

E6893 Big Data Analytics:

Delving into the Q&A network – text mining and graph analysis

Zhen Liang, Xinli Wang
ZI2406, xw2341



December 17, 2015

- **Project Overview**
- **Motivation**
- **Tools and Data Sets**
- **Algorithms**
- **Performance**
- **Future Work**

- Stack Exchange is a Q&A platform where software engineers, scientists, students share knowledge and get questions answered.
- As users, we are interested in:
 - What are heated discussed topics
 - How to filtering best answers among all the given answers

As developers, we are interested in:

-The problems users are facing and how they can take such information to improve their products and documentation.



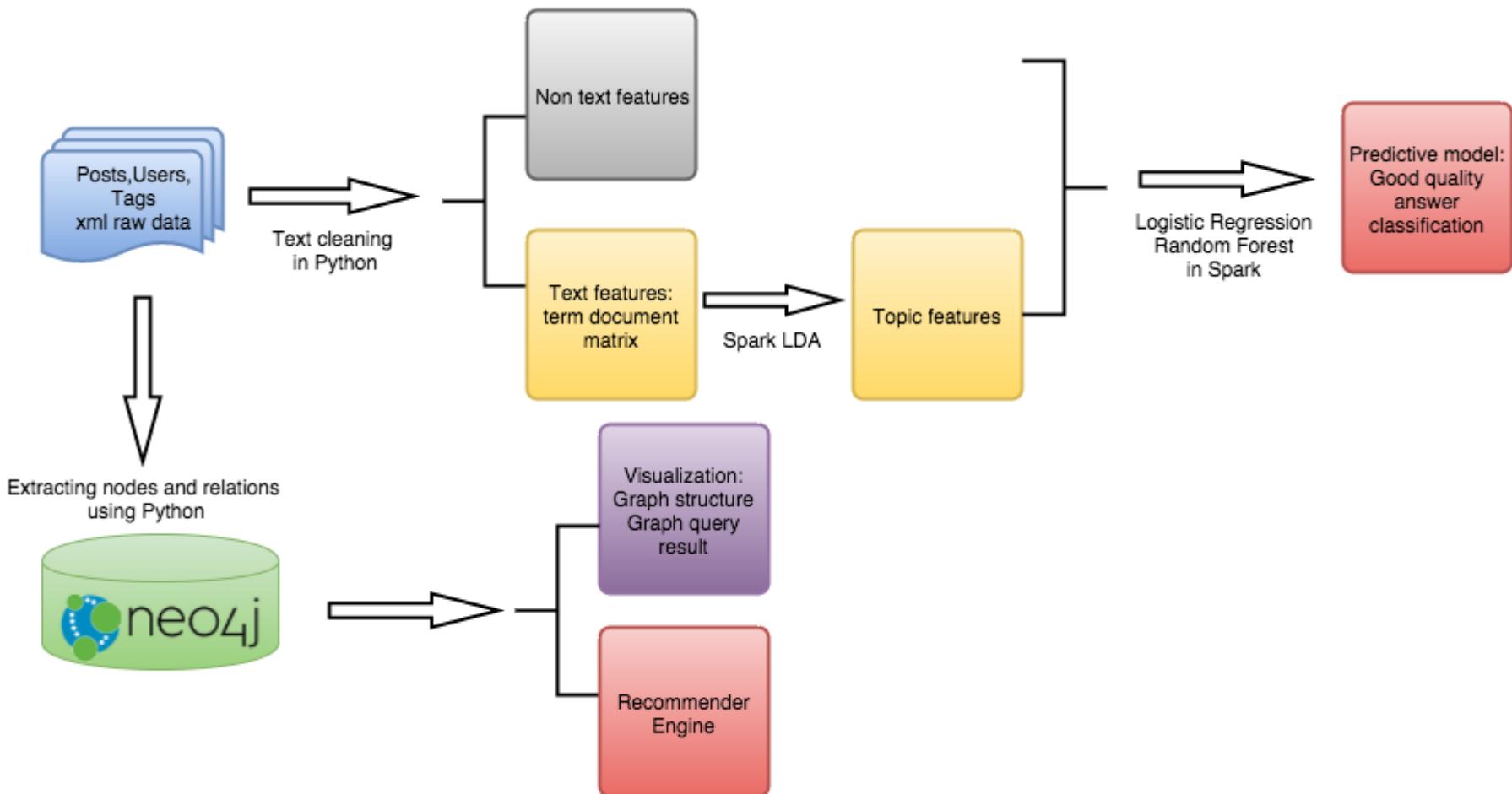
StackExchange

Our project addresses such problems by

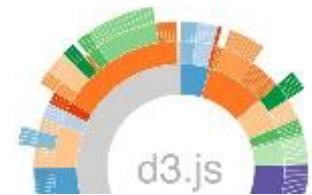
- Extracting topics out of large amount of posts and the topic distribution of each document
- Predicting the best answers by building a classification model
- Visualizing the “network” of questions, to know what’s the trends and relationships among discussed topics

