**Team NC-227027**
**Team members**
- Adarsh S
- Narendran Omprakash

# No Code ML

**24ᵗʰ January 2022**

## DATASET OVERVIEW AND IMPORT (Provided in Dataset description)

- We import the dataset into Knowledge Studio using the excel import node
- The shape and size of the data are given below
  - 7043 customer details are provided with their own unique and random customerID
  - Demographical characteristics such as Gender, Senior Citizenship, and presence of partner or dependents are provided.
  - The telecom company provides two basic services: Phone Service and Internet service. Certain other bundled services can also be availed by the customer if they use each of these services. The phone service allows you to get a multiple-line collection.
  - Details of the payment duration, method, and amount for customers have been provided.
  - Finally, the important (output) feature has been provided- whether or not the customer discontinued (stopped using) the service in the past few months.

## Size and shape of data

| # | Field Name | Field Label | Data Type | Cardinality | Mean | andard Deviatic |
|---|---|---|---|---|---|---|
| 1 | customerID | customerID | String | 7043 | | |
| 2 | gender | gender | String | 2 | | |
| 3 | SeniorCitizen | SeniorCitizen | Number | 2 | 0.16 | 0.37 |
| 4 | Partner | Partner | String | 2 | | |
| 5 | Dependents | Dependents | String | 2 | | |
| 6 | tenure | tenure | Number | 73 | 32.37 | 24.56 |
| 7 | PhoneService | PhoneService | String | 2 | | |
| 8 | MultipleLines | MultipleLines | String | 3 | | |
| 9 | InternetService | InternetService | String | 3 | | |
| 10 | OnlineSecurity | OnlineSecurity | String | 3 | | |
| 11 | OnlineBackup | OnlineBackup | String | 3 | | |
| 12 | DeviceProtection | DeviceProtection | String | 3 | | |
| 13 | TechSupport | TechSupport | String | 3 | | |
| 14 | StreamingTV | StreamingTV | String | 3 | | |
| 15 | StreamingMovies | StreamingMovies | String | 3 | | |
| 16 | Contract | Contract | String | 3 | | |
| 17 | PaperlessBilling | PaperlessBilling | String | 2 | | |
| 18 | PaymentMethod | PaymentMethod | String | 4 | | |
| 19 | MonthlyCharges | MonthlyCharges | Number | 1585 | 64.76 | 30.09 |
| 20 | TotalCharges | TotalCharges | Number | 6531 | 2,283.30 | 2,266.77 |
| 21 | Discontinued | Discontinued | String | 2 | | |

## DATA CLEANUP

### Missing value treatment

- Missing value substitution is done by the following criteria
  - Discrete Variable - Substitute Mode
  - Continuous Variable - Substitute Mean
- 11 entries in TotalCharges are missing. This is treated by substituting the mean.
- This can be done using the Variable transformations node in Knowledge studio using the following expression

      CASE WHEN [TotalCharges] IS NULL THEN Avg([TotalCharges]) ELSE [TotalCharges] END

### Outliers Treatment

The monthly and total charges columns don't contain any outliers, they lie in the range that was previously specified in the data documentation.

### Variable Selection

To highlight potential predictors, we use a variety of different measures such as Entropy, Gini, Information Value, likelihood test statistic ratio, etc. By using the union option, we select the top 2 variables for each measure.

## DATA ANALYSIS

### Preliminary observations

- The dataset has an equal number of male (50.48%) and female (49.52%) records.
- 16.2% of the customers in the dataset are senior citizens.
- Almost half of the customers in the dataset have a partner (48.30%).
- 30% of the customers have dependants.
- The tenure of the customers seems to follow a roughly symmetrical bimodal distribution with peaks at 0-8 months or 65-72 month periods.
- Only about 10% of customers do not avail the phone service.
- 43% of customers avail fiber optic internet service while 34% use DSL and 21% have not opted for internet service.
- 55% of customers are on a month-to-month contract while 20% are on a yearly contract and 24% are on a biannual contract.
- 26.54% of the customers in the dataset have discontinued the service.

