

Taming Topic Instability For Reasoning Over Unstructured Data in Software Engineering

Amritanshu Agrawal
Computer Science, NCSU
Raleigh, North Carolina
aagrawa8@ncsu.edu

Tim Menzies
Computer Science, NCSU
Raleigh, North Carolina
tim.menzies@gmail.com

Wei Fu
Computer Science, NCSU
Raleigh, North Carolina
wfu@ncsu.edu

Abstract—Topic Modeling has been widely used to identify patterns in a corpus across various fields. The Latent Dirichlet Allocation (LDA) algorithm is a widely used topic modeling algorithm which in practice is called using its defaults. We report that these defaults can lead to highly misleading results in Software Engineering (SE). It is very important to generate stable topics using the Non-Deterministic LDA method. But one of the curse is setting the right parameters that control the LDA. We ran differential evolution (DE) as an optimizer to explore the tuning space (as a first step) then tested it against the default parameter settings. We found that these tunings were remarkably simple, and took hundreds, not thousands of evaluations to obtain very good results. Since the improvements are so large, and (2) the tuning is so simple, that many industries, who are using LDA, should change the standard methods for better throughput and accuracy of their products. The implication for other kinds of topic modeling algorithms is now an open and pressing issue.

Keywords—Topic modeling, Stability, LDA, tuning, differential evolution.

I. INTRODUCTION

In the 21st century, we have now access amount of data about Software Engineering. Most of these data is in an unstructured format [52]. According to International Data Corporation (IDC), shown in figure 1, this number is about 1600 Exabytes. Unstructured data does not have a pre-defined data model and is typically text-heavy. Finding insights among unstructured text is extremely difficult unless we can search, characterize, and classify the text data in a meaningful way. One of the common techniques for finding related topics within unstructured text (an area called topic modeling) is Latent Dirichlet Allocation (LDA).

Topic models built by LDA are non-deterministic in behaviour. Due to non-determinism, topics generated by LDA are not stable [56], [7]. It learns the various distributions which is a problem of Bayesian inference [11]. We are talking about model instability itself rather than results instability. Many industries, working on a product development in text analytics have been widely using it. They want better results in their product. Menzies et al. [47] showed how model instability can give rise to large effort estimations across various projects. And we do need to avoid these large effort estimations.

LDA takes certain set of parameters which try to take us closer to cluster distributions. But it is an art to find the right parameters setting that control the choices within a data miner. But it is very impractical to get the right parameters settings due to larger search space of configurations. Nevertheless, we

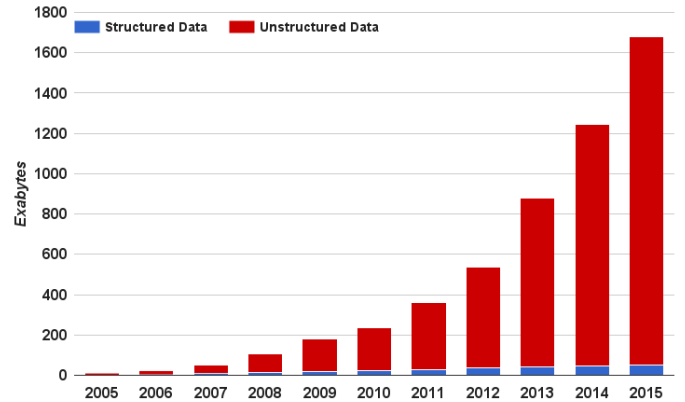


Figure 1: Data Growth 2005-2015

rarely tuned LDA since we reasoned that a LDA's default tunings have been well-explored by the developers of those algorithms (in which case tuning would not lead to stable topic models). Due to larger search space of configurations, it is not even feasible to find such tunings. And that is why, LDA is usually used with off-the-shell default parameters. Many researchers have agreed to the idea of using different configurations in LDA. But, we have seen rare work done when it comes to tune the important parameters of LDA.

Wei et al. [18] have shown that one of the simplest *multiobjective optimizers (differential evolution* [64]) works very well for tuning in certain kinds of software analytics. There has been success stories where DE has helped in tuning the parameters [8], [14]. Prior to this work, our intuition was that tuning would change the behavior of LDA, to some degree. Also, we suspected that tuning would take so long time and be so CPU intensive that the benefits gained would not be worth the effort. We have now seen large instability in 2 data miners, namely, defect predictors and LDA. Results have shown that we can contain the instability to a large extent using DE. So, this makes us wonder Search-Based Software Engineering (SBSE) techniques can help us.

The results of this paper show that the above points are false since:

- Tuning LDA is remarkably simple;
- And can dramatically improve the stability of LDA model.

Those results were found by exploring six research questions:

- **RQ1:** *Do the default settings lead to misleading results?* We will show that topics generated with the default settings are misleading as stability scores start to drop after top 5 terms, which is the basis to define a particular topic generated from LDA.
- **RQ2:** *Does tuning improve the stability scores of LDA?* In our work, we will show dramatic improvement in the stability scores with the tuned parameters rather than going for default parameters.
- **RQ3:** *Does different data need different configurations to make LDA stable? Does it change some predefined parameter values of LDA?* We will show that different sets of configurations are found out by DE. For one of the results, we showed that we need topic size of less than 30 rather than 67 [21].
- **RQ4:** *Is tuning easy?* We have seen that one of the simpler multiobjective optimizers (differential evolution [64]) works very well for tuning in major software analytics.
- **RQ5:** *Is tuning impractically slow?* We achieved dramatic improvements in the stability scores of LDA in less than 300 evaluations (!); i.e., very quickly.
- **RQ6:** *Should data miners be used “off-the-shelf” with their default tunings?* Wei et al. [18] showed that for defect prediction, their answer is an emphatic “no”. And we also conclude that for LDA, it is a “no”.

Based on answers to these questions, we can surely mitigate the instability of non-deterministic LDA. LDA should not be used with “off-the-shelf” default tunings. Any future paper on LDA should include a tuning study. Tuning needs to be repeated whenever data or goals are changed.

Topic Modeling have been used in various spheres of SE. Sun et al. reported a survey of LDA’s usage to support various SE tasks between 2003 and 2015 [66]. The number of papers published across top major conferences and journals are listed in Table I. The survey in figure 2 show that there is an increasing concern in this area.

Venue	Full Name	Count
ICSE	International Conference on Software Engineering	4
CSMR-WCRE / SANER	International Conference on Software Maintenance, Reengineering, and Reverse Engineering / International Conference on Software Analysis, Evolution, and Reengineering	3
ICSM / ICSE	International Conference on Software Maintenance / International Conference on Software Maintenance and Evolution	3
ICPC	International Conference on Program Comprehension	3
ASE	International Conference on Automated Software Engineering	3
ISSRE	International Symposium on Software Reliability Engineering	2
MSR	International Working Conference on Mining Software Repositories	8
OOPSLA	International Conference on Object-Oriented Programming, Systems, Languages, and Applications	1
FSE/ESEC	International Symposium on the Foundations of Software Engineering / European Software Engineering Conference	1
TSE	IEEE Transaction on Software Engineering	1
IST	Information and Software Technology	3
SCP	Science of Computer Programming	2
ESE	Empirical Software Engineering	4

Table I: Survey Venue and Statistics

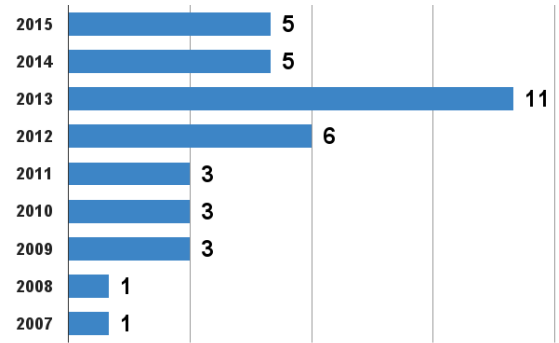


Figure 2: The distribution of the papers published in different years

- Topic modeling is applied in various SE tasks, including source code comprehension, feature location, software defects prediction, developer recommendation, traceability link recovery, re-factoring, software history comprehension, software testing and social software engineering.
- There are works in Requirements Engineering where it was necessary to analyze the text and come up with the important topics [5], [68], [45].
- People have used topic modeling in prioritizing test cases, & identifying the importance of test cases [27], [79], [76].
- Increasingly, it has become very important to have automated tools to do a Systematic Literature Review (SLR) [71]. We found these papers [63], [2], [42] who have used clustering algorithms (topic modeling) to do SLR.

