



# CHURN CASE

---

AB testing

Customer Segmentation

Decision Tree

Logistic Regression

Model Evaluation

Customer Lifetime Value

Rui Li

2020/5/15

<https://github.com/rui629/MarketingAnalyticsAssignment>

0

# Preprocessing

<https://github.com/ruili629/MarketingAnalytics>  
Assignment

# Drop some cols or rows

- Reason:
- 1. One single value (Language, Country)
- 2. Abnormal value (Leave only those age<100 or eg. change 1997 into  $2020 - 1997 = 23$ )
- 3. Uncontrollable or useless factor in segmentation (eg. subid)

# Age

# change those > 1900 to 2020 - year

```
: subscribers[subscribers.age>100]
```

intended_use	weekly_consumption_hour	num_ideal_streaming_services	retarget_TF	age	male_TF	...	creation_until_cancel_days	cancel_before_trial_end	trial_
access to exclusive content	27.301448	NaN	False	1955.0	False	...	7.0	False	20
access to exclusive content	25.851492	NaN	False	1950.0	False	...	30.0	True	20
expand regional access	27.301448	1.958566	False	1957.0	False	...	NaN	True	20
supplement OTT	31.651317	1.874212	False	1969.0	False	...	4.0	False	20

# Fill in the NA

```
values = {'num_weekly_services_utilized': 0, 'weekly_consumption_hour': 0, 'num_ideal_streaming_services': 0,  
          'op_sys': 'unknownOp', 'creation_until_cancel_days' : 0, 'revenue_net' : 0, 'join_fee' : 0,  
          'payment_type' : 'not_paid'}  
df = df.fillna(value=values)
```

# Change Data Type

```
# change data type  
df['preferred_genre'] = df['preferred_genre'].astype(str)  
df['package_type'] = df['package_type'].astype(str)  
df['intended_use'] = df['intended_use'].astype(str)  
df['payment_type'] = df['payment_type'].astype(str)  
df['op_sys'] = df['op_sys'].astype(str)
```

# 1

## AB Testing

单击此处添加副标题内容

# AB Testing - paid

	paid	all	conversion_rate
base_uae_no_trial_7_day_guarantee	1.0	1	1.000000
high_jpy_14_day_trial	1.0	1	1.000000
low_eur_no_trial	1.0	1	1.000000
low_sar_no_trial	1.0	1	1.000000
low_uae_no_trial	134.0	167	0.802395
base_eur_14_day_trial	11.0	18	0.611111
high_aud_14_day_trial	1.0	2	0.500000
base_uae_14_day_trial	91776.0	227096	0.404129
high_uae_14_day_trial	119.0	325	0.366154
high_sar_14_day_trial	4.0	12	0.333333
low_gbp_14_day_trial	1.0	4	0.250000

**\*\*Hypothesis setup:\*\***

Variant B: high\_uae\_14\_day\_trial

Variant A: base\_uae\_14\_day\_trial

H0: Variant B and Variant A had the same conversion rates.

HA: Variant B had a lower conversion rate than Variant A

**\*\*Assumptions:\*\***

Variant A represents the population and we can treat the population mean as known and equal to the mean of Variant A.

Result:

$$|Z| = 2.2 > 1.64$$



# AB Testing - cancel\_before\_trial\_end

	cancel	all	churn_rate
base_uae_no_trial_7_day_guarantee	1.0	1	1.000000
high_jpy_14_day_trial	1.0	1	1.000000
low_eur_no_trial	1.0	1	1.000000
low_sar_no_trial	1.0	1	1.000000
low_uae_no_trial	141.0	167	0.844311
base_eur_14_day_trial	11.0	18	0.611111
high_aud_14_day_trial	1.0	2	0.500000
high_sar_14_day_trial	6.0	12	0.500000
base_uae_14_day_trial	103253.0	227096	0.454667
high_uae_14_day_trial	140.0	325	0.430769
low_gbp_14_day_trial	1.0	4	0.250000

**\*\*Hypothesis setup:\*\***

Variant B: high\_uae\_14\_day\_trial

Variant A: base\_uae\_14\_day\_trial

H0: Variant B and Variant A had the same churn(cancel) rates.

HA: Variant B had a lower churn rate than Variant A

**\*\*Assumptions:\*\***

Variant A represents the population and we can treat the population mean as known and equal to the mean of Variant A.