### Graphs and charts

**Tenure**



No. of customers vs. tenure

**Additional services**

## Online security

Yes
28.7%

No internet
21.7%

No
49.7%

## Online Backup

Yes
34.5%

No internet
21.7%

No
43.8%

## Device protection

Yes
34.4%

No internet
21.7%

No
43.9%

## Tech Support

Yes
29.0%

No internet
21.7%

No
49.3%

## Streaming TV

Yes
38.4%

No internet
21.7%

No
39.9%

## Streaming Movies

Yes
38.8%

No internet
21.7%

No
39.5%

**Payment methods**

## Payment methods



Mailed check
22.9%

Bank transfer
21.9%

Credit card
21.6%

Electronic check
33.6%

## Monthly charges

## Monthly Charges



No of customers (y-axis)

MonthlyCharges (x-axis)

Bins: [18.25, 28.30), [28.30, 38.35), [38.35, 48.40), [48.40, 58.45), [58.45, 68.50), [68.50, 78.55), [78.55, 88.60), [88.60, 98.65), [98.65, 108.70), [108.70, 118.75]

## Segment Analysis
We take every value and compare it against the target variable (Discontinued)

| R | Field Name | Whole Dataset | Discontinued=No | Discontinued=Yes |
|---|---|---|---|---|
| 1 | gender |  |  |  |
| 2 | SeniorCitizen |  |  |  |
| 3 | Partner |  |  |  |
| 4 | Dependents |  |  |  |

| R | Field Name | Whole Dataset | Discontinued=No | Discontinued=Yes |
|---|---|---|---|---|
| 5 | tenure |  |  |  |

## Observations

- Gender has no Information value and hence can be eliminated. Partner, Senior citizen and Dependants have an Information value of [0.08, 0.16].
- Customers with less tenure seem to be more probable to discontinue than customers with large tenures.

| R | Field Name | Whole Dataset | Discontinued=No | Discontinued=Yes |
|---|---|---|---|---|
| 6 | PhoneService | | | |
| 7 | MultipleLines | | | |
| 8 | InternetService | | | |
| 9 | OnlineSecurity | | | |
| 10 | OnlineBackup | | | |
| 11 | DeviceProtection | | | |

**Observations**

- Phone service, multiple lines have no Information value.
- Customers who opt for fiber optic internet seem more likely to discontinue the service
- Customers who opt for additional services (Online backup, Online security, Device protection, Tech support, Streaming TV, Streaming Movies) are generally less likely to discontinue the service.

| R | Field Name | Whole Dataset | Discontinued=No | Discontinued=Yes |
|---|---|---|---|---|
| 12 | TechSupport | | | |
| 13 | StreamingTV | | | |
| 14 | StreamingMovies | | | |
| 15 | PaperlessBilling | | | |
| 16 | PaymentMethod | | | |

## Observations

- Customers who opt for automatic payments are less likely to discontinue the service. Customers who send electronic checks discontinue the most.

| R | Field Name | Whole Dataset | Discontinued=No | Discontinued=Yes |
|---|---|---|---|---|
| 17 | MonthlyCharges |  |  |  |
| 18 | TotalCharges |  |  |  |
| 19 | Contract |  |  |  |

## Observations

- Customers who discontinue are mostly on a monthly contract (88%).
- Only 3% of customers with two-year contracts and 8% of customers with annual contracts have discontinued, while 43% of customers on monthly contracts have discontinued.
- Most customers (90%) with low monthly charges have not discontinued.

# Characteristic Analysis

Characteristic analysis charts

## Contract



## Monthly charges

## Payment method



DV - Discontinued

## Tenure



DV - Discontinued

## DATA PRE-PROCESSING

### Target attribute: Discontinued=Yes

### Train-test split

**The dataset is split into the following partitions:**

- Train - 60%
- Test - 30%
- Fit - 10%

This can be done using Knowledge studio's partition node.

