

IBM Supervised Machine Learning Classification

Course Final Project:
Major Projects Classification by Cost Range

By Javier A. Jaime-Serrano
June 16, 2021

Abstract

For this Project we used the Major Project dataset from the province of Alberta in Canada [1], it contains more than 700 currently active projects on the province, this Data set is filtered for Projects valued at \$5 million or greater.

This Dataset contains a lot of valuable information on the Major Projects. The Estimated Cost, Sector & Type and the Location data will be used.

We want to be able to classify projects by type, cost and capacity and compare them with similar projects from the same dataset.

We will use Decision Tree and Random Forest Classifiers from Scikit-learn [2].

Objectives

From the Major Projects in the dataset, we will like to know if the cost of the project is within a reliable range for the project type and size and/or capacity.

If the classifier results are projects “too cheap” or “too expensive”, or there is not enough data available in the data set, further analysis will be required.

For the type of projects that we have enough data within the range, we will be able to estimate the cost of similar projects, with a rough order of magnitude (ROOM).

This quick ROOM Cost Estimate can be used as benchmark in multiple industries.

Data Cleaning

In order to prepare and clean the dataset:

- We drop the projects where there is no estimated cost.
- Made some assumptions about schedule completion and status.
- Drop not required columns and renamed the remaining columns.
- Made corrections on project types and sectors.
- A problem encountered was how to extract the location coordinates (Longitude & Latitude) from a GeoJSON column. The problem was solved with Python code that loop over all rows and extract the start (first) locations by a type condition.

Data Cleaning

Figure 1. Cleaned Dataset

ProjectId	Project Name	Estimated Cost (millions)	Municipality	Forecasted Completion	Sector	Type	Stage	Developer	Start Latitude	Start Longitude
7	StoneGate Landing	3000.0	Calgary	2021	Mixed-Use	Mixed-Use	Started	WAM Development Group / AIMCo	51.172501	-113.975800
11	Shepard Station Suburban Office Campus Building 1	22.0	Calgary	2020	Commercial	Office: Low-Rise	Started	Shepard Development Corp.	50.931721	-113.970596
22	Barron Building Renovation	100.0	Calgary	2021	Residential	Apartment: Mid-Rise	Proposed	Strategic Group	51.046070	-114.076614
26	Quarry Crossing II Office Building	72.8	Calgary	2027	Commercial	Office: Low-Rise	Proposed	Remington Development Corp.	50.966900	-114.002899
32	Nolan Hill TownHomes	5.0	Calgary	2027	Residential	Townhouses	Proposed	Jayman Modus	51.162041	-114.160912

Data Exploration

We explored the data, first with descriptive statistics and bar charts (see Figure 2).

Second, we used box plots for the cost estimate ranges by type (see Figure 3).

Third, we used folium library to create a map using latitude and longitude values (See Figure 4).