Since topic modeling has been widely used in various spheres of SE, and it is really important to generate stable topics. Prior Papers who reports LDA may have been severely compromised. This gives all of us the motivation of generating stable topics in software engineering which will make product development faster, cheaper and accurate.

II. RELATED WORK

Are the impacts of instability in LDA addressed in the text analytics? To answer that question, in April 2016, we searched scholar.google.com for the conjunction of “lda” and “topics” or “stable” or “unstable” or “coherence”. At the time of this writing (July 2016), we got 475 results. We selected for papers since 2012, we got to 189 results. After discarding the non-SE papers, we read over this sample of 57 highly-cited SE text analytics papers. 28 papers talk about instability in LDA (more details can be found at <https://goo.gl/Bpc6Vb>). We found in that sample was that only two authors acknowledged the impact of tunings. About 35% papers in our sample did not adjust the “off-the-shelf” configuration of the LDA. Among the other 65%, 83% papers, explored the configurations space with brute force. Their was not an effective and faster way of choosing the right set of parameters. Those papers, which accepted the problem of instability in LDA, came to some conclusions which can be found in Table II. Conclusions of remaining papers are categorized in the following subsections.

A. Problems with LDA

- In LDA, learning the various distributions is a problem of Bayesian inference. It uses the alternative inference techniques like Gibbs sampling [74], [25] and variational expectation maximization (VEM) [50]. These techniques

Reference ID	Year	Citations	Venues	Mentions instability in LDA?	Talks about using-off-the-shell parameters?	Does tuning?	Conclusion	Tasks / Use cases
[62]	2011	112	WCRE	Y	Y	N	Explored Configurations without any explanation.	Bug Localisation
[56]	2010	108	MSR	Y	Y	N	Explored Configurations without any explanation. Reported their results using multiple experiments.	Traceability Link recovery
[7]	2014	96	ESE	Y	Y	N	Explored Configurations without any explanation. Choosing right set of parameters is a difficult task.	Stackoverflow Q&A data analysis
[60]	2013	75	ICSE	Y	Y	Y	Uses GA to tune parameters. They determine the near-optimal configuration for LDA in the context of only some important SE tasks.	Finding near-optimal configurations
[20]	2013	61	ICSE	Y	Y	N	Explored Configurations without any explanation.	Software Requirements Analysis
[29]	2011	45	MSR	Y	Y	N	They validated the topic labelling techniques using multiple experiments.	Software Artifacts Analysis
[26]	2014	44	RE	Y	Y	N	Explored Configurations without any explanation.	Requirements Engineering
[67]	2011	44	ICSE	Y	Y	N	Open issue to choose optimal parameters.	A review on LDA
[13]	2012	35	MSR	Y	Y	N	Choosing the optimal number of topics is difficult.	Software Defects Prediction
[68]	2014	35	SCP	Y	Y	N	Explored Configurations without any explanation.	Software Artifacts Analysis
[69]	2014	31	ESE	Y	Y	N	Choosing right set of parameters is a difficult task.	Software Testing
[6]	2009	29	MSR	Y	Y	N	Explored Configurations without any explanation and accepted to the fact their results were better because of the corpus they used.	Software History Comprehension
[39]	2013	27	ESEC/FSE	Y	Y	Y	Explored Configurations using LDA-GA.	Traceability Link recovery
[38]	2013	20	MSR	Y	Y	N	In Future, they planned to use LDA-GA.	Stackoverflow Q A data analysis
[32]	2014	15	WebSci	Y	Y	N	Explored Configurations without any explanation.	Social Software Engineering
[28]	2012	13	ICSM	Y	Y	N	Explored Configurations without any explanation (Just with no. of topics).	Software Requirements Analysis
[24]	2013	13	SCP	Y	Y	N	Their work focused on optimizing LDA's topic count parameter.	Source Code Comprehension
[19]	2015	6	IST	Y	Y	N	Explored Configurations without any explanation. Choosing right set of parameters is a difficult task.	Software re-factoring
[21]	2016	5	Computer Science Review	Y	Y	N	Explored Configurations without any explanation.	Bibliometric and citation analysis
[37]	2014	5	ISSRE	N	Y	N	Explored Configurations without any explanation.	Bug Localisation
[55]	2015	3	JIS	Y	Y	N	They improvised LDA into ISLDA which gave stability across different runs.	Social Software Engineering
[65]	2015	2	IST	Y	Y	Y	Explored Configurations using LDA-GA.	Software Artifacts Analysis
[12]	2016	0	JSS	N	Y	N	Explored Configurations without any explanation. Choosing right set of parameters is a difficult task.	Software Defects Prediction

Table II: Literature Survey

in LDA try to maximize the resulting lower bound on the log likelihood [11]. The expected complete log likelihood of the data have many local maximas, which leads to different distributions and in turn leads to instability in the output. The distributions which are found with the help of LDA, resulting from the same dataset with the same vocabulary and model parameters, any differences between them are entirely due to the randomness in inference techniques [32]. This randomness affects word and document ratios. There is a problem of finding the optimal number of clusters, but different configuration parameters may lead us to the stable topics.

- The topics learned by LDA sometimes are difficult to interpret by end users [77], [60]. We want topics to be finer grained (more number of topics) or coarse grained (less number of topics) according to the use case. Rationality about the topics instability is important if an industry is using topic modeling in all their use cases [34], [58]. We do not want vaguely defined topics.
- Online learning [30] for LDA uses constant memory requirements and empirically converges faster than batch collapsed Gibbs sampling. However, it still requires a full pass through the entire corpus each iteration. It can therefore be slow to apply to very large datasets, and is not naturally suited to settings where new data is constantly arriving. New data arriving can induce more randomness in the inference techniques.
- We observed that this is just not happening with only particular programming language or particular library. This instability is happening irrespective of the tools in which LDA has been implemented. Some of the papers mentioned using GibbsLDA++ written in C++ and they

observed the same instability [40], [70], [26]. Some papers mentioned about using Python implementation of LDA and observed the same instability [26]. LDA implemented in JAVA language produced the same incoherent topics [43], [29]. And, we observed this problem is in the Scikit-Learn version of LDA, implemented in Python, as well as Spark Mllib library.

B. Tuning: Important and Ignored

This section argues that tuning is an under-explored software analytics. The impact of tuning is well understood [10]. Yet issues of tuning are poorly addressed in the data analytics literature. When we tune a data miner, what we are really doing is changing how a learner applies its heuristics. This means tuned data miners use different heuristics, which means they ignore different possible models, which means they return different models; i.e. how we learn changes what we learn.

Wei et al. [18] shows that only few authors acknowledged the impact of tunings. A few papers did acknowledge that one data miner may not be appropriate for all data sets. Those papers tested different “off-the-shelf” data miners on the same data set. For example, Elish et al. [16] compared support vector machines to other data miners for the purposes of defect prediction. Those papers tested different “off-the-shelf” data miners on the same data set. Wei et al. [18] showed that finding useful tunings is very easy. This result is both novel and unexpected. Learners are very amenable to tuning. Hence, they are also very amenable to significant performance improvements. Given the low number of evaluations required, they asserted that tuning should be standard practice for anyone using data miners.

