

---

# Surrounding Entities Impacting Sales

Identify the most important surrounding entities impacting POS



# TOC

Problem Statement

Feature Engineering & Target Creation

ML Methods & Models

Performance Results

Results & Conclusion

# Problem Statement

Every company wants to succeed by achieving maximum sales. Many companies distribute their goods at physical Point Of Sales (POSs). For all of them the challenge is to devise a strategy that will drive the sales at POSs. Possible solution could be to place the product in the most convenient location for consumers.

Given:

- 1) Sales information about the sales volumes of a product at particular POS
- 2) Surroundings information about 90 different amenities (restaurants, shops, beauty salons etc.) about surroundings of each POS.

The goal is to create a model that identifies important attributes in the surroundings that impact sales.



Fig : POS Sales



# Feature Engineering & Target Creation

# Target Variable - Histogram Analysis (HA)

Classification Problem : ( High Sales Shop = 1, Not High Sales = 0 )

- 1) Separated data into train and test set 70:30 ratio
- 2) Created a histogram of log of sales values of the training set.
- 3) Based on visual inspection of the histogram choose a cut-point ( $x=8.35$ ) , such that you have almost equal data-points on both sides
- 4) Binarized the target variable (train & test) based on the cut point.
- 5) Ensured that the cut-point is chosen, such that ratio of classes in the train and test data are similar .

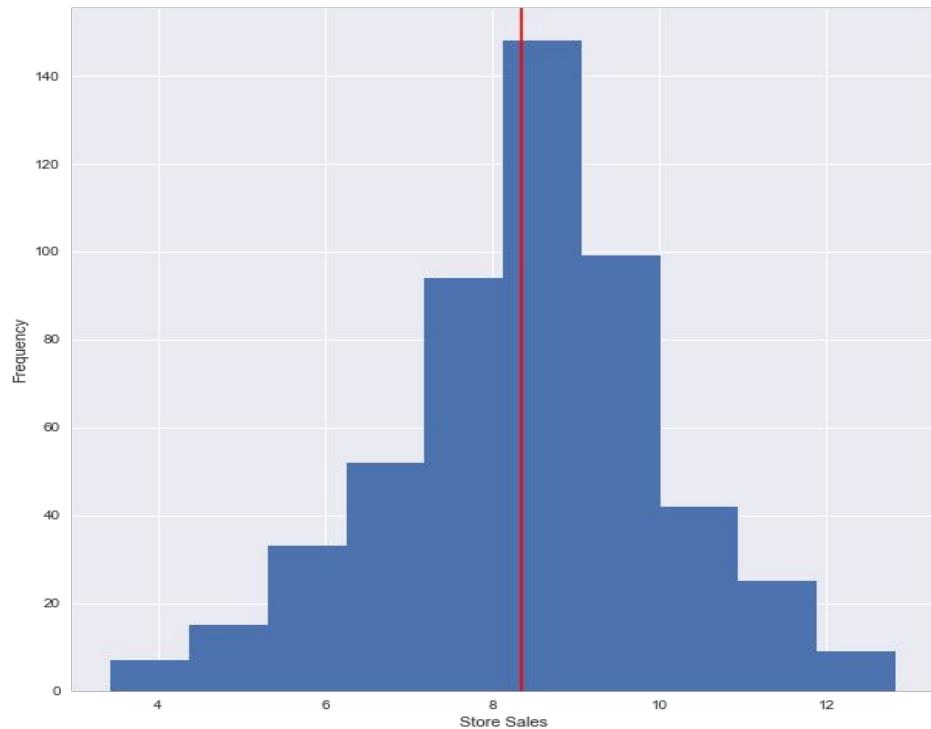


Fig : Log Total Sales



# Creation of Features

Assumptions :

Stores are selling cigarettes.

No.of ratings is a proxy for no.of people visiting the place .

Overall rating indicates the attractiveness of the place (People would go / return to places with high rating ).

A few natural features already present , in the json document (a bit of parsing required ). eg no.of entity\_types around the store (bars, train station etc ). Following custom built features were created .

1

Equal entity share assigned to each store, added over all entities of a particular type .

2

Average rating per entity type for each store .

3

No.of stores around a particular store calculated on post-code level .ie stores in the same surrounding .  
(Competition or Similarity effect )

4

Average no.of reviews of an entity (per entity type ) shared by each store as calculated by share factor in point 1.



# Experiment Design

# Experiment Design



## Data Cleaning Target & Feature Creation

Targets were created in two different ways and various features were created too, focussing on interaction between store and entity on a micro level .