Project Overview

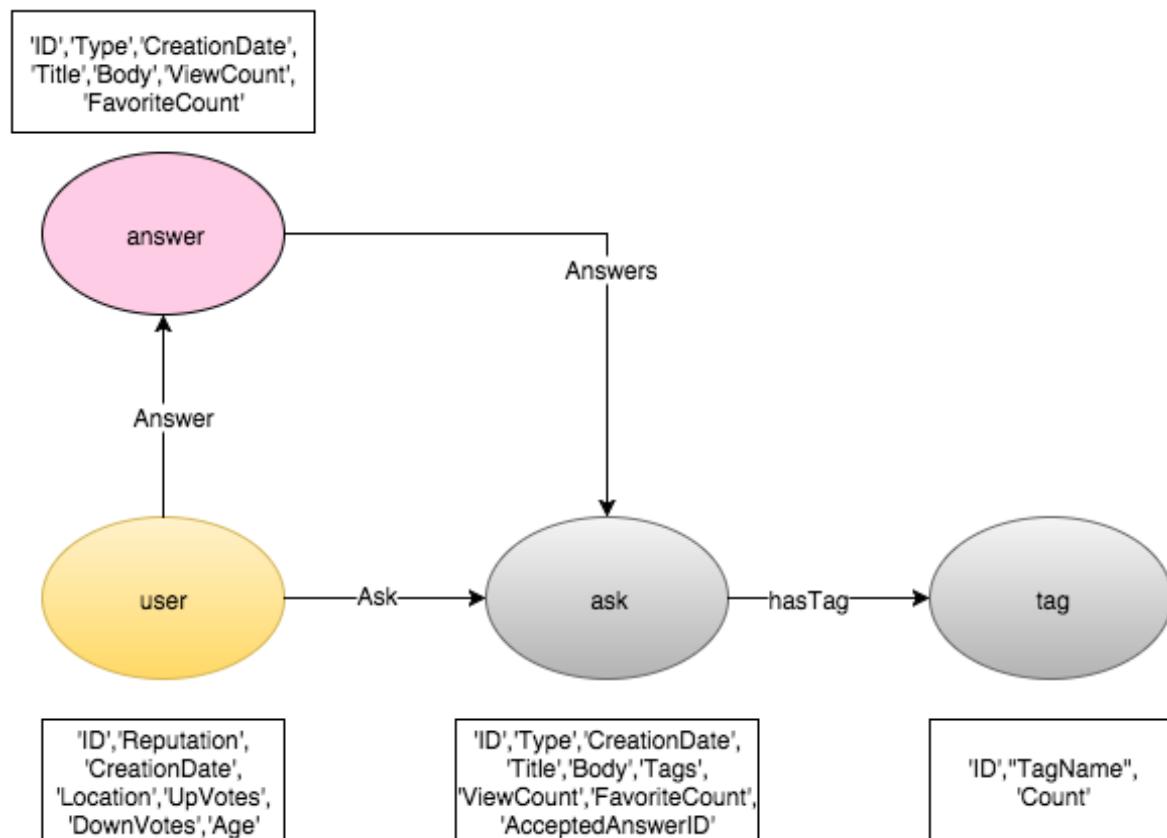
- Text mining
- Graph database and visualization