C. WorkAround to LDA Instability Problem

1) Unsupervised:

Solution 1: Most papers [44], [43], [26] manually accessed the topics and then used for further experiments. Some made use of Amazon Mechanical Turk to create gold-standard coherence judgements [34]. This solution is related to results stability rather than doing model stability

This workaround will take extensive manual effort and time which is not feasible for many. In the age of automation, we need an efficient way of making LDA stable.

Solution 2: Users have external knowledge regarding word correlation, which can be taken into account to improve the semantic coherence of topic modeling methods. SCLDA [77] can handle different kinds of knowledge such as word correlation, document correlation, document label and so on. One advantage of SCLDA over existing methods is that it is very fast to converge.

User knowledge is valuable but highly subjective and we seek fully automated methods.

Solution 3: The most common parameters in the LDA are number of clusters (k), number of iterations (n), document topic prior (α), and word topic prior (β). The other solution is tuning these parameters of LDA. Some researchers [60], [39], [65] used Genetic Algorithm (GA) to tune the parameters and then using different configurations for LDA. Some people [20], [70] achieved higher stability by just increasing the number of cluster size.

We seek stability and LDA-GA uses stochastics. But we found tuning combined with Wei et al. [18] work to have better advantage than the other solutions. And in our work, we have used this solution.

2) Supervised:

Solution 3: We found so many improvised frameworks of LDA namely, *Labeled LDA* [61], *Dirichlet Forest LDA* [3], *Logic LDA* [4], *Quad-LDA* [54], *NMF-LDA* [75], *Interactive Topic Modeling* (ITM) [31], *MRTF* [15] and many more, helping in generating better stable topics.

We will not be able to use any of the above supervised solutions as our datasets are not labeled. And all these frameworks work best with the classified datasets. This is going to be our future work which we will be addressing next.

D. Evaluation

To evaluate the topics coherence or the cluster stability using LDA, there has been number of evaluation measures proposed. There is a direct approach, by asking people about topics, and an indirect approach by evaluating *pointwise mutual information (PMI)* [34], [58] between the topic words. PMI being an automatic method, we are not sure of exact

details of it which made us to not use this measure. We could not use direct approach, due to resource limitation to ask an expert for the type of datasets we used.

Perplexity is the inverse of the geometric mean per-word likelihood. monotonically decreasing in the likelihood of the data. It shows how well topic-word and word-document distributions predict new test samples. The smaller the perplexity, the better (less uniform) is the LDA model. The usual trend is that as the value of perplexity drops, the number of topics should grow [32]. Researchers caution that the value of perplexity doesn't remain constant with different topic size and with dictionary sizes [32], [80]. A lot depend on the code implementation of perplexity and the type of datasets used. Since, we are using different implementations of LDA as well as different datasets, we are not using Perplexity as evaluation measure.

Another approach is that we use Jaccard Similarity for evaluating stability [58], [20]. We define our measure similar to it. We say topics are stable, when there are x times $n\%$ of terms overlap in a topic across m runs. We can see an example of 5 terms overlap in figure 3. We take median of all the scores generated by evaluating topic. In this example, figure 3, median score is 1.0 for each run. So overall median score comes out to be 1.0. We always select top 10 terms within a topic suggested by some researchers [60], [41]. It makes sense to us that to have a successful topic model, each topic can be identified using top few terms. *If we get different words then it can mean different topics which can put any product based on topic modeling at risk.*

To understand our results, we have defined some terminologies. We call the above median score as **Raw (\mathbb{R}) score** which is 1.0. The \mathbb{R} score can come from tuning and untuned results defined similarly. We also define a **Delta Score** which is the absolute difference between **Raw Score_{tuned}** - **Raw Score_{untuned}**.

RUN 1:		\mathbb{R} SCORES
Topic 0:	glori telemetri command spacecraft trace smrd tim spec pip parent	2/4 = 0.50
Topic 1:	spec smrd parent child glori artifact referenc verif matrix data	3/4 = 0.75
Topic 2:	test case accuraci roll document glori plan pitch yaw valu	1/4 = 0.25
Topic 3:	command specifi ground softwar telemetri initi pip data configur band	4/4 = 1.00
RUN 2:		
Topic 0:	command specifi ground softwar initi telemetri data pip configur band	4/4 = 1.00
Topic 1:	document test glori initi plan roll case pitch yaw point	1/4 = 0.25
Topic 2:	spec smrd parent child glori artifact verif point matrix spacecraft	3/4 = 0.75
Topic 3:	pip capabl glori ground command smrd spec tim list spacecraft	2/4 = 0.50
RUN 3:		
Topic 0:	spec smrd parent child glori referenc artifact verif matrix spacecraft	3/4 = 0.75
Topic 1:	command telemetri pip ground spacecraft modesoftware document spec glori	2/4 = 0.50
Topic 2:	command specifi ground softwar initi telemetri data configur pip band	4/4 = 1.00
Topic 3:	test initi accuraci glori roll plan matrix yaw pitch valu	1/4 = 0.25
RUN 4:		
Topic 0:	command specifi tim pip ground telemetri glori includ softwar document	2/4 = 0.50
Topic 1:	test initi glori plan roll accuraci case pip point yaw	1/4 = 0.25
Topic 2:	command specifi ground softwar telemetri initi data pip configur band	4/4 = 1.00
Topic 3:	spec smrd child glori artifact referenc parent verif matrix data	3/4 = 0.75

Figure 3: Example of 5 terms overlap across 4 runs and its stability score

E. Topic Modeling with LDA

Latent Dirichlet Allocation (LDA) is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar. It learns the various distributions (the set of topics, their associated word probabilities, the topic of each word, and the particular topic mixture of each document) and is a problem of Bayesian inference. We ran 2 different inference techniques like Gibbs sampling and variational expectation maximization (VEM) on python platform. We also ran VEM onto Spark platform. We chose 2 different inference techniques on the recommendation from Layman [35], [36] which happened over personal communication.

The pseudocode for untuned LDA is shown in Algorithm 1 with the default set of parameters. In the following description, superscript numbers denote lines in the pseudocode. The data is shuffled everytime LDA is run^{L5}. Data is in the form of term frequency scores of each word per document. Shuffling is done in order to induce maximum variance among different runs of LDA. Topics^{L6} is a list of list which contains topics from all the different runs. A stability score is evaluated on every 10 runs, and this process is continued 10 times. At the end, median score is selected as the untuned final score^{L3-L11}.

Algorithm 1 Pseudocode for untuned LDA with Default Parameters

Input: *no.oftersms, Data*
Output: *Final_Score*

```

1: function LDAScore( no.oftersms, Data)
2:   Score  $\leftarrow \emptyset$ 
3:   for j = 0 to 10 do
4:     for i = 0 to 10 do
5:       data  $\leftarrow$  shuffle(Data)
6:       Topics.append(lda(k = 10,  $\alpha$  = none,  $\beta$  = none))
7:     end for
8:     Score.append(Overlap(Topics, no.oftersms, l[0]))
9:   end for
10:  Final_Score  $\leftarrow$  median(Score)
11:  return Final_Score
12: end function
```

F. Tuning Algorithms

In section II-C, we talked about our reasoning why are we choosing tuning the parameters of LDA as the solution. Here is a list of optimizers used widely in research: simulated annealing [17], [48]; various genetic algorithms [23] augmented by techniques such as differential evolution [64], tabu search and scatter search [22], [9], [51], [53]; particle swarm optimization [59]; numerous decomposition approaches that use heuristics to decompose the total space into small problems, then apply a response surface methods [33], [81]. Of these, the simplest ones are simulated annealing (SA) and differential evolution (DE), each of which can be coded in less than a page of some highlevel scripting language. Our reading of the current literature is that there are more advocates for differential evolution than SA. For example, Vesterstrom and Thomsen [73] found DE to be competitive with particle swarm optimization and other GAs. DEs have been applied before for parameter tuning (e.g. see [57], [14], [18]) but this is the first time they have been applied to tune the configurations of LDA.

