



# Coursera Final Project

## **Predicting the best location with more potential revenue**

RAPHAEL FERREIRA



# 1. Introduction

PREDICTING THE BEST LOCATION WITH MORE POTENTIAL REVENUE

# 1.1 Background

- ▶ A start-up based on Toronto wants to expand to New York. Its core business is delivery: customers want to receive something from coffee shops, drugstores, pizza places, groceries, for example.
- ▶ Through the app customer demands some product, which it can indicate a specific place or not, and where it wants that product to be delivered. Most of the products are delivered by walk or by bike.
- ▶ So, spatial conditions are very important for such start-up, specially the category of venues that exist in each place, since it represents places where it could have more delivery transactions, and therefore more revenue.

## 1.2 Problem

- ▶ In order to begin its operations in NY, this start-up, our client, wants to know the best location, which might have more potential to bring deliveries orders to them, and therefore more revenue.
- ▶ They already know its best location in Toronto, but they need to estimate where to start in NY.
- ▶ By the most common type of venue categories in its best location in Toronto, we need to predict the best location in NY, for this first move of our client.

# Summary

- ▶ This work is divided in two main parts:
  - ▶ Unsupervised/spatial: where we get Foursquare Data and create clusters for Toronto and NY neighborhoods. In Toronto we find our target cluster: Cluster 3 - location with venues that are customers of our client
  - ▶ Supervised: where we work with data wrangling (categories spatially sorted and transform data into numerical), visualization (threshold for classification) and prediction (LR and Decision Tree Models)



## 2. Data acquisition and cleaning

PREDICTING THE BEST LOCATION WITH MORE POTENTIAL REVENUE

## 2.1 Data Sources

- ▶ The dataset that we are going to use come from Foursquare, mainly the type of venues and its quantities for each neighborhood in Toronto and NY, so spatial data.
- ▶ The acquisition is made by API that gets the data from Foursquare. Through such data we first create the clusters for Toronto, recognize the most important cluster, then we have data frames that have the 10 most common types of venues for each neighborhood.
- ▶ The same procedure is done for NY. After that we transform such categorical data in numerical in order to predict which cluster in NY might have more revenue potential for our client.

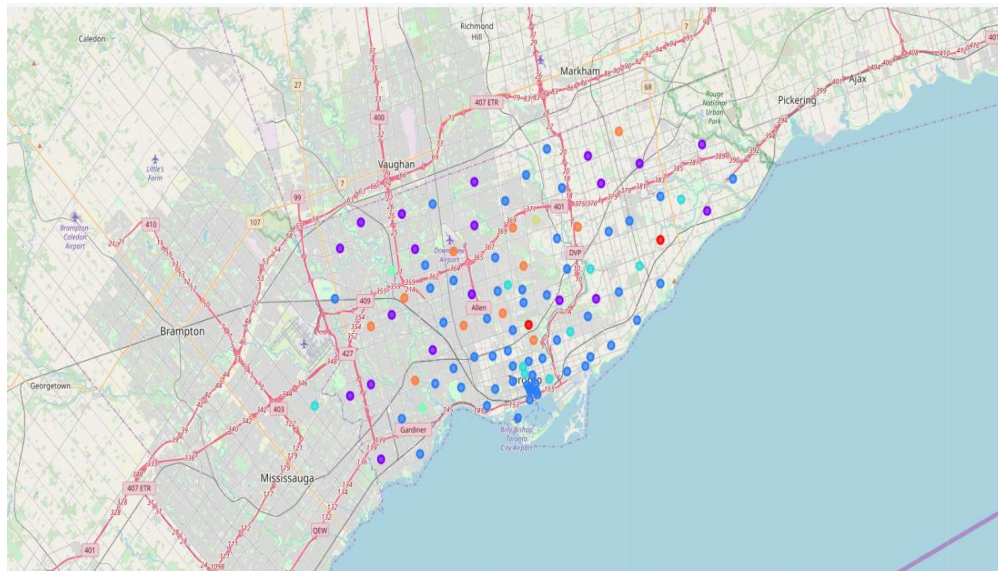
## 2.2 Data Preparation

- ▶ The data acquisition come as JSON from Foursquare. So, we need to get the venues information for each neighborhood in Toronto and NY, and it is done using latitude and longitude for each neighborhood. After that, it is created data frames that have venues names and their categories.
- ▶ Another data frame that contains venue categories for each neighborhood is created, however such data is categorical, and it is transformed in numerical with dummy variables, and the mean of their quantities is calculated.
- ▶ After clustering, each neighborhood is labeled with the cluster that it belongs to. So, it is created another data frame for each cluster, containing the 10 most common types of venues (categories) for each neighborhood in that cluster.
- ▶ The final data frame is the number of categories (rows) by most common (columns), in order words, the quantity of categories is spatial ranked, since this was made during the clustering process.
- ▶ In order to train the models, some data preparation is made, splitting the data into training set and test set.



