
Multi-Factor Duplicate Question Detection in Stack Overflow

— **Team Name:** Hyperplane
Scavengers —

DataSet Curation

- Official MSR Dataset: Using techniques in paper, number of duplicates detected were only 2
- Finally resorted to using the Stack Exchange Data dump which contained data till 2021
- Postlinks.xml contained relationship for duplicate detection
- Postgres database for easy handling, sorting of data, visualizations and error analysis

Preprocessing

- Everything lowercase and the phrase Duplicate was removed from the title
- **<blockquote>** tag removed (since it contained the links to “duplicate”)
- All HTML tags removed but the data was preserved except **<code>**
- All whitespace converted to single space
- Punctuation except # removed, special case - C#
- Stemming
- Stopword Removal
- Tokenization

DupPredictor Model

4 Components: Title, Body, Topics (from Title + Body) and Tags

- Bag of Words and Cosine Similarity
 - Bag of Words Embedding for each component
 - Cosine Similarity for each component
- LDA
 - Get trend of the conversation, try to extract some meaning from the text
 - Topic modelling on Title + Body
 - Gensim Library with $K = 100$
- Greedy Composer
 - Random Restart Approach
 - Greedy Approach with params increasing in steps of 0.05
 - Weighted Sum and top K ($K = 10, 20$) are stored

Paper Results

Evaluation Metric

- Recall-Rate@K

$$\text{recall} - \text{rate}@k = \frac{N_{\text{detected}}}{N_{\text{total}}}$$

- N_{detected} is the number of duplicate questions whose master was in top K
- N_{total} is the total number of duplicate questions

Research Question 1: How good is DupPredictor over its 4 individual components?

- Is our model good?
- Measure recall-rate@K, K = 10 and 20 of DupPredictor with its single components
- DupPredictor performs very good as compared to its other individual components and by a large margin

Algorithm	Recall Rate@10	Improvement (%)
DupPredictor	0.56	0.0
Title similarity	0.33	69.69
Description similarity	0.2	180
Topic similarity	0.08	600
Tag similarity	0.28	100

Algorithm	Recall Rate@20	Improvement (%)
DupPredictor	0.65	0.0
Title similarity	0.42	54.76
Description similarity	0.26	150
Topic similarity	0.14	364.28
Tag similarity	0.34	91.17

Research Question 2: Effect of Varying Number of Training Questions

- How many questions is enough for training due to the high compute?
- Trained the model over 2 training set sizes, the first 100 and 300 questions. Tested on the same testing set (300-400) for uniformity. Measure the recall-rate@10 and recall-rate@20
- General parameter trend remains the same

$$\alpha \geq \beta \geq \delta \geq \gamma$$

Trained on first	Alpha (Title)	Beta (Body)	Gamma (Topic)	Delta (Tags)
K = 10				
100	0.5	0.5	0.0454	0.45
300	1.0	0.6	0.3	0.6
K = 20				
100	0.4	0.5	0.13738	0.1992
300	0.9	0.85	0.3153	0.3

Trained on first	Recall@10 Training	Recall@10 Testing	Recall@20 Training	Recall@20 Testing
100	0.74	0.55	0.78	0.62
300	0.67	0.56	0.74	0.65

- Model trained on 300 performs better on test set so it must be generalizing better
- However performance difference is negligible
- Implies: Model is able to learn and tune its parameters even on a lesser amount of data and hence we can restrict training to 300 questions

Research Question 3: Does DupPredictor estimate the 4 weights of its 4 constituent components well?

