

Application of AI on IT Service Management



Arun singh Sivaprakash
Pramod kumar Nagaraj
Van Tien NGUYEN
Alexander POPPE

EPITA Graduate School of Computer Science

Advisor
Olivier Berthet

Master of Science in Artificial Intelligence Systems

16-07-2021

Contents

1	Introduction	1
1.1	Motivations	1
1.2	Context of the Project	1
1.3	Objectives and Contributions	2
1.4	Overview of thesis	2
2	State of Art	3
2.1	NLP pipeline	3
2.1.1	Text Processing	3
2.1.2	Feature Extraction	4
2.1.3	Modeling	5
2.2	Text Data Visualization	7
2.2.1	Importance of Data Visualization	8
2.2.2	Tools used for Text Data Visualization	8
3	System's Artchitecture	10
3.1	Insights Dashboard Architecture	10
3.2	Clustering System Artchitecture	11
4	Methodology	12
4.1	Data Visualization – ELK Dashboard	12
4.2	Clustering	12
5	Results	13
6	Conclusion	14
7	References	15

Abstract

Large enterprise Information technology infrastructure components generate large volume of alerts and incident tickets. These are manually screened, but it is otherwise difficult to extract information automatically from them to gain insight in order to improve operational efficiency. We propose a framework to cluster alters and incident tickets based on the text in them, using unsupervised machine learning. This would be a step towards eliminating manual classification of the alters and incidents, which is very labour intense and costly.

The challenges that arrive at analysing huge amount text data starts with the data quality. Data quality metrics concerning the interpretation of text are too difficult to compute automatically and too pragmatic because they heavily depend on the linguistic content and the specific situation in which the text is being used. Several statistical approaches have been used to determine some level of data quality in this project.

The visual representation of text data helps to bring more insights from the data. Building dashboard with various filters will help the users to customize the search. We have implemented state of art visualization dashboards with different open-source tools.

We proposed the techniques to visualize the clusters that make human readable in an easy and understandable way compared to the traditional methods, we find a simple way to define prototype of cluster for easy interpretation. This framework for clustering and visualization will enable enterprises to prioritize the issue in their IT infrastructure and improve the reliability and availability of their services.

Chapter 1

Introduction

1.1 Motivations

Organization of all sizes struggles with a common problem in Infrastructure and Operational Management. Thousands of automated alerts with semi-structure text are generated every day from hundreds of infrastructure tools. There are thousands of incidents tickets with manually entered unstructured text are created daily by support personnel. After collecting all the details of the incidents maintaining of those tickets is very hard and costly and difficult to analysis each incident without grouping them.

To solve this issue, we proposed a machine learning technique which do clustering alters with semi-structured text and clustering incidents with unstructured text. With the help of Natural language process, we have pre-processed the data i.e, data tokenization and removed all the common words i.e, text mining methods. Before clustering we also be working in data quality such that we can able check how good the data is for the clustering them.

1.2 Context of the Project

Organization operating in IT enables business segment over a decade. On an average, the company receive more than 25k IT incidence or tickets. Which are handle to best practice by certain framework which are done by researchers with incident management, problem management etc.

We are developing the framework which help the IT tech companies to visualize the data quality and data clustering. By this company can understand the problem of the user in systematic manner within no time, and the companies auto robots can answer for the user problem with more effectively.

1.3 Objectives and Contributions

1.4 Overview of thesis

Chapter 2

State of Art

Natural language processing (NLP) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data. The goal is a computer capable of "understanding" the contents of documents, including the contextual nuances of the language within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themselves. In our project, NLP is used to extract information tickets and evaluate the quality of data.

2.1 NLP pipeline

In general, NLP contains three main stages: text processing, feature extraction, modeling.



Figure 2.1: NLP pipeline

2.1.1 Text Processing

In this step, we take text raw as an input. Then we process the input into a form that is the best for next step: feature extraction. More specifically,

raw data needs to be cleaned that means removing special characters such as HTML tags, etc. In other words, input data does not contain any info for the model to learn and are irrelevant or noisy data. After that, we do the data normalization which involves the case normalization, punctuation removal so that the text is in a single format for the machine to learn. In case normalization, we convert all capitalization to lower to bring a common case, and in punctuation removal, we replace all punctuations with space. The next process is breaking up text documents into individual words called tokens and remove the stop words. Stop word removal means removing the non-important words like "a", "is", "the", "and", "an", etc. After that, we need to reduce words to its normalized form, and this process can be done by stemming or lemmatization technique. Stemming is a process of reducing a word to its root form, for example: "branching", "branched" and "branches" can all be reduced to "branch", meanwhile lemmatization is the algorithmic process of determining the canonical form (lemma) of a word. Unlike stemming, it is not only a word reduction but depends on the meaning of the word in a sentence and needs to consider a language's full vocabulary, for instance: "was" will be transformed into "be", "meeting" will be "meet" or "meeting" depending on the context.