## Feature Selection

Of all the features created, a few of them were selected based on Logistic Regression with L1 penalty Select KBest (f\_classify methods )



## Classification

The features selected in the feature selection phase were passed on to DT,RF,ET classifiers. The results of the classification give us an confidence , of the experiment setup .

## Interpretation

All the classifiers used ,give us an indication of feature importance . We try to compare the output of various classifiers and select the most common features..



# Getting Ready- Common Settings

---

## Training & Testing datasets

01

Test data ratio : 0.3

## Model Evaluation

02

K Fold cross validation  
AUC score .

## Grid Search

03

RandomizedSearchCV



# Feature Selection

---

## Logistic Regression with L1 Penalty

### O1

Logistic regression with L1 penalty was used to preselect features .

Benefit of L1 is that it can push feature coefficients to 0, creating a method for feature selection.



# Feature Selection

## Select KBest Features

02

Scikit learn library provides out of the box methods for feature selection . eg Kbestfeatures

Compute the ANOVA F-value between each feature and the target.

Use the selected features as input to other classifiers.

Ref : [http://scikit-learn.org/stable/modules/generated/sklearn.feature\\_selection.SelectKBest.html](http://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBest.html)



# Classification

## Decision Tree

O1

Though Decision tree classifiers suffer from the disadvantages such as growing complex and still unable to generalise the data well.

It is simple to understand and trees can be visualised .

Training / Testing Scores : 0.74 / 0.69

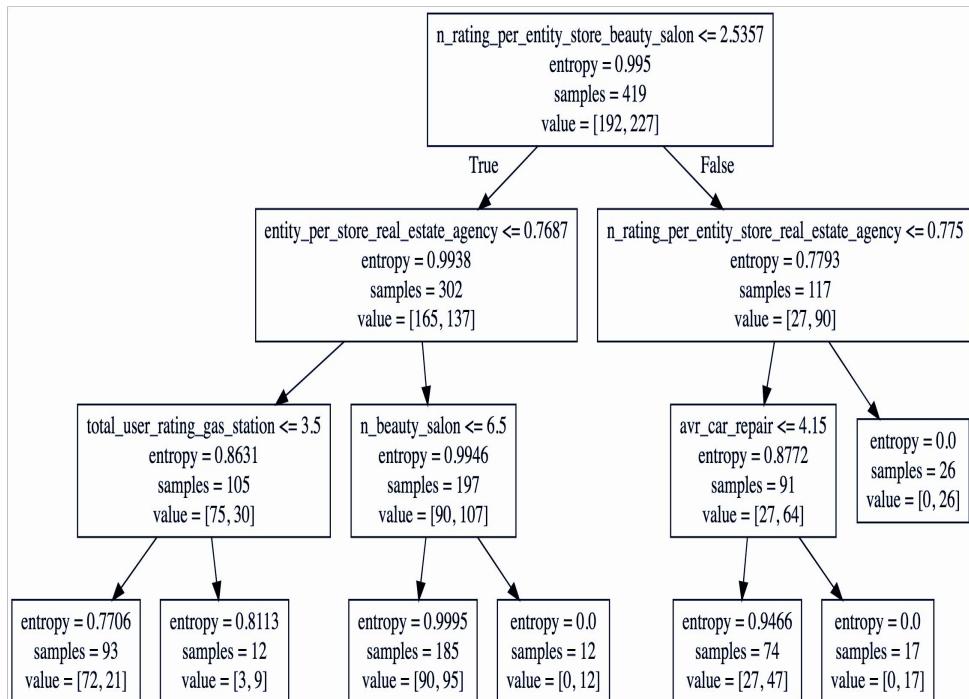


Fig : Decision Tree

# Classification

## Random Forest & Extra Trees Classifier

02

One of the most robust set of classifiers, inherently reduces overfitting .

Provides a list of feature importance, which can help in identifying entity attributes.