Result:

$$|Z| = 12.22 > 1.64$$



# AB Testing - renew

	renew	all	conversion_rate
base_eur_14_day_trial	13.0	23	0.565217
high_jpy_14_day_trial	1.0	2	0.500000
high_aud_14_day_trial	2.0	5	0.400000
base_uae_14_day_trial	72971.0	209383	0.348505
high_uae_14_day_trial	129.0	375	0.344000
high_sar_14_day_trial	5.0	15	0.333333
low_uae_no_trial	25.0	88	0.284091
base_uae_no_trial_7_day_guarantee	0.0	2	0.000000
low_gbp_14_day_trial	0.0	1	0.000000

**\*\*Hypothesis setup:\*\***

Variant B: high\_uae\_14\_day\_trial

Variant A: base\_uae\_14\_day\_trial

H0: Variant B and Variant A had the same renew rates.

HA: Variant B had a lower renew rate than Variant A

**\*\*Assumptions:\*\***

Variant A represents the population and we can treat the population mean as known and equal to the mean of Variant A.

Result:

$$|Z| = 9.71 > 1.64$$

# Optimal Sample Size and Sequential Testing - (paid)

The optimal sample size for each segment is 2616

The challenger wins 0.0% of the time.

Assume  $P(X_i=1)$  under  $H_0$  = p-varA and  $P(X_i=1)$  under  $H_1$  = p-treatment.

Set desired type 1 error = 5% and type 2 error = 20%.

Success rate is 80.0%

# 2

---

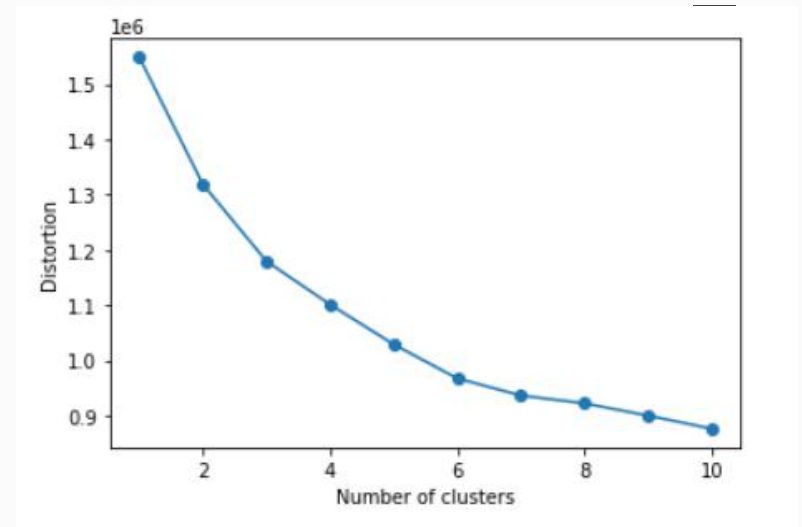
## Customer Segmentation

单击此处添加副标题内容

# Customer segmentation

## Cluster 0 (High value) customer description:

1. They completed less videos but they rated much more often.
2. Customer internet package: base mainly.
3. Preferred genre: comedy and drama mainly.
4. Main Attribution\_technical\_brand: "sem intent google", "email blast".
5. Plan\_type: high\_uae\_14\_day\_trial
6. Main Payment\_type:Paypal
7. Main Payment\_type\_Standard:Charter
8. High Value because they pay.



Cluster	cust_service_msgs	num_videos_completed	num_videos Rated	retarget_IF	ancel_before_trial_eninitial_credit_card_declined	revenue_net	paid_IF	refund_after_trial_IF	
0	0.3726	2.6224	0.0684	0.0438	1.0000	0.0025	4.0263	1.0000	0.2587
1	0.5936	3.0025	0.0005	0.0359	0.0840	0.0827	-0.0026	0.0000	0.0000
2	0.5655	2.9997	0.0008	0.0231	0.0954	0.0952	0.0000	0.0000	0.0000
3	0.5812	3.0539	0.0005	0.0317	0.0726	0.0711	-0.0014	0.0000	0.0000
Cluster	genre_comedy	genre_drama	genre_international	genre_regional	ccess to exclusive co	use_education	expand international acse_expand regional acces	package_base	
0	0.3424	0.1510	0.0298	0.0549	0.2403	0.0571	0.1391	0.1827	0.3205
1	0.6850	0.2034	0.0291	0.0590	0.4172	0.0000	0.0007	0.0007	0.0000
2	0.0000	0.0007	0.0002	0.0002	0.0000	0.1451	0.3806	0.4208	0.0005
3	0.6585	0.2407	0.0290	0.0479	0.4911	0.0004	0.0009	0.0008	1.0000
Cluster	affiliate	bing	bing_organic	brand sem intent binbrand sem intent googl	direct_mail	discovery	display	search engine	
0	0.0387	0.0029	0.0014	0.0174	0.0923	0.0012	0.0171	0.0042	0.0707
1	0.0381	0.0033	0.0012	0.0080	0.0761	0.0005	0.0098	0.0033	0.0247
2	0.0330	0.0036	0.0024	0.0071	0.0733	0.0000	0.0000	0.0011	0.0542
3	0.0362	0.0037	0.0018	0.0100	0.0772	0.0007	0.0111	0.0037	0.0259
Cluster	facebook	facebook_organic	google_organic	internal	organic	podcast	referral	vod	an high_uae_14_ trial
0	0.2464	0.0181	0.0640	0.0153	0.1022	0.0027	0.0452	0.0019	0.0039
1	0.3414	0.0166	0.0505	0.0049	0.0933	0.0053	0.0243	0.0013	0.0024
2	0.3719	0.0149	0.0517	0.0104	0.0899	0.0014	0.0252	0.0000	0.0000
3	0.3592	0.0145	0.0520	0.0047	0.0928	0.0055	0.0245	0.0013	0.0024