2.1.2 Feature Extraction

Feature Extraction is a way of extracting feature vectors from the text after the text processing step so that it can be used in the machine learning model as input. Word Embedding is one such technique where we can represent the text using vectors. The more popular forms of word embeddings are: Bag-of-Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF). BoW is a way of extracting features from the text for use in modeling, it treats each document as a collection or bag of word. This means a representation of text that describes the occurrence of words within a document. It involves two things: a vocabulary of known words in the corpus (set of documents) and a measure of the presence of known words. For example: we have two documents: Document1 is "Xavier likes to play football. Eric likes football too." and Document2 is "Eric prefers tennis to football", we can have the representation below:

	Xavier	likes	to	play	football	Eric	too	prefers	tennis
Document1	1	2	1	1	2	1	1	0	0
Document2	0	0	1	0	1	1	0	1	1

Figure 2.2: Example of BoW

Term Frequency-Inverse Document Frequency is a ponderation method used in information retrieval. This statistical measure gives an evaluation of how important is a word to a document, depending on the corpus considered. TF-IDF is calculated as: $tfidf(t, d, D) = tf(t, d) * idf(t, D)$ Inverse document frequency: $idf(t, D) = \log(\frac{N}{|\{d \in D: t \in d\}|})$ with N the number of documents in corpus D. We continuously calculate tf-idf for the above example:

	Xavier	likes	to	play	football	Eric	too	prefers	tennis
Document1	1	2	1	1	2	1	1	0	0
Document2	0	0	1	0	1	1	0	1	1
idf	Log(2)	Log(2)	0	Log(2)	0	0	Log(2)	Log(2)	Log(2)
tf_doc1	1/9	2/9	1/9	1/9	2/9	1/9	1/9	0	0
tf_doc2	0	0	1/5	0	1/5	1/5	0	1/5	1/5
tf-idf_doc1	0,033	0,067	0	0,033	0	0	0,033	0	0
tf-idf_doc2	0	0	0	0	0	0	0	0,06	0,06

Figure 2.3: Example of TF-IDF calculation

2.1.3 Modeling

The final stage of the NLP pipeline is modeling, which includes designing a statistical or machine learning model, fitting its parameters to training data, using an optimization procedure, and then using it to make predictions about unseen data. In our project, we use clustering machine learning algorithms. Here is an overview of this algorithm. Clustering is an unsupervised machine learning task. It involves automatically discovering natural grouping in data. Unlike supervised learning, clustering algorithms only interpret the input data and find natural groups or clusters in features in feature space. A cluster is often an area of density in the feature space where examples from the domain (observations or rows of data) are closer to the cluster than other clusters. The cluster may have a center (the centroid) that is a sample or a point feature space and may have a boundary or extent. Clustering can be helpful as a data analysis activity in order to learn more about the problem domain. It also can be useful as a type of feature engineering, where existing and new examples can be mapped and labeled as belonging to one of the identified clusters in the data. Evaluation of identified clusters is subjective and may require a domain expert, although many clustering-specific quantitative measures do exist. Typically, clustering algorithms are compared academically on synthetic datasets with pre-defined clusters which an algorithm is expected to discover. There are two common clustering

algorithms: K-Means and Gaussian Mixture Model.

K-Means:

This is one of the simplest and frequently used unsupervised learning algorithms, especially in data mining and statistics. Being a partitioning algorithm, its goal is to form groups of data points based on the number of clusters, represented by the variable k . K needs to be predefined before the execution. K-means uses an iterative refinement method to produce its final clustering based on the number of clusters defined by the user and the dataset. Initially, k-means randomly chooses k as the mean values of k clusters, called centroids, and find the nearest data points of the chosen centroids to form k clusters. Then, it iteratively recalculates the new centroids for each cluster until the algorithm converges to one optimum value. K-means clustering would be suited with the numerical data with a low dimensionality because numerical data is used to compute the mean value. The type of data best suited for K-means clustering would be numerical data with a relatively lower number of dimensions.

Problem description: A set of observations (x_1, x_2, \dots, x_n) where each observation is a d -dimensional real vector, this algorithm aims to partition the n observations into $k (<= n)$ sets $S = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within cluster sum of squares.

Algorithm: initial randomly set of k means (m_1, m_2, \dots, m_k) , the algorithm proceeds by altering between two steps. The first step is assignment step: assign each observation to the cluster with the nearest mean that with the least squared Euclidean distance: $S_i = \arg \min_k ||x_i - m_k||^2$. The second step is update step which recalculates the means (centroids) for observation assigned to each cluster: $m_k = \frac{\sum_{i:k_i=k} x_i}{|i:k_i=k|}$. The algorithm has converged when the assignments no longer change.

K-means has some advantages which are simple clustering algorithm so that it can be implemented easily. It also faster computationally than other clustering algorithms because K-means has only a few computations. Additionally, this algorithm can scale up to large dataset and easily adapt to new data samples. Besides, there are some disadvantages: K-means has difficulty with clustering data set of varying sizes and density. Its result can also vary depending on initial values and the number of clusters has to be specified manually.

Gaussian Mixture Model (GMM):

Probabilistic models and use the soft clustering approach for distributing the points in different clusters, it assumes that there are a certain number of Gaussian distributions and each of these distributions represent a cluster. therefore, a GMM tends to group the data points belonging to a single distribution together.

Gaussian distribution (or Normal distribution): is a type of continuous probability distribution for a real-valued random variable. The general form of its probability density function: $f(x) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$, with parameters: μ is mean of the distribution and σ is standard deviation.

To implement GMM, we use Expectation-Maximization algorithm (EM). EM is a statistical algorithm for finding the right model parameters. It typically used for handling missing values or in other words latent variables. This algorithm has two steps: E-step, the available data is used to estimate the values of the missing variables and M-step, based on the estimated values generated in the e-step, the complete data is used to update the parameters. EM in GMM: We need to assign k number of clusters. This means that there are k Gaussian distributions, with the mean and covariance values to be $\mu_1, \mu_2, \dots, \mu_k$ and $\sigma_1, \sigma_2, \dots, \sigma_k$. Additionally, there is another parameter for the distribution that defines the number of points for the distribution. Or in other words, the density of the distribution is represented with Π

2.2 Text Data Visualization

Data visualization is essential to assist businesses in quickly identifying data trends, which would otherwise be a hassle. The pictorial representation of data sets allows analysts to visualize concepts and new patterns. With the increasing surge in data every day, making sense of the quintillion bytes of data is impossible without Data Proliferation, which includes data visualization.

Every professional industry benefits from understanding their data, so [data visualization](#) is branching out to all fields where data exists. For every business, information is their most significant leverage. Through visualization, one can prolifically convey their points and take advantage of that information

2.2.1 Importance of Data Visualization

A dashboard, graph, infographics, map, chart, video, slide, etc. all these mediums can be used for visualizing and understanding data. Visualizing the data enable decision-makers to interrelate the data to find better insights and reap the importance of data visualization, which are:

Analysing the Data in a Better Way

Analysing reports helps business stakeholders focus on the areas that require attention. The visual mediums help analysts understand the key points needed for their business. Whether it is a sales report or a marketing strategy, a visual representation of data helps companies increase their profits through better analysis and better business decisions.

Faster Decision Making

Humans process visuals better than any tedious tabular forms or reports. If the data communicates well, decision-makers can quickly take action based on the new data insights, accelerating decision-making, and business growth simultaneously.

Making Sense of Complicated Data

Data visualization allows business users to gain insight into their vast amounts of data. It benefits them to recognize new patterns and errors in the data. Making sense of these patterns helps the users pay attention to areas that indicate red flags or progress. This process, in turn, drives the business ahead.

2.2.2 Tools used for Text Data Visualization

Data visualization tool provides users with an easier way to create visual representations of large data sets. When dealing with data sets that include hundreds of thousands or millions of data points, automating the process of creating a visualization, at least in part, makes a designer's job significantly easier.