Random Forest :

Training / Testing scores : 0.68/0.63

Extra Trees Classifier :

Training / Testing scores : 0.74/0.64

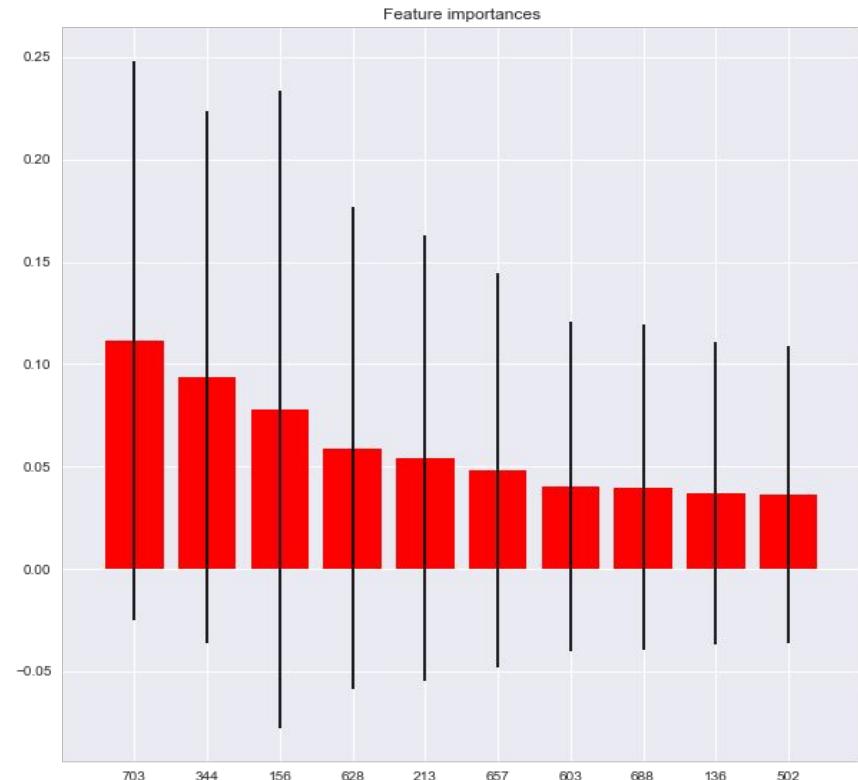


Fig : Extra Trees Classifier Feature Importance



# Results



# Performance Analysis



A run is considered good enough for usage if :

- Train & Test ROC-AUC is  $\geq 0.60$
- Difference between Training AUC and Testing AUC  $\leq 10$  percentage points

Classifier	Decision Tree Clf (1)	Extra Tree Clf (2)	Decision Tree Clf (3)	Random Forest Clf (4)	Extra Tree Clf (5)
Feature Preselection Method	NA	NA	KBest(F_classif)	KBest(F_classif)	Logistic Regression(L1 Penalty)
Training AUC	0.74	0.74	0.68	0.68	0.69
Test AUC	0.69	0.64	0.61	0.63	0.61

# Entities Affecting Sales

- Features related to the same entity type have the same color code.



(1)Decision Tree Clf NA HA (0.74/0.69)	(3)Decision Tree Clf KBest(F_classif) HA (0.68/0.61)	(2)Extra Tree Clf NA HA (0.74/0.64)	(4)Random Forest Clf KBest(F_classif) HA (0.68/0.63)	(5)Extra Tree Clf Logistic Regression(L1 Penalty) HA (0.69/0.61)
n_rating_per_entity_store_beauty_salon	total_user_rating_doctor	rating_per_entity_store_physiotherapist (0.087518)	rating_per_entity_store_hair_care (0.112934)	rating_per_entity_store_hair_care (0.056445)
entity_per_store_real_estate_agency	total_user_rating_beauty_salon	city_GR (0.071863)	total_user_rating_beauty_salon (0.110141)	rating_per_entity_store_insurance_agency (0.054403)
total_user_rating_gas_station		rating_dentist (0.060138)	total_user_rating_hair_care (0.103838)	n_rating_per_entity_store_pharmacy (0.035163)
n_beauty_salon		entity_per_store_jewelry_store (0.054807)	total_user_rating_doctor (0.091116)	n_rating_per_entity_store_dentist (0.034688)
n_rating_per_entity_store_real_estate_agency		rating_home_goods_store (0.054117)	n_rating_per_entity_store_doctor (0.079437)	n_store (0.031897)
		entity_per_store_bar (0.050996)	n_review (0.079117)	rating_per_entity_store_spas (0.031538)
		avr_gym (0.040799)	avr_gym (0.064840)	avr_post_office (0.031239)
		avr_electrician (0.039407)	rating_per_entity_store_dentist (0.062976)	n_rating_per_entity_store_gym (0.028741)
		n_store (0.038358)	n_store (0.059283)	rating_per_entity_store_furniture_store (0.027427)

# Entities Affecting Sales

---

- Ratings of Beauty salon,Haircare and Spa
- Entities related to medical facilities (doctor,physiotherapist and dentist)
- No.of entities of type store
- No.of Gym





---

Thank you.

