



# Classifying Question Quality of StackOverflow

---

MIDS 266 - Fall 2022  
Paco Valdez



# The Problem

---

There are multiple forums where coding questions can be asked, an interesting problem is to be able to prioritize “good” questions and do the opposite to “bad” questions. The applications to this solution would help these forums to provide feedback to users while they are asking the question, also “good” questions could be moved to the next step of the workflow that could be an automated tagging system or to a manual review.

# Previous Efforts

---

- Piyush Arora, Debasis Ganguly, and Gareth J.F. Jones. 2016. **Nearest Neighbour based Transformation Functions for Text Classification: A Case Study with StackOverflow**. In Proceedings of the 2016 ACM International Conference on the Theory of Information Retrieval (ICTIR '16). Association for Computing Machinery, New York, NY, USA, 299–302. <https://doi.org/10.1145/2970398.2970426>
- László Tóth, Balázs Nagy, Tibor Gyimóthy, and László Vidács. 2020. **Why will my question be closed? NLP-based pre-submission predictions of question closing reasons on stack overflow**. In Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER '20). Association for Computing Machinery, New York, NY, USA, 45–48. <https://doi.org/10.1145/3377816.3381733>
- Antoaneta Baltadzhieva, Grzegorz Chrupała. 2015. **Predicting the Quality of Questions on Stackoverflow**. Proceedings of Recent Advances in Natural Language Processing, pages 32–40, Hissar, Bulgaria, Sep 7–9 2015. <https://aclanthology.org/R15-1005.pdf>

# The dataset

Stack Overflow releases all the questions that have been asked and also it annotates the questions with badges. This can be translated into the following labels to be used for training.

Score-Based	View-Based	Bookmark-Based	Label
Great Question (Score $\geq 100$ )	Famous Question (Views $\geq 10,000$ )	Stellar Question (Bookmarks $\geq 100$ )	0
Good Question (Score $\geq 25$ )	Notable Question (Views $\geq 2,500$ )	Favorite Question (Bookmarks $\geq 25$ )	1
Nice Question (Score $\geq 10$ )	Popular Question (Views $\geq 1,000$ )	*Missing Category (Bookmarks $\geq 10$ )	2
No Badge	No Badge	No Badge	3