- Is greedy random restart method accurate and well predicted?
- Randomly generated 50 sets of α , β , γ , δ
- Simplistic approach outperforms all randomly generated weights.
- If the random generated params follows the trend $\alpha \geq \beta \geq \delta \geq \gamma$ and $\alpha \approx 1$, then model performs best
- Worst performing give to little attention to title, or too much attention to topic

alpha	beta	gamma	delta	recall-rate@10	recall-rate@20
0.9	0.85	0.3152	0.3	0.56	0.64
0.89	0.0687	0.29	0.2593	0.53	0.62
0.8118	0.5148	0.6062	0.3224	0.53	0.61
0.2605	0.88	0.6911	0.08	0.36	0.45
0.1749	0.4582	0.7276	0.0998	0.36	0.45

Error Analysis And Basic Stats

Pruning the search space improves performance on multiple fronts

- 94/100 questions in the test set had at least one tag in common with their duplicates.
- The 6 cases where tags were not common were very subtle misses.

- For eg:

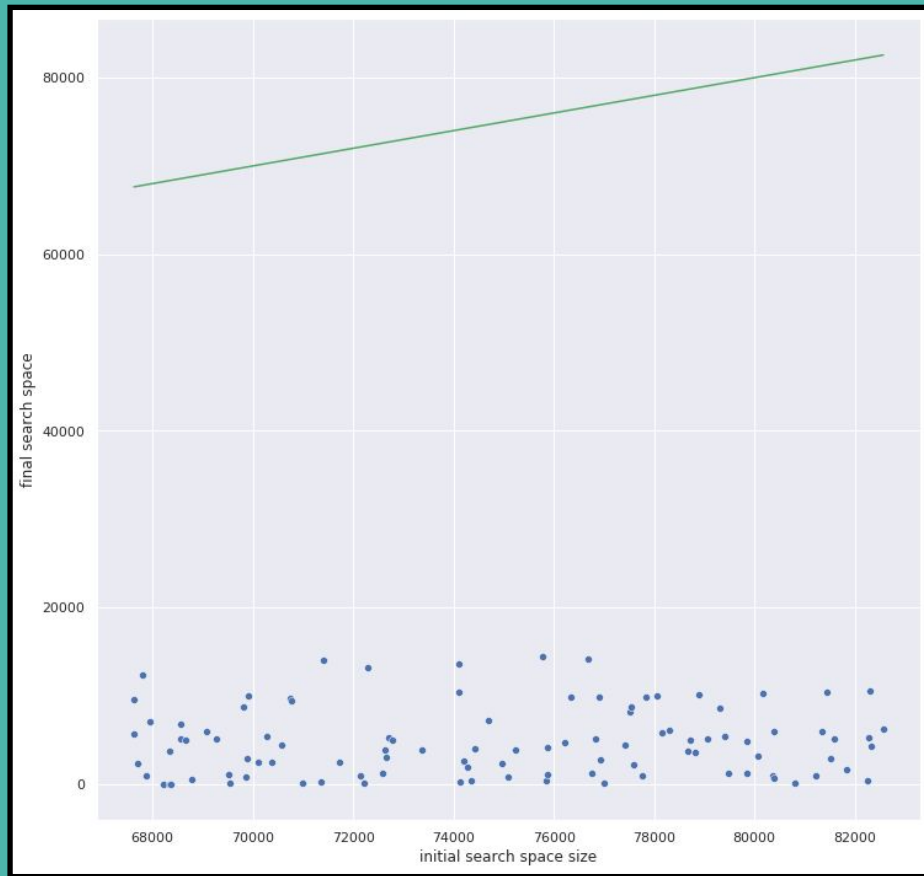
Tags of a question with ID 441910:

internet-explorer image firefox rendering

Tags of its duplicate ie 130161:

html css internet-explorer-6 png

- Figure on the right shows the massive reduction in search space size with 2 benefits:
 - Less time to evaluate candidates in search space
 - Higher rank for an actual duplicate