3

---

# Decision Tree

Preprocess our input DF  
Feature Selection  
Process  
Visualization

# Which features to select?

- After we finish the cleaning data and use get dummies to get our input dataframe, we get a problem:

Too many variables!

- **Recursive Feature Ranking with recursive feature elimination and cross-validation**
- We choose to use the RFECV, which performs RFE in a cross-validation loop to find the optimal number or the best number of features. Hereafter a recursive feature elimination applied on decision tree classifier or logistic regression with automatic tuning of the number of features selected with cross-validation.

- Classifier: decision tree; Dependent variable: renew;

```
: len(df_input.columns)
: 101

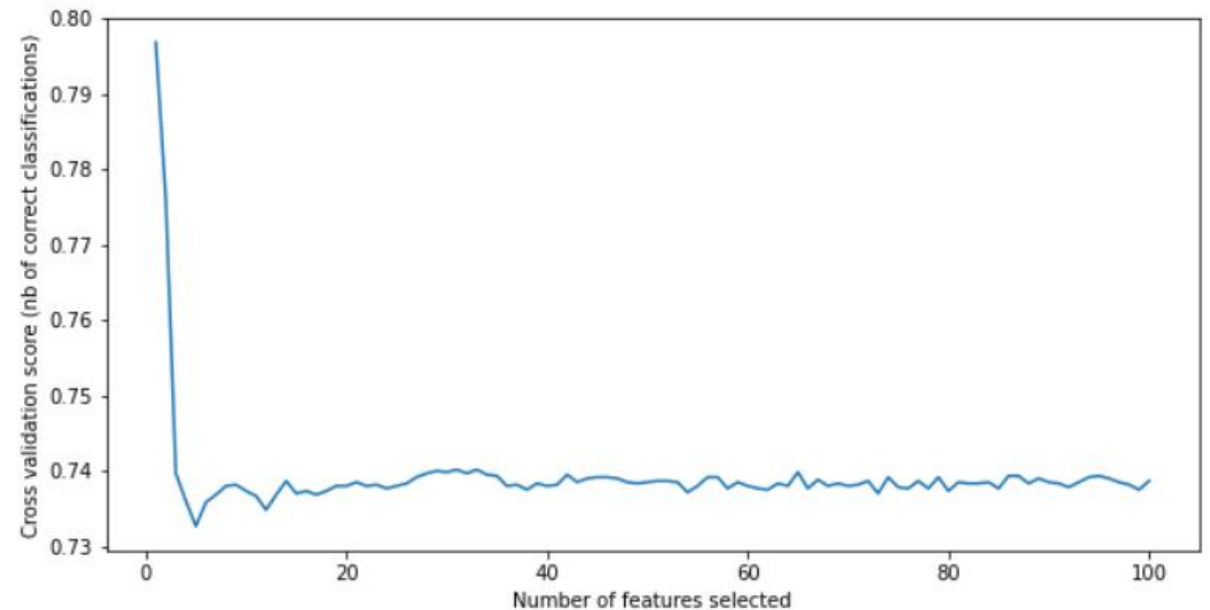
: df_input.columns
: Index(['num_weekly_services_utilized', 'weekly_consumption_hour',
        'num_ideal_streaming_services', 'retarget_TF', 'age',
        'months_per_bill_period', 'monthly_price', 'discount_price',
        'creation_until_cancel_days', 'cancel_before_trial_end',
        ...
        'plan_type_high_sar_14_day_trial', 'plan_type_high_uae_14_day_trial',
        'plan_type_low_gbp_14_day_trial', 'payment_type_Apple Pay',
        'payment_type_CBD', 'payment_type_Najim', 'payment_type_Paypal',
        'payment_type_RAKBANK', 'payment_type_Standard Charter',
        'payment_type_not_paid'],
        dtype='object', length=101)
```