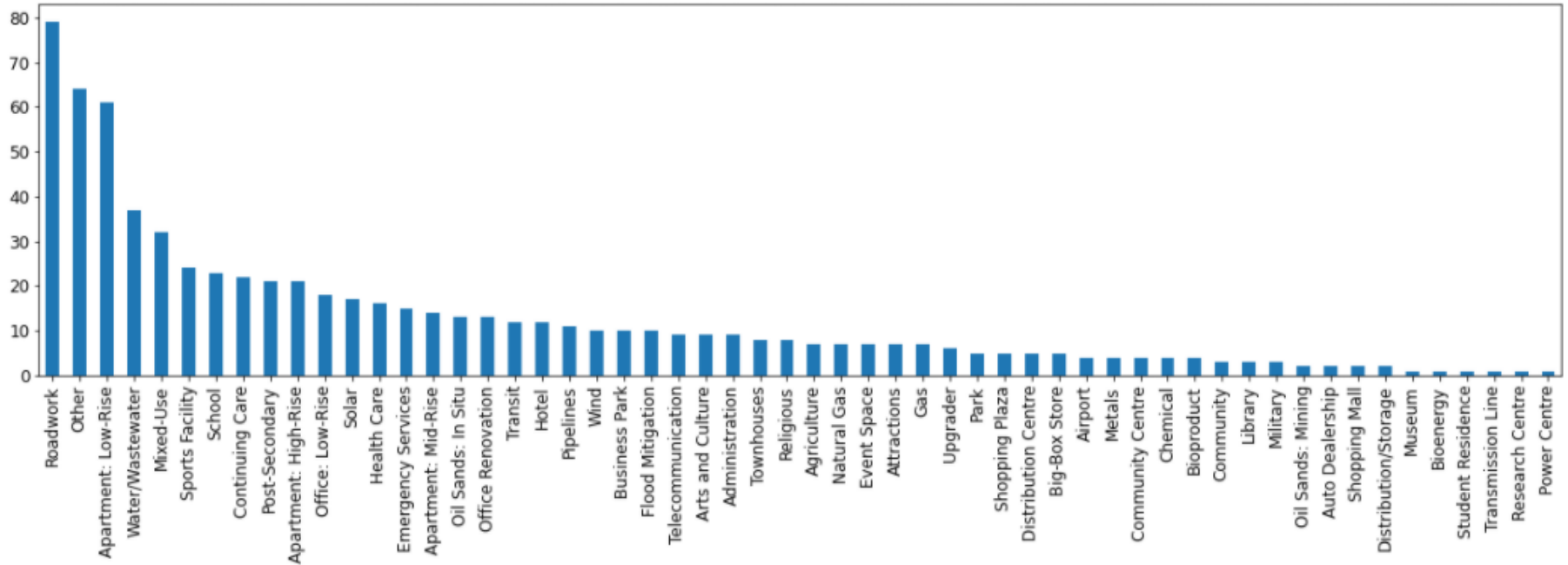
We also used KMeans [2] to Cluster the projects by geographical region using the location coordinates (See Figure 5).

And then we used the Classifiers from Scikit-learn [2] to predict one of 4 labels:

1. Too Cheap 2. In the Range. 3. Too Expensive 4. Not Enough Data

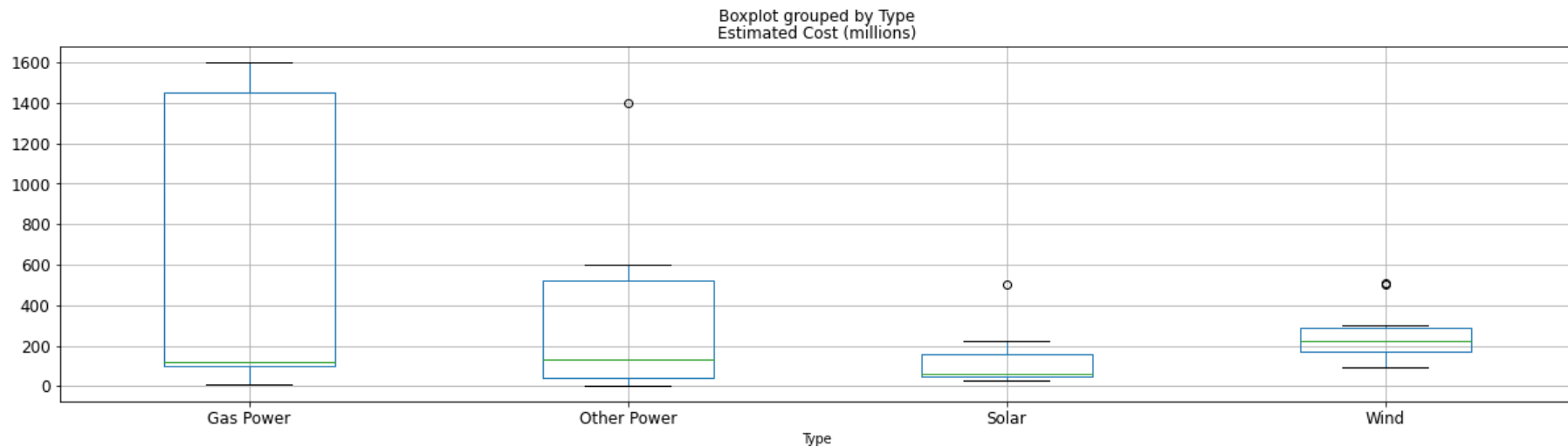
Data Exploration

Figure 2. Project types bar chart



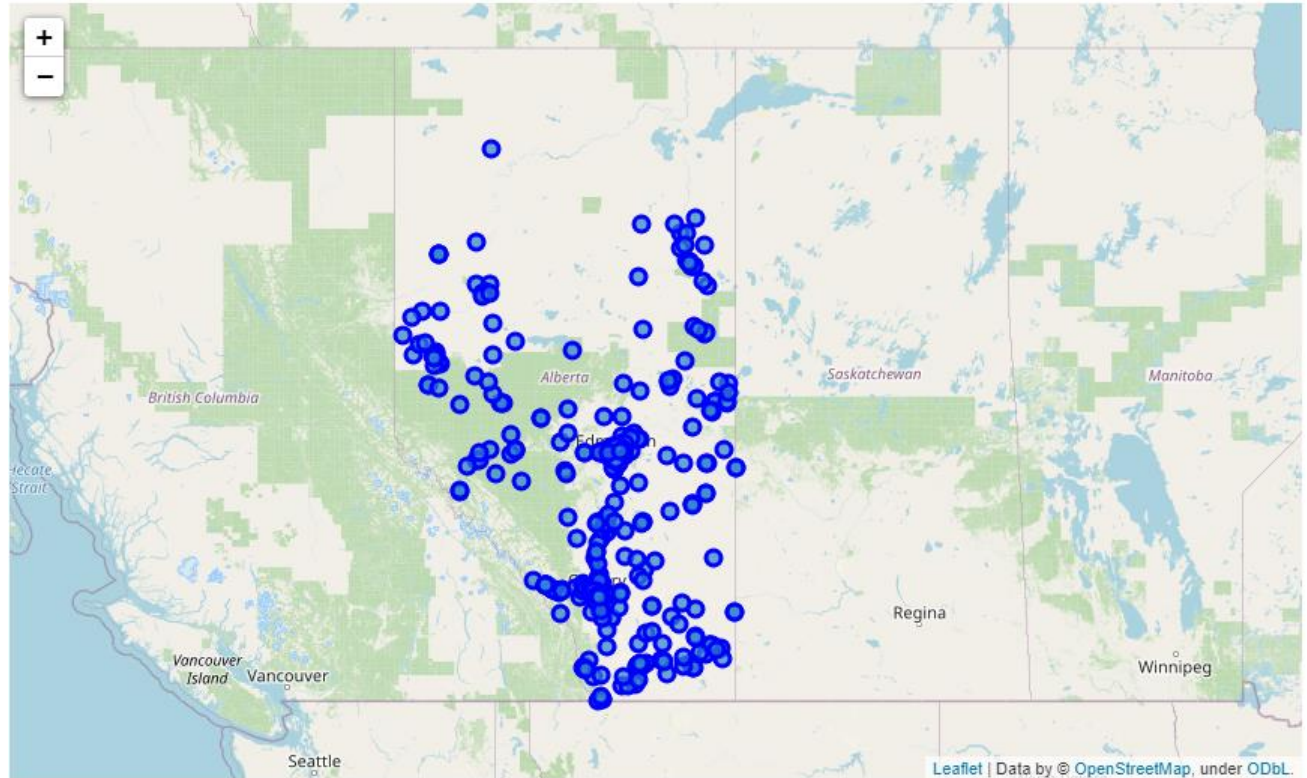
Data Exploration

Figure 3. Box Plot Estimated Cost for Power Sector



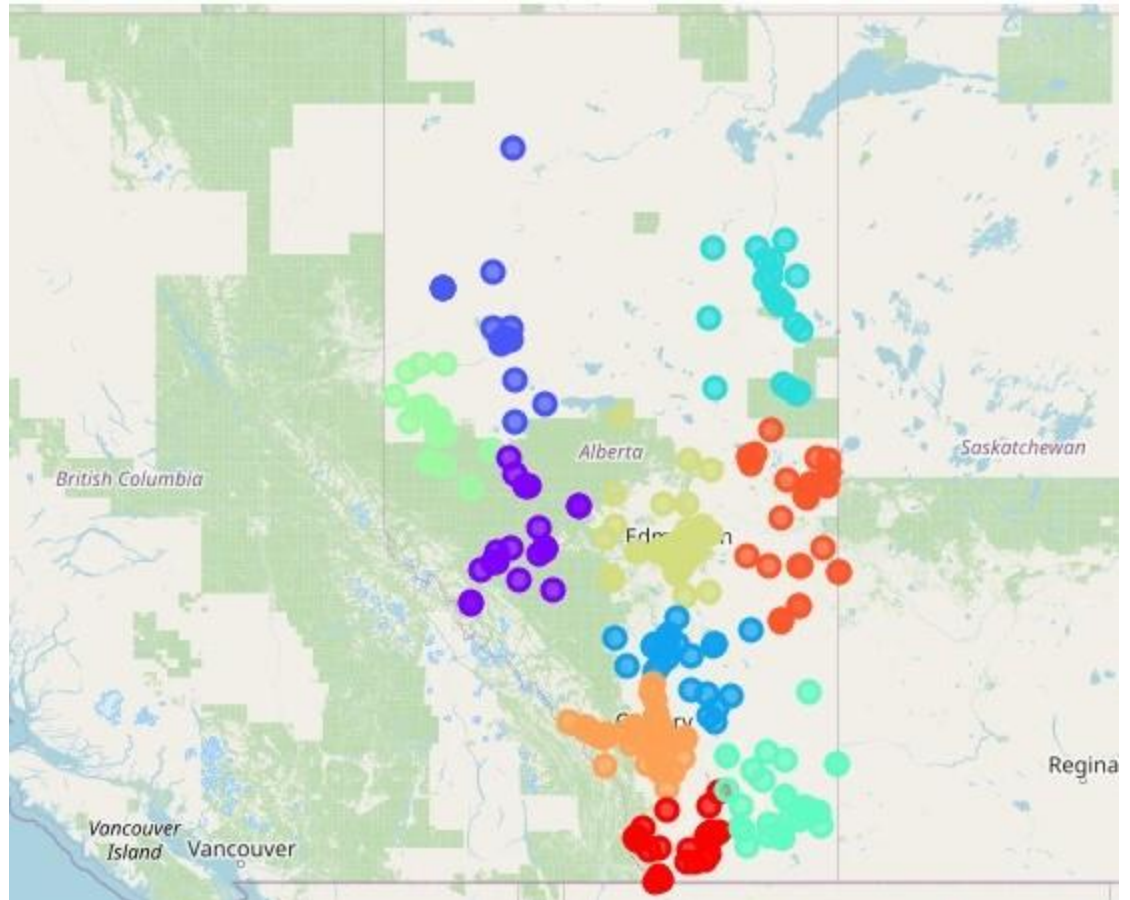
Data Exploration

Figure 4.
Projects Map



Data Exploration

Figure 5. Projects Map
KMeans Clusters by
Region with Project type



Findings

After corrections, we ended with 58 unique types in 9 sectors (see figure 6).

For the feature engineering, we extracted size and capacity data by type from the dropped project details column, adding units and cleaning it manually in excel.

The target variable was set from a percentage range in relation to average cost, divided by capacity per type of projects, labeled with conditional formulas in excel.

The datasets were merged back, and one hot encoding was used again to transform the categorical features into numerical dummy features (89 Columns).

Findings

Figure 6.
Unique Sectors
and Types

		Project Name						
Sector	Type		Institutional	Administration	9	Residential	Apartment: High-Rise	18