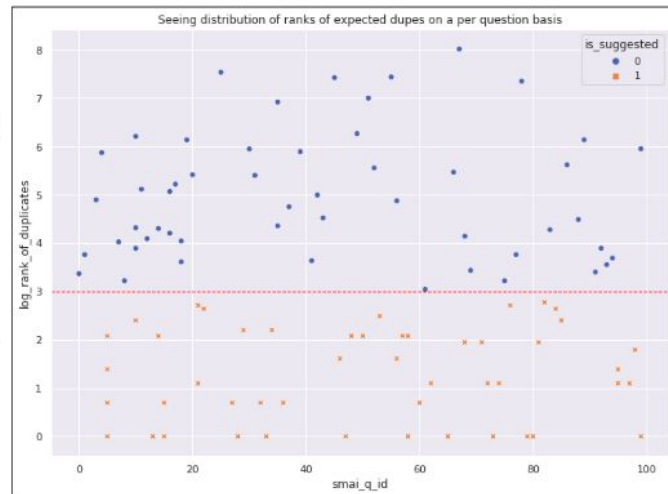
SUMMARY OF ABLATION STUDY SCORES (at Recall@20)



Model description	Training Recall@20	Test Recall@20		Mean rank of highest ranked duplicate		Best rank of a duplicate		Worst rank of a duplicate		Parameters
		Without Pruning	With pruning	Without Pruning	With pruning	Without Pruning	With pruning	Without Pruning	With pruning	
K=20 with all components	0.74	0.65	0.68	175.84	99.9	1	1	2365	1359	[0.9, 0.85, 0.31, 0.3, 0.0]
K =20 with all components except Title Similarity	0.61	0.53	0.53	454	287	1	1	12366	4048	[0, 0.85, 0.35, 0.64, 0.0]
K =20 with all components except Topic Similarity	0.733	0.61	0.61	229	160	1	1	3691	3613	[0.2, 0.19, 0, 0.07, 0.0]
K =20 with all components except Body Similarity	0.707	0.61	0.62	214	120	1	1	4058	1972	[1.0, 0, 0.45, 0.6, 0.0]
K =20 with all components except Tags Similarity	0.67	0.51	0.62	540	122	1	1	25818	3177	[0.9, 0.8, 0.1, 0, 0.0]
K=20 with all components + Jaccard Coeff included	0.763	0.65	0.67	105	78	1	1	2054	1197	[0.45, 0.35, 0.25, 0.29, 1.0]

Some insights based on the table above and the detailed analysis in the next section

- As per the ablation study, the most significant parts for detection seem to be **“title” and the “tags”**.
- Pruning the search space is able to improve results and is quicker computationally
- The most significant increase in results after pruning was in the case where tag similarity was discarded. This was as pruning based on tags does a job ALMOST similar to tag similarity.
- With moderators strictly enforcing titles to be precise and accurate, “title similarity” seems to be a great indicator.
- TAGS become important as it helps in **restricting search to relevant languages and technologies**. Eg: a question regarding retrieval queries on MySQL is likely to be similar in body and title when compared to the same question wrt MongoDB. However, the two questions being different tagged is what is allowing the model to successfully differentiate between the two.
- In the adjoining figure, x: question id ; y : rank of one of the expected duplicates. Since a question can have more than one expected duplicate, there may be multiple dots on the same vertical line. The red line depicts the threshold (20) beyond which questions are not suggested.

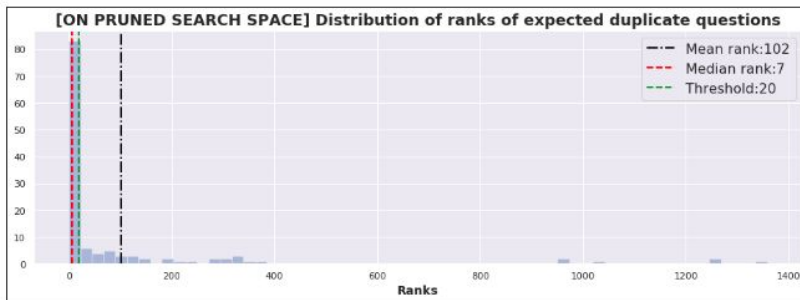


ALL 4 components (title + body + tags + topics) involved

In general, the 0.68% accuracy seems decent.