# EDA - Label distribution

Label

Distribution

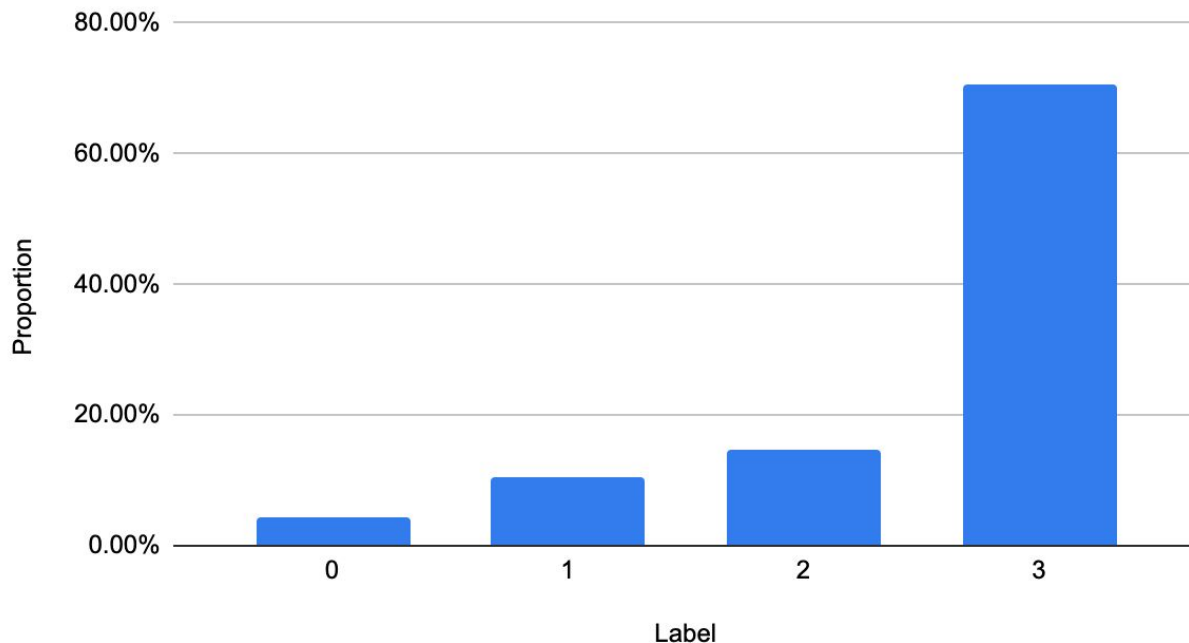
Great - 0 - 4.25%

Good - 1 - 10.43%

Nice - 2 - 14.85%

Bad - 3 - 70.47%

Label histogram



# The dataset - How questions are represented?

---

**Title:** Error with Redshift jdbc connector

**Body:** <p>Hey, I have this code to connect to Redshift using jdbc<p><code>

```
var jdbc = new ( require('jdbc') );
var config = {
  libpath: 'C:/Users/ABCD/Desktop/jar/RedshiftJDBC41-1.1.6.1006.jar',