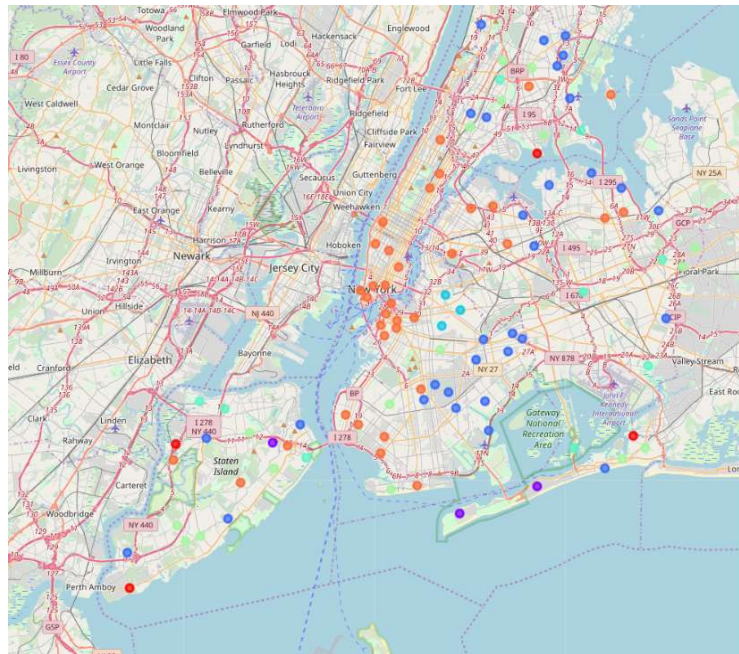
### 3. Unsupervised: clustering

- Toronto's clusters: we used K-means in this task. 7 clusters were created. And the cluster 3 is our client best cluster.



### 3. Unsupervised: clustering

- NY's clusters: we used K-means in this task as well. 8 clusters were created. But just 3 clusters (2,4 and 8) have a good number of neighborhoods (around 340 in average)





# 4. Exploratory Analysis

PREDICTING THE BEST LOCATION WITH MORE POTENTIAL REVENUE

## 4. Exploratory Analysis

- Counting the numbers of categories in the clusters give us the kind of venues that bring more delivery transactions.