We will tune hyperparameters of LDA using DE. DE also comes up with parameters like *number of iterations*, *mutation*,

crossover rate, *population size*, and *different selection strategies*. The objective of LDA is to maximize the stability scores. Based on the literature review, we found that we don't need to have topics more than 10 terms. The parameters which affect the most in LDA and that can be tuned can be referred in Table III. We also used online variational Bayesian approximation as well as gibbs sampling for LDA, with random input row order of data to make sure that in each run the learning is different. This gave us more concrete results as few tuned set of parameters which we found performed good across each run.

The pseudocode for DE is shown in Algorithm 2. DE evolves a *NewGeneration* of candidates from a current Population. Our DE will terminate after a certain number of evaluations. We ran initial experiments to define this constant. We will talk about this in section V. Each candidate solution in the Population is a pair of (Tunings, Scores). Tunings are selected from Table III and Scores come from the stability score after running our LDA experiment^{L23-L30}.

The premise of DE is that the best way to mutate the existing tunings is to Extrapolate^{L31} between current solutions. Three solutions a, b, c are selected at random. For each tuning parameter i, at some probability cr, we replace the old tuning x_i with y_i . For booleans, we use $y_i = x_i$ (see line 39). For numerics, $y_i = a_i + f \times (b_i - c_i)$ where f is a parameter controlling crossover. The trim function^{L41} limits the new value to the legal range min..max of that parameter.

The main loop of DE^{L9} runs over the Population, replacing old items with new Candidates (if new candidate is better). This means that, as the loop progresses, the Population is full of increasingly more valuable solutions. This, in turn, also improves the candidates, which are Extrapolated from the Population.

III. EXPERIMENTATION

A. Data Sets

To answer our research questions, we will be working on couple of open source datasets to verify that tuning does help the stability scores. All these datasets have been preprocessed using the normal steps of stopwords removal, stemming, and then considering the top 5% of tf-idf scores. More details about the statistics of datasets can be found in the table IV. The datasets used are:-

- **PITS:** This is a text mining data set generated from project and issue tracking system (PITS) reports [46], [49]. This data is about the bugs and changes in the code, to submit and review patches, to manage quality assurance, to support communication between developers, etc. The aim is to identify the top topics which can identify each severity separately. The dataset can be downloaded from this repository [1].
- **StackOverflow:** StackOverflow is a privately held website, the flagship site of the Stack Exchange Network which features questions and answers on a wide range of topics in computer programming. The data can be downloaded from various sources^{1 2}. This dataset is so

¹<http://data.stackexchange.com>

²<http://blog.stackoverflow.com/tags/cc-wiki-dump/>

Parameters	Default	Tuning Range	Description
n_topics (k)	10	[10,100]	Number of topics or cluster size
doc_topic_prior (α)	None	[0,1]	Prior of document topic distribution. This is called alpha
topic_word_prior (β)	None	[0,1]	Prior of topic word distribution. This is called beta

Table III: List of parameters tuned by this paper

Algorithm 2 Pseudocode for DE with a constant number of evaluations

```

Input: np= 10, f= 0.7, cr= 0.3, iter= 3, Goal  $\in$  Finding maximum score
Output:  $S_{best}$ ,  $final\_generation$ 
1: function DE(np, f, cr, iter, Goal)
2:   Cur_Gen  $\leftarrow \emptyset$ 
3:   Population  $\leftarrow$  InitializePopulation(np)
4:   for i = 0 to np - 1 do
5:     Cur_Gen.append(Population[i], ldscore(Population[i],
term, Data))
6:   end for
7:   for i = 0 to iter do
8:     NewGeneration  $\leftarrow \emptyset$ 
9:     for j = 0 to np - 1 do
10:       $S_i \leftarrow$  Extrapolate(Population[j], Population, cr, f, np)
11:      if ldscore( $S_i$ )  $\geq$  Cur_Gen[j][1] then
12:        NewGeneration.append( $S_i$ , ldscore( $S_i$ , term, Data))
13:      else
14:        NewGeneration.append(Cur_Gen[j])
15:      end if
16:    Cur_Gen  $\leftarrow$  NewGeneration
17:  end for
18:  end for
19:   $S_{best} \leftarrow$  GetBestSolution(Cur_Gen)
20:   $final\_generation \leftarrow$  Cur_Gen
21:  return  $S_{best}$ ,  $final\_generation$ 
22: end function
23: function LDAScore(l, term, Data)
24:   Topics  $\leftarrow \emptyset$ 
25:   for i = 0 to 10 do
26:     data  $\leftarrow$  shuffle(Data)
27:     Topics.append(lda(k = l[0],  $\alpha$  = l[1],  $\beta$  = l[2]))
28:   end for
29:   return Overlap(Topics, term, l[0])
30: end function
31: function EXTRAPOLATE(old, pop, cr, f, np)
32:   a, b, c  $\leftarrow$  threeOthers(pop, old)
33:   newf  $\leftarrow \emptyset$ 
34:   for i = 0 to np - 1 do
35:     if cr  $\leq$  random() then
36:       newf.append(old[i])
37:     else
38:       if typeof(old[i]) == bool then
39:         newf.append(not old[i])
40:       else
41:         newf.append(trim(i, (a[i] + f * (b[i] - c[i]))))
42:       end if
43:     end if
44:   end for
45:   return newf
46: end function

```

big that it can only be run on Spark with a cluster of nodes to reduce the runtime. Results of this dataset is only coming from Spark mllib implementation.

- **Citemap:** Citemap dataset contains titles and abstracts of papers from a database of 11 conferences from 1993-2013. This data was obtained in the form of a sql dump from the work of Vasilescu et al. [72]. Dataset can be downloaded from ³.

Datasets	Size	
	Before Preprocessing	After Preprocessing
PitsA	1.2MB	292 KB
PitsB	704KB	188KB
PitsC	143KB	37KB
PitsD	107KB	26KB
PitsE	650KB	216KB
PitsF	549KB	217KB
Citemap	8.6MB	3.7MB
Stackoverflow	7GB	527MB

Table IV: Statistics on the datasets

IV. EXPERIMENTAL RESULTS

Note in the following, we report these results after running on Local machine just for PITS and Citemap datasets. These results are implemented using VEM method of LDA in Python (Scikit-learn). The StackOverflow dataset is so large that it can only be run on Spark cluster. We have specifically mentioned if any of the other results are coming from Gibbs Implementation or Spark platform.

A. RQ1: Do the default settings lead to misleading results?

In figure 4, X-axis represents number of terms overlap. Y-axis represents the \mathcal{R} scores. For definition of \mathcal{R} and *Delta Scores*, you can refer to Section II-D. Different lines are for the different datasets. From the figure 4, it can be clearly seen that instability occurs in the untuned results. Researchers [60], [41] said that, 10 terms per topic are ideal to define a topic of interest but we can see \mathcal{R} scores start dropping after 5 terms. With default settings, it can surely lead to misleading results.

Result 1

There exists systematic errors in prior LDA results.

B. RQ2: Does tuning improve the stability scores of LDA?

In figure 5, x-axis represents number of terms overlap and y-axis represents the \mathcal{R} scores. Different lines are for the different datasets. We can see that tuning results are either always above the default settings (untuned) or stayed the same. Since figure 5 is hard to read, we are reporting delta improvement in figure 6. In figure 6, x-axis represents number of terms overlap and y-axis represents the delta improvement which is (*tuning* - *untuned*) results. If the lines are above x-axis, that means tuning performed better than the untuned. Different lines are for the different datasets. It clearly shows that the answer to RQ2 is “yes” - tuning has a positive effect on stability scores. Tuning either helped it or remained the same but it never had any negative effect. The dramatic improvement started showing after 5 terms overlap. The most improvement

³<https://github.com/ai-se/citemap/blob/master/citemap.csv>

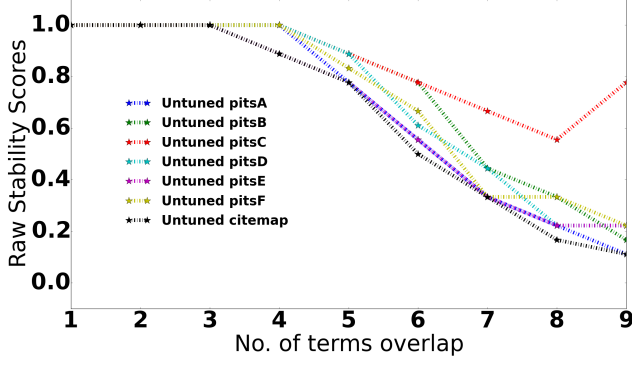


Figure 4: Terms vs Raw Scores between tuning and default parameters

which we observed was in PitsD dataset for 8 terms overlap of about 80%. Delta improvement remained 0 till 5 terms overlap for most number of datasets. That is why we recommend that we shouldn't list more than 5 terms per topic, as stability starts dropping.