ELK STACK

The ELK stack is an acronym used to describe a stack that comprises of three popular open-source projects: Elasticsearch, Logstash, and Kibana. ELK

stack gives us the ability to aggregate logs from all systems and applications, analyze these logs, and create visualizations for application and infrastructure monitoring, faster troubleshooting, security analytics, and more.

E = Elasticsearch

Elasticsearch is an open-source, RESTful, distributed search and analytics engine built on Apache Lucene. Support for various languages, high performance, and schema-free JSON documents makes Elasticsearch an ideal choice for various log analytics and search use cases.

L = Logstash

Logstash is an open-source data ingestion tool that allows you to collect data from a variety of sources, transform it, and send it to your desired destination. With pre-built filters and support for over 200 plugins, Logstash allows users to easily ingest data regardless of the data source or type.

K = Kibana

Kibana is an open-source data visualization and exploration tool for reviewing logs and events. Kibana offers easy-to-use, interactive charts, pre-built aggregations and filters, and geospatial support and making it the preferred choice for visualizing data stored in Elasticsearch.

We have implemented the dashboard with ELK to visualize the ticket data and generate as much insights from it.

Chapter 3

System's Architecture

3.1 Insights Dashboard Architecture

Raw Text data can be sent into ELK through Elasticsearch or Logstash. Logstash parses and the transforms the data and push it into Elasticsearch where the data is indexed. Kibana is used to create charts, graphs and dashboards. ELK is scalable and easier to deploy in both cloud and on-premise architecture. Below is the pictorial representation of implementation.

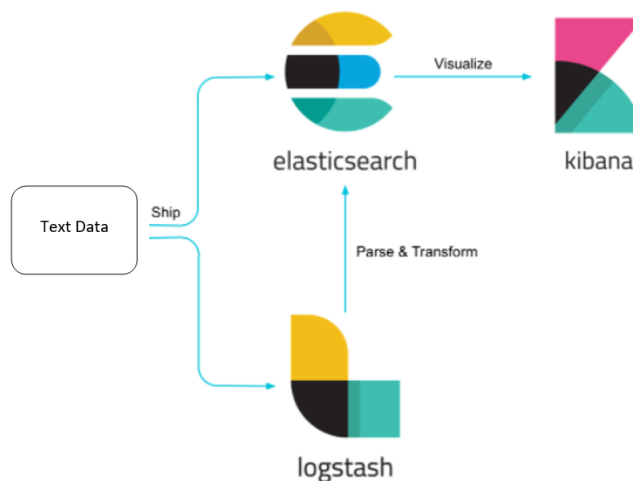


Figure 3.1: Dashboard Architecture

3.2 Clustering System Architecture

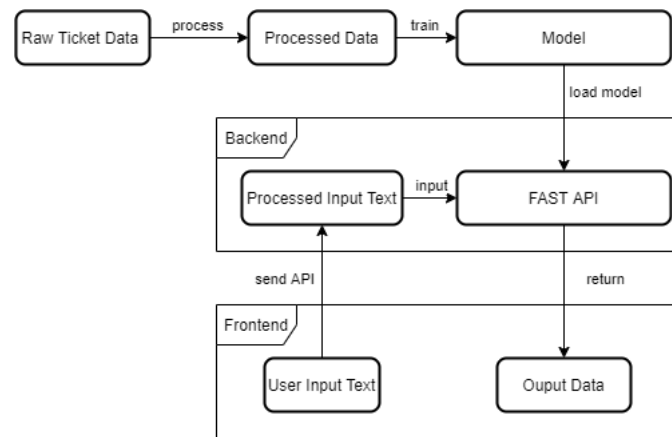


Figure 3.2: Clustering System Architecture

Chapter 4

Methodology

4.1 Data Visualization – ELK Dashboard

Data visualization is essential to assist businesses in quickly identifying data trends, which would otherwise be a hassle. The pictorial representation of data sets allows analysts to visualize concepts and new patterns. With the increasing surge in data every day, making sense of the quintillion bytes of data is impossible without Data Proliferation, which includes data visualization.

Every professional industry benefits from understanding their data, so data visualization is branching out to all fields where data exists. For every business, information is their most significant leverage. Through visualization, one can prolifically convey their points and take advantage of that information.

4.2 Clustering

Chapter 5

Results

Chapter 6

Conclusion

Chapter 7

References

Bibliography

- [1] [Wikipedia: Natural language processing](#)
- [2] [Article: iClass Gyansetu](#)