				Continuing Care	22		Apartment: Low-Rise	61
Commercial	Business Park	10		Emergency Services	15		Apartment: Mid-Rise	14
	Distribution Centre	5		Health Care	16		Community	6
	Office Renovation	13		Library	3		Other Residential	4
	Offices	18		Military	3		Townhouses	8
	Other Commercial	4		Other Institutional	7	Retail	Auto Dealership	2
Industrial	Agriculture	7		Post-Secondary	21		Big-Box Store	5
	Bioproduct	4		Religious	8		Mixed-Use	33
	Chemical	4		School	23		Other Retail	4
	Metals	4	Oil and Gas	Distribution/Storage	2	Shopping Mall	2	
Infrastructure	Other Industrial	14		Gas	7		Shopping Plaza	5
	Telecommunication	9		Oil Sands: In Situ	13	Tourism	Arts and Culture	9
	Airport	4		Oil Sands: Mining	2		Attractions	8
	Flood Mitigation	10		Other Oil and Gas	2		Community Centre	4
	Other Infrastructure	16		Pipelines	11		Event Space	7
		Roadwork	79		Upgrader	6	Hotel	12
		Transit	12	Power	Gas Power	7	Other Tourism	9
		Water/Wastewater	37		Other Power	8	Park	5
					Solar	17	Sports Facility	24
					Wind	10		

Results

The Decision Tree Classifier was trained with 70% of the Dataset, reserving the remaining 30% for Testing.