### Checking statistical characteristics of train and test data

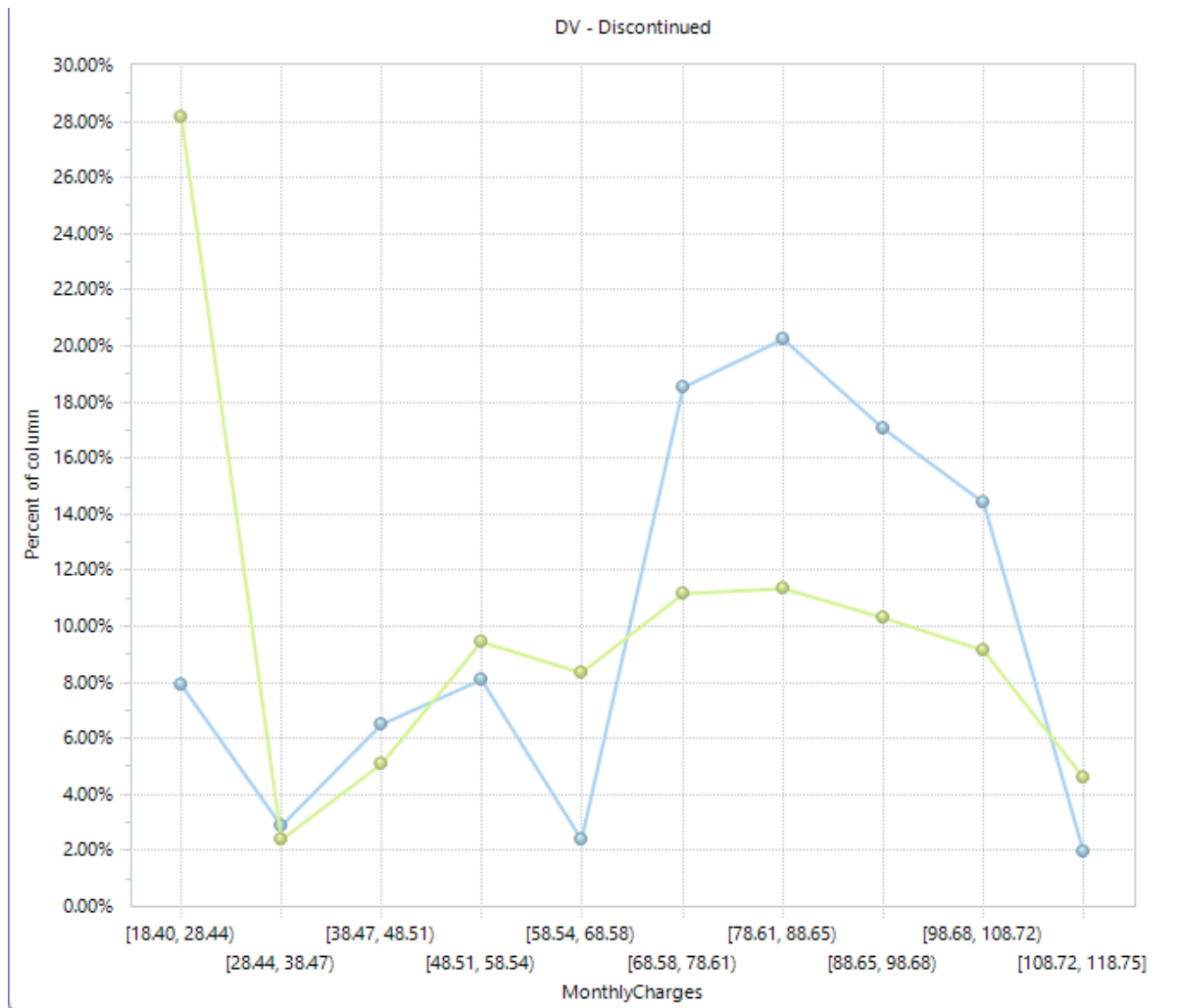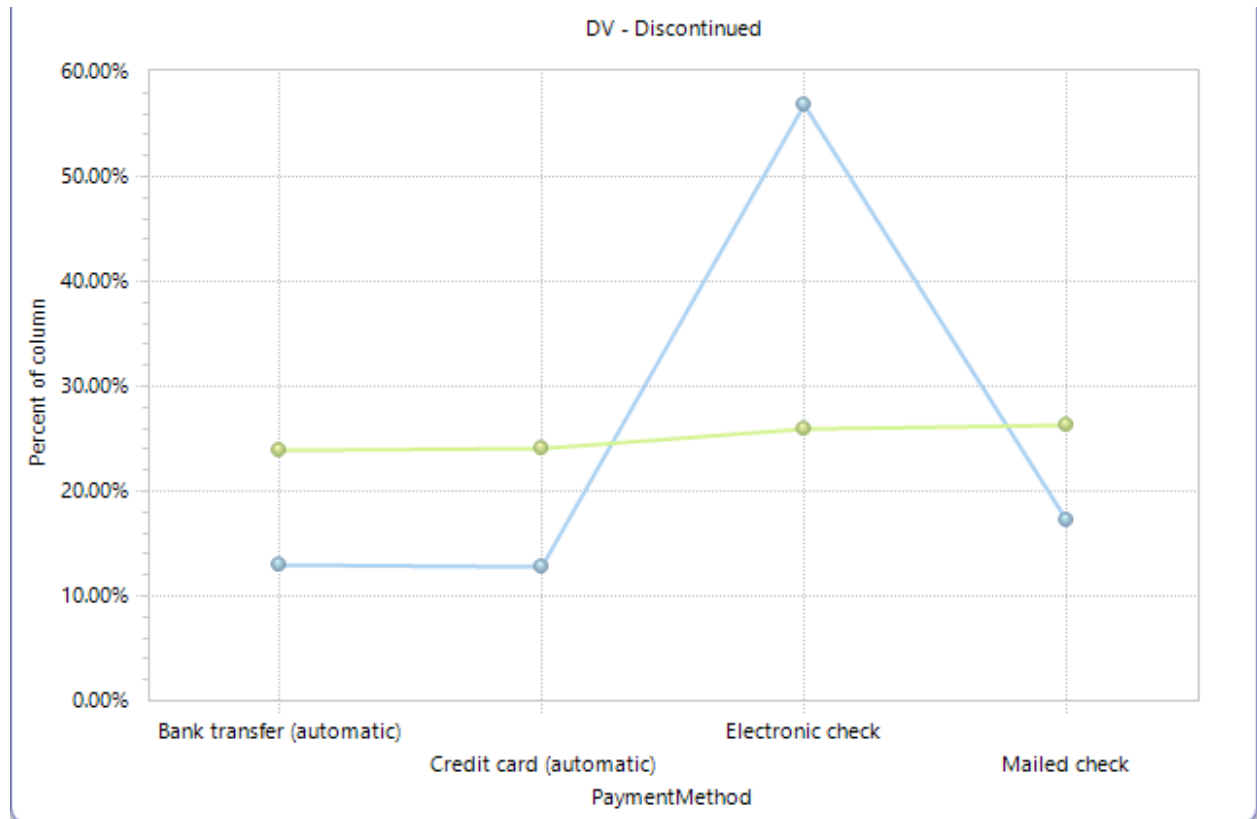