Table 1. Neighborhood and Venues in Toronto's cluster 3

Neighbourhood	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue	6th Most Common Venue	7th Most Common Venue	8th Most Common Venue	9th Most Common Venue	10th Most Common Venue
1	Highland Creek/Rouge Hill/Port Union	Construction & Landscaping	Bar	Yoga Studio	Dessert Shop	Event Space	Empire Restaurant	Electronics Store	Eastern European Restaurant	Drugstore
2	Guildwood/Morningside/West Hill	Electronics Store	Mexican Restaurant	Breakfast Spot	Pizza Place	Medical Center	Intersection	Rental Car Location	Spa	Tech Startup
4	Cedarbrae	Caribbean Restaurant	Lounge	Bakery	Hakka Restaurant	Fried Chicken Joint	Thai Restaurant	Athletics & Sports	Bank	Diner
7	Charm, Golden Mile/Oakridge	Bus Line	Bakery	Park	Intersection	Fast Food Restaurant	Metro Station	Bus Station	Soccer Field	Crepes
8	Cliffcrest/Cliffside/Scarborough Village West	Motel	American Restaurant	Yoga Studio	Empire Restaurant	Electronics Store	Eastern European Restaurant	Drugstore	Dog Run	Discount Store
9	Beck-CO/Cliffside West	College Stadium	General Entertainment	Skating Rink	Cafe	Deli / Bodega	Electronics Store	Eastern European Restaurant	Drugstore	Discount Store
10	Dorset Park/Scarborough Town Centre/Westford Hill	Indian Restaurant	Latin American Restaurant	Pet Store	Vietnamese Restaurant	Chinese Restaurant	Deli / Bodega	Electronics Store	Eastern European Restaurant	Drugstore
11	Maryvale/Westford	Middle Eastern Restaurant	Shopping Mall	Sandwich Place	Bakery	Auto Garage	Breakfast Spot	Yoga Studio	Empire Restaurant	Electronics Store
12	Agincourt	Breakfast Spot	Chinese Restaurant	Sandwich Place	Lounge	Department Store	Empire Restaurant	Electronics Store	Eastern European Restaurant	Drugstore
17	Willowdale	Golf Course	Mediterranean Restaurant	Pool	Dog Run	Yoga Studio	Dance Studio	Eastern European Restaurant	Drugstore	Discount Store
18	Fairview/Hwy 7/Fairview	Ten Room	Movie Theater	Smoothie Shop	Shopping Mall	Julius Bar	Burger Joint	Fast Food Restaurant	Bakery	Department Store
19	Bayview Village	Cafe	Chinese Restaurant	Bank	Japanese Restaurant	Department Store	Empire Restaurant	Electronics Store	Eastern European Restaurant	Drugstore
22	Willowdale South	Ramen Restaurant	Cafe	Pet Store	Indonesian Restaurant	Movie Theater	Plaza	Shopping Mall	Fast Food Restaurant	Steakhouse
26	Dan Mills North	Cafe	Gym / Fitness Center	Baseball Field	Japanese Restaurant	Caribbean Restaurant	Department Store	Empire Restaurant	Electronics Store	Eastern European Restaurant
27	Flemington Park/Dan Mills South	Gym	Donut Shop	Italian Restaurant	Japanese Restaurant	Discount Store	Beer Store	Bike Shop	Asian Restaurant	Sporting Goods Shop
32	Downsview Central	Food Truck	Korean Restaurant	Baseball Field	Home Service	Yoga Studio	Empire Restaurant	Electronics Store	Eastern European Restaurant	Drugstore
36	Woodbine Heights	Pharmacy	Athletics & Sports	Cosmetics Shop	Curling Ice	Bus Stop	Skating Rink	Beer Store	Furniture / Home Store	Dance Studio
37	The Beaches	Other Great Outdoors	Health Food Store	Tail	Pub	Yoga Studio	Dance Studio	Eastern European Restaurant	Drugstore	Dog Run
38	Leslieville	Gym	Sushi Restaurant	Grocery Store	Fish & Chips Shop	Liquor Store	Coffee Shop	Clothing Store	Restaurant	Burger Joint
39	Thorncliffe Park	Indian Restaurant	Yoga Studio	Housing Development	Pharmacy	Pizza Place	Discount Store	Sandwich Place	Burger Joint	Supermarket
41	The Danforth West/Riverside	Greek Restaurant	Ice Cream Shop	Italian Restaurant	Brewery	Fruit & Vegetable Store	Dessert Shop	Cosmetics Shop	Pub	Yoga Studio
42	The Beaches West/Innis Bazaar	Gym	Sushi Restaurant	Ice Cream Shop	Fish & Chips Shop	Italian Restaurant	Fast Food Restaurant	Liquor Store	Movie Theater	Park
45	Danforth North	Gym	Breakfast Spot	Hotel	Food & Drink Shop	Park	Clothing Store	Grocery Store	Art Gallery	Pub
47	Danforth	Dessert Shop	Italian Restaurant	Sushi Restaurant	Park	Coffee Shop	Pizza Place	Seafood Restaurant	Sandwich Place	Cafe
51	Cabbagetown/St. James Town	Cafe	Pet Store	Taiwanese Restaurant	Gastropub	Indian Restaurant	Italian Restaurant	Japanese Restaurant	Jewelry Store	Diner
52	Church and Wilesey	Tea Room	Themed Restaurant	Breakfast Spot	Bookstore	Juice Bar	Diner	Salon / Barber Shop	Restaurant	Ramen Restaurant
54	Ryerson Garden District	Cafe	Burger Joint	Ramen Restaurant	Tea Room	Buncho	Thai Restaurant	Theater	Plaza	Pizza Place
55	St. James Town	Coffee Shop	Gastropub	Japanese Restaurant	BBQ Joint	Food Truck	Italian Restaurant	Middle Eastern Restaurant	Crepes	Cosmetics Shop
56	Stacey Park	Farmers Market	Coffee Shop	Museum	Liquor Store	Seafood Restaurant	Bookstore	Thai Restaurant	French Restaurant	Breakfast Spot
58	Addicks/King Richmond	Steakhouse	Coffee Shop	Pizza Place	Spa/Salon	Bar	Hotel	Plaza	Asian Restaurant	Food Court
59	Harbourfront East/Toronto Islands Union Station	Hotel	Plaza	Bubble Tea Shop	Sporting Goods Shop	Supermarket	Solid Place	Deli / Bodega	Bakery	Skating Rink
60	Design Exchange/Toronto Dominion Centre	Coffee Shop	Cafe	Gym	Hotel	Beer Bar	Japanese Restaurant	Restaurant	Bakery	Pub

