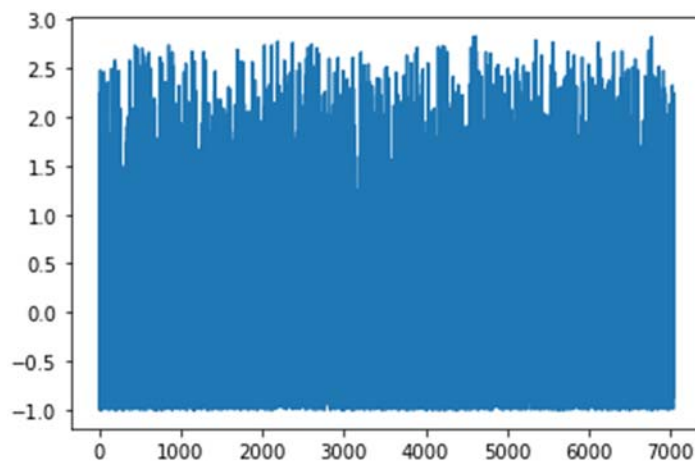
4. Data Transformation and Data Discretization:

Data Normalization using Z-score:



1. Z- Score for Monthly Charges:

2. Z- Score for Total Charges:



3. Z- Score for tenue:

