# Multi-Document Text Summarization for Competitor Intelligence: A Framework

A THESIS
SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE **FELLOW PROGRAMME IN MANAGEMENT (INDUSTRY)**,
INDIAN INSTITUTE OF MANAGEMENT, INDORE

By

Swapnajit Chakraborti
FPM (Industry) 2013
IIM Indore, India
MARCH, 2017

## Thesis Advisory Committee

1.    ------------------------------------------------------ [Chair]
            [Prof. Shubhamoy Dey]

2.    ------------------------------------------------ [Member]
            [Prof. D. L. Sunder]

3.    ------------------------------------------ [Member]
            [Prof. Rajhans Mishra]

# Abstract

Proliferation of web as an easily accessible information resource has led many business organizations to gather competitor intelligence (CI) directly from various resources on the internet, namely, news, blogs, reports, reviews etc. While collection of such information is easy from the internet, the collation and structuring of them for perusal of business decision makers, is not very straight-forward. This work is an endeavor towards exploring the effectiveness and feasibility of text clustering followed by text summarization as a technique for comprehending competitor intelligence and thereby use it as an input for strategic decision making by the managers. Although applications of various text clustering and text summarization techniques are abundant in state-of-the-art literature, there have been very few authentic research on understanding their impact in organizing competitor intelligence for managerial decision making. This research contributes towards developing an integrated framework which supports analysis of competitor information by clustering the corpus collected from internet using a variation of K-means, called Multi-level K-means clustering, and then generating the extractive summaries of these topical clusters by applying global optimization techniques, such as, Artificial Bee Colony (ABC) optimization, Differential Evolution (DE) etc. The results of analysis of the quality of these system generated summaries, measured w.r.t. human generated golden summaries, are also presented in this thesis. As effectiveness of the system generated summaries towards competitor intelligence analysis, can be best judged by the practicing managers at various organizations, a study is also undertaken to obtain managerial feedback on these automatically generated system summaries and corresponding findings are reported.

_____

# List of Figures

# List of Tables

# Glossary of Terms

| Term | Meaning | Chapter Reference |
|---|---|---|
| LSA (Latent Semantic Analysis) | Finding out the hidden/latent topics within a corpus by applying matrix factorization | 2 |
| Centroid | Central theme of a text corpus represented by important words in the corpus | 2 |
| Topic Identification | Finding the topics from a corpus; done by clustering as well as probabilistic modelling | 5 |
| Corpus | Collection of text documents | 4 |
| TF-IDF | Term Frequency*Inverse Document Frequency | 4 |
| S-KM | Standard K-means (K-means++) clustering technique used as building block for Multi-level K-Means | 5 |
| APS | Average Pairwise Similarity; used as metric for cluster quality | 5 |
| AAPS | Average of APS of all clusters generated from a corpus | 5 |
| Extractive Summary | Summary generated by extracting sentences as it is from the corpus | 6 |
| System Summary | Summary generated automatically by applying optimization techniques | 6 |
| Golden Summary | Summary generated by human volunteer from a corpus | 6 |
| Total Penalty | Loss of similarity with central theme of a cluster computed by removing one sentence at a time from candidate system summary | 6 |
| ABC Optimization | Artificial Bee Colony Optimization; latest swarm optimization based on foraging behaviour of honey bees | 6 |
| DE Optimization | Differential Evolution based optimization | 6 |
| ROUGE | Recall Oriented Understudy of Gisting Evaluation; summary quality measurement tool | 6 |
| Recall Score | Fraction of golden summary that matches with system summary; measured using sequence overlap | 6 |
| Precision Score | Fraction of system summary that matches with golden summary; measured using sequence overlap | 6 |
| ML-KM | Multi-level K-means Clustering technique introduced for clustering CI corpus | 5 |
| CI | Competitor Intelligence | 1 |
| ATS | Automatic Text Summarization | 1 |

| TSS | Total Summary Score; Optimization function used to evaluate a candidate summary generated during optimization process | 6 |

# Table of Contents

_____