Params: [0.9, 0.85, 0.3125, 0.3] i.e. high importance to title and body text

- qid for whom we want a duplicate: 441529
- Link to qid: [LINK](#)
- Actual best candidate:** 263400
 - scores: [0.51, 0.11, 0.5, 0.4]
 - final score: 0.85 | rank 68
- Chosen best candidate:** 34505
 - Scores: [0.5, 0.5, 0.7, 1]
 - final score: 1.4287 | rank : 1
- We see that the q_chosen beats q_expected in all 4 criteria.
- Conclusion:** The model is misled by the use of common words like “objects” and “hash function” in q_asked and q_chosen. Even though q_expected is clearly a better match; but the model fails to capture the fact that the q_asked is actually discussing an algorithm and naively allots q_chosen as the higher score (despite it being very different) just because of common words. This can be improved by considering synonyms or word embedding such that phrases like “intelligently built” etc can be linked to words like “algorithm”.



question asked	Expected best	Chosen best
<div><h3>How do you implement GetHashCode() on objects? [duplicate]</h3><p>Asked 12 years, 10 months ago · Active 12 years, 10 months ago · Viewed 2k times</p><p>This question already has answers here: Closed 12 years ago.</p><p>Duplicate: What is the best algorithm for an overridden System.Object.GetHashCode?</p><p>If you've written an object with a variety of data-members, how do you intelligently implement GetHashCode?</p><p>One developer told me he just XORs (^ operator) the Hash of relevant data-fields, but I am unconvinced this is a "best-practices" implementation.</p><p>If I had my way, there would be functions <code>Object.CombineHashes(Object[])</code>, and/or <code>Object.CombineHashes(int[])</code> to help intelligently build hashes of complex objects.</p><p>How would you write these functions?</p><p>cs .net</p></div>	<div><h3>What is the best algorithm for overriding GetHashCode?</h3><p>Asked 13 years, 1 month ago · Active 1 month ago · Viewed 252k times</p><p>1567</p><p>In .NET, the <code>gethashcode</code> method is used in a lot of places throughout the .NET base class libraries. Implementing it properly is especially important to find items quickly in a collection or when determining equality.</p><p>Is there a standard algorithm or best practice on how to implement <code>GetHashCode</code> for my custom classes so I don't degrade performance?</p><p>.net algorithm hashcode gethashcode</p><p>Share Edit Follow</p><p>edited Sep 18 '19 at 17:55</p><p>poke 325k ● 64 ● 498 ● 558</p><p>asked Nov 4 '08 at 20:53</p><p>bitbonk 46.7k ● 35 ● 177 ● 271</p></div>	<div><h3>Is Object.GetHashCode() unique to a reference or a value?</h3><p>Asked 13 years, 3 months ago · Active 6 years, 2 months ago · Viewed 8k times</p><p>25</p><p>The MSDN documentation on <code>Object.GetHashCode()</code> describes 3 contradicting rules for how the method should work.</p><ol style="list-style-type: none">1. If two objects of the same type represent the same value, the hash function must return the same constant value for either object.2. For the best performance, a hash function must generate a random distribution for all input.3. The hash function must return exactly the same value regardless of any changes that are made to the object.<p>Rules 1 & 3 are contradictory to me.</p><p>Does <code>Object.GetHashCode()</code> return a unique number based on the value of an object, or the reference to the object. If I override the method I can choose what to use, but I'd like to know what is used internally if anyone knows.</p><p>cs .net</p></div>

ALL components (Except Title Similarity) involved

In general, the 0.53% accuracy WITHOUT using the title is reasonable.

Params: [0, 0.85, 0.35, 0.64] i.e. without the title, the content now shifts attention to TAGS.

