# Modeling Community Question-Answering Archives

Zainab Zolaktaf, Fatemeh Riahi, Mahdi Shafiei and Evangelos Milios
Email : z.zolaktaf@gmail.com
Computer Science Department, Dalhousie University

NSERC Business Intelligence Network

DALHOUSIE UNIVERSITY
*Inspiring Minds*

## 1 Community Question-Answering Archives

### Abstract

Community Question Answering (CQA) services contain large archives of previously asked questions and their answers. We present a statistical topic model for modeling Question-Answering archives. We refer to our model as Question-Answering Topic Model or QATM. The model explicitly captures relationships between questions and their answers by modeling topical dependencies. To analyze the model, we use it for Question Answering and Automatic Tagging. Our model achieves improved performance in retrieving the correct answer for a query question compared to the LDA model. In addition, for large numbers of topics, our model achieves better clustering performance in terms of F-measure.

### Example Question-Answer Pair

**Question:**
How do I replace all occurrences of a word in a document with another word? Any solution is welcome!

**Answer:**
Regular expressions(regex) allow you to search and manipulate text based on patterns. In some languages, standard string manipulation functions are available e.g. replaceAll() method from Java's string class. In other languages, such as Perl, regex's are integrated into the language itself. Utilities such as grep, vi, etc can also perform this type of pattern matching and replacement.

### Properties of Question-Answers

❖ Answer topics influenced by question topics

❖ Answer topics more technical and specific

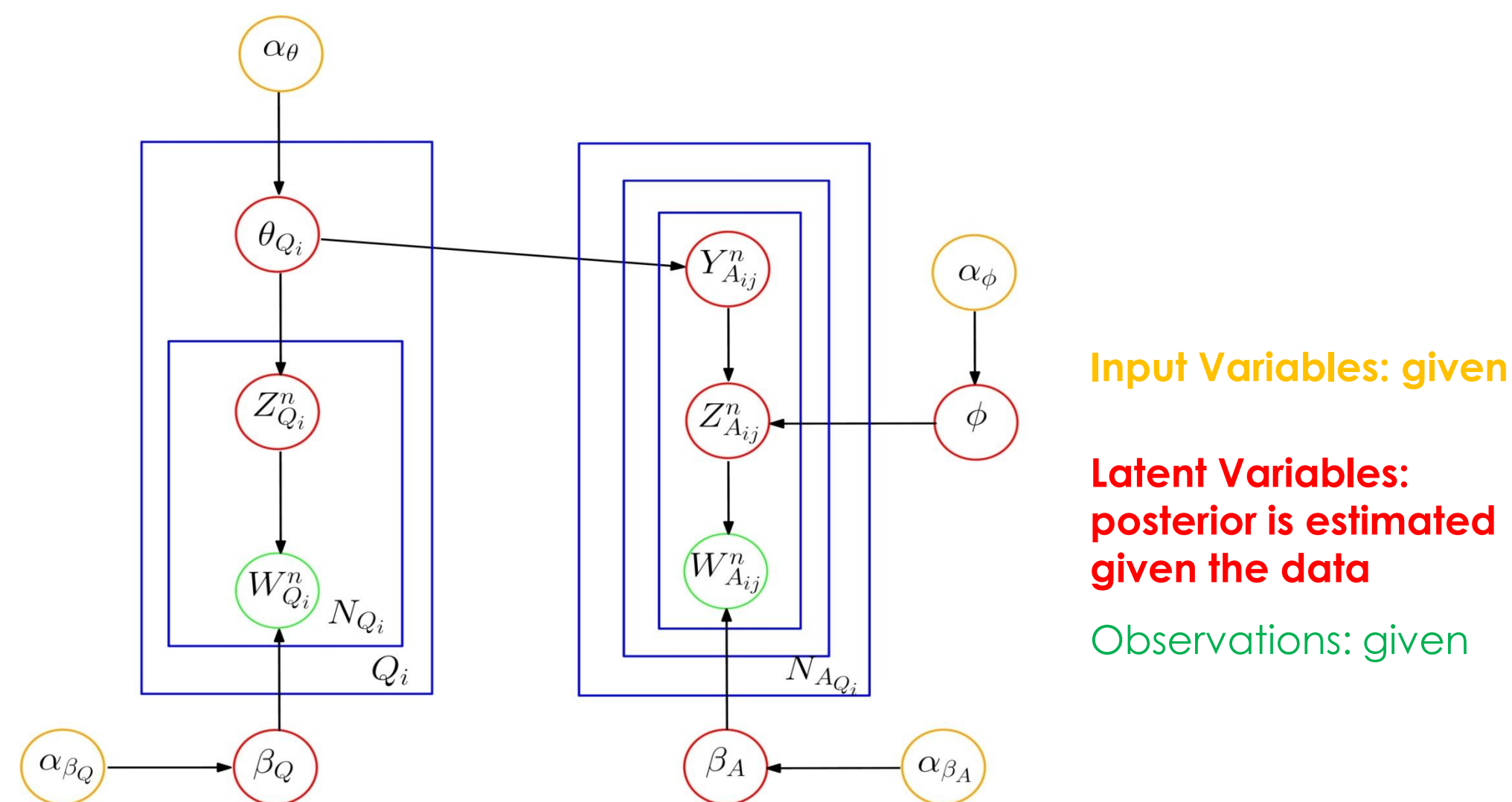❖ Answers may contain additional topics that are correlated with question topics