- Python
- Spark
- Cypher
- D3 and JavaScript

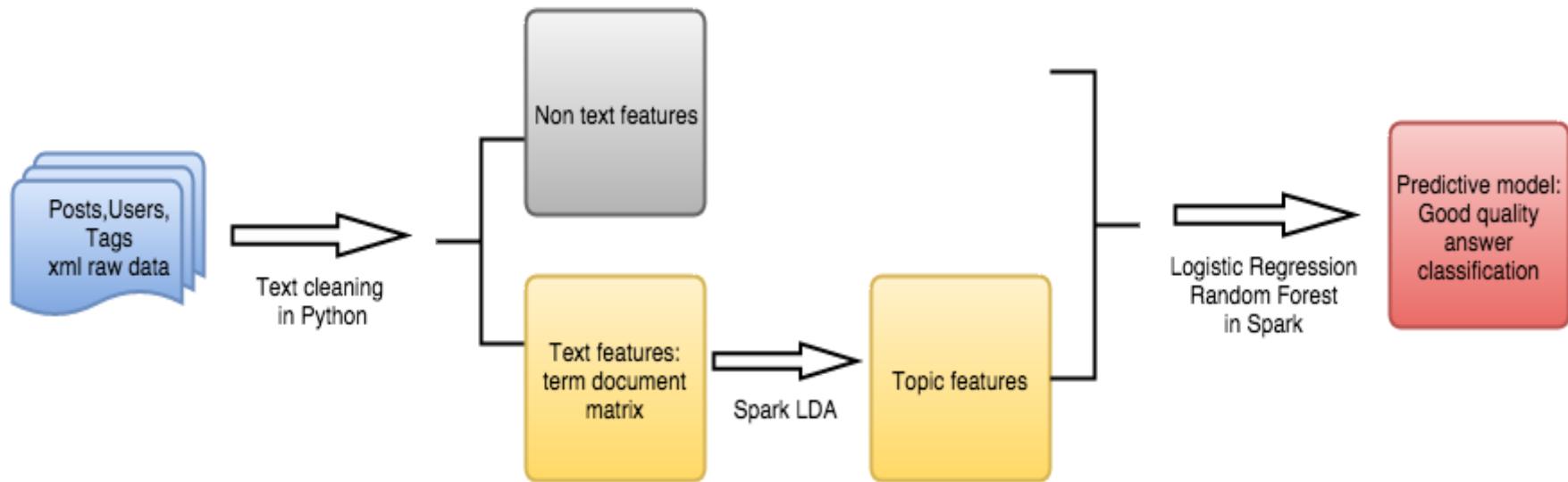


- Q&A data from Stack Exchange Data Dump, we choose the data science categories to conduct our analysis.
- Containing information of users, posts, and users-posts interaction.
- Specifically, information of posts, their links, votes tags, history, users information

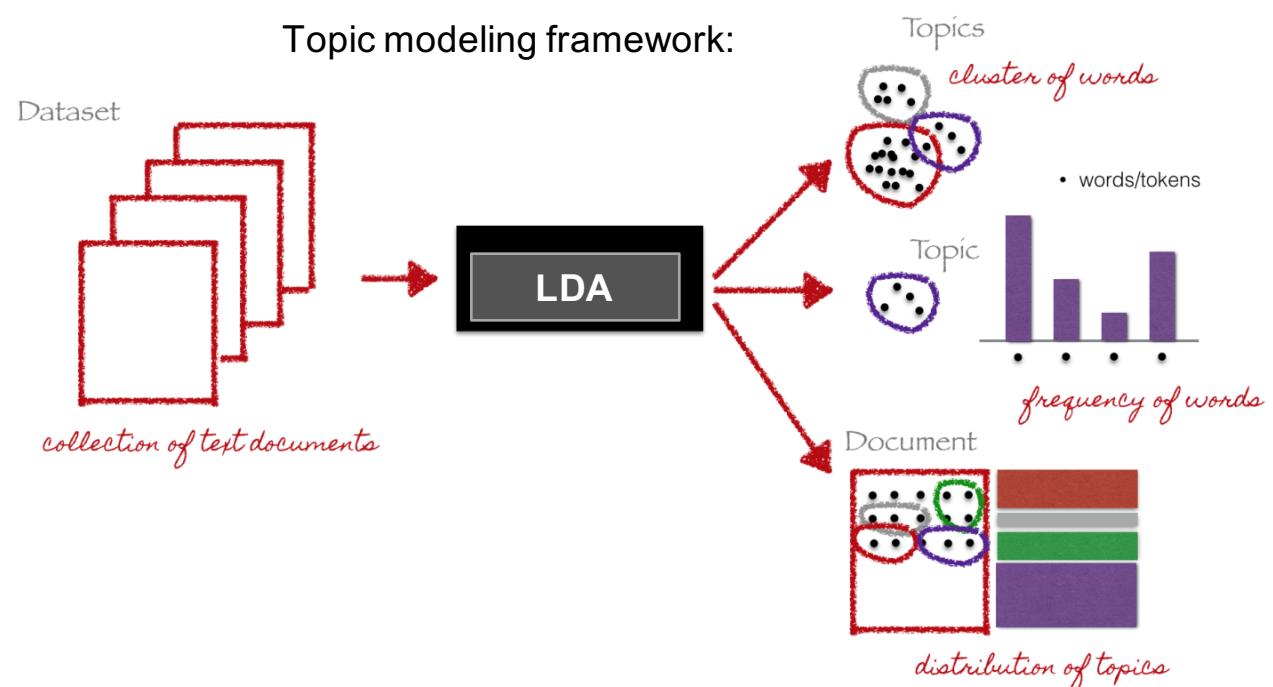


Textmining and Classification:

- Algorithms for data and Video pre-processing.
- Event-related field/response analysis.
- Classification and Collaborative Filtering.

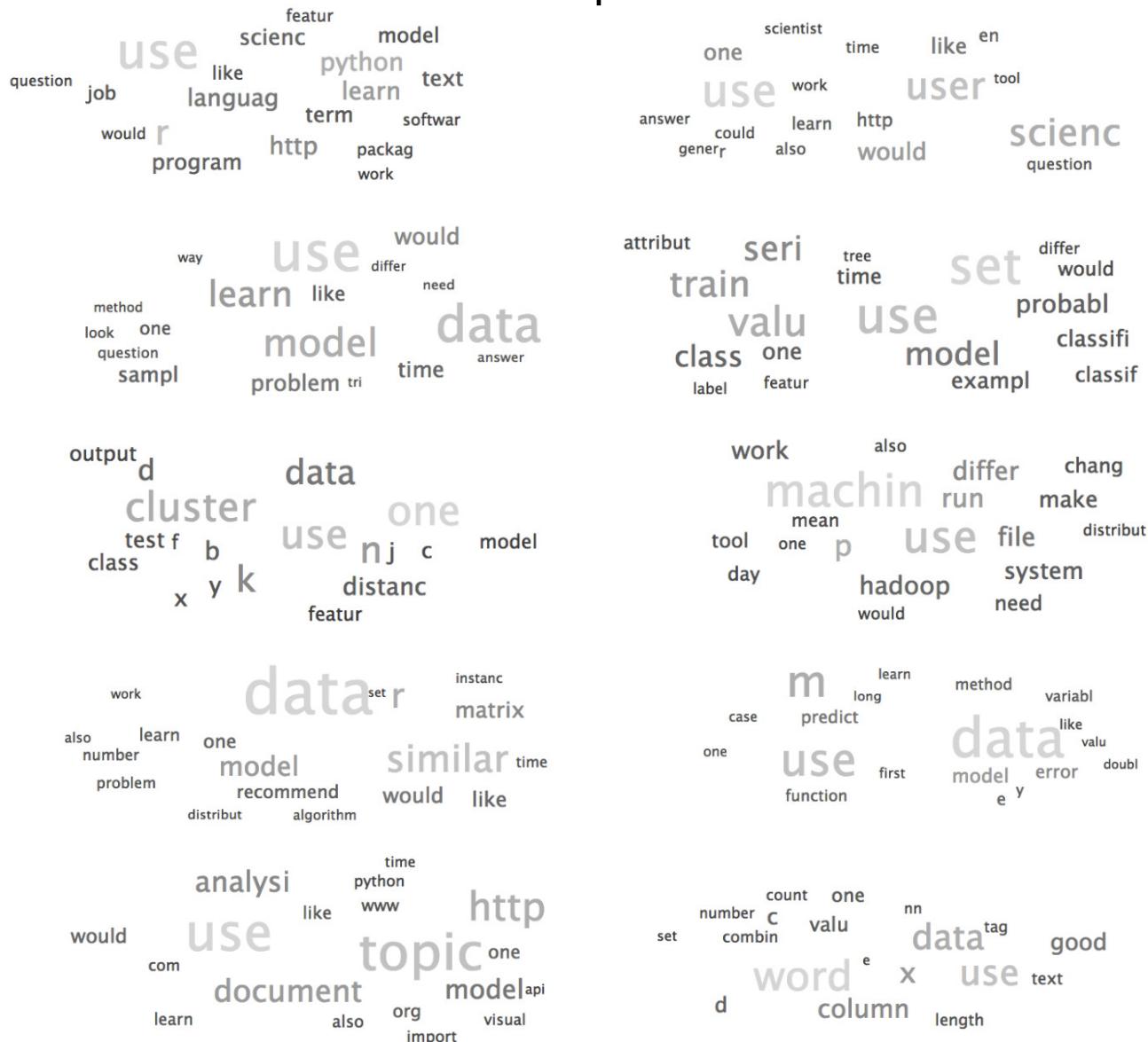


- Clean text
- Term document matrix
- Run Spark LDA using Python API to get an arbitrary number of topics (e.g. 10 topics)
- Run Spark LDA using Scala API to get the topic distribution on each document



Text Feature Generation: LDA result

Top 10 word distribution on each topic:



Some good quality answers are accepted by the people who raises the questions. Building a model to predict what answers are going to be accepted, so as to predict the good quality answers:

Objective:

Predicting good quality answers is useful for

1. ranking the answers
2. understanding the successful features that a good quality answer have so as to provide the direction to improve the overall answers quality.