  drivename: 'com.amazon.redshift.jdbc41.Driver',
  url: 'jdbc:redshift://examplecluster.abc123xyz789.us-west-2.redshift.amazonaws.com:5439/dev',
  user: 'xxxx',
  password: 'xxxxx'
};
```

```
jdbc.initialize(config, function(err, res) {
  if (err) {
    console.log(err);
  }
});
```

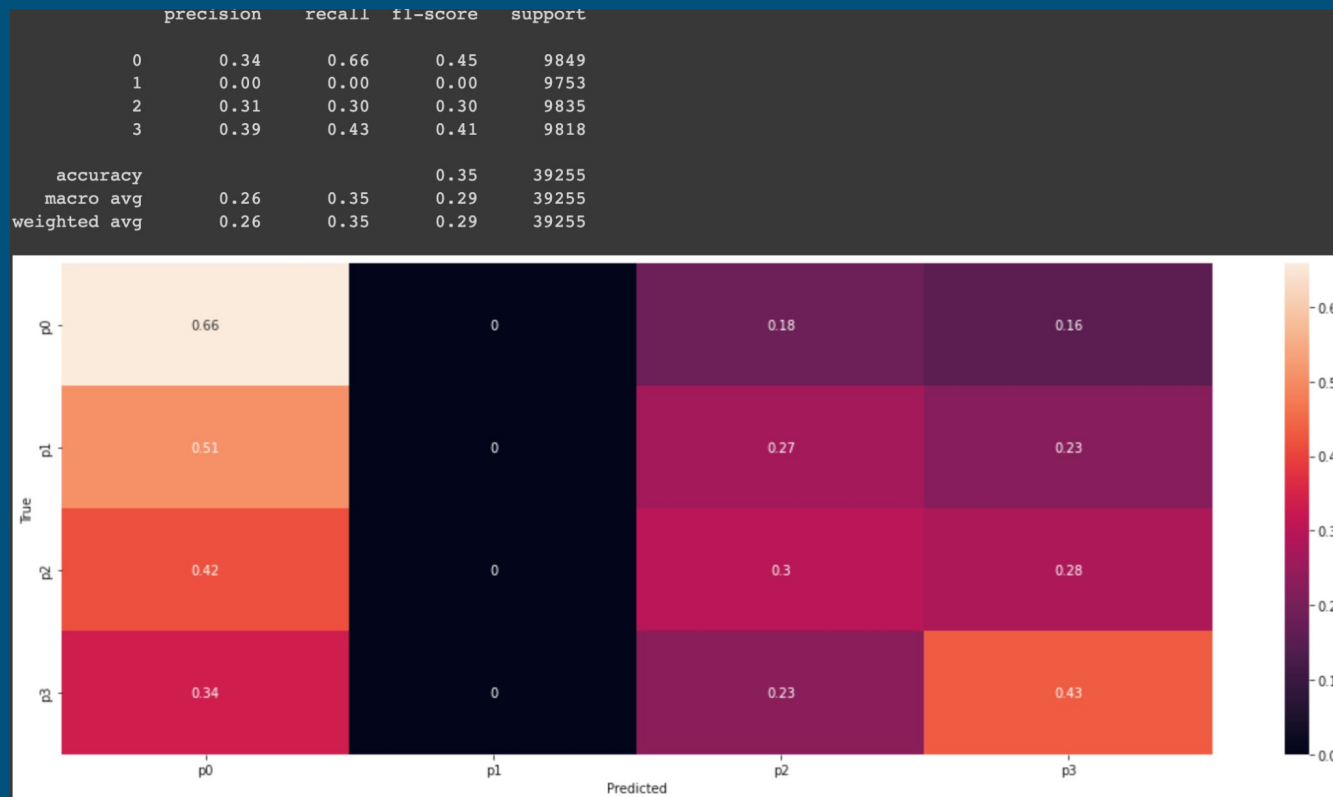
</code><p>The error is the following:<p><code>  
^ TypeError: undefined is not a function

