

# **PREDICTING CHURN TO MINIMIZE** **BUSINESS LOSS**

---

## **CONTENTS**

Predicting Churn to Minimize Business Loss .....	0
1 - Description of the Problem: .....	1
What is churn and why is it of value to businesses: .....	1
Why Is It Necessary? .....	1
Where Is It Used? .....	2
2 - Description of the Data: .....	2
Data before reading: .....	3
Data After reading into Dataframe: .....	3
Size of Data: .....	4
Feature Engineering: .....	5
Adding location from FourSquare API for Chicago area: .....	5

## **1 - DESCRIPTION OF THE PROBLEM:**

### **WHAT IS CHURN AND WHY IS IT OF VALUE TO BUSINESSES:**

Customer churn, also known as customer attrition, occurs when customers stop doing business with a company. The companies are interested in identifying segments of these customers because the price for acquiring a new customer is usually higher than retaining the old one. For example, if Netflix knew a segment of customers who were at risk of churning they could proactively engage them with special offers instead of simply losing them.

Churn can also be defined as “The use of customer data and/or feedback to forecast the likelihood of a customer or group of customers discontinuing their subscription in the future.”

### **WHY IS IT NECESSARY?**

Having the ability to accurately predict future churn rates is necessary because it helps your business gain a better understanding of future expected revenue.

In addition, when you're able to use churn prediction to forecast the potential churn rate of a particular customer, it allows you to target that individual in an attempt to prevent them from discontinuing their subscription with you.

And, since the cost of acquiring a new customer is 5x higher than keeping an existing one, there's plenty of revenue-based reason to do everything in your power to keep those existing customers.

Predicting churn rates can also help your business identify and improve upon areas where customer service is lacking. And, by

making those improvements, you can decrease churn and improve revenue numbers.

In the end, the bottom line is that churn prediction is essential because it helps you understand what preventative steps are necessary to ensure lost revenue is minimized.

## **WHERE IS IT USED?**

Churn prediction is used in a variety of different industries and types of businesses.

It is, however, most relevant to SaaS companies and membership based businesses that charge an ongoing monthly, quarterly, or annual fee for their software or services.

As far as how churn prediction can be used within your business, it's one of the key components of determining the lifetime value of customers. And, armed with accurate, real-time data about the lifetime value of your customers, your company will be in a much better position to ensure that you're making decisions that keep you moving forward.


So, this project is aimed to leverage customer's data of sales of Chicago metropolitan area and predict whether he/she will churn in a given time or will remain loyal to the business!

## **2 - DESCRIPTION OF THE DATA:**

The data I will be using in this project is of a Mall located in Chicago's multiple cities consisting of customer's sales data. First let's have a look at how our data looks like:

Our data is in text form but it is not text data rather, it is numerical data stored in text file, so we first load data into jupyter environment and save a copy of it to a dataframe in order to manipulate it.

## **DATA BEFORE READING:**

 Customer\_Sales\_Transactional\_Data\_CSV.txt - Notepad

```
File Edit Format View Help
SALES_date,CUSTOMER_ID,SALES_AMOUNT
10/18/2014,34810920,205.44000
9/22/2014,1026037818,51.36000
9/30/2014,1095693062,222.56000
10/25/2014,15142688,171.20000
10/18/2014,1022245368,171.20000
9/29/2014,3020948,85.60000
9/26/2014,53524262,171.20000
10/4/2014,9820168,428.00000
10/7/2014,949544616,51.36000
10/13/2014,1089626432,256.80000
9/26/2014,103416852,77.04000
10/24/2014,1097139564,171.20000
9/25/2014,32063112,342.40000
10/1/2014,1001021338,102.72000
10/24/2014,830317874,85.60000
```

## **DATA AFTER READING INTO DATAFRAME:**

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import datetime
```

```
data=pd.read_csv('Customer_Sales_Transactional_Data_CSV.txt')
```

```
df=data.copy()
```

```
df.head()
```

	SALES_dATE	CUSTOMER_ID	SALES_AMOUNT
0	10/18/2014	34810920	205.44
1	9/22/2014	1026037818	51.36
2	9/30/2014	1095693062	222.56
3	10/25/2014	15142688	171.20
4	10/18/2014	1022245368	171.20

So, we can see that our data initially contains 03 columns: 'SALES\_dATE', 'CUSTOMER\_ID' and 'SALES\_AMOUNT'. So given a customer\_id, we know his/her sales date and sales amount and there are more than one rows against one customer id.

## **SIZE OF DATA:**

Our data contains more than one million rows.