Definition: Accepted answers are defined as best answers

There are two classes of answers in our dataset, accepted vs. not accepted.

Pre-processing data: delete accepted answers that are given and accepted by the same users or that have been edited after being accepted.

A. Presentation Quality

- Length (count of words in answers' posts)
- URL count
- Score = upvotes - downvotes
- Comment Count

A. Attitude

Scores are returned to describe intensity emotional style of a post. The results are between -1 and 1.

alchemyAPI is used for sentiment analysis

A. Time

The waiting time between a question posted and an answer is given (in seconds)

A. Reputation

Answerer's Reputation

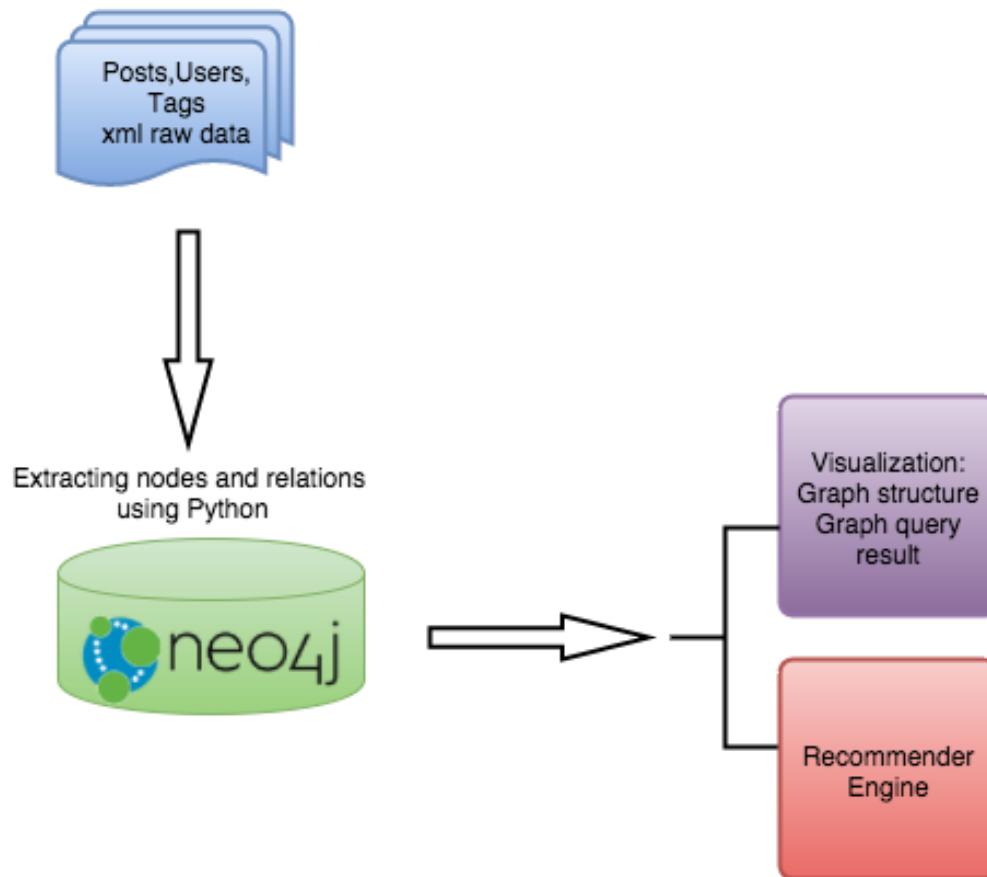
E. Topic features

1. Splitting dataset into train and test (4: 1) randomly
2. Applying Random Forest on all features
3. Classification is applied on subsets of training and test sets by removing one of features at a time.