Result 2

Based on Figure 6, we highly recommend tuning in future LDA.

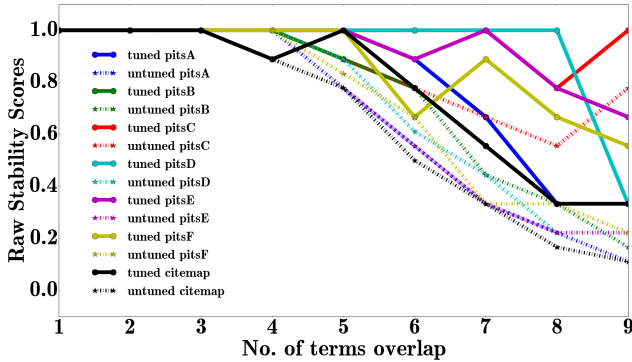


Figure 5: Terms vs Raw Scores between tuning and default parameters

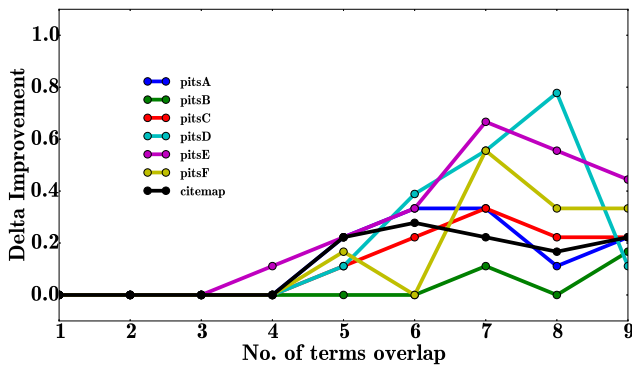


Figure 6: Terms vs Delta Improvement between tuning and default parameters

C. RQ3: Does different data need different configurations to make LDA stable? Does it change some predefined parameter values of lda?

Figures 7, 8, and 9, x-axis represents different datasets and y-axis represents values of median and IQR. Each figure shows the variation in the values of number of topics (k), alpha (α) and beta (β). These results are for 5 terms overlap. From these figures, it clearly shows how tuning selects the different ranges of values of parameters. These figures represent a very high Median Score and as well as very high interquartile range (IQR) for different datasets. This means that tuning helped us to arrive at different sets of configurations making the topic generation more constant. In figure 7 for pitsB dataset, we are seeing 0 IQR, that means tuning helped us to find only 1 set of configurations which got us to that Stable score. We found many such similar numbers throughout the results. (Refer this link at <https://goo.gl/Nin5pV> for elaborate results).

The other result which we found is in the citemap dataset. If we look in the spreadsheet, from Row 121-129, we see that we are getting good stable scores within 30 topic size. But in the [21] paper, they found that we will need a topic size of 67. Definitely, it can change predefined parameters of LDA. We also observed the similar trend for other term overlaps as well as other sampling methods.

Result 3

Tuning the parameters of LDA needs to be repeated whenever data or goals are changed.

D. RQ4: Is tuning easy?

In terms of the search space explored via tuning, it is much smaller. To see this, recall from Algorithm 2 that DE explores a Population of size $np = 10$. This is a very small population size since Rainer Storn (one of the inventors of DE) recommends setting np to be ten times larger than the number of parameters being optimized. We also ran experiments with different sets of F , CR , and population size (np), and we see that it didn't change much.

We have shown results only for Citemap dataset which can be referred in Figure 12. For other datasets, we can refer at <https://goo.gl/HQNASF>. In Figure 12, different lines show different combinations of F , CR and Population size. F is selected either 0.3 or 0.7, and similarly CR is selected 0.3 or 0.7, and population size is selected either 10 or 30 (which is 10 times the number of parameters being optimized). These numbers are according to the original Rainer Storn [64] recommended settings.

After reviewing the results from all the datasets, we can say that there isn't much of an improvement by using different F , CR , and Population size. So our all other experiments used $F = 0.7$, $CR = 0.3$ and $Popsiz = 10$.

E. RQ5: Is tuning impractically slow?

Figure 10 and 11, x-axis represents different datasets and y-axis represents the runtimes in seconds (*Log scale*). From the table V, we will show that we just need 300 evaluations to