```
: df.shape
```

```
: (1058198, 3)
```

This is a six week sales data of different customers. The problem of churn prediction will be solved by this data in such a way that first five weeks data will be used to train the model and sixth week data will be used to test or predict churn against a given customer.

## **FEATURE ENGINEERING:**

Since our data consists of only 03 columns. We need to do feature engineering to get better insight from our data.

### **ADDING LOCATION FROM FOURSQUARE API FOR CHICAGO AREA:**

I used FourSquare API to add location to given data and then predict churn with respect to area.

This will help in focusing on areas with greater churn and hence that particular mall.

```
df["area"] = np.random.choice(chicago_areas, size=len(df))
```

```
df.head()
```

	SALES_date	CUSTOMER_ID	SALES_AMOUNT	area
0	10/18/2014	34810920	205.44	Genoa City, WI-IL
1	9/22/2014	1026037818	51.36	Michigan City-LaPorte, IN-MI
2	9/30/2014	1095693062	222.56	Michigan City-LaPorte, IN-MI
3	10/25/2014	15142688	171.20	Chicago-Aurora-Elgin-Joliet-Waukegan
4	10/18/2014	1022245368	171.20	Round Lake Beach-McHenry-Grayslake

Now, we need to split date column into month, day and date plus week number of the year [0-52] respectively like:

```
weeks=[]
for i in range(1058198):
    weeks.append(datetime.date(df.loc[i, 'Year'],df.loc[i, 'Month'],df.loc[i, 'Day']).isocalendar()[1])
```

```
df['weeks']=weeks
```

```
df.head()
```

	SalesDate	CoustomerId	SalesAmount	Year	Day	Month	weeks
0	10/18/2014	34810920	205.44	2014	18	10	42
1	9/22/2014	1026037818	51.36	2014	22	9	39
2	9/30/2014	1095693062	222.56	2014	30	9	40
3	10/25/2014	15142688	171.20	2014	25	10	43
4	10/18/2014	1022245368	171.20	2014	18	10	42

Next thing to do is to find maximum, minimum and average amount of sale against each customer.

```
Week_0_5.columns=['CustomerId', 'Mean_1', 'Max_1', 'Min_1', 'Count_1', 'StandardDeviation_1', 'Revenue_1', 'Mean_2', 'Max_2', 'Min_2',
```

```
Week_0_5.fillna(0,inplace=True)
```

```
Week_0_5['TotalRevenue']=Week_0_5['Revenue_1'] + Week_0_5['Revenue_2'] + Week_0_5['Revenue_3'] + Week_0_5['Revenue_4'] + Wee
```

```
Week_0_5.head()
```

	CustomerId	Mean_1	Max_1	Min_1	Count_1	StandardDeviation_1	Revenue_1	Mean_2	Max_2	Min_2	...	Count_4	StandardDeviation_4	Re
0	7260	171.20	171.20	171.20	1.0	0.0	171.20	128.4	85.60	171.20	...	3.0	19.768473	
1	11024	68.48	68.48	68.48	1.0	0.0	68.48	85.6	51.36	119.84	...	1.0	0.000000	
2	12476	51.36	51.36	51.36	1.0	0.0	51.36	0.0	0.00	0.00	...	1.0	0.000000	
3	14156	171.20	171.20	171.20	1.0	0.0	171.20	171.2	171.20	171.20	...	1.0	0.000000	
4	17606	171.20	171.20	171.20	1.0	0.0	171.20	0.0	0.00	0.00	...	0.0	0.000000	

5 rows × 32 columns

And similarly other columns are also introduced which is all present in jupyter notebook