# Apply it right away?

- If we apply the RFECV without further preprocessing, we would get:  
**Optimal number of features: 1**  
**Selected features: ['paid\_TF']**
- This is not satisfying since paid\_TF is only a result and it is highly correlated with “renew”

Optimal number of features: 1  
Selected features: ['paid\_TF']



# Drop those “result” variable

- Dropping the variables that are highly correlated with “renew”:

```
['paid_TF','payment_type_not_paid','revenue_net','cancel_before_trial_end','renew']
```

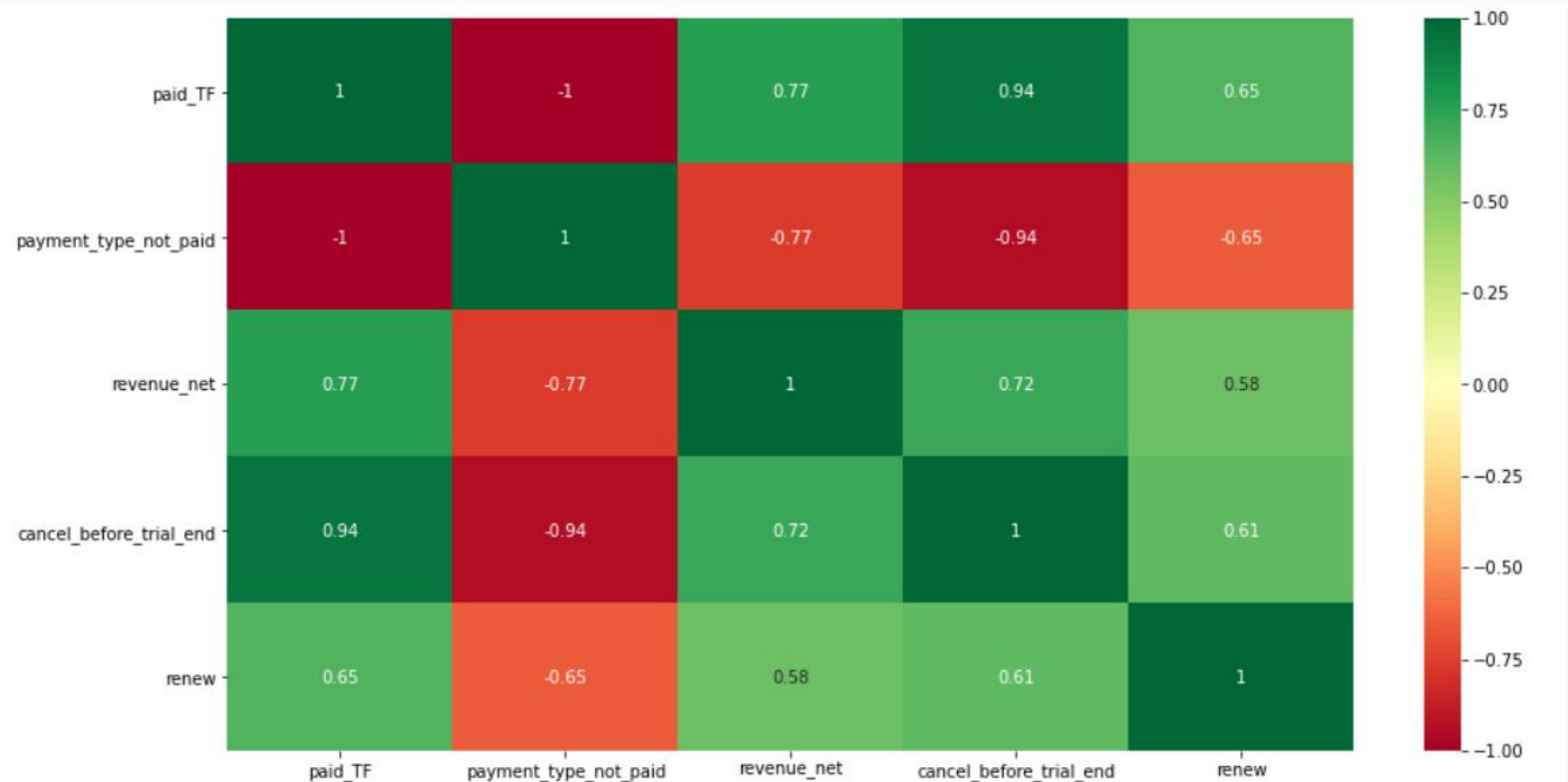
```
a = df_input.corr()['renew']
```

```
a[a>0.5]
```

```
cancel_before_trial_end    0.612260  
revenue_net                0.577446  
paid_TF                   0.652731  
renew                     1.000000  
Name: renew, dtype: float64
```

```
a[a<-0.5]
```

```
payment_type_not_paid    -0.652731  
Name: renew, dtype: float64
```



```

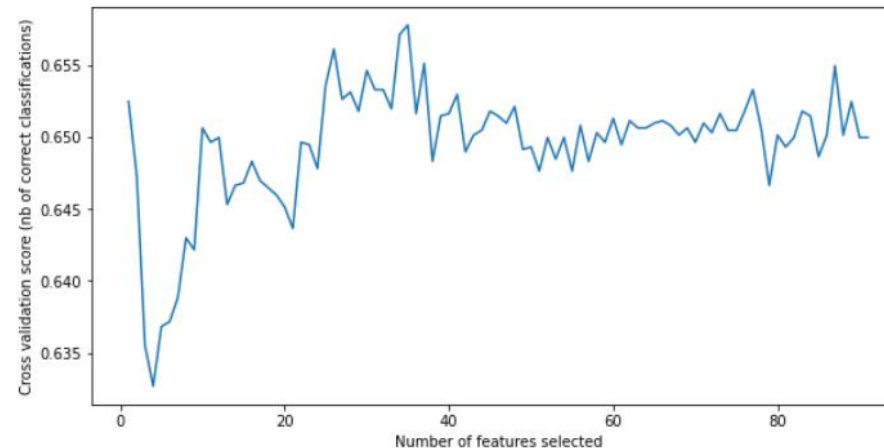
: feature_selection_class(tree.DecisionTreeClassifier(),
                        X.drop(columns = ['paid_TF', 'payment_type_not_paid', 'revenue_net'
                        , 'cancel_before_trial_end', 'creation_until_cancel_days',
                        'payment_type_Paypal', 'payment_type_RAKBANK',
                        'payment_type_Standard Charter', 'num_weekly_services_utilized']),
                        y)