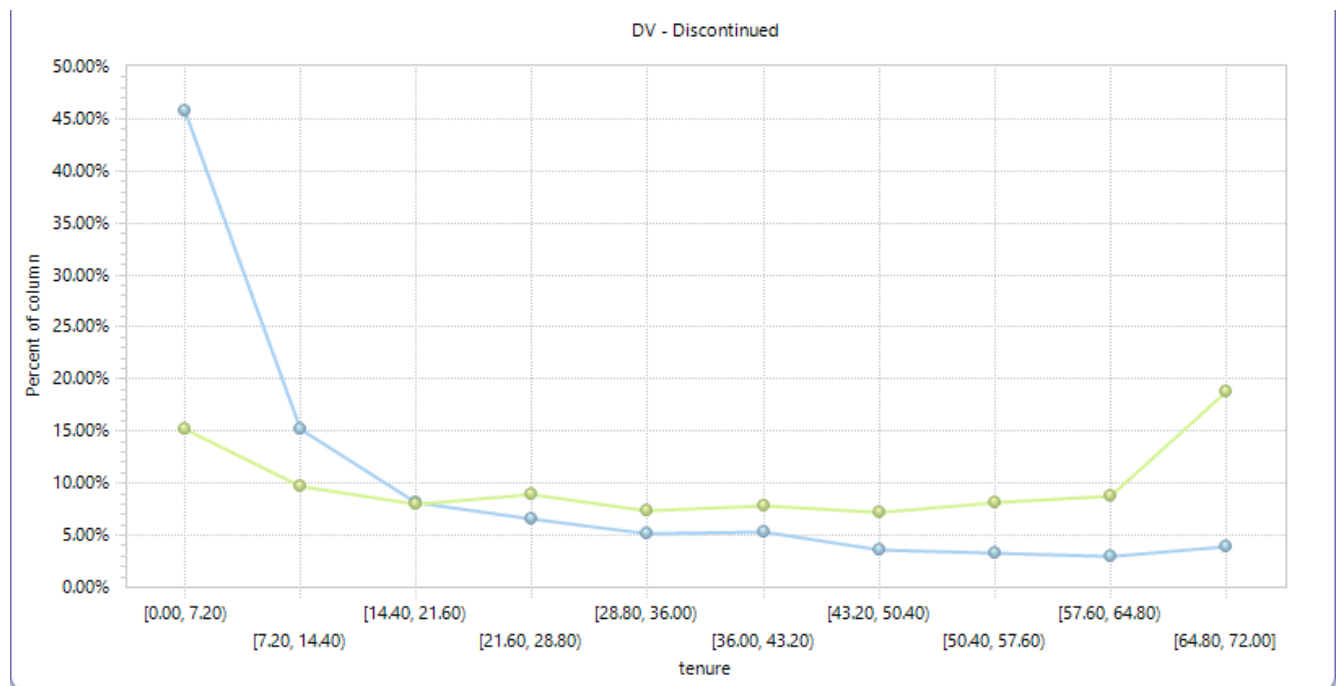We check for some of the comparable attributes of the train, test and fit data to ensure they have similar characteristics

