

Knowledge Graph For Researchers

Kyle Vincent and Emma Meno CS 4624: Multimedia, Hypertext, and Information Access Dr. Edward Fox Virginia Tech, Blacksburg, VA 24061 April 28, 2020



OUTLINE

| | Project Deliverables | [3] |
|---|------------------------|------|
| • | Accomplishments | [5] |
| • | Testing and Evaluation | [20] |
| • | Lessons Learned | [22] |
| | Future Work | [26] |
| | Acknowledgements | [28] |
| | References | [29] |



Deliverables



Project Deliverables

- Build ontology for Twitter-based project info an SME might want to query, answer or calculate
- Build knowledge graph based on created ontology
- Build set of APIs to trigger set of network algorithms based on info queried to graph





Accomplishments

Literature Review: Twitter Data Source

A VTechWorks Home / Student Works / CS5604: Information Retrieval

2016 🐼
2017 🐼

CS5604: Information Retrieval

| 3ROWSE BY | | | | | | | | |
|---------------|---------|--------|----------|--|--|--|--|--|
| By Issue Date | Authors | Titles | Subjects | | | | | |

Search within this collection:

Go

This collection contains the final projects of the students in various offerings of the course Computer Science 5604: Information Retrieval. This course is taught by Professor Ed Fox. Analyzing, indexing, representing, storing, searching, retrieving, processing and presenting information and documents using fully automatic systems. The information may be in the form of text, hypertext, multimedia, or hypermedia. The systems are based on various models, e.g., Boolean logic, fuzzy logic, probability theory, etc., and they are implemented using inverted files, relational thesauri, special hardware, and other approaches. Evaluation of the systems' efficiency and effectiveness.

Built Ontology

| | A | В | С | D | E | F | G | н | 1 | J | ĸ | L | М |
|----|------------|--|--|--|----------------------|--|--------------|--------------------|---|-------------------|--------------------|-----------------|--|
| 1 | Task ID | Task Name | Input File(s) | Output File(s) | Dependent Task ID | Functions / Libraries | Manual (Y/N) | Data Collection | Command Line Parameters | Project | | | Report Names |
| 2 | 201 | 16 Reports (CS 5604 | 4: Information Retrieva | I) | | | | | | | | | |
| 3 | | Remove stop words and 0 lemmatize | Cleantext from HBase (RDD of tweets) | Cleaned Lemmatized Tweets | | StanfordNLP | Y | | table name data type for each row-key desired columns | https://vtechworl | ks.lib.vt.edu/hand | lle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report |
| 4 | | 1 Manually Classify Tweets | Cleaned Lemmatized Tweets AND Set of Events Classes | Tweet Training Data | c |) N/A | Y | | collection ID batch size | https://vtechwor | ks.lib.vt.edu/hand | lle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report |
| 5 | | 2 Feature Selection | Tweet Training Data | Word Feature Representation | 1 | Word2Vec Model | N | | | https://vtechwor | ks.lib.vt.edu/hand | lle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report |
| 6 | : | 3 Feature Selection | Tweet Training Data | Word Feature Representation | 1 | Association Rules Based Classifier (provided by Dr. Pereira) | N | | path_to_jar_file support threshold block size cluster node base_dir_on-HDFS path_to_training_file path_to_test_file stopwords_file output_directory | https://vtechwori | ks.lib.vt.edu/hand | lle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report |
| 7 | | 4 Tweet Classification | Word Feature Representation | Predicted Class of the Given Tweet // Associated Probabilities of Each Class | 2,3 | Multi-class logistic regression classifier | N | | collection number batch size retrain metric <training file=""> <test file=""></test></training> | https://vtechwor | ks.lib.vt.edu/hand | lle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report |
| 8 | | 5 Convert to CSV | Cleaned Tweets in HBase | CSV Format Tweets | 37 | Pig Script | N | | | https://vtechworl | ks.lib.vt.edu/hand | lle/10919/73739 | CS5604: Information and Storage Retrieval Fall 2016 - CMT (Collection Management Tweets) |
| 9 | | Record User, URL, and 6 Tweet Relationships | CSV Format Tweets | Social Network Matrix | 5 | i | Y | | | https://vtechwori | ks.lib.vt.edu/hand | lle/10919/73739 | CS5604: Information and Storage Retrieval Fall 2016 - CMT (Collection Management Tweets) |
| 10 | | Calculate Importance 7 Factors for Nodes | Social Network Matrix | Weighted Social Network Matrix | 6 | 5 | Y | | | https://vtechwori | ks.lib.vt.edu/hand | lle/10919/73739 | CS5604: Information and Storage Retrieval Fall 2016 - CMT (Collection Management Tweets) |