```

**After dropping the highly correlated 'result-type' variables like cancel, revenue\_net, we ran the RFECV again to get the best number and types of features.**

Optimal number of features: 35

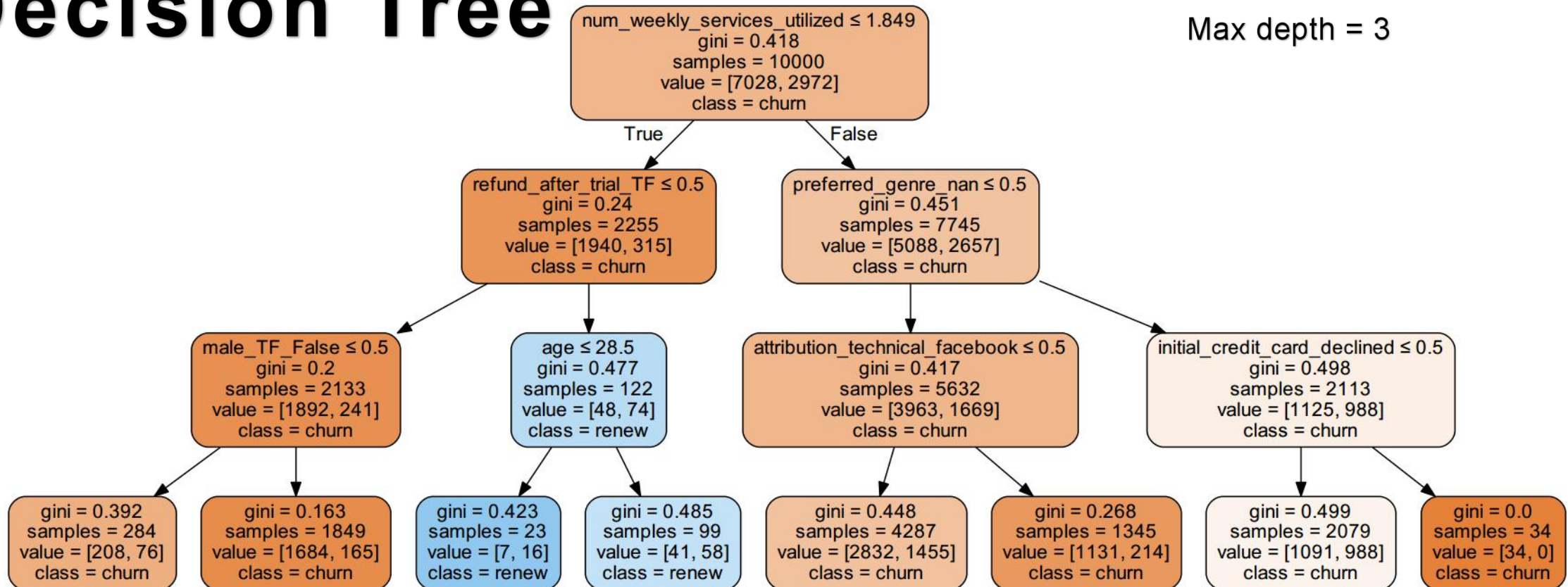
Selected features: ['weekly\_consumption\_hour', 'num\_ideal\_streaming\_services', 'retarget\_TF', 'age', 'initial\_credit\_card\_declined', 'join\_fee', 'refund\_after\_trial\_TF', 'package\_type\_base', 'package\_type\_economy', 'package\_type\_enhanced', 'preferred\_genre\_comedy', 'preferred\_genre\_nan', 'intended\_use\_access to exclusive content', 'intended\_use\_expand international access', 'intended\_use\_other', 'intended\_use\_replace OTT', 'intended\_use\_supplement OTT', 'male\_TF\_True', 'attribution\_technical\_brand sem intent google', 'attribution\_technical\_email', 'attribution\_technical\_email blast', 'attribution\_technical\_facebook', 'attribution\_technical\_google\_organic', 'attribution\_technical\_organic', 'attribution\_technical\_pinterest', 'attribution\_technical\_referral', 'attribution\_technical\_search', 'attribution\_survey\_facebook', 'attribution\_survey\_referral', 'attribution\_survey\_search', 'attribution\_survey\_tv', 'op\_sys\_Android', 'op\_sys\_ios', 'payment\_type\_CBD', 'payment\_type\_Najim']