at JDBCConn.initialize  
(C:\Users\ABCD\node\_modules\jdbc\lib\jdbc.js:62:20)  
</code>

# Baseline - Uncased-Bert Fine Tuned with a 20% random sample of all questions from 2012 to 2022

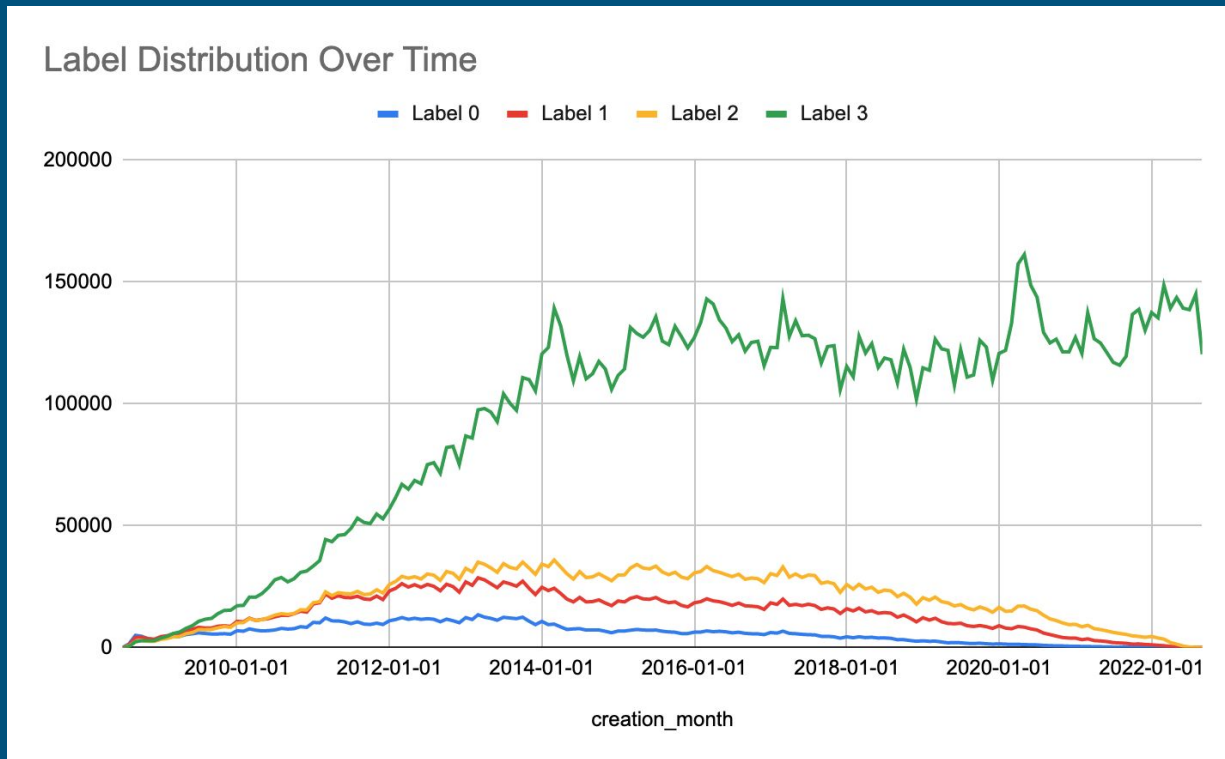


# Bertweet Fine Tuned with a 20% random sample of all questions from 2012 to 2022

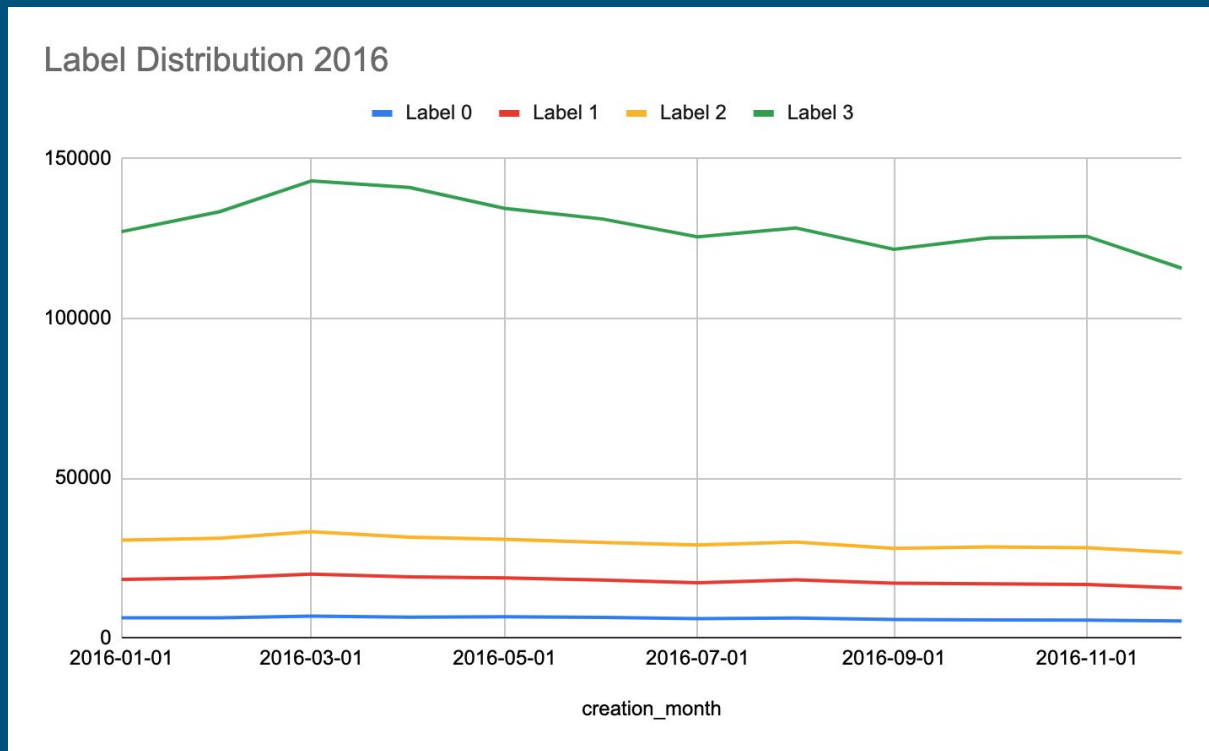




# Exploratory Analysis -Label distribution over time



# Exploratory Analysis - Label Distribution 2016



# Add special token - Unprocessed Text

---

Vuex getter not updating

+++++

<p>I have the below getter:</p>

<pre class="lang-js prettyprint-override">

```
<code>  withEarmarks: state => {  
    var count = 0;  
    for (let l of state.laptops) {  
      if (l.earmarks.length > 0) {  
        count++;  
      }  
    }  
    return count;  
  }  
</code>
```