| Attribute | Train | Test | Fit |
|---|---|---|---|
| Cardinality | 4225 | 2113 | 705 |
| Gender | Male: 50.37%<br>Female: 49.63% | Male: 48.84%<br>Female: 51.16% | Male: 56.03%<br>Female: 43.97% |
| Senior citizen | Yes: 16.73% | Yes: 16.31% | Yes: 15.14% |
| Phone service | Yes: 90.65% | Yes: 89.65% | Yes: 89.87% |
| Internet service | DSL: 34.27%<br>Fiber optic: 44.00%<br>No: 21.73% | DSL: 33.19%<br>Fiber optic: 45.67%<br>No: 21.13% | DSL: 34.97%<br>Fiber optic: 43.30%<br>No: 21.72% |
| Contract | Monthly: 55.72%<br>One year: 20.47%<br>Two year: 23.81% | Monthly: 54.09%<br>One year: 21.39%<br>Two year: 24.51% | Monthly: 53.62%<br>One year: 22.13%<br>Two year: 24.26% |
| Tenure | Mean: 31.73<br>Std. deviation: 24.42 | Mean: 32.84<br>Std. deviation: 24.83 | Mean: 34.80<br>Std. deviation: 24.44 |
| Monthly charges | Mean: 64.61<br>Std. deviation: 29.90 | Mean: 64.41<br>Std. deviation: 30.02 | Mean: 66.72<br>Std. deviation: 31.38 |
| Total charges | Mean: 2218.25<br>Std. deviation: 2237.5 | Mean: 2319.13<br>Std. deviation: 2278.4 | Mean: 2565.79<br>Std. deviation: 2381.8 |

## MODEL TRAINING AND TESTING

### Decision Trees

Decision trees are good at identifying potential predictors that can be used in other models. Their aim is to resolve a dependent variable from a series of independent variables.