### Goal

Discovering the topical structure of CQA archives

❖ Automatic Tagging: Annotating the content with topics discovered to help browse and understand the archive.

❖ Question Answering: Given a newly submitted question, use estimated topics to retrieve relevant question-answers from the archive.

## 2 QATM

### Graphical Model



**Input Variables: given**

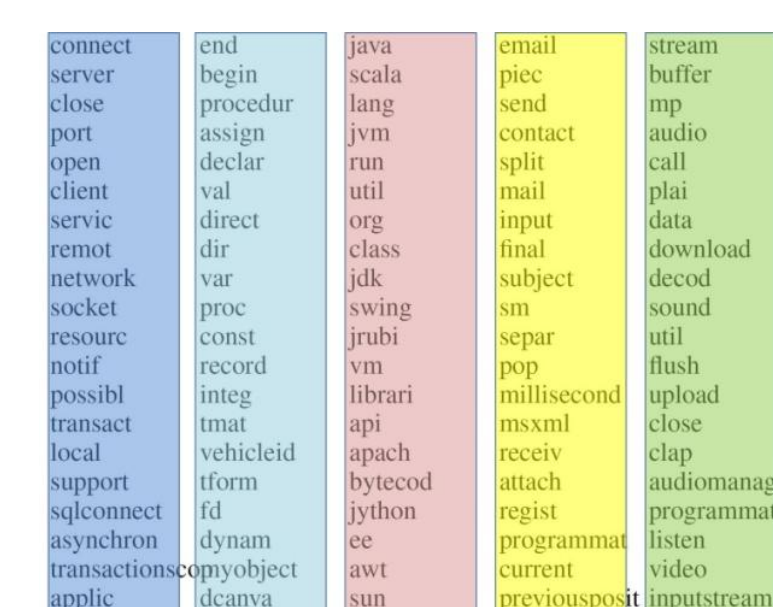**Latent Variables: posterior is estimated given the data**

Observations: given

**Simple intuition**

❖ Documents exhibit multiple topics

❖ Topics in answers are influenced by topics in question

### Generative Process

❖ Documents are generated by realizations of random variables that are sampled from probability distributions

❖ Corpus contains two types of topics, Question topics (Q-topics) and Answer topics (A-topics)

❖ Each question is a mixture of Q-topics

❖ Each answer is a mixture of A-topics

❖ Model captures the dependency between question topics and answer topics by conditioning each A-topic in an answer on Q-topics drawn from the topic distribution of the corresponding question.

### Inference and Parameter Estimation

❖ Essentially, an optimization problem: Model specifies a joint probability distribution with some parameters. What are the optimal values of the model parameters?

❖ To fit the model to data, the generative process is inverted and parameter values are generated from given observations(the words)

❖ Exact inference is intractable. Gibbs sampling, an approximate inference algorithm is used.
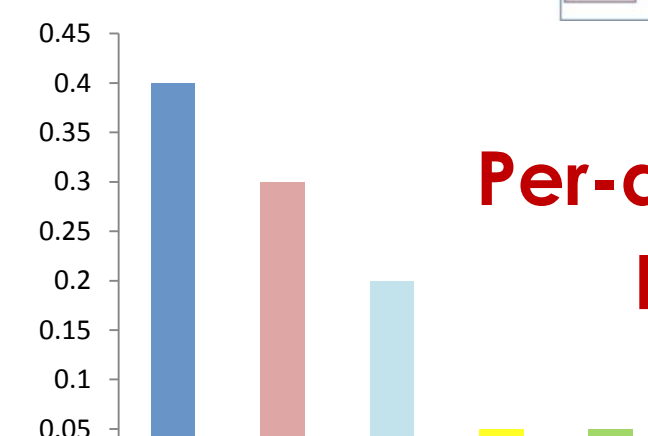
❖ Inference Output