- qid for whom we want a duplicate: 508191 | Link to qid: [LINK](#)
- **Actual best candidate:** 200090
 - scores: [0.67, 0, 0.95, 0.5]
 - final score: 0.65 | rank 62
- **Chosen best candidate:** 275733
 - Scores: [0.4, 0.129, 0.88, 0.7]
 - final score: 0.87 | rank : 1
- q_chosen beats q_expected in "body_score" which has a high weightage of 0.85.
- **Conclusion:** Due to removal of title, q_expected loses out. Also, because the body in q_asked is almost entirely empty (since we remove the code before inputting to model), there are hardly any words in the "filtered body" of q_asked to give any meaningful insights for matching. However, since the coefficient for the body is high, q_chosen overtakes q_expected. The model where title similarity exists is able to CORRECTLY detect this test case.

question asked	Expected best	Chosen best
<div><h3>How can I convert a string to an integer in C++ [d]</h3><p>Asked 12 years, 10 months ago · Active 12 years, 10 months ago · Viewed 2k times</p><div><p>4</p><p>This question already has answers here: Closed 12 years ago.</p></div><p>I am trying to copy the value in bar into the integer foo.</p><p>This is what I have so far. When I run it I got a different hex value. Any help would be great.</p><pre>int main() { string bar = "0x00EB0C62"; int foo = (int)bar; cout << hex << foo; ChangeNewWal("pinball.exe", (void*) foo, "100000", 4); return 0; }</pre><p>So the output should be 0x00EB0C62.</p><p>c++ string int casting</p></div>	<div><h3>How do you convert a C++ string to an int? [d]</h3><p>Asked 13 years, 1 month ago · Active 5 years, 1 month ago · Viewed 137k times</p><div><p>46</p><p>This question already has answers here: Closed 10 years ago.</p></div><p>Possible Duplicate: How to parse a string to an int in C++?</p><p>How do you convert a C++ string to an int?</p><p>Assume you are expecting the string to have actual numbers in it ("1", "345", "38944", for example).</p><p>Also, let's assume you don't have boost, and you really want to do it the C++ way, not the cruffy old C way.</p><p>c++ parsing int stdstring</p></div>	<div><h3>How can I tokenize a C++ string? [closed]</h3><p>Duplicate: How do I tokenize a string in C++?</p><p>I have a character array in C++:</p><pre>arr="abc def ghi"</pre><p>I want to get the strings "abc" "def" "ghi" out of the string. Are there any built in functions to do this?</p><p>stackoverflow</p><p>edited Mar 5 at 4:12 Rich B 6,723 ●2 ●12 ●31</p><p>asked Nov 9 at 8:08 santhosh 35 ●2</p><p>comments (1)</p></div>

ALL components (Except Body Similarity) involved

In general, the 0.62% accuracy WITHOUT using the body shows that BODY does not play a major role in helping the detector..

Params: [1, 0, 0.5, 0.6] with emphasis on TITLE and TAGS.

- qid for whom we want a duplicate: 490973 | Link to qid: [LINK](#)
- **Actual best candidate:** 296020
 - scores: [0.15, 0.6, 0.7, 0.4]
 - final score: 0.8 | rank 42
- **Chosen best candidate:** 150606
 - Scores: [0.36, 0.26, 0.65, 0.7]
 - final score: 1.11 | rank : 1
- Conclusion: The question wants an answer on how to FREEZE cells in EXCEL using HTML/CSS. Q_chosen gives an answer which is **NOT EVEN REMOTELY related** to EXCEL. Clearly, the model fails to capture the "excel" factor and hence, blunders. The BODY of q_expected and q_ask match to a large extent to give a score of 0.6 . But since coeff for body score is ZERO in this model, the test case gets misdetected by the model.