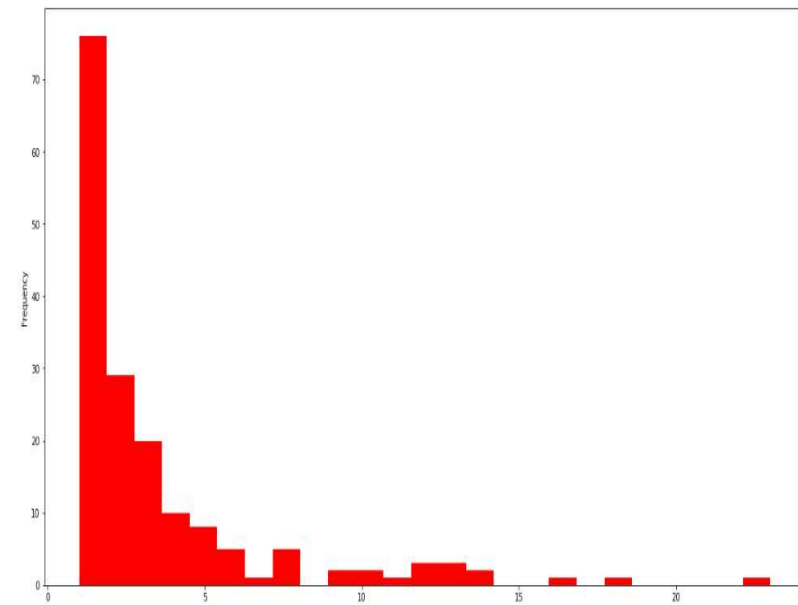
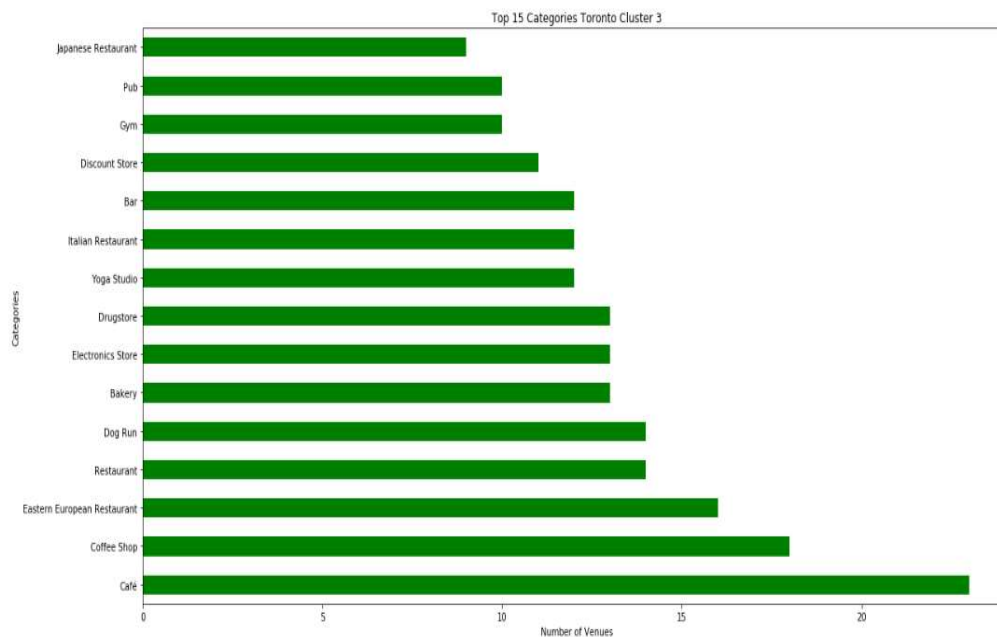