**Per-corpus topic distributions**

**Per-word topic assignments**

**Per-document topic proportions**



## 3 Dataset

### Stackoverflow

**Stackoverflow Statistics**

| | |
|---|---|
| #Questions | 1,188,585 |
| #Answers | 2,939,767 |
| #Tags | 27,659 |
| #Users | 440,672 |
| Avg #answer per question | 2.4818 |
| Avg #tags per question | 2.9254 |
| Avg score of questions | 1.4967 |
| Avg score of answers | 1.8478 |

### Dataset Creation



**Selected Tags**

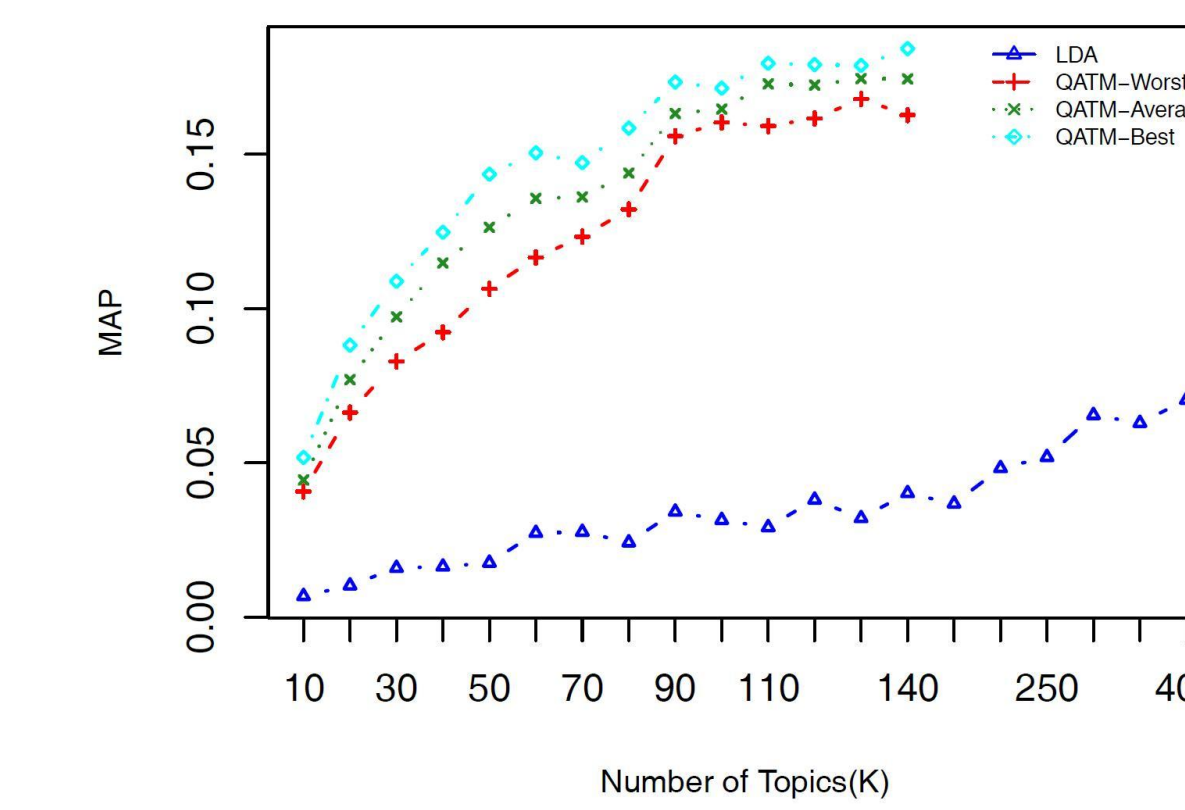| C# | .net | Sql |
|---|---|---|
| Sql-server | Java-script | css |
| Java | Ruby | Ruby-on-rails |
| Wpf | Iphone | Web-development |
| Android | Windows | Delphi |
| Django | Python | C |
| Bash | Linux | Homework |

**Tag Distribution in Training Set and Test Set**



## 4 Experiments – Question Answering

### Results on the Stackoverflow dataset

**QATM retrieval performance compared to LDA model in terms of Mean Average Precision and TopN measures**



| | Top1 | Top2 | Top3 | Top4 | Top5 |
|---|---|---|---|---|---|
| LDA | 0.023 | 0.026 | 0.029 | 0.03 | 0.032 |
| QATM-Worst | 0.108 | 0.127 | 0.138 | 0.144 | 0.148 |
| QATM-Average | 0.122 | 0.142 | 0.151 | 0.156 | 0.16 |
| QATM-Best | 0.131 | 0.152 | 0.161 | 0.164 | 0.168 |

**Topical dependencies captured by QATM with examples of Q-topics and A-topics represented by their first 20 most probable words**



## 4 Experiments – Automatic Tagging

### Results on the Stackoverflow dataset

**Clustering results in terms of Precision, Recall and F-measure.**