</pre>

<p>And in a component, this computed property derived from that getter:</p>

<pre class="lang-js prettyprint-override">

```
<code>    withEarmarks() { return this.$store.getters.withEarmarks; },  
</code>
```

</pre>

<p>The value returned is correct, until I change an element within the laptops array, and then the getter doesn't update.</p>

# Add special token - processed Text

---

vuex getter not updating i have the below getter [CODE] and in a component this computed property derived from that getter [CODE] the value returned is correct until i change an element within the laptops array and then the getter doesnt update

```
from transformers import AutoTokenizer, TFAutoModel

tokenizer = AutoTokenizer.from_pretrained(model_checkpoint)
model = TFAutoModel.from_pretrained(model_checkpoint)

tokenizer.add_tokens(['[CODE]'], special_tokens=True)

model.resize_token_embeddings(len(tokenizer))
```

# Uncased Bert with Special Token Fine Tuned with 2016 data



# Bertweet with Special Token Fine Tuned with 2016 data



# Fine Tuned T5 (2016 data)

---

Text:<p>I have the below getter:</p><pre class="lang-js prettyprint-override">

<code> withEarmarks: state =&gt; { var count = 0;

```
    for (let l of state.laptops) {
      if (l.earmarks.length &gt; 0) {
        count++;
      }
    }
```

```
    return count;
  }
}
```

</code>

</pre>

<p>And in a component, this computed property derived from that getter:</p>

<pre class="lang-js prettyprint-override">

<code> withEarmarks() { return this.\$store.getters.withEarmarks; },

</code>

</pre>

<p>The value returned is correct, until I change an element within the laptops array, and then the getter doesn't update.</p>

Label: good</s>

Categories: ['great', 'good', 'nice', 'bad']

# T5 Fine Tuned with 2016 data





# Analysis - Real Special Token Frequency

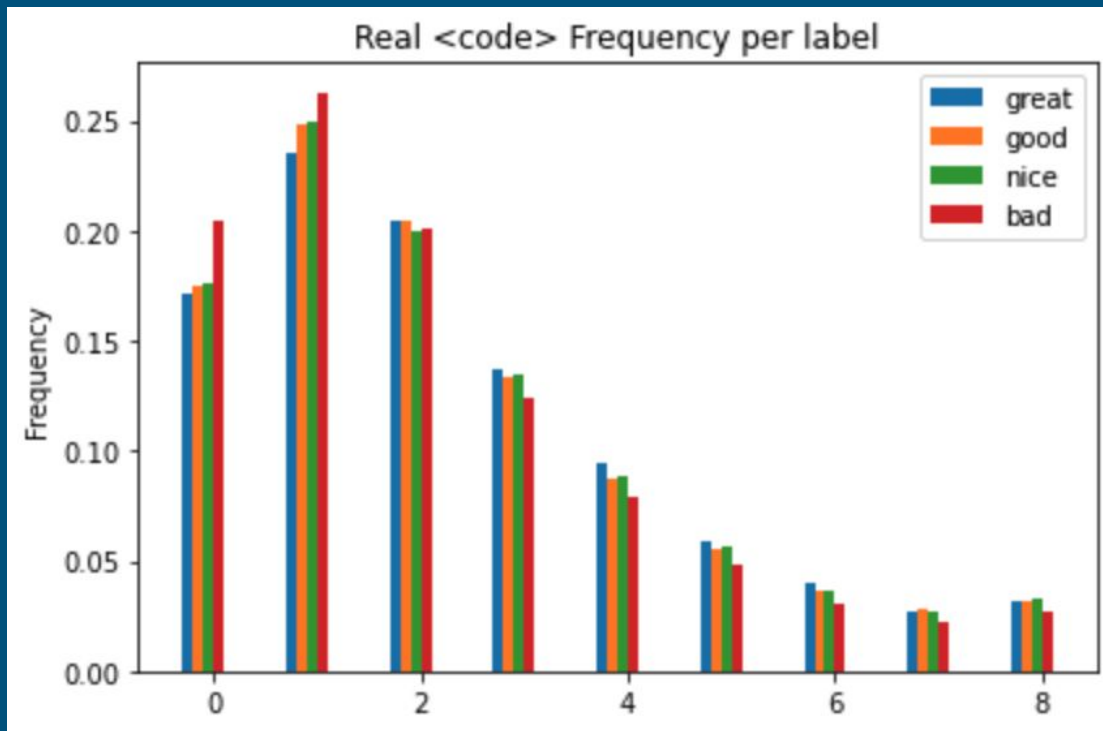
Special token frequency:

Category: **great** Mean: **2.845**

Category: **good** Mean: **2.846**

Category: **nice** Mean: **2.854**

Category: **bad** Mean: **2.557**



# Analysis - Predicted Special Token Frequency

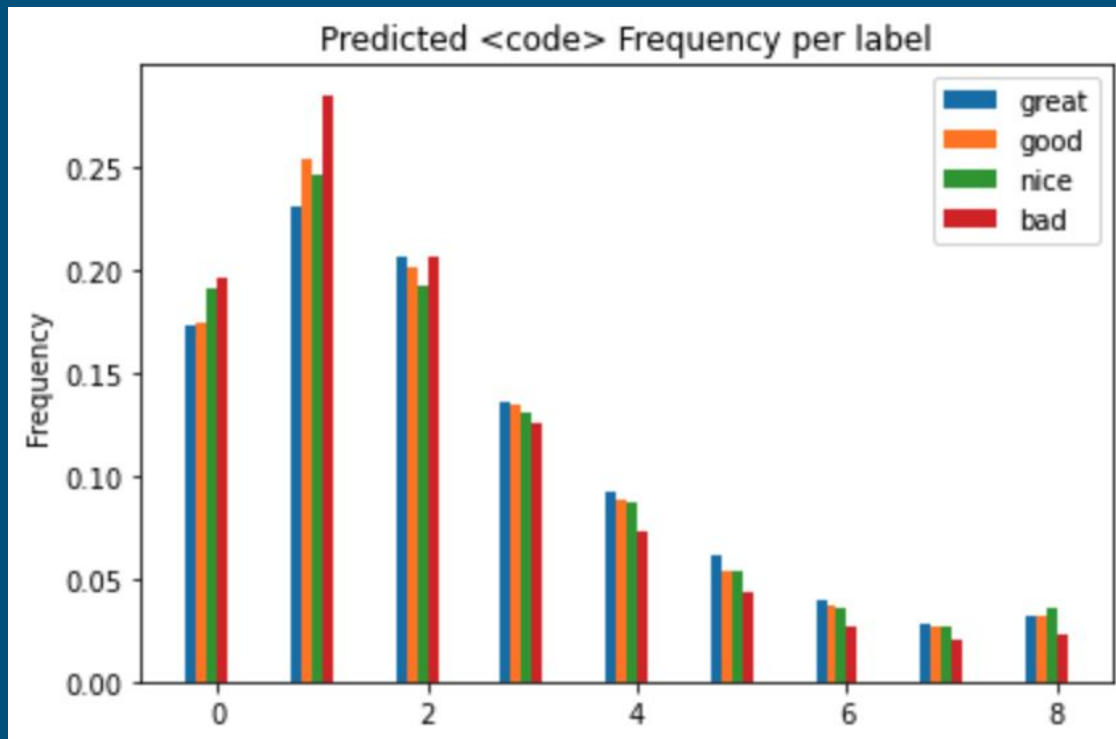
Special token frequency:

Category: **great** Mean: **2.856**

Category: **good** Mean: **2.839**

Category: **nice** Mean: **2.923**

Category: **bad** Mean: **2.376**



# Next Steps

---

- Use only a subset of the questions. (i.e. only python questions, or exclude esoteric languages)
- Adjust the score of a question based on when was voted on.
- Use a model trained on code to analyze the code blocks.
- More complex models on top of T5 or Bertweet.
- Do binary classification.

# Questions

---

Thank you!