Table 2. Toronto's cluster 3 category frequencies

	categ	clot3
33	Café	23.0
42	Coffee Shop	18.0
62	Eastern European Restaurant	16.0
138	Restaurant	14.0
60	Dog Run	14.0
15	Bakery	13.0
63	Electronics Store	13.0
61	Drugstore	13.0
169	Yoga Studio	12.0
104	Italian Restaurant	12.0
17	Bar	12.0
59	Discount Store	11.0
90	Gym	10.0
135	Pub	10.0
105	Japanese Restaurant	9.0

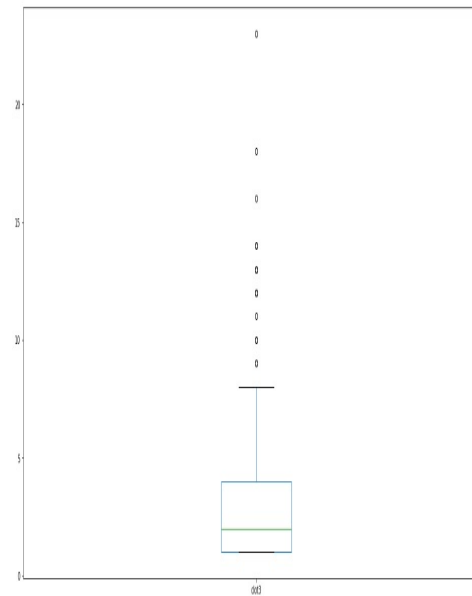
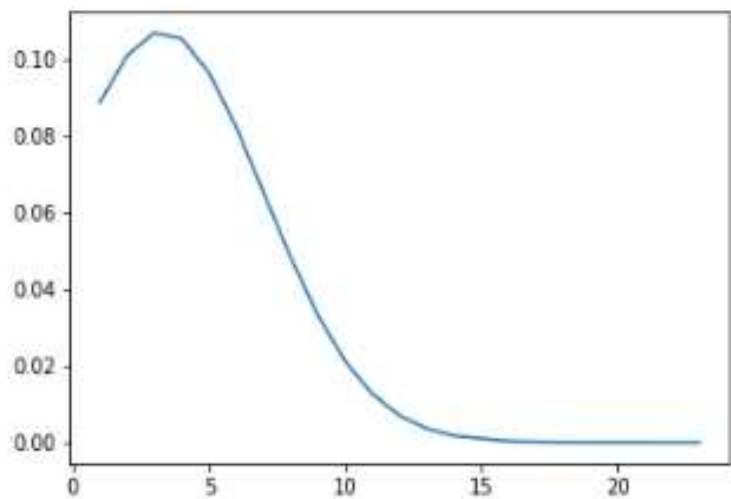
## 4. Exploratory Analysis

- Data visualization: Café, Coffee Shop, Restaurants, Bakery – food delivery. But other kind of venues are also important: Gym, Drug