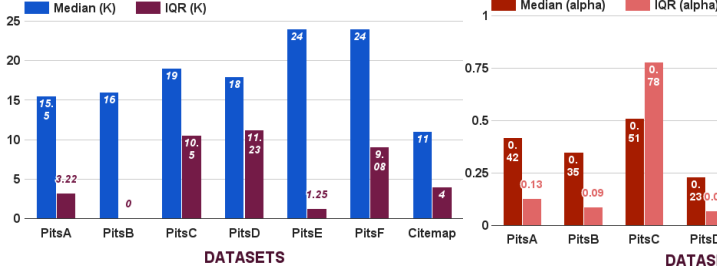


Figure 7: Datasets vs Parameter (k) variation

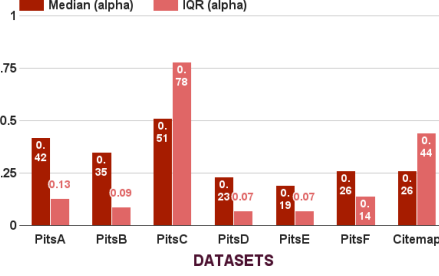


Figure 8: Datasets vs Parameter (α) variation

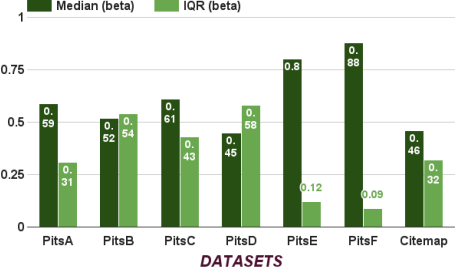


Figure 9: Datasets vs Parameter (β) variation

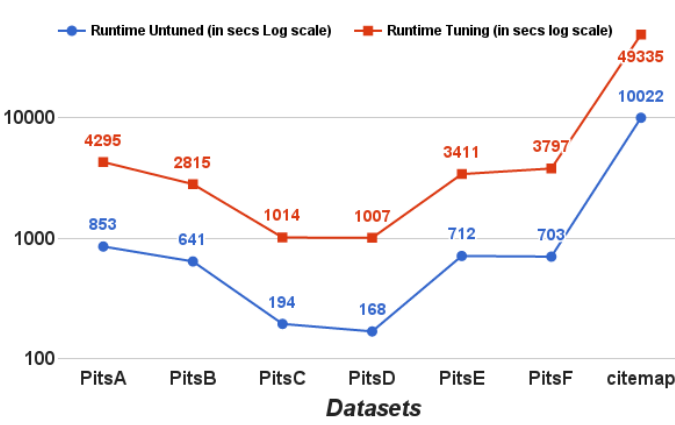


Figure 10: Gibbs: Datasets vs Runtimes

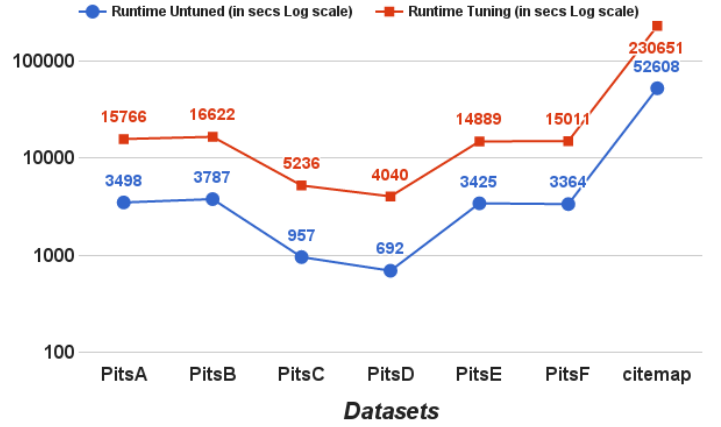


Figure 11: VEM: Datasets vs Runtimes

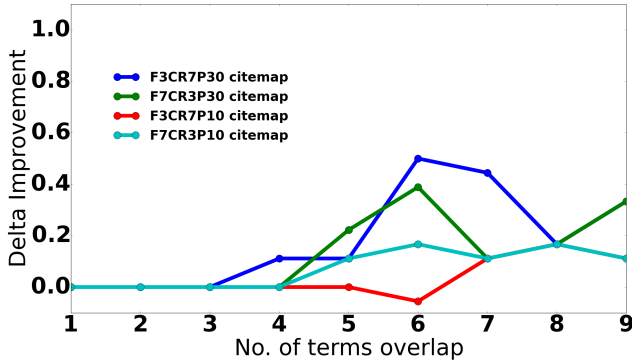


Figure 12: Terms vs Delta Improvement using Different settings of DE

do tuning. Using this criteria, we can see that tuning runtimes is only about 5 times the runtimes without tuning. Figure 10 is with Gibbs implemented in Python, and Figure 11 is with VEM implemented in Python.

Result 4

These tuning runtimes are consistent with Wei et al. [18] which makes tuning feasible.

F. RQ6: Should data miners be used “off-the-shelf” with their default tunings?

From figures 7, 8, and 9, we can easily see that for different datasets, we need different parameters. Such large IQRs show that we get quite varied ranges of parameters. Hence, we answer RQ6 as “no” since, to achieve the improvements seen in this paper, tuning has to be repeated whenever the goals or data sets are changed. Given this requirement to repeatedly run tuning, it is fortunate that (as shown above) tuning is so easy and so fast.

V. THREATS TO VALIDITY

Is it a quirk of the implementation? This instability problem remained the same across any implementation of LDA. Figure 13, is with VEM method implemented in Spark. In this figure, x-axis represents different term overlaps and y-axis represents the \mathcal{R} scores. Dashed lines represents the untuned results and the solid lines represents the tuned results. We can see that tuning helped in achieving stable topic generation. The results of other datasets, can be referred at <https://goo.gl/UVaql1>. So it is not a quirk of implementation.

Is it a quirk of the sampling method used? In Figure 14, dashed lines represents the Gibbs implementation and solid line represents the tradition VEM method. The y-axis represents the delta improvement between tuning and untuned results. We can see that even after 5 terms, the stability

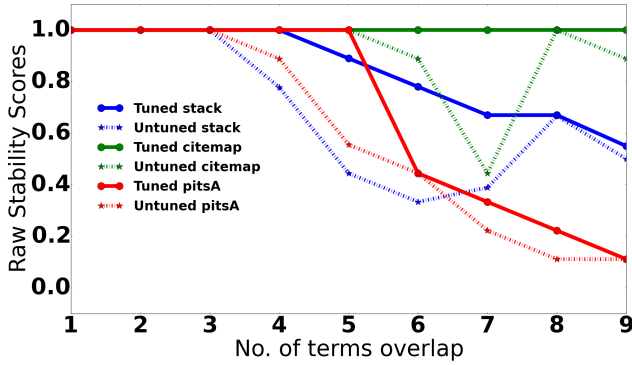


Figure 13: Spark Results

score went down in both Gibbs Sampling and VEM. We also saw the same magnitude of improvements. So it is not a quirk of a particular sampling method. Here results of only 3 datasets are shown. For other datasets, results can be referred at <https://goo.gl/faYAcg>. So it is not a quirk of implementation.

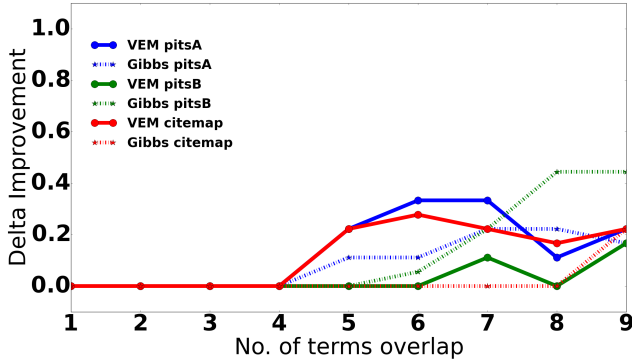


Figure 14: GIBBS vs VEM

Terminating criteria for DE From table III, we can see that after 300 evaluations, there was no improvement in the stability scores across all the datasets and it stayed the same.

Datasets\Evaluations	100	200	300	400
PitsA	0.9	0.9	1.0	1.0
PitsB	0.9	0.9	0.9	1.0
PitsC	0.9	1.0	1.0	1.0
PitsD	0.9	1.0	1.0	1.0
PitsE	0.9	0.9	1.0	1.0
PitsF	0.9	0.9	0.9	0.9
Ciemap	0.67	0.67	0.77	0.77
Stackoverflow	0.6	0.7	0.8	0.8

Table V: Evaluations vs Stability Scores

VI. CONCLUSION AND FUTURE WORK

Our exploration of the six research questions listed in the introduction shows that when doing topic modeling, analytics

without parameter tuning are considered harmful and misleading. As more data getting generated day by day, it is really necessary to do accurate topic modeling. Now, with the help of tuning, we can generate stable topics to quite a good extent. Now we can delve into prior papers to improve their work. We can use the more stable LDA to find features and fed into a classifier to get improved precision and recall [12], [63]

This paper showed that tuning improved the stability scores of LDA, sometimes the improvement is quite dramatic (about 80%). This paper also highlighted that we can now select the right set of parameters for different datasets to get stable topics. This paper combined with Wei et al. [18] suggests that data miners should not be used off-the-shell with their default tunings.

As to future work, it is now important to explore the implications of these stable topics generated by LDA in a product development. We can now actually work with other workarounds mentioned in the section II-C2 with classified datasets. To improve stability, we can even try to improve the actual clusters to get better results.

This paper just investigated on stability of LDA using one optimizer. Hence, we can make no claim that DE is the best optimizer for all data miners. Rather, our point is that there exists at least some learners whose performance can be dramatically improved by at least one simple optimization scheme. And there are already claims of other unstable data miners like Decision Tree Learning, Neural Networks, and Bayesian Learning [78] where optimizers can mitigate instability into all these miners. We also hope that this work inspires much future work as this community develops and debugs best practices for tuning software analytics.