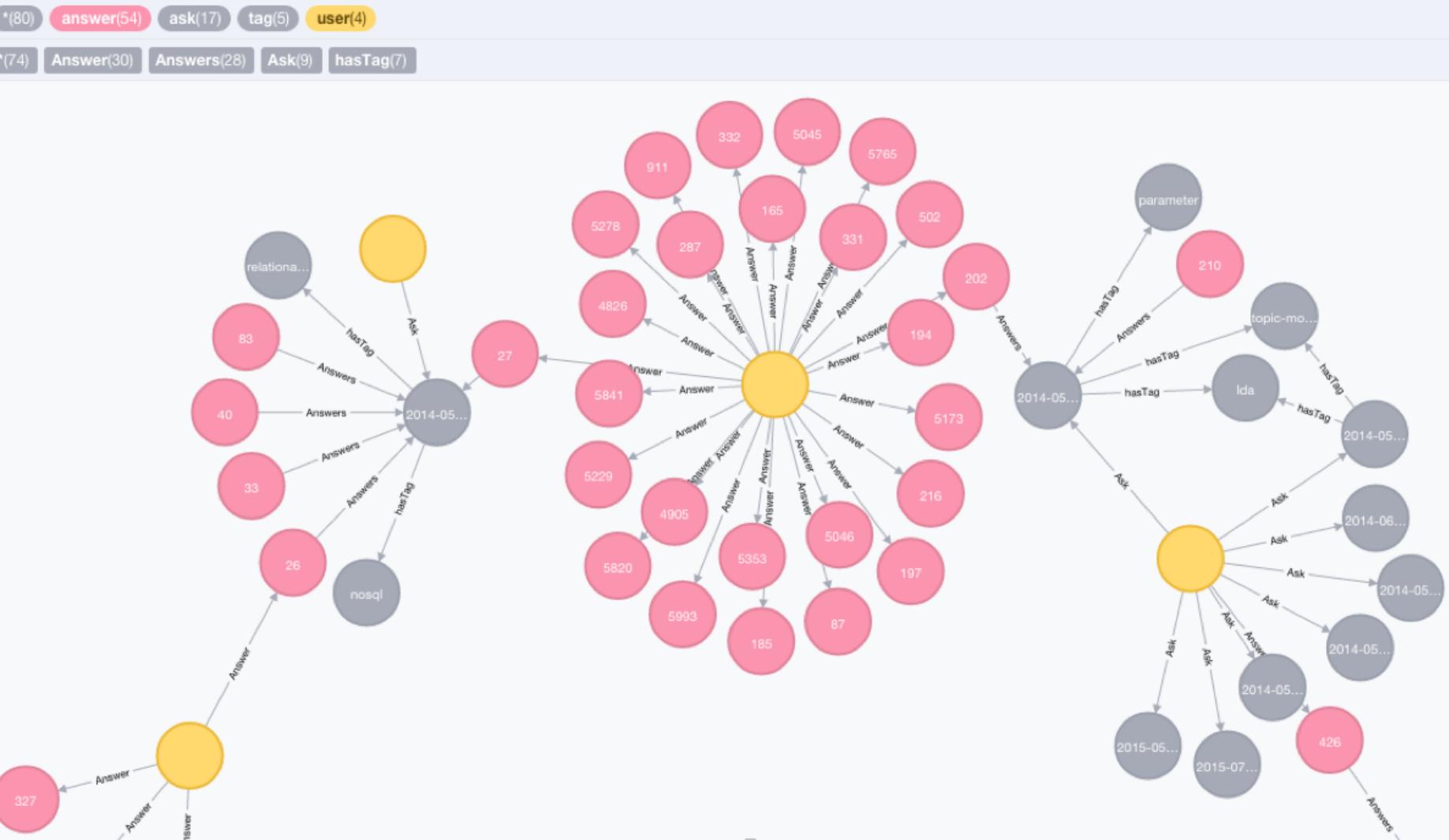
	AUC	F	PREC	REC
All Factors	0.60	0.77	0.80	0.81
w/o Reputation	0.62	0.79	0.82	0.83
w/o Quality	0.58	0.76	0.76	0.80
w/o Time	0.52	0.70	0.73	0.78
w/o Topic Features	0.73	0.85	0.85	0.86
w/o Attitude	0.64	0.79	0.80	0.82

Table 1 Prediction Results for Random Forest Models Using Different Sets of Features

- Graph database and visualization



Neo4j Graph Database - a user-post network



next we are going to extract data from Neo4j and draw a network graph of tags using Force-Directed Graph

In [17]:

Slide Type -

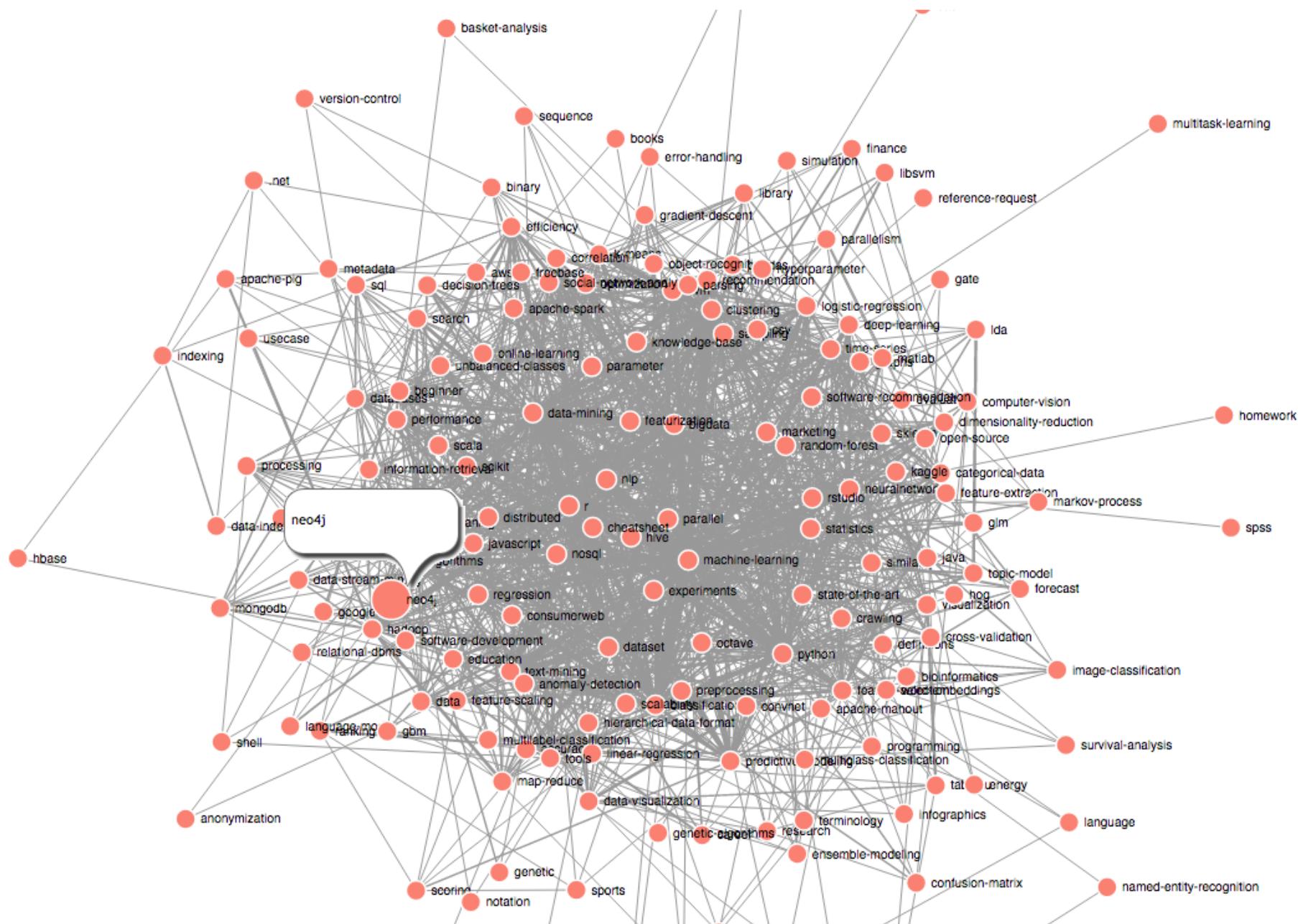
```
%load_ext cypher
# get all pairs of tags what appear in the same documents
query="""
MATCH (u:tag)<-[ :hasTag]-(p:ask)-[:hasTag]->(q:tag)
Where toInt(q.ID)>toInt(u.ID)
RETURN u.TagName, qTagName, count(*) AS tag_cooccurrence
ORDER BY tag_cooccurrence DESC
"""

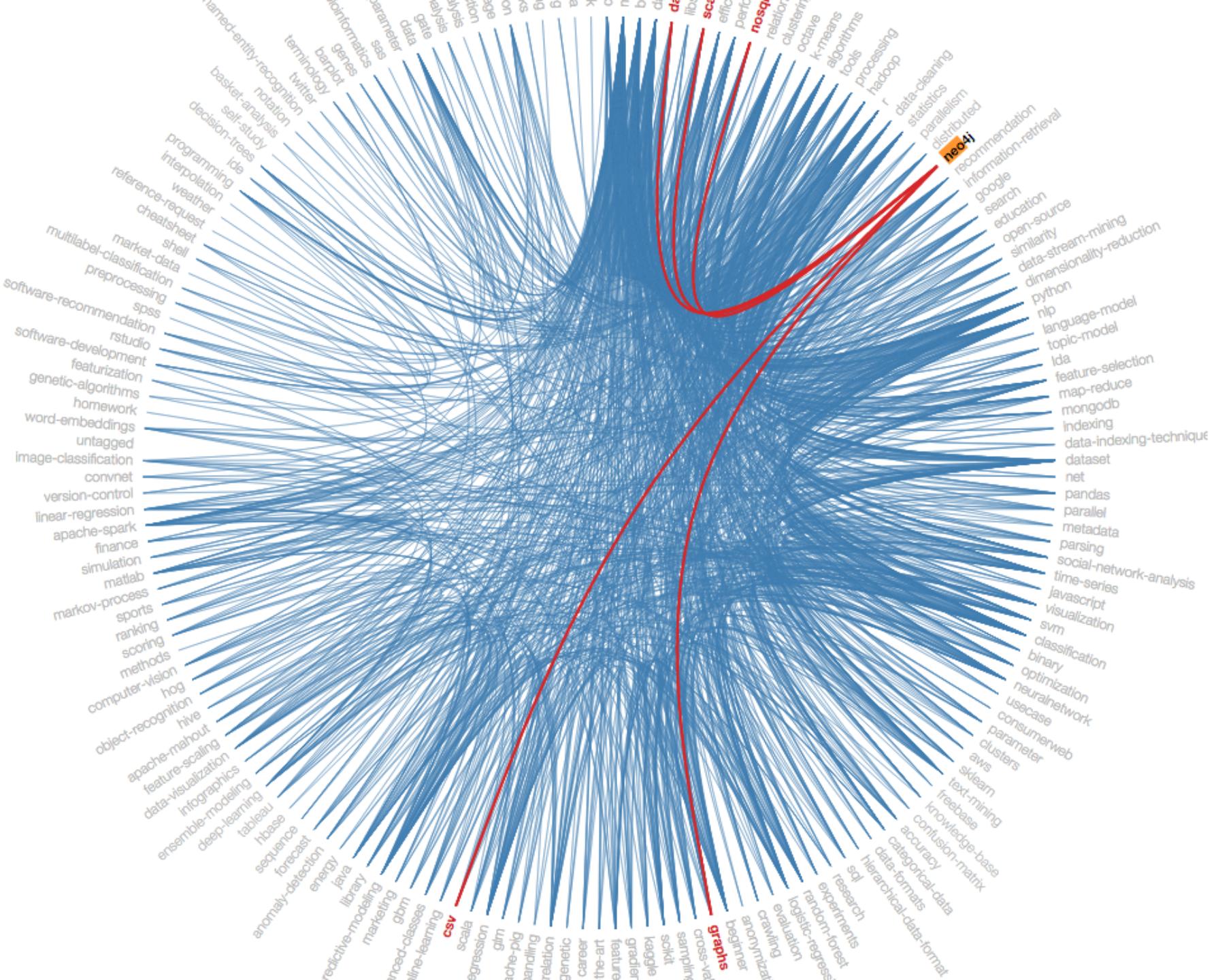
tag_cooccurrence = graph.cypher.execute(query)
tag_cooccurrence
```

