# Spreadsheet Information Retrieval using NLP

Deepak Gami (170101020)
Devaishi Tiwari (170101021)
Soumik Paul (170101066)
Utkarsh Mishra (170101083)

Group 4

## 1   Problem Statement

Spreadsheet software is one of the most important tool for any organisation. Microsoft's Excel alone has an estimated 750 Million users. They are used to perform data-operations ranging from simple ones like search and filter to very complex ones like pivot tables and linear regression.

However, the plethora of operations provided by these software, makes spreadsheets complicated, especially for novice users. One has to have thorough knowledge about the menus and options within each menu, to be able to execute all set of operations. Apart from the learning complexity, the process of navigating menus is also time consuming.

### 1.1   Idea Description

We aim to solve this problem by creating an intelligent layer on top of a spreadsheet, which provides an alternative interaction method via Natural Language. Hence, instead of having to navigate menus, a user can simply describe his/her query in natural language and get the required results faster.

An open-sourced NLP enabled powerful spreadsheet can substantially democratize usage of complex statistical operations over spreadsheets.

### 1.2   Advantages

The advantages of our work could be:

- Faster & easier workflows for existing users of spreadsheet

- Reduced knowledge barrier to perform statistical analysis over spreadsheets

- By combining this with speech processing, spreadsheet software can be made usable by people with disabilities who have difficulties in handling mouse and keyboard

## 2   Major Challenges

As discussed in Microsoft's Research Paper [GM14], the complexity is enhanced due to two reasons:

1. Variations in Language on Same Task Variations in Task and Composition

2. Variations in Task and Composition

Apart from this, all queries are highly contextual with multiple references to column names, table name and/or data entries. Our methods have to take this into consideration.

## 3   Existing Works

- NLyze [GM14]: Researchers at Microsoft built an interface using a Programming by Natural Language (PBNL) methodology by developing a Domain specific language(DSL) which expresses the desired task and translation algorithm for converting natural language specifications into a ranked set of likely programs in the DSL.

- NLP-SIR[FMM09] : This research has outlined the effectiveness and efficiency of spreadsheet information retrieval using a Natural language interface which is uses NLP techniques to interpret the commands.

Figure 1: Examples of sample sentences that showcase the variations

# 4 Proposed Direction

After some preliminary readings, we have designed the following method:

1. Generate a set of prototypical Natural Language Queries to perform the operations. These NLQ's have to be abstract, so that they can be replaced with contextual parameters later. Eg: For operation "FILTER", a prototypical Natural Language Query can "List all <table_name> that have <col_name> greater than <data_entry>". Terms in <> will be replaced as per the queries context. Let the set of NLQ's be $S$. Our set will be tuple with each entry being {NLQ, operation}.

2. Now given a query(say $q$), get the contextual information and replace all the abstractions ( ie <> terms) in $S$ to get $S_q$

3. Next, find the NLQ from $S_q$ which is semantically closest to $q$. This will tell us what operation needs to be performed.

4. Given the operation & contextual information, we can execute and return the results.

This is just the basic idea, and we need to test it before finalizing. Our next step would be to pick a small subset of operations ( 1 or 2 simple ones like FILTER, ADD ) and test this method. After we have a good method to perform basic operations, we will try to enhance it to handle complex operations (eg: ADD + FILTER).

# References

[FMM09]   Derek Flood, Kevin McDaid, and Fergal McCaffery. "NLP-SIR: A Natural Language Approach for Spreadsheet Information Retrieval". In: *CoRR* abs/0908.1193 (2009). arXiv: 0908.1193. http://arxiv.org/abs/0908.1193.

[GM14]    Sumit Gulwani and Mark Marron. "NLyze: Interactive Programming by Natural Language for SpreadSheet Data Analysis and Manipulation". In: *SIGMOD'14, June 22-27, 2014, Snowbird, UT, USA*. June 2014. https://www.microsoft.com/en-us/research/publication/nlyze-interactive-programming-natural-language-spreadsheet-data-analysis-manipulation/.