REFERENCES

- [1] The promise repository of empirical software engineering data, 2015.
- [2] A. Al-Zubidy and J. C. Carver. Review of systematic literature review tools.
- [3] D. Andrzejewski, X. Zhu, and M. Craven. Incorporating domain knowledge into topic modeling via dirichlet forest priors. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 25–32. ACM, 2009.
- [4] D. Andrzejewski, X. Zhu, M. Craven, and B. Recht. A framework for incorporating general domain knowledge into latent dirichlet allocation using first-order logic. In *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, volume 22, page 1171, 2011.
- [5] H. U. Asuncion, A. U. Asuncion, and R. N. Taylor. Software traceability with topic modeling. In *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering-Volume 1*, pages 95–104. ACM, 2010.
- [6] S. K. Bajracharya and C. V. Lopes. Mining search topics from a code search engine usage log. In *MSR*, pages 111–120. Citeseer, 2009.
- [7] A. Barua, S. W. Thomas, and A. E. Hassan. What are developers talking about? an analysis of topics and trends in stack overflow. *Empirical Software Engineering*, 19(3):619–654, 2014.
- [8] Y. Bazi, N. Alajlan, F. Melgani, H. AlHichri, S. Malek, and R. R. Yager. Differential evolution extreme learning machine for the classification of hyperspectral images. *Geoscience and Remote Sensing Letters, IEEE*, 11(6):1066–1070, 2014.
- [9] R. P. Beausoleil. “moss” multiobjective scatter search applied to non-linear multiple criteria optimization. *European Journal of Operational Research*, 169(2):426–449, 2006.
- [10] J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(Feb):281–305, 2012.

- [11] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. volume 3, pages 993–1022. JMLR. org, 2003.
- [12] T.-H. Chen, W. Shang, M. Nagappan, A. E. Hassan, and S. W. Thomas. Topic-based software defect explanation. *Journal of Systems and Software*, 2016.
- [13] T.-H. Chen, S. W. Thomas, M. Nagappan, and A. E. Hassan. Explaining software defects using topic models. In *Mining Software Repositories (MSR), 2012 9th IEEE Working Conference on*, pages 189–198. IEEE, 2012.
- [14] I. Chiha, J. Ghabi, and N. Liouane. Tuning pid controller with multi-objective differential evolution. In *Communications Control and Signal Processing (ISCCSP), 2012 5th International Symposium on*, pages 1–4. IEEE, 2012.
- [15] H. Daumé III. Markov random topic fields. In *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, pages 293–296. Association for Computational Linguistics, 2009.
- [16] K. O. Elish and M. O. Elish. Predicting defect-prone software modules using support vector machines. *Journal of Systems and Software*, 81(5):649–660, 2008.
- [17] M. S. Feather and T. Menzies. Converging on the optimal attainment of requirements. In *Requirements Engineering, 2002. Proceedings. IEEE Joint International Conference on*, pages 263–270. IEEE, 2002.
- [18] W. Fu, T. Menzies, and X. Shen. Tuning for software analytics: Is it really necessary? *Information and Software Technology*, 76:135–146, 2016.
- [19] Y. Fu, M. Yan, X. Zhang, L. Xu, D. Yang, and J. D. Kymer. Automated classification of software change messages by semi-supervised latent dirichlet allocation. *Information and Software Technology*, 57:369–377, 2015.
- [20] L. V. Galvis Carreño and K. Winbladh. Analysis of user comments: an approach for software requirements evolution. In *Proceedings of the 2013 International Conference on Software Engineering*, pages 582–591. IEEE Press, 2013.
- [21] V. Garousi and M. V. Mäntylä. Citations, research topics and active countries in software engineering: A bibliometrics study. *Computer Science Review*, 19:56–77, 2016.
- [22] F. Glover and C. McMillan. The general employee scheduling problem. an integration of ms and ai. *Computers & operations research*, 13(5):563–573, 1986.
- [23] A. T. Goldberg. On the complexity of the satisfiability problem. 1979.
- [24] S. Grant, J. R. Cordy, and D. B. Skillicorn. Using heuristics to estimate an appropriate number of latent topics in source code analysis. *Science of Computer Programming*, 78(9):1663–1678, 2013.
- [25] T. L. Griffiths and M. Steyvers. Finding scientific topics. *Proceedings of the National academy of Sciences*, 101(suppl 1):5228–5235, 2004.
- [26] E. Guzman and W. Maalej. How do users like this feature? a fine grained sentiment analysis of app reviews. In *Requirements Engineering Conference (RE), 2014 IEEE 22nd International*, pages 153–162. IEEE, 2014.
- [27] H. Hemmati, Z. Fang, and M. V. Mantyla. Prioritizing manual test cases in traditional and rapid release environments. In *Software Testing, Verification and Validation (ICST), 2015 IEEE 8th International Conference on*, pages 1–10. IEEE, 2015.
- [28] A. Hindle, C. Bird, T. Zimmermann, and N. Nagappan. Relating requirements to implementation via topic analysis: Do topics extracted from requirements make sense to managers and developers? In *Software Maintenance (ICSM), 2012 28th IEEE International Conference on*, pages 243–252. IEEE, 2012.
- [29] A. Hindle, N. A. Ernst, M. W. Godfrey, and J. Mylopoulos. Automated topic naming to support cross-project analysis of software maintenance activities. In *Proceedings of the 8th Working Conference on Mining Software Repositories*, pages 163–172. ACM, 2011.
- [30] M. Hoffman, F. R. Bach, and D. M. Blei. Online learning for latent dirichlet allocation. In *advances in neural information processing systems*, pages 856–864, 2010.
- [31] Y. Hu, J. Boyd-Graber, B. Satinoff, and A. Smith. Interactive topic modeling. *Machine learning*, 95(3):423–469, 2014.
- [32] S. Koltcov, O. Koltsova, and S. Nikolenko. Latent dirichlet allocation: stability and applications to studies of user-generated content. In *Proceedings of the 2014 ACM conference on Web science*, pages 161–165. ACM, 2014.
- [33] J. Krall, T. Menzies, and M. Davies. Gale: Geometric active learning for search-based software engineering. *Software Engineering, IEEE Transactions on*, 41(10):1001–1018, 2015.
- [34] J. H. Lau, D. Newman, and T. Baldwin. Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality. In *EACL*, pages 530–539, 2014.
- [35] L. Layman, 2016.
- [36] L. Layman, A. P. Nikora, J. Meek, and T. Menzies. Topic modeling of nasa space system problem reports: research in practice. In *Proceedings of the 13th International Workshop on Mining Software Repositories*, pages 303–314. ACM, 2016.
- [37] T.-D. B. Le, F. Thung, and D. Lo. Predicting effectiveness of ir-based bug localization techniques. In *2014 IEEE 25th International Symposium on Software Reliability Engineering*, pages 335–345. IEEE, 2014.
- [38] M. Linares-Vásquez, B. Dit, and D. Poshyvanyk. An exploratory analysis of mobile development issues using stack overflow. In *Proceedings of the 10th Working Conference on Mining Software Repositories*, pages 93–96. IEEE Press, 2013.
- [39] S. Lohar, S. Amornborvornwong, A. Zisman, and J. Cleland-Huang. Improving trace accuracy through data-driven configuration and composition of tracing features. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, pages 378–388. ACM, 2013.
- [40] S. K. Lukins, N. A. Kraft, and L. H. Etzkorn. Source code retrieval for bug localization using latent dirichlet allocation. In *Reverse Engineering, 2008. WCRE’08. 15th Working Conference on*, pages 155–164. IEEE, 2008.
- [41] S. K. Lukins, N. A. Kraft, and L. H. Etzkorn. Bug localization using latent dirichlet allocation. *Information and Software Technology*, 52(9):972–990, 2010.
- [42] C. Marshall and P. Brereton. Tools to support systematic literature reviews in software engineering: A mapping study. In *2013 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, pages 296–299. IEEE, 2013.
- [43] W. Martin, M. Harman, Y. Jia, F. Sarro, and Y. Zhang. The app sampling problem for app store mining. In *Mining Software Repositories (MSR), 2015 IEEE/ACM 12th Working Conference on*, pages 123–133. IEEE, 2015.
- [44] G. Maskeri, S. Sarkar, and K. Heafield. Mining business topics in source code using latent dirichlet allocation. In *Proceedings of the 1st India software engineering conference*, pages 113–120. ACM, 2008.
- [45] A. K. Massey, J. Eisenstein, A. I. Antón, and P. P. Swire. Automated text mining for requirements analysis of policy documents. In *Requirements Engineering Conference (RE), 2013 21st IEEE International*, pages 4–13. IEEE, 2013.
- [46] T. Menzies. Improving iv&v techniques through the analysis of project anomalies: Text mining pits issue reports-final report. 2008.
- [47] T. Menzies, A. Butcher, D. Cok, A. Marcus, L. Layman, F. Shull, B. Turhan, and T. Zimmermann. Local vs. global lessons for defect prediction and effort estimation,” ieee trans. software eng., preprint, published online dec. 2012; <http://goo.gl/k6qno>. tim menzies is a full professor in computer science at the lane. In *Department of Computer Science and Electrical Engineering, West Virginia University. Contact. Citeseer*.
- [48] T. Menzies, O. Elrawas, J. Hihn, M. Feather, R. Madachy, and B. Boehm. The business case for automated software engineering. In *Proceedings of the twenty-second IEEE/ACM international conference on Automated software engineering*, pages 303–312. ACM, 2007.
- [49] T. Menzies and A. Marcus. Automated severity assessment of software defect reports. In *Software Maintenance, 2008. ICSM 2008. IEEE International Conference on*, pages 346–355. IEEE, 2008.
- [50] T. Minka and J. Lafferty. Expectation-propagation for the generative aspect model. In *Proceedings of the Eighteenth conference on Uncertainty in artificial intelligence*, pages 352–359. Morgan Kaufmann Publishers Inc., 2002.
- [51] J. Molina, M. Laguna, R. Martí, and R. Caballero. Sspmo: A scatter tabu search procedure for non-linear multiobjective optimization. *INFORMS Journal on Computing*, 19(1):91–100, 2007.