To better visualize the tree, I used the max depth = 3, yet the best depth tested for now is 10

# Decision Tree

Max depth = 3



# 4

---

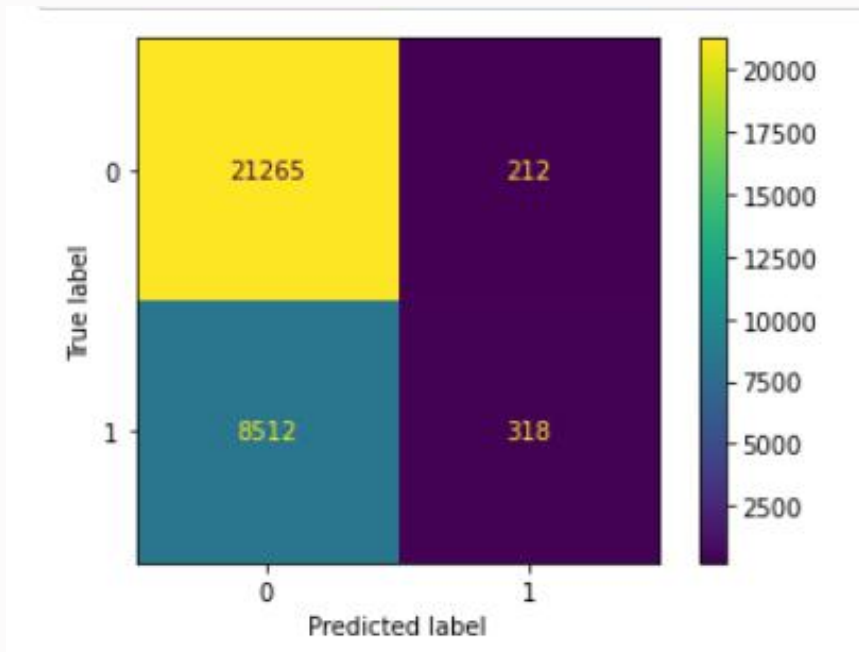
## Model Evaluation

1. Under the same process of building decision tree, the logistic regression was also built
2. Comparing the result of Logistic Regression