Out[17]:

	u.TagName	q.TagName	tag_cooccurrence
1	machine-learning	classification	76
2	machine-learning	data-mining	69
3	machine-learning	r	48
4	machine-learning	python	34
5	machine-learning	neuralnetwork	34
6	machine-learning	statistics	30
7	machine-learning	algorithms	27
8	bigdata	data-mining	26
9	machine-learning	bigdata	26
10	machine-learning	predictive-modeling	26
11	machine-learning	dataset	25
12	machine-learning	nlp	25
13	machine-learning	clustering	23
14	bigdata	hadoop	22
15	data-mining	classification	22
16	data-mining	clustering	21
17	data-mining	dataset	21

Tag Force-directed graph





For those questions which do not have answer yet, return the users who answered questions with the same tags before.

```
query = """MATCH (q:ask)-[:hasTag]->(t:tag)<-[ :hasTag]-(:ask)<-[ :Answers]-()<-[ :Answer]-(u:user)
Where not (q)<-[ :Answer]-()
RETURN q.ID, u.ID, t.TagName"""
recommend = graph.cypher.execute(query)
recommend
```

Out[37]:

	q.ID	u.ID	t.TagName
1	16	51	machine-learning
2	61	51	machine-learning
3	115	51	machine-learning
4	116	51	machine-learning
5	159	51	machine-learning
6	169	51	machine-learning
7	196	51	machine-learning
8	211	51	machine-learning
9	215	51	machine-learning
10	265	51	machine-learning
11	266	51	machine-learning
12	310	51	machine-learning
13	319	51	machine-learning
14	323	51	machine-learning
15	326	51	machine-learning
16	369	51	machine-learning
17	406	51	machine-learning
18	410	51	machine-learning

Future Work

- Improve topic models, try models that can incorporate the name entity and meta data of each documents.
- Implement more classification algorithms and do a comparison between different methods
- Improve the recommender system

Questions?