- [52] A. Nadkarni and N. Yezhkova. Structured versus unstructured data: The balance of power continues to shift. *IDC (Industry Development and Models)* Mar, 2014.
- [53] A. J. Nebro, F. Luna, E. Alba, B. Dorronsoro, J. J. Durillo, and A. Beham. Abyss: Adapting scatter search to multiobjective optimization. *Evolutionary Computation, IEEE Transactions on*, 12(4):439–457, 2008.
- [54] D. Newman, E. V. Bonilla, and W. Buntine. Improving topic coherence with regularized topic models. In *Advances in neural information processing systems*, pages 496–504, 2011.
- [55] S. I. Nikolenko, S. Koltcov, and O. Koltsova. Topic modelling for qualitative studies. *Journal of Information Science*, page 0165551515617393, 2015.
- [56] R. Oliveto, M. Gethers, D. Poshyvanyk, and A. De Lucia. On the equivalence of information retrieval methods for automated traceability link recovery. In *Program Comprehension (ICPC), 2010 IEEE 18th International Conference on*, pages 68–71. IEEE, 2010.
- [57] M. G. Omran, A. P. Engelbrecht, and A. Salman. Differential evolution methods for unsupervised image classification. In *Evolutionary Computation, 2005. The 2005 IEEE Congress on*, volume 2, pages 966–973. IEEE, 2005.
- [58] D. O’Callaghan, D. Greene, J. Carthy, and P. Cunningham. An analysis of the coherence of descriptors in topic modeling. *Expert Systems with Applications*, 42(13):5645–5657, 2015.
- [59] H. Pan, M. Zheng, and X. Han. Particle swarm-simulated annealing fusion algorithm and its application in function optimization. In *Computer Science and Software Engineering, 2008 International Conference on*, volume 1, pages 78–81. IEEE, 2008.
- [60] A. Panichella, B. Dit, R. Oliveto, M. Di Penta, D. Poshyvanyk, and A. De Lucia. How to effectively use topic models for software engineering tasks? an approach based on genetic algorithms. In *Proceedings of the 2013 International Conference on Software Engineering*, pages 522–531. IEEE Press, 2013.
- [61] D. Ramage, D. Hall, R. Nallapati, and C. D. Manning. Labeled lda: A supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1-Volume 1*, pages 248–256. Association for Computational Linguistics, 2009.
- [62] S. Rao and A. Kak. Retrieval from software libraries for bug localization: a comparative study of generic and composite text models. In *Proceedings of the 8th Working Conference on Mining Software Repositories*, pages 43–52. ACM, 2011.
- [63] A. Restificar and S. Ananiadou. Inferring appropriate eligibility criteria in clinical trial protocols without labeled data. In *Proceedings of the ACM sixth international workshop on Data and text mining in biomedical informatics*, pages 21–28. ACM, 2012.
- [64] R. Storn and K. Price. Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces. *Journal of global optimization*, 11(4):341–359, 1997.
- [65] X. Sun, B. Li, H. Leung, B. Li, and Y. Li. Msr4sm: Using topic models to effectively mining software repositories for software maintenance tasks. *Information and Software Technology*, 66:1–12, 2015.
- [66] X. Sun, X. Liu, B. Li, Y. Duan, H. Yang, and J. Hu. Exploring topic models in software engineering data analysis: A survey. In *2016 17th IEEE/ACIS International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD)*, pages 357–362. IEEE, 2016.
- [67] S. W. Thomas. Mining software repositories using topic models. In *Proceedings of the 33rd International Conference on Software Engineering*, pages 1138–1139. ACM, 2011.
- [68] S. W. Thomas, B. Adams, A. E. Hassan, and D. Blostein. Studying software evolution using topic models. *Science of Computer Programming*, 80:457–479, 2014.
- [69] S. W. Thomas, H. Hemmati, A. E. Hassan, and D. Blostein. Static test case prioritization using topic models. *Empirical Software Engineering*, 19(1):182–212, 2014.
- [70] K. Tian, M. Reville, and D. Poshyvanyk. Using latent dirichlet allocation for automatic categorization of software. In *Mining Software Repositories, 2009. MSR’09. 6th IEEE International Working Conference on*, pages 163–166. IEEE, 2009.
- [71] G. Tsafnat, P. Glasziou, M. K. Choong, A. Dunn, F. Galgani, and E. Coiera. Systematic review automation technologies. *Systematic reviews*, 3(1):1, 2014.
- [72] B. Vasilescu, A. Serebrenik, and T. Mens. A historical dataset of software engineering conferences. In *Proceedings of the 10th Working Conference on Mining Software Repositories*, pages 373–376. IEEE Press, 2013.
- [73] J. Vesterstrøm and R. Thomsen. A comparative study of differential evolution, particle swarm optimization, and evolutionary algorithms on numerical benchmark problems. In *Evolutionary Computation, 2004. CEC2004. Congress on*, volume 2, pages 1980–1987. IEEE, 2004.
- [74] X. Wei and W. B. Croft. Lda-based document models for ad-hoc retrieval. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 178–185. ACM, 2006.
- [75] P. Xie, D. Yang, and E. P. Xing. Incorporating word correlation knowledge into topic modeling. In *Conference of the North American Chapter of the Association for Computational Linguistics*, 2015.
- [76] G. Yang and B. Lee. Predicting bug severity by utilizing topic model and bug report meta-field. *KIISE Transactions on Computing Practices*, 21(9):616–621, 2015.
- [77] Y. Yang. *Improving the Usability of Topic Models*. PhD thesis, NORTHWESTERN UNIVERSITY, 2015.
- [78] D. Zhang and J. J. Tsai. *Machine learning applications in software engineering*, volume 16. World Scientific, 2005.
- [79] Y. Zhang, M. Harman, Y. Jia, and F. Sarro. Inferring test models from kate’s bug reports using multi-objective search. In *Search-Based Software Engineering*, pages 301–307. Springer, 2015.
- [80] W. Zhao, J. J. Chen, R. Perkins, Z. Liu, W. Ge, Y. Ding, and W. Zou. A heuristic approach to determine an appropriate number of topics in topic modeling. *BMC bioinformatics*, 16(13):1, 2015.
- [81] M. Zuluaga, G. Sergeant, A. Krause, and M. Püschel. Active learning for multi-objective optimization. In *Proceedings of the 30th International Conference on Machine Learning*, pages 462–470, 2013.