Nodes.csv

| | A | В |
|---|--------|--|
| 1 | fileId | name |
| 2 | 1 | Cleantext from HBase (RDD of tweets) |
| 3 | 2 | Cleaned Lemmatized Tweets |
| 4 | 3 | Set of Events Classes |
| 5 | 4 | Tweet Training Data |
| 6 | 5 | Word Feature Representation |
| 7 | 6 | Predicted Class of the Given Tweet |
| 8 | 7 | Associated Probabilities of Each Class |
| 9 | 8 | Cleaned Tweets in HBase |



Edges.csv

| rds and 1 / Tweets 2 | outputId 2 4 | functionsAndLibraries StanfordNLP N/A | Y | commandLineParameters table name; data type for each row-key; desired columns | reportUrl https://vtechworks.lib.vt.edu/handle/10919/73713 | reportName CS5604 Fall 2016 Classification Team Final Report | domainCollection Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur, Hurricane Sandy, Hurricane Isaac, China Factory Explosion, Texas Plant Explosion, Manhattan Explosion Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur. |
|------------------------------|--|---|----------------|---|--|--|--|
| rds and 1 1 / Tweets 2 | 2 | StanfordNLP | Y | table name; data type for each row-key; desired columns | https://vtechworks.lib.vt.edu/handle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report | Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur, Hurricane Sandy, Hurricane Isaac, China Factory Explosion, Texas Plant Explosion, Manhattan Explosion Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur. |
| / Tweets 2 | 4 | N/A | | | | | Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur, |
| | | | Y | collection ID; batch size | https://vtechworks.lib.vt.edu/handle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report | Hurricane Sandy, Hurricane Isaac, China Factory Explosion, Texas Plant Explosion, Manhattan Explosion |
| Tweets 3 | 4 | N/A | Y | collection ID; batch size | https://vtechworks.lib.vt.edu/handle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report | Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur, Hurricane Sandy, Hurricane Isaac, China Factory Explosion, Texas Plant Explosion, Manhattan Explosion |
| n 4 | 5 | Word2Vec Model | N | | https://vtechworks.lib.vt.edu/handle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report | Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur, Hurricane Sandy, Hurricane Isaac, China Factory Explosion, Texas Plant Explosion, Manhattan Explosion |
| n 4 | 5 | Association Rules Based Classifier (provided by Dr. Pereira) | N | path_to_jar_file; support threshold; block size; cluster node; base_dir_on-HDFS; path_to_training_file; path_to_test_file; stopwords_file; output_directory | https://vtechworks.lib.vt.edu/handle/10919/73713 | CS5604 Fall 2016 Classification Team Final Report | Kentucky Shooting, Newton Shooting, Firefighter Shooting, Hurricane Arthur, Hurricane Sandy, Hurricane Isaac, China Factory Explosion, Texas Plant Explosion, Manhattan Explosion |
| | Tweets 3 1 4 1 4 | Tweets 3 4 | Tweets 3 4 N/A | Tweets 3 4 N/A Y A 4 5 Word2Vec Model N Association Rules Based Classifier (provided by Dr. (provided by Dr. (provided by Dr. N | Tweets 3 4 N/A Y collection ID; batch size 1 4 5 Word2Vec Model N 1 4 5 Word2Vec Model N 1 4 5 Word2Vec Model N 1 Association Rules Based Classifier (provided by Dr. 4 5 path_to_jar_file; support threshold; block size; cluster node; base_dir_on-HDFS; path_to_training_file; path_to_test_file; stopwords_file; output_directory | Tweets 3 4 N/A Y collection ID; batch size https://vtechworks.lib.vt.edu/handle/10919/73713 1 4 5 Word2Vec Model N https://vtechworks.lib.vt.edu/handle/10919/73713 1 4 5 Word2Vec Model N https://vtechworks.lib.vt.edu/handle/10919/73713 1 4 5 Word2Vec Model N path_to_jar_file; support threshold; block size; cluster node; base_dir_on-HDFS; path_to_test_file; atopwords_file; output_directory https://vtechworks.lib.vt.edu/handle/10919/73713 | Tweets 3 4 N/A Y collection ID; batch size https://techworks.lib.vt.edu/handle/10919/73713 CS5604 Fall 2016 Classification Team Final Report 1 4 5 Word2Vec Model N https://techworks.lib.vt.edu/handle/10919/73713 CS5604 Fall 2016 Classification Team Final Report 1 4 5 Word2Vec Model N https://techworks.lib.vt.edu/handle/10919/73713 CS5604 Fall 2016 Classification Team Final Report 1 Association Rules Based Classifier (provided by Dr. 5 path_to_jar_file; support threshold; block size; cluster node; base_dir_on-HDFS; path_to_itest_file; output_directory https://techworks.lib.vt.edu/handle/10919/73713 CS5604 Fall 2016 Classification Team Final Report |

Dataset Algo Dataset

Path Metrics



Interim Project Toolkit - The GRANDstack



Implemented Project Toolkit - Grakn

- Database storage/interface
- Graql queries for reads, writes, & analysis
- Three Client APIs:
 - Node.js
 - Python
 - o Java





The And/Or Problem

- With traditional node-edge-node relationships, we weren't sure how to differentiate between two scenarios present in our data:
 - A certain output requires Input 1
 AND Input 2 to be produced
 - You can use either Input 1 OR Input
- This means the difference between returning 2 paths and 1 path



Solution: The Hypergraph

- Recommended by our Client
- A hyperedge can connect any number of entities
- Now we can distinguish between hyperedges and normal edges



Grakn Database Schema

- Our client recommended a tool called Grakn.
 - Designed for Knowledge Graphs specifically
 - Open Source
 - Allows Hyper-Relationships
- Sample of our Schema:

defin

#Our addaset nodes
file sub entity,
 key fileId,
 has name,
 plays inputfile,
 plays outputfile;

task sub relation, abstract,

key taskId, has name, has inputId, has outputId, has functionsAndLibraries, has functionsAndLibraries, has commandLineParameters, has reportUrl, has reportVame, has domainCollection, relates inputfile, plays ancestor, plays decendent;

#points from input to outpu produces sub task, relates inputfile, relates outputfile;

need-two-inputs sub task, relates outputfile, relates inputfile;

Migration Scripts

- With schema established, starting migrating the data
- Grakn has three clients:
 - Python
 - o Java
 - Node.js
- Wrote data migration scripts that could use these clients and migrate from CSV \rightarrow Grakn DB
- This proved way more challenging than it had been with Neo4j



Possible User Interface: Console

- One possible way to query the data is the Grakn Console
- It is included with the download of the software
- Uses the Graql language.



Grakn Workbase UI





Testing and Evaluation

Testing and Evaluation

- First test to be conducted by client shortly after project submission
- Next test will be to have a subject matter expert use both the console and the Workbase UI





Lessons Learned

Challenges Faced

- Waiting on network science data
- Gathering sufficient material from project reports
- Selecting an appropriate project toolkit
- Transition to virtual learning



Solutions



- Focus on different data \rightarrow Twitter
- Including reports from multiple years → 2016 & 2017
- Shift to Grakn
- Use of Zoom and other online platforms







Overall Lessons

- Asking more questions earlier in the design process
- Communicating with the professor, client, and team
- Documentation and project tools





Future Work

Future Work

- Upload the rest of the data to the graph
- Build upon schema to aid better path queries
- Use Node.js Cient to establish a website as front end
- Automate process of parsing through reports



Acknowledgements

Client: Prashant Chandrasekar

- Fifth-year PhD student
- Works with Dr. Fox
- Our Work: Proof of Concept for PhD



Professor: Dr. Edward Fox

 Director of Digital Library Research Laboratory



References

"Grakn.AI – Text Mined Knowledge Graphs",

https://www.stockwerk.co.at/event/grakn-ai-textmined-knowledge-graphs/

"Grakn Schema Overview", <u>https://dev.grakn.ai/docs/schema/overview</u>

Directed Hypergraph and Applications: https://www.sciencedirect.com/science/article/pii/0166218X9390045P

Grakn Documentation: <u>https://grakn.ai/</u>



Questions?