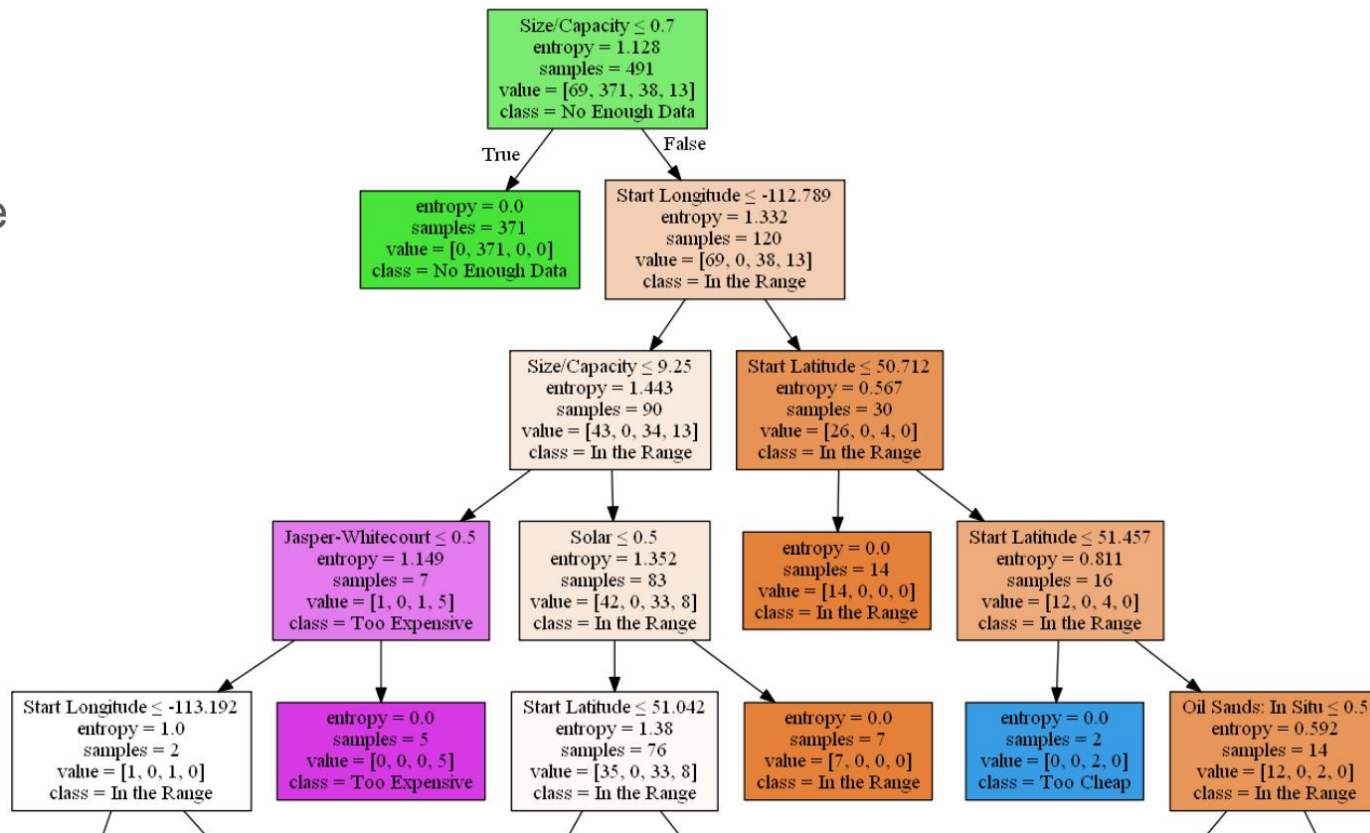
We used the Entropy (information gain) as the classifier criterion, and a maximum depth of 20 nodes (these are the 20 questions of the Project Name).