Looking at the split report given by the decision tree, we can see that the variable "Contract" plays a vital role in predicting the dependent variable as suggested by its entropy value. The higher the entropy of a particular variable, the lower the information gained from that variable.
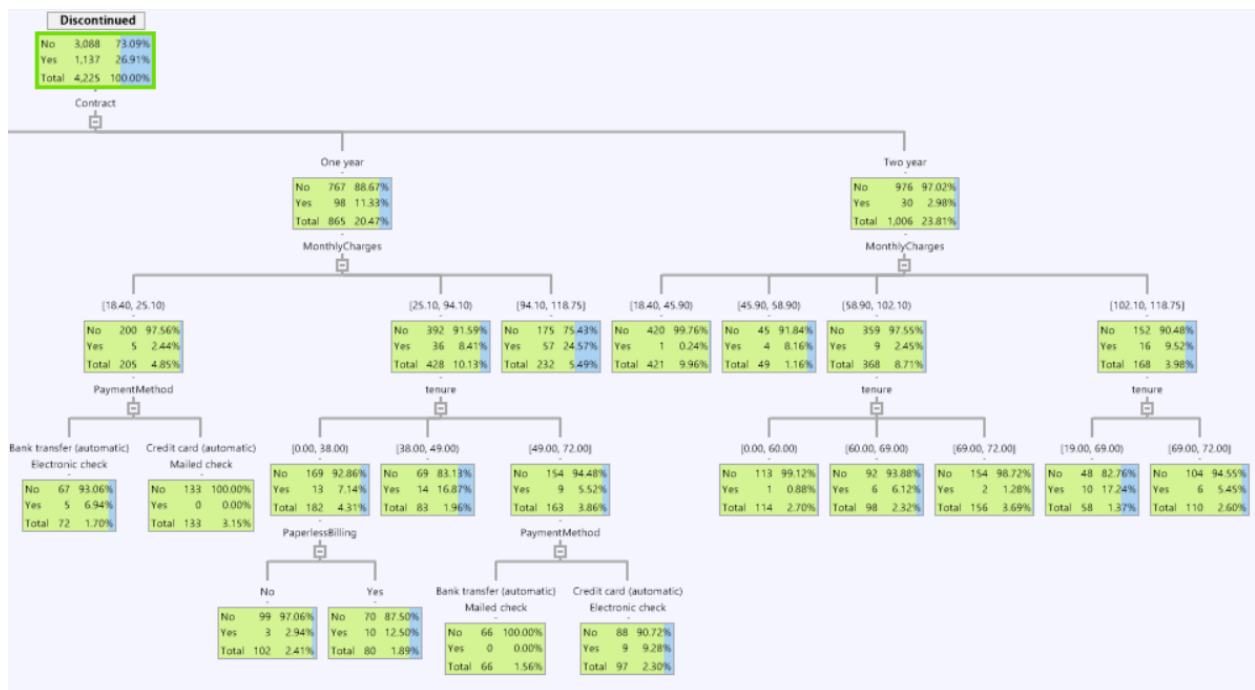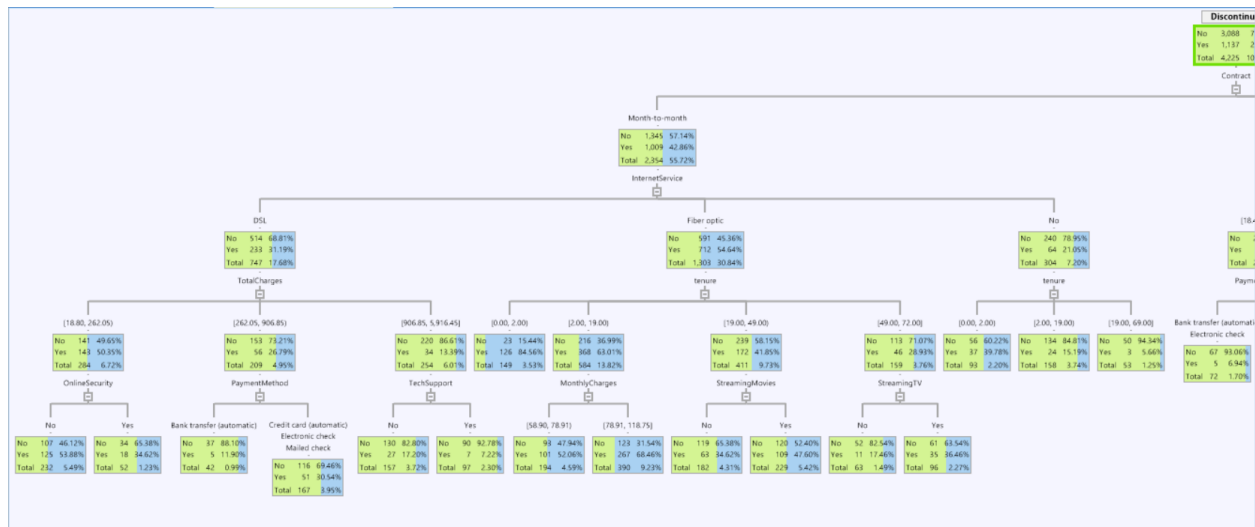
Insert Decision Tree - Split Report

