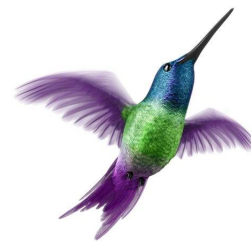# SOW Risk Scoring Model

Hafiza Umair, Khaled AlNajjar, Kiranpal Sohi, Ramona Torabian

April 7, 2020

# Introduction

**Objective:**

- Using NLP techniques, extract textual contents and clauses from SOW (Statement of Work) with TELUS suppliers/vendors
- Identify risk and label each SOW based on its risk level.
- Use unsupervised algorithm, namely clustering model to identify risk.

**Method:**

- Extract data from publicly available PDF SOWs downloaded from http://www.pdfsearchengine.net/
- Specific Data as required by Telus includes: 1) Work Description, 2) ETA or delivery date, 3) Headcount & rates, 4) Nature of resources, 5) Location of work, 6) Committed Funds, 7) Committed volume, 8) Currency (USD or CAD), 9) SOW Category (IT, fleet, field services, consulting, cleaning).
- Use extracted data to cluster the individual SOWs as High, Medium or Low risk.

**Limitations:**

- Documents were not of high quality as they were downloaded from internet. The content and subject of SOWs varied significantly depending on what was available freely online.
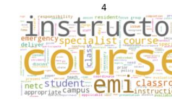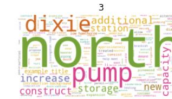
# Preparation

- A total of 48 PDF Statement of Works were downloaded from the Internet.
- The content varied from software to hospital system solution to Construction.
- A total of 30 PDFs were readable and uploaded to Google Colab notebook for code development.
- Used NLTK, and spaCy:
  - Clean up data - removing stopwords, lower case, punctuation removal
  - Lemmatization applied
- Built Data Corpus  - bag of words
- Built Document Term Matrix - DTM using CountVectorizer
  - Numerical representation of the text so that machine can read and process
- Analyzed common words in all SOWs, added more stopwords and rebuilt DTM

# Word Cloud Analysis

- Using spaCy library, the Word Clouds were created for analysis.
- Prominent words (high frequency) in Word Cloud provide quick understanding of each individual PDF SOW.
- The subject of SOWs varied from including: medical, furniture, education, technology.

# spaCy Word Label

- Using spaCy library, standard word labels were used to identify the following:
  - ETA or Delivery Date
  - Committed Funds
  - Location of Work
  - Committed Volume
- For the following key items the code is still under development:
  - Work Description
  - Headcount & Rates
  - Nature of Resources
  - Currency

# K-Mean Cluster

- The text content of PDFs was used to develop a clustering algorithm.
- The high risk SOWs based on Committed Funds ($s) is Cluster 0: Construction Industry based SOWs
- The value of Medium and Low risk SOWs were difficult to assess as the value of Committed Funds were either not in the publicly available SOWs or were difficult to extract for analysis.



FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead. [_testing.py:19]

**Clusters**

```
0 : water, aquifer, floridan, well, reclaim, mgd, phase, supply, wellfield, construction
1 : contractor, shall, task, work, service, eir, client, project, ltch, sfdc
2 : child, caf, vgisc, injury, care, mct, city, wcb, process, service
```

# Conclusion

- Word Cloud provided a visual analysis of each individual PDF Statement of Work. It can be used by Telus to quickly determine the keywords and subject matter of SOW.

- spaCy library provided a good means to clean up the data and allowed for extraction of approximately half of the key data from SOWs. This process can be used to extract data from Telus SOWs.

- K-Mean Clustering can be used to group SOWs into High, Medium and Low risk. With Telus, the volume of PDFs used for clustering can help enrich the clustering learning and provide a more robust clusters analysis.

# Improvement Opportunities

- Identifying the confidence intervals
- Using Custom Name Entity Recognition - NER to train spaCy model and better clean the data for further analysis and increase accuracy.
- K-Mean Cluster analysis requires cleaner and greater volume of data to ensure robust results.
- Additional PDFs are required and if possible with similar structure. Telus PDFs may be more consistent and therefore easier to extract data from.