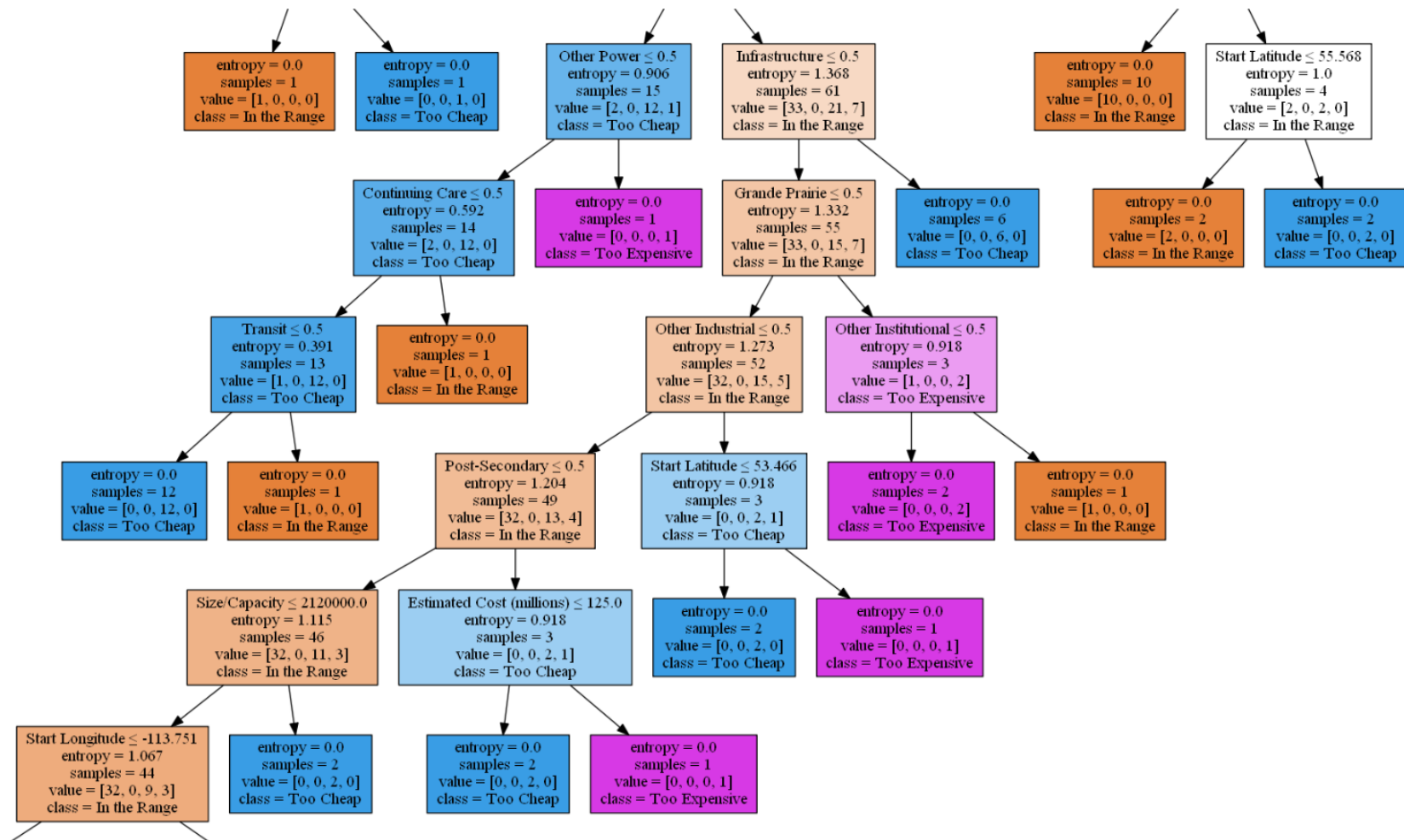
We used Grid Search to find the best hyperparameters of number of nodes and maximum depth of the tree.