## 4. Exploratory Analysis

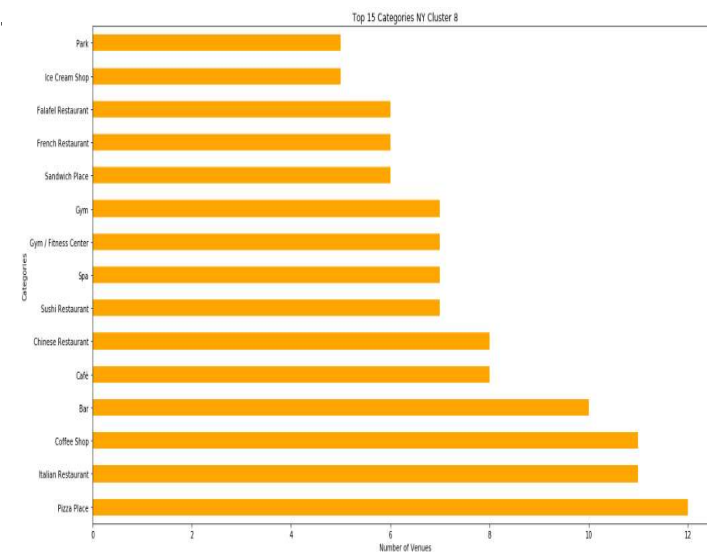
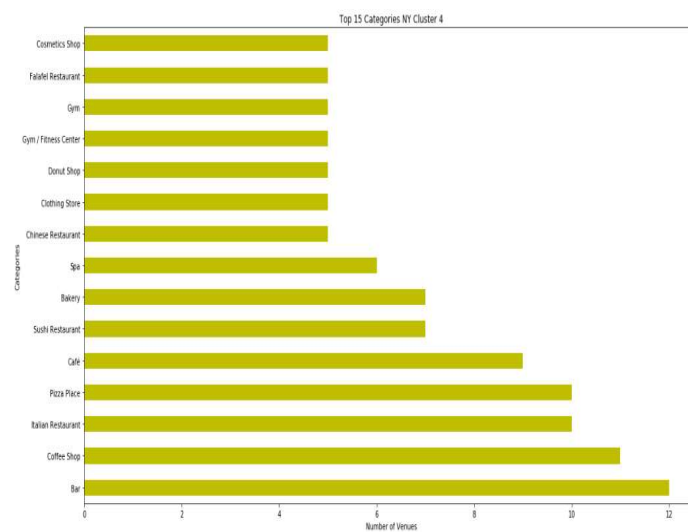
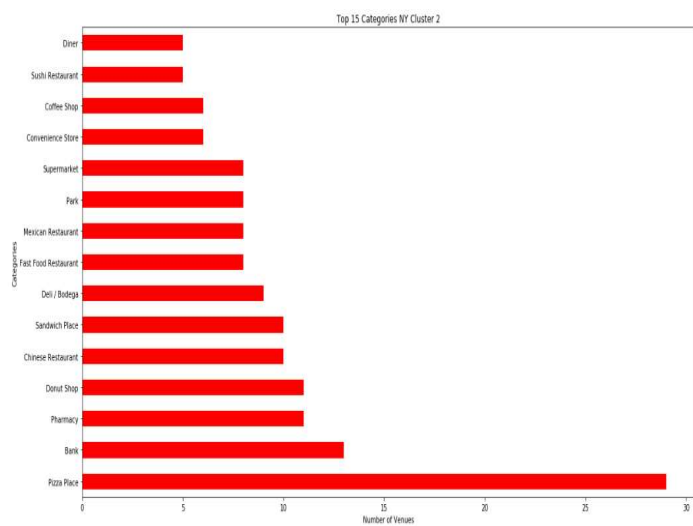
- Descriptive Statistics: data is concentrated around 3 (venues per category), but there are many outliers such as Café, Coffee Shops and Restaurants



clot3	
count	170.000000
mean	3.294118
std	3.738351
min	1.000000
25%	1.000000
50%	2.000000
75%	4.000000
max	23.000000

## 4. Exploratory Analysis

- Data visualization: NY cluster 2 seems to be more similar with Toronto cluster 3





# 5. Methodology

PREDICTING THE BEST LOCATION WITH MORE POTENTIAL REVENUE



## 5.1 Models

- ▶ We decided to use two kind of models (machine learning): logistic regression and decision tree.
- ▶ We created the “y” variable for Toronto's cluster 3 with threshold in 3: if a category venue has a frequency higher than 3, that it is a potential client to our client (a good proxy for delivery transactions: places with certain type of venues have more delivery transactions than others).
- ▶ Therefore, we have a classification problem. Two type of models will help us in compare results from the training and testing part, as well the prediction part, when we classify the NY 3 clusters in order to verify which has more potential to our client.

## 5.2 Model Evaluation

- ▶ Logistic regression presented better performance than the decision tree with the testing data set
- ▶ Jaccard index (how many classifications were correct) for LR was 94.11% and for Decision tree was 78,31%

Table 4. Logistic Regression performance metrics

	precision	recall	f1-score	support
0.0	0.92	1.00	0.96	23
1.0	1.00	0.82	0.90	11
micro avg	0.94	0.94	0.94	34
macro avg	0.96	0.91	0.93	34
weighted avg	0.95	0.94	0.94	34

