

TASK 1 DATA MINING CAPSTONE PROJECT

1. TASK 1.1

In the first task of the Data Mining Casptone Project, it was used a LDA algorithm applied to a 100,000 reviews sample provided by the python script in order to get exactly 10 topics.

I first used in python the following line on the terminal:

```
python py27_processYelpRestaurants.py -sample
```

to get the 100,000 restaurant reviews from 10 different cuisines. After that, the command

```
python py27_ldaTopicModeling.py -o "sample_topics10.txt" -K 10
```

provided me with 10 different topics with its respective weights from the LDA Topic Modelling Algorithm.

In the following visualisation, I depicted with *D3.js* using the Radial Reingold-Tilford Tree (for reference see <http://bl.ocks.org/mbostock/4063550>) the distribution from the topics with its terms: Each size of the nodes describes the weight of the term in the topic.

The visualisation (next page) is convenient and very appropriate and shows in a direct way how is the distribution of the topics from the given sample. Although if one reviews terms, they are not entirely consistent within each topic suggesting a tune in the LDA model.

2. TASK 1.2

Due to lack of time, I was not able to complete this task. Hopefully I will have more time for Task 2 :).