question asked	Expected best	Chosen best
<div><h3>Emulating Excel's "freeze cells" in an HTML table [duplicate]</h3><p>Asked 12 years, 10 months ago · Active 5 years, 2 months ago · Viewed 16k times</p><p>This question already has answers here: Closed 10 years ago.</p><p>Possible Duplicate: How can I lock the first row and first column of a table when scrolling, possibly using javascript and CSS?</p><p>I have an HTML table that contains a large number of rows and columns. The top row contains headers, and the first cell in every row below that contains a header.</p><p>I need to allow the user to scroll the table whilst keeping the top row and left column visible at all times (similar to what can be achieved using the "freeze cells" option in Excel).</p><p>The solution only needs to work for IE7+ as this is for an internal application.</p><p>Ideally the solution should degrade gracefully in other browsers, but that's not essential.</p><p>My hunch is that there's a JavaScript solution to this...</p><p>javascript html user-interface html-table</p><p>Share · Edit · Follow</p><p>edited May 23 '17 at 12:03</p><p>asked Jan 29 '09 at 9:36</p></div>	<div><h3>How can I lock the first row and first column of a table when scrolling, possibly using JavaScript and CSS?</h3><p>Asked 13 years ago · Active 1 year ago · Viewed 108k times</p><p>How can I create a table that has its first row and first column both locked, as in Excel, when you activate "freeze panes"? I need the table to both scroll horizontally and vertically (a lot of solutions for this exist, but only allow vertical scrolling).</p><p>So, when you scroll down in the table, the first row will stay put, since it will have the column headings. This may end up being in a <code><thead></code>, or it may not, whatever makes the solution easier.</p><p>When you scroll right, the first column stays put, since it holds the labels for the rows.</p><p>I'm pretty certain this is impossible with CSS alone, but can anyone point me toward a JavaScript solution? It needs to work in all major browsers.</p><p>javascript css excel scroll css-tables</p><p>Share · Edit · Follow</p><p>edited Sep 15 '16 at 8:22</p><p>Brian Tompsett - 止藤 6,384 · 57 · 52 · 129</p><p>edited Nov 17 '18 at 16:17</p><p>pkwending 13,9k · 30 · 98 · 137</p><p>Hi, I know this is a bit old, but did you get a working solution to this? The answer marked as correct now has broken links. I'm trying to find out the same thing here: stackoverflow.com/questions/743663/... - Damovise Apr 13 '09 at 12:40</p><p>I just tried the links in the accepted answer, and they worked for me. Are you still having trouble? - pkwending Apr 14 '09 at 17:00</p></div>	<div><h3>JavaScript highlight table cell on tab in field</h3><p>Asked 13 years, 2 months ago · Active 13 years, 2 months ago · Viewed 6k times</p><p>I have a website laid out in tables. (a long mortgage form)</p><p>1 in each table cell is one HTML object. (text box, radio buttons, etc)</p><p>What can I do so when each table cell is "tabbed" into it highlights the cell with a very light red (not to be obtrusive, but tell the user where they are)?</p><p>javascript html</p><p>Share · Edit · Follow</p><p>asked Sep 25 '08 at 23:52</p><p>Jason</p></div>

ALL components (Except TAGS Similarity) involved

In general, the accuracy without using tags is 0.51%. But in the pruned search space, tags are automatically taken care of and hence, leads to an accuracy of 0.62% (very minor drop) when tags are not considered.

Params: [0.9 , 0.8 , 0.1 , 0] with emphasis on TITLE and BODY.

- **qid for whom we want a duplicate:** 521687 | Link to qid: [LINK](#)
- **Actual best candidate:** 43021
 - scores: [0.5, 0.08, 0.35 , 1]
 - final score: 0.77 | rank 28
- **Chosen best candidate:** 141108
 - Scores: [0.8, 0.25, 0.64, 0.4]
 - final score: 0.999 | rank : 1
- **Conclusion:** q_ask is concerned with foreach in C#. Q_expected is related to foreach in C# but q_chosen is related to foreach in PHP. Both of them are not filtered as both of them have the “foreach” tag common with q_ask. Q_expected has exactly the same tags as q_asked but q_chosen is predicted (even though it **deals with a totally different language**). This is due to the zero value of coefficients of TAGS_SCORE.