Table 5. Decision Tree performance metrics

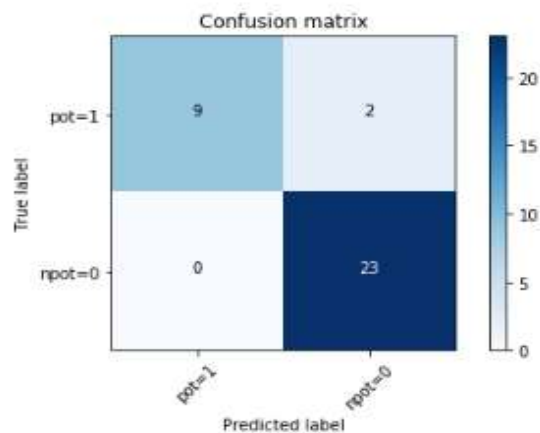
	precision	recall	f1-score	support
0.0	0.72	1.00	0.84	28
1.0	1.00	0.52	0.69	23
micro avg	0.78	0.78	0.78	51
macro avg	0.86	0.76	0.76	51
weighted avg	0.85	0.78	0.77	51

## 5.2 Model Evaluation

- ▶ we can realize that both classification models are good in classifying non potential clients instead of potential clients.

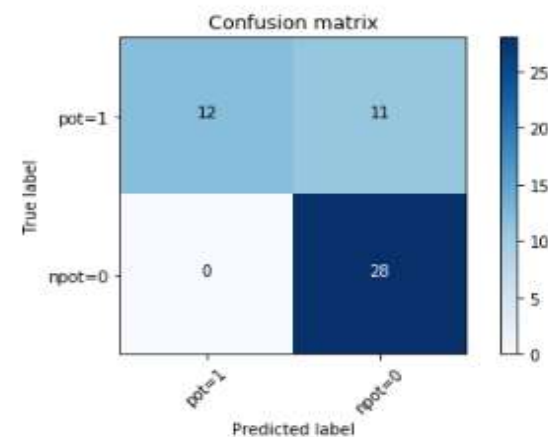
Logistic Regression Confusion Matrix.

Confusion matrix, without normalization  
[[ 9 2]  
[ 0 23]]



Decision Tree Confusion Matrix.

Confusion matrix, without normalization  
[[12 11]  
[ 0 28]]





# 6. Results

PREDICTING THE BEST LOCATION WITH MORE POTENTIAL REVENUE

## 6. Results

- ▶ the models classified, using each cluster data set, the categories venues that are interesting in terms of clients to our client – more potential to bring delivery transactions, and therefore revenue. The results were not so different between models and between NY clusters.
- ▶ Their proportion of adherent venue categories (classified as good for our client,  $y = 1$ ) are similar, around 25%, just for cluster 8 in the Decision Tree model we see a number quite distant from the others

Table 5. Predictions for NY Clusters using Logistic Regression Model

	NY Cluster2	NY Cluster4	NY Cluster8
Proportion adherent	29.4	26.0	23.0
Cluster Size	340.0	340.0	350.0
Potential clients	100.0	88.0	82.0

Table 6. Predictions for NY Clusters using Decision Tree Model

	NY Cluster2	NY Cluster4	NY Cluster8
Proportion adherent	21.2	23.1	17.1
Cluster Size	340.0	340.0	350.0
Potential clients	72.0	79.0	82.0



# 7. Conclusion

PREDICTING THE BEST LOCATION WITH MORE POTENTIAL REVENUE

## 7. Conclusion

- ▶ We analyzed spatial data from Foursquare about neighborhood and venue categories in order to identify the best location for a start-up (delivery business) based on Toronto to begin its expansion in NY.
- ▶ We created clusters for Toronto and for NY (unsupervised), and from the best cluster for our client in Toronto we predict the best cluster in NY for our client (supervised).
- ▶ From the results, cluster 2 seems to be the best choice for our client begin its expansion in NY. This cluster had the best results in the best model.
- ▶ But cluster 4 is a good choice as well, since its results are not so far from those of cluster 2. Perhaps it could be a second phase of our client's expansions in NY.



# 8. Future Directions

PREDICTING THE BEST LOCATION WITH MORE POTENTIAL REVENUE



## 8. Future Directions

- ▶ We used just Foursquare data.
- ▶ A good improvement could be gathering other kind of data, like socioeconomic data for example.
- ▶ Such information can change the results that we found, bring more information.
- ▶ Other improvement could be testing SVM and Neural Network models in order to compare results.
- ▶ Moreover, the value of the threshold is something that could change results if it is changed, as well as the number of clusters for each city.