| Variable Name | Rank | Info | Include | info value | input entropy | output entropy Δ | ratio entropy |
|---|---|---|---|---|---|---|---|
| Contract | 1 | 16.75868 | ☑ | 1.21256 | 0.84019 | 0.69938 | 0.16759 |
| tenure | 2 | 12.59267 | ☑ | 0.80747 | 0.84019 | 0.73439 | 0.12593 |
| TechSupport | 3 | 11.0053 | ☑ | 0.70289 | 0.84019 | 0.74772 | 0.11005 |
| OnlineSecurity | 4 | 10.66243 | ☑ | 0.68257 | 0.84019 | 0.7506 | 0.10662 |
| InternetService | 5 | 9.8985 | ☑ | 0.63153 | 0.84019 | 0.75702 | 0.09898 |
| OnlineBackup | 6 | 8.36715 | ☑ | 0.54133 | 0.84019 | 0.76989 | 0.08367 |
| DeviceProtection | 7 | 7.66751 | ☑ | 0.50056 | 0.84019 | 0.77577 | 0.07668 |
| PaymentMethod | 8 | 7.00115 | ☑ | 0.41379 | 0.84019 | 0.78137 | 0.07001 |
| MonthlyCharges | 9 | 5.67902 | ☑ | 0.36337 | 0.84019 | 0.79247 | 0.05679 |
| StreamingTV | 10 | 5.48697 | ☑ | 0.37525 | 0.84019 | 0.79409 | 0.05487 |
| TotalCharges | 11 | 5.37844 | ☑ | 0.31858 | 0.84019 | 0.795 | 0.05378 |
| StreamingMovies | 12 | 5.35491 | ☑ | 0.36776 | 0.84019 | 0.7952 | 0.05355 |
| PaperlessBilling | 13 | 3.1693 | ☑ | 0.19309 | 0.84019 | 0.81356 | 0.03169 |
| Dependents | 14 | 2.66133 | ☑ | 0.1654 | 0.84019 | 0.81783 | 0.02661 |
| Partner | 15 | 1.90953 | ☑ | 0.1143 | 0.84019 | 0.82414 | 0.0191 |
| SeniorCitizen | 16 | 1.53783 | ☑ | 0.08883 | 0.84019 | 0.82727 | 0.01538 |

Analyze

Cancel    < Back    Next >    Save    Run    Help

Entropy variance is used as the measure for the decision tree splits at each level. The split search method is cluster. Automatic tree generation is selected for building the tree.

# Explaining the decision Tree





The first split is done using the Contract variable. It can be observed that 97% of the two-year contracts and 89% of the one-year contracts are not discontinued. But, 43% of the month-month contracts are discontinued.

In the 43% of the contracts that are discontinued, the majority of the users having no internet connection seem to be satisfied with the service having just a 21% rate of discontinuation. Users using fiber optics, on the other hand, have a higher rate of discontinuation with 54%. Within the first two months, 45% of the users who used fiber optics discontinued. Users using DSL, have a

comparatively lower rate of discontinuation at 31% but the ones who used it the least with very low total charges discontinued their service.

## Bagging, Boosting and Random Forests - Ensemble methods

Boosting works by recursively training a particular tree and correcting the errors made by that tree. When correcting the errors made by the initial tree, we create another tree that can classify the records correctly.

