

Data Mining mini project

Katayoon Sadeghi mehr Feb 2020

Logistic Regression/Decision Tree: Churn Costumers

The dataset of churn contains 14 columns. **customerID**, **tenure**, **PhoneService**, **Contract**, **Paperless**, **PaymentMethod**, **MonthlyCharge**, **MontlyGroceryspent**, **Revenue**, **gender**, **Age**, **Partner**, **N_dependent** and **Churn** as dependent variable.

Monthly Charges has 12, churn 52 and revenue 44 missing values.

median is used to replace revenue and monthly charges missing values, and "Undecided" for the 52 missing values for $y=\text{churn}$.

After cleaning we have **7042** rows of data.

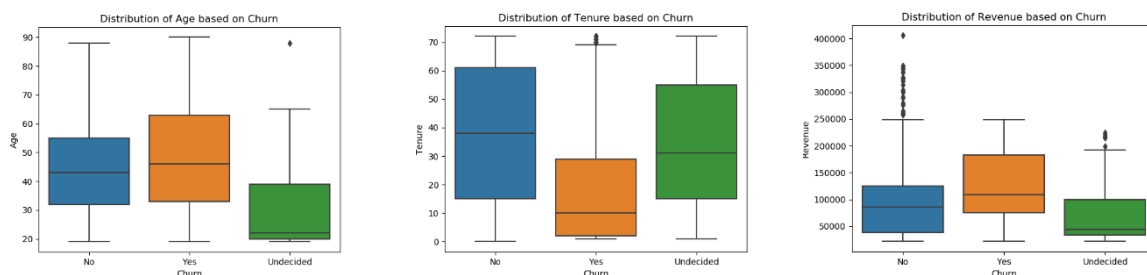
Encoding on categorical variables, Phone Service with 2 levels yes=1, no=0, Contract with 3 levels month-to-month=1, one-year=2 and two-year=3, Paperless with 2 levels yes=1, no=0, PaymentMethod with 4 levels Bank transfer (automatic)=1, Credit card (automatic)=2, Electronic check=3 and Mailed check=4, gender with 2 levels male=1 and female=2, Partner with 2 levels yes=1, no=0, churn with 3 levels undecided=1, no=2 and yes=3 is done.

customerID	tenure	Phone Service	Contract	Paperless	Payment Method	Monthly Charges	Monthly Grocery spent	Churn	Revenue	gender	Age	Partner	Dependents	N_ dependent
6823-SIDFQ	28	Yes	One year	No	Credit card (automatic)	18.25	259	No	22024	Male	53	No	No	0

Distribution of variables:

Average of age, tenure and revenue for different levels of churn. Apparently, the undecided ones are from younger ages. Those who decided to churn have lower tenure and those who do not want to churn are mostly from lower income with some outliers of high incomes.

Churn	Age	tenure	Revenue
Yes	47.65	17.95	117474.57
No	44.48	37.61	95520.22
Undecided	29.19	34.49	73489.66



Correlation and association between variables:

There is pretty high correlation between $y=\text{churn}$ and some dependent variables. Also some of the dependent variables are highly correlated with each other.

