

## Documentation of Machine Learning Software

- Users of software documentation are not always software experts.
  - Scientists and engineers of varying expertise and depth of knowledge on data science and machine learning
  - automated generation and adaptation of machine learning software documents for users with varying levels of expertise.
- stack
- Investigation of ^ overflow Q&A and classify<sup>ing</sup> the documentation related to Q/A's within machine learning domain.
  - Find trigger of problems and potential change requests to the documentation
  - Build techniques for automatic documentation generation and extending to adoption, summarization, and explanation of software functionalities.

## Introduction

- Documentation is any artifact which its purpose is to communicate information about the software system to which it belongs, to individuals involved in the production of that software.
- Problem: Existing hyper-technicality of software engineering causes users to do extra searching in other resources.
- develop expertise aware software documentation for ML products
- Research Questions: 1) What type of problems do users face when using documentation of ML software?  
 2) How does the expertise level of users impact their understandings of the documentation of ML Software?  
 3) How does documentation evolve in relation to the SO questions?

## Methodology

### Research Question 1:

- Use soft stack overflow platform to analyse problems and painpoints of ML software tools.
- Mining following questions:- what types of problems do users have with ML software, and in what areas?
  - which types of documentations are used and questioned?
  - why these documentations are referred to, and in which area?
  - How the documentation referrals happen?
- Gathered SO Q&As tagged as ML or popular ML Libraries. i.e. Tensorflow or Pytorch
- Labeling found SO Q&A with mining questions mentioned above.
- Used two annotators for thematic analysis for identification by manual analysis
- the latter method helped in defining pitfalls of SO Q&A.

### Research Question 2:

- ultimately adapt software documentation for different user groups and information needs
- developing understanding of software and domain (ML)

- Movashowitz et al.: The number of upvoted questions, number of up voted answers, number of accepted answers, number of down voted answers.
- consider # of posts tagged in relation to ML
- compare groups with different levels of expertise and reason out the extent of difference between their needs from software documentation.

### Research Question 3:

- relationship between*
- analyse the questions to the opened issues in the respective software repository and the commit messages while changing the documentation
  - associate each question to a release of a software assuming that a valid question that is asked before Release<sub>j</sub> can't be responded in Release<sub>i</sub>, where  $j < i$
  - cosine similarity between the text of the SO question and issue in the repository
  - for questions not matching issues, similarity of commit messages → were matched.

### Summary and Future Work.

- For RQ1 →
- understanding user's behaviour in using ML documentation and understanding the process of evolution of the software documents in ML domain
  - 500 questions with tensor flow tag and manually categorizing questions, 16.6% of questions are related to documentation
    - SO questions are concerning parameter tuning, Model creation, and error/exception → 10.7%
    - Most questions were triggered as users weren't able to replicate example in documentation and lack of description
      - Hilroy 24.8%
      - 14.2%
      - 12.9%

- 60.7% of documentations are related to official Tensorflow documentations, while others refer to third-party materials like tutorials, videos, books, scientific papers
- 72.8% hyperlinked to mentioned documentation and others used a screenshot or mentioned the name of the documentation.

## What do developers know about machine learning: a study of ML discussions on Stack Overflow

- Machine Learning can be used for bug prediction, and software development effort estimation.
- Investigation to understand what educators should focus on, and how different online programming discussion communities can be more helpful.
- Some machine learning topics are significantly more discussed
- Latent Dirichlet Allocation (LDA) can suggest more appropriate tags that make machine learning tags more visible.
- The last point helps getting faster response on sites like SO

### Introduction

- Need to find specific gaps (in those interested) knowledge in Machine learning
- Topic modeling approach for making machine learning questions more searchable through tagging system.
  - RQ1: what machine learning topics are discussed on SO?
  - RQ2: what exactly do the developers discuss about those machine learning topics?
  - RQ3: what are the characteristics of machine learning posts considering their popularity and efficiency?
  - RQ4: Do the developers tag machine learning posts correctly, and can we improve such tagging system with topic modeling?

Hilroy

## Methodology

### Data collection:

- Used java SAX<sup>to</sup> parse through dataset
- found 28010 machine learning posts with machine learning tag.
- Entire analysis in this paper is based on STorrent

### Topic Modeling using LDA:

- LDA is a method of topic modelling
- java based package for topic modelling "Mallet"
- Requires Set of documents and parameters  $\alpha, \beta$ , and # of topics (K)
- Used results from data collections stage
- topic word matrix displayed top words
- document word matrix displayed described topic weights for each input document.
- LDA don't have a label but are rather probabilistic distribution of words only

## Topic Categorization

2022-05-10

- LDA was used to identify the topics of machine learning post (28,010)
- Tried categorization with 20 topics and 50 topics ( $K=20$  &  $K=50$ )
- Set  $\alpha, \beta$  to 0.01, 0.25, 1.00 for each K value
- when  $\alpha/\beta = 0.01$ : hard to identify leading topics
  - important topics remain hidden
  - many topics get associated with document when they aren't related

2022-05-11