question asked	Expected best	Chosen best
<div><div>foreach with index [duplicate]</div><div>Asked 12 years, 10 months ago · Active 2 years, 3 months ago · Viewed 229k times</div><div><div>167</div><div>This question already has answers here: How do you get the index of the current iteration of a foreach loop? (25 answers) Closed 8 years ago.</div><div>Is there a C# equivalent of Python's <code>enumerate()</code> and Ruby's <code>each_with_index</code>?</div><div>cf · foreach</div><div><div>edited Jul 12 '18 at 1:18</div><div> 86.5k · 43 · 226 · 265</div><div><div>asked Feb 8 '09 at 19:01</div><div> 4,824 · 6 · 27 · 25</div></div><div>If you're using LINQ, there are overrides of the various functions that allow for enumeration. Otherwise, you're usually stuck using a variable that you increment yourself. – GWJLosa Feb 6 '09 at 19:03</div><div>1 Can we update this question to C# 7.0 where there are tuples now? I wonder how a solution would look using tuples. – Hordvik Wiese Mar 13 '17 at 8:57</div><div>Isn't that a feature of <code>foreach</code> that you can process each element absolute decoupled from its position in the list? – Konstantin A. Moog Apr 9 '18 at 11:21</div></div></div></div>	<div><div>How do you get the index of the current iteration of a foreach loop?</div><div>Asked 13 years, 3 months ago · Active 1 year, 2 months ago · Viewed 1.2m times</div><div><div>1109</div><div>Is there some rare language construct I haven't encountered (like the few I've learned recently, some on Stack Overflow) in C# to get a value representing the current iteration of a foreach loop?</div><div>For instance, I currently do something like this depending on the circumstances:</div><div><pre>int i = 0; foreach (Object o in collection) { // ... i++; }</pre></div><div>cf · foreach</div></div></div>	<div><div>How to find the foreach index?</div><div>Asked 13 years, 2 months ago · Active 6 months ago · Viewed 814k times</div><div><div>575</div><div>Is it possible to find the <code>foreach</code> index?</div><div>In a <code>for</code> loop as follows:</div><div><pre>for (\$i = 0; \$i < 10; ++\$i) { echo \$i . ' ' ; }</pre></div><div><code>\$i</code> will give you the index.</div><div>Do I have to use the <code>for</code> loop or is there some way to get the index in the <code>foreach</code> loop?</div><div>php · loops · foreach</div></div></div>

Further Explorations

- Attempted to further improve the best results obtained.
 - Mainly experiment with 2 techniques
 - Increasing the number of features while calculating the composer score
 - Further improving the predictions with the help of NLP
-

Jaccard Coefficient

- Composer is treating each of its similarities as separate entities
- Use Jaccard to get a holistic view
- combined the processed title, body and tags as BOW

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- Included a fifth parameter, Jaccard similarity
- Model predicted jaccard to be very important
- Though training recall increased from 0.74 to 0.76, testing recall didn't increase.
- Conclusion: Title, body and tags are disjoint

BERT embeddings

- Reranking smaller pruned rank list (by DupPredictor)
- Top 30 predictions
- Actual masters not in top 30 appended at the end
- Computed bert embeddings for title and body text

$$BertScore = \alpha \times CosineSim(BERT(Q1_{Title}), BERT(Q2_{Title})) + \beta \times CosineSim(BERT(Q1_{Body}), BERT(Q2_{Body}))$$

- Improved recall-rate@20 - from **0.65 to 0.88**
- For **29** duplicate, improved richness by including **32** more correct candidate questions
- Able to include **2** questions in top 20 - initially **not even in top 30**

Website

Live Demo

Limitations

- Size of training and testing set
 - Huge number of duplicate questions
 - Had to compute large number of pairwise scores
 - Training set - 60,000 candidate questions per duplicate - 8 hours time
 - Test set - 80,000 candidate questions per duplicate - 5 hours time
 - Including more duplicate became tough
- Estimation of parameters
 - Able to compute majority of parameters, but greedy approach costs 4 hours compute time per training.
- Bert Training
 - Embedding creation for each post was compute heavy
 - Approx 20 Hours with 14GB RAM. Laptops started hanging
 - Resorted to the re-ranking approach - reduced search space

Thanks