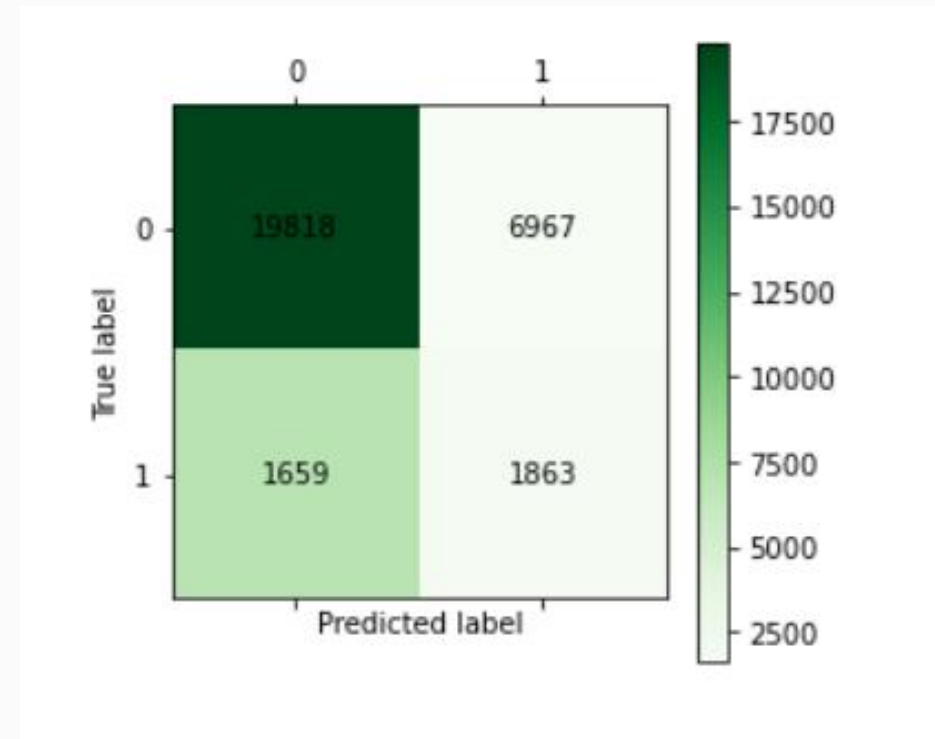
# Confusion Matrix

- Comparing Confusion Matrix of two models, Logistic Regression model has a higher true positive rate but also higher false negative rate.