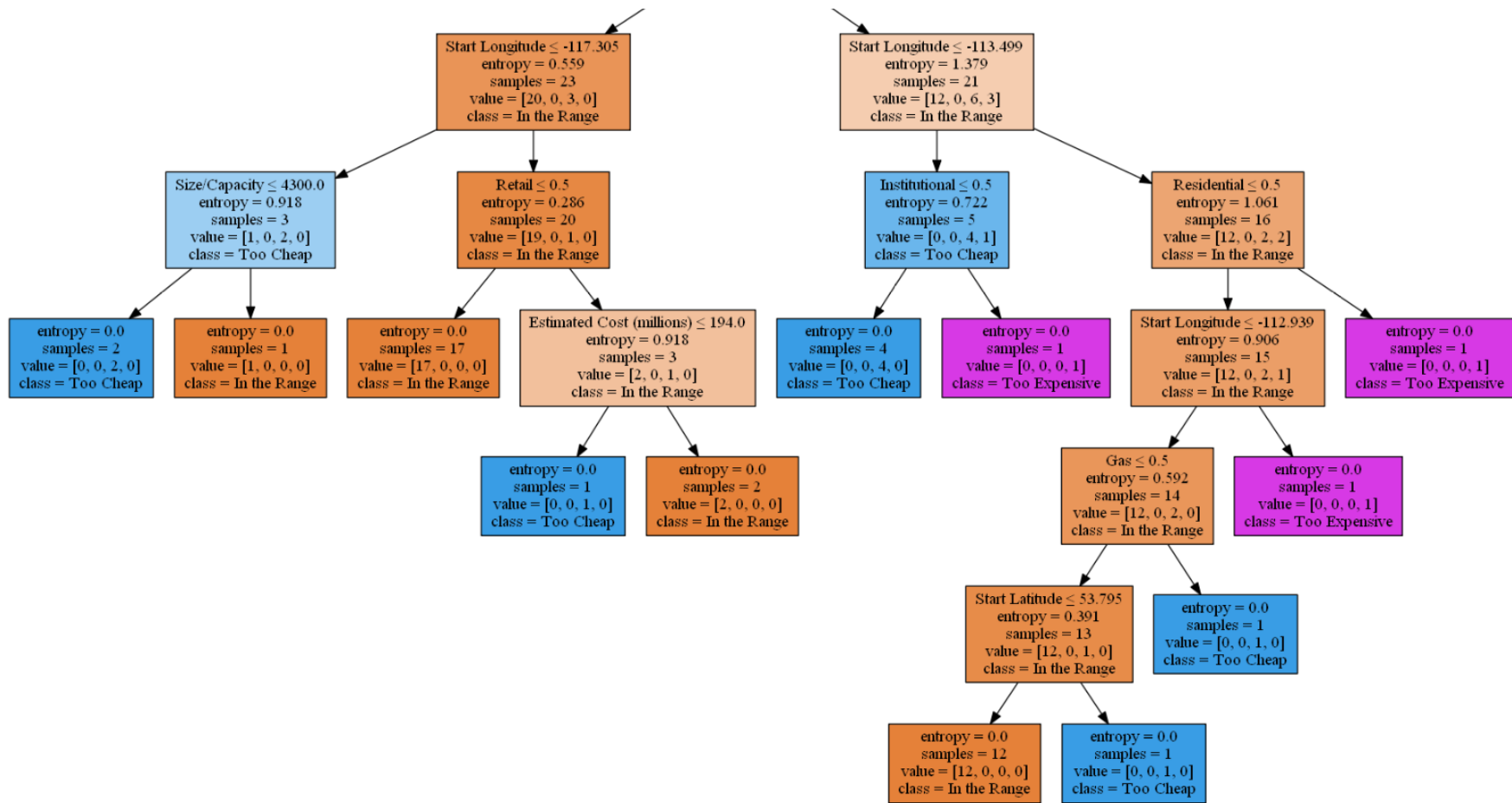
The Decision Tree Accuracy is 89% on the test set, the maximum number of nodes, depth is 16 and number of nodes is 67 (see figure 7).

Results

Figure 7.
Decision Tree







Limitations

The results of Decision Tree Classifier show the limitations of the Dataset, the first split take most of the projects out, classifying them as “Not Enough Data”.

Only 97 projects were classified as “In the Range”, and Only 3 types has more than 10 Projects in this valid range.

This work was done with only one Province Dataset, and other Canadian Provinces has Datasets available for Future work in similar formats.

Conclusions

We were able to classify project by type, cost and capacity with the use of a Decision Tree Classifier with 89% accuracy on the training set.

Most of the projects has not enough data for the engineered features of size and/or capacity. Further work is required to search for capacity data from other sources (company websites, industry associations, etc.).

We found additional uses for the classifier, as the detected outliers were in a couple of cases, big projects announcements (as publicity stunts) and later withdrawn. And projects that were too expensive were never completed.

References

[1] Alberta Major Projects: <https://majorprojects.alberta.ca/>

[2] Scikit-learn library: <https://scikit-learn.org/stable/>

[3] Jupiter Notebook: <https://github.com/javier-jaime/IBM-Machine-Learning-Capstone/>