Bagging iteratively takes a subset of the dataset and creates a large number of trees. During prediction, they take an average of the predictions made by each of the trees.

Random forest is similar to bagging but it also takes a subset of variables for each of the trees.

For boosting, similar parameters to whatever was given for the decision tree algorithm is given. That is, the split is taken into consideration by Entropy variance and the split search method is cluster. The number of iterations for the model to improve upon its previous tree's failures is set to 250.

Similar properties were set for the bagging model, on top of that, sampling with replacement is checked. Sampling with replacement allows us to reduce the variance without increasing the bias.



As previously mentioned, in random forests, along with the training data, the features are sampled as well. The percentage of features that are sampled is given as 33% for our random forest model.

## Deep Learning

While using the deep learning node in Knowledge studio, we use the test dataset - Dataset_fit, to guard against the model overfitting on the training dataset. Overfitting happens when the model trains too much on the training data to a point where it memorizes it which leads to poor generalization. A model that has overfitted on a training dataset will perform poorly on the testing dataset.

We select all the variables to train the deep learning model. A singular value threshold is set by default, which helps in finding collinearity and excluding variables that are collinear.

The link function is set to logistic which helps in non-linear transformation. The optimizer is set to conjugate gradient, I have no idea what it is.

## Regularized Regression

To overcome overfitting, a regularization parameter is added to a model's loss function which penalizes weights and drives them towards zero. L1 regularization can even drive weights towards zero.

We are going to use a ridge regression model, ridge regression uses L2 regularization which squares the weights and adds them to the loss function thereby penalizing them.

**Regularization - Variable Selection**                           — ☐ ✕

**Dependent Variable**

Link function:        Logit (Logistic Regression)           ⌄

Dependent variable:   Discontinued                      ⌄   👓

Target category:      Yes                                    ⌄

**Regularization Type**

Method:               Ridge Regression                       ⌄

🔘 Lagrangian Form (additive terms appended to deviance)

⚪ Constrained Form (minimize deviance subject to constraint)

Standardization options:   Report coefficients for un-standardized data   ⌄

Standardization type:      Center                                         ⌄

☐ Supress intercept        ☐ Scale $\lambda$ by number of records

☐ Regularize intercept

**Parameters**

Minimize:             Deviance + $\lambda \| \beta \|_z^2$

$\lambda =$           1

**Independent Variable Selection**

Not selected:  👓                    Selected:  👓

|  |  |
|---|---|
|  | Contract |
|  | Dependents |
|  | DeviceProtection |
| > | InternetService |
| < | MonthlyCharges |
|  | OnlineBackup |
| >> | OnlineSecurity |
| << | PaperlessBilling |
|  | Partner |
|  | PaymentMethod |
|  | SeniorCitizen |
|  | StreamingMovies |

Attribute Editor

Cancel    < Back    Next >    Save    Run    Help

## Train and Test accuracies

## Model Selection

For selecting the best model from the previously explained set of models in which we trained our dataset, we use the Model Selector node in Knowledge Studio. We choose the evaluation metric as AUC (Area under the ROC curve). AUC is commonly used to determine the performance of classification algorithms.

| Model name | AUC |
|---|---|
| Dataset_RndForest | 0.85279 |
| Dataset_Regul | 0.8519 |
| Dataset_Reg | 0.85161 |
| Dataset_DL | 0.85149 |
| Dataset_Bagging | 0.84812 |
| Dataset_Boosting | 0.83855 |
| Dataset_Tree_Inst | 0.83512 |

From the table above, we can see that the Random Forest model outperformed the other models but the difference in model performance is pretty low.

## CONCLUSION AND IMPROVISATION

We have implemented all different models in Altair Knowledge Studio to predict the customers who are likely to discontinue the Telecom company's services. We have achieved a model with 85% accuracy. Knowledge Studio's AutoML features came in very handy and helped us implement all the models without any hassle.

To provide better results, the company can do one or more of the following

- Provide more data points to research
- Survey customers on why they are satisfied or dissatisfied

From the study, we can conclude that

- Providing additional services help in retaining customers
- A longer term contract up front might with incentives can help retain customers