Decision Tree



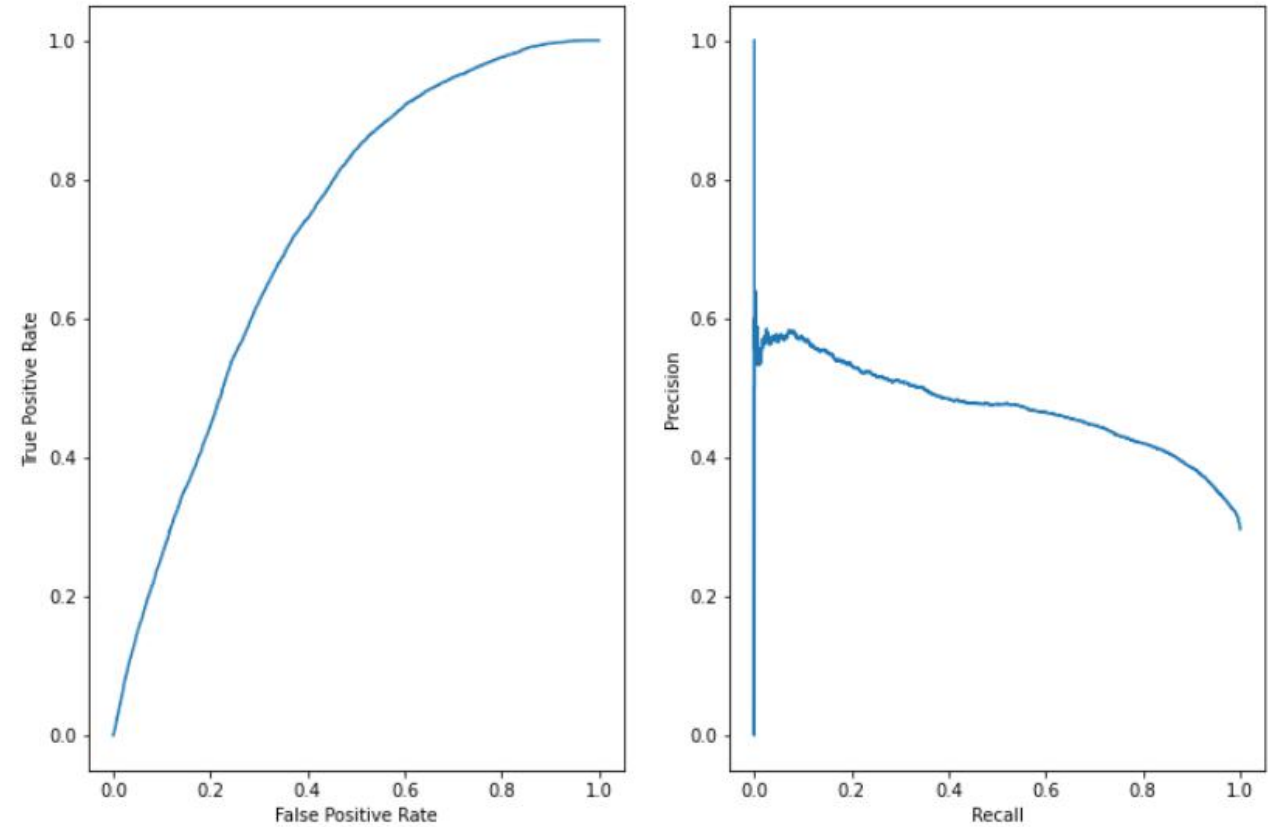
Logistic Regression





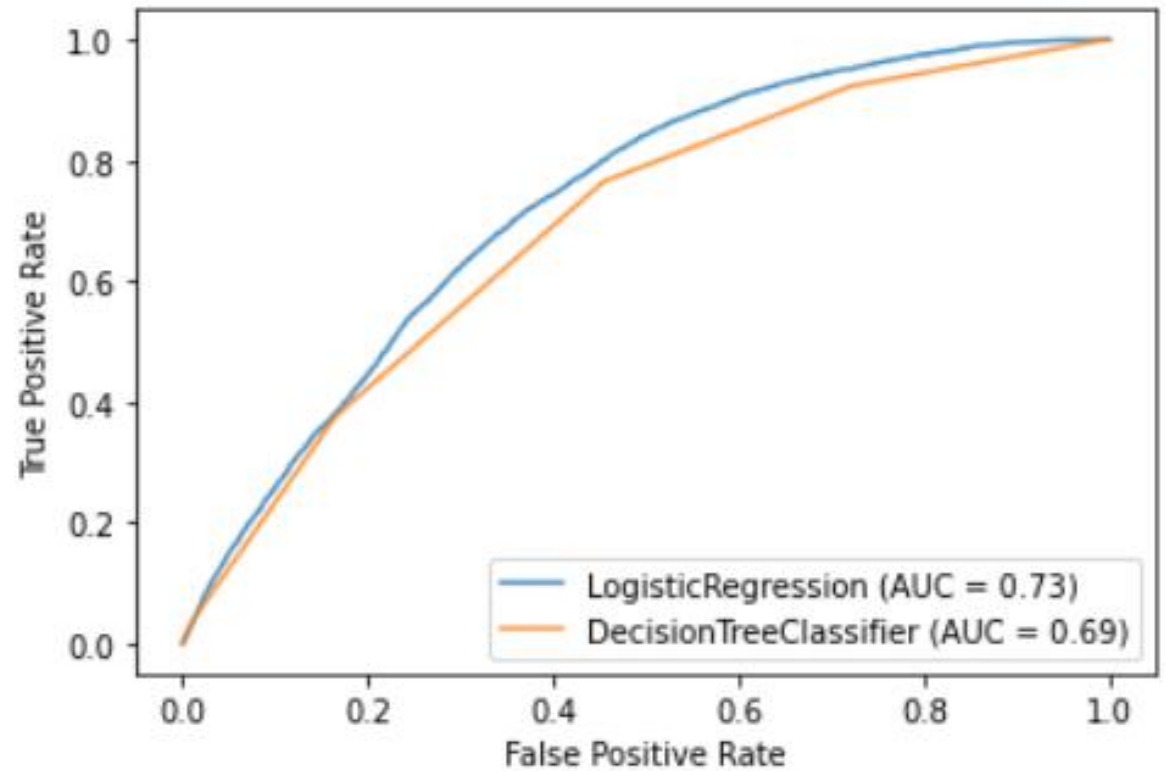
# Confusion Matrix

- Since it seems that the logistic regression has a bit better performance than decision tree, and the result of its recall and precision is just satisfying.



# ROC Curve

- From the ROC curve, it appears that the logistic regression has a bit better performance



5

---

# Customer Lifetime Value

Using logistic regression to predict

# CLV caculation

- As it appears that the logistic regression has a bit better performance, we shall use the logistic model to predict how likely this customer is sill going to renew in the nex period, so as to simulate and get the lifetime value.

```
prob_list = [renew for renew, churn in log.predict_proba(X.drop(columns = drop_cols))]
```

```
prob_list[1:10]
```

```
[0.7261671245082209,  
 0.7261671245082209,  
 0.8123141328142311,  
 0.8123141328142311,  
 0.9270975978280358,  
 0.8862471829100999,  
 0.6174895556199471,  
 0.7911038113703474,  
 0.9223666593475623]
```

# CLV Caculation Method

```
import random

# using r_year = 0.1 as default
r_year = 0.1
# effective discount rate of each month
r = pow(1+r_year,1/12)-1

revenue_list = list(df['revenue_net'])
# 0 if churned
all_churn_list = list(df['revenue_net'])

t=1
while sum(all_churn_list) != 0:
    for i in range(len(prob_list)):
        if all_churn_list[i] != 0:
            if random.random() < prob_list[i]:
                revenue_list[i] = revenue_list[i]+revenue_list[i]/pow(1+r,t)
            else:
                all_churn_list[i] = 0

    t += 1
```

- **Discounted**

$$CLV = \frac{\sum_{i=1}^n \sum_{t=1}^T (gross\ margin_{it})^{-(1+r)^t}}{n}$$

$r$  = discount rate,  $t$  indexes time periods