**分领域：新教师论坛**

# 面向API文档的数据分析研究

## 张静宣
## 2019年11月24日

南京航空航天大学
Nanjing University of Aeronautics and Astronautics

# 自我介绍

## 张静宣(Jingxuan Zhang)

助理教授，硕士生导师
南京航空航天大学 计算机科学与技术学院
Email: jxzhang@nuaa.edu.cn
主页: https://csjxzhang.github.io
研究方向：软件数据分析，软件仓库挖掘

| 时间 | 单位 | 经历 | 导师 |
|---|---|---|---|
| 2008.9-2012.6 | 大连理工大学 | 本科 | -- |
| 2012.9-2018.6 | 大连理工大学 | 博士 | 江贺 |
| 2014.5-2014.8 | 新加坡管理大学 | 研究助理 | David Lo |
| 2018.7-现在 | 南京航空航天大学 | 助理教授 | -- |

# 汇报提纲

# API和API文档

## API衍生新的商业模式



**API文档**-帮助开发者学习API的正确使用

# API文档的分类



API 规范

API 教程

博客

论坛

Maalej W, Robillard M P. Patterns of knowledge in API reference documentation, IEEE Transactions on Software Engineering, 2013, 39(9): 1264-1282.

# API文档的内容

既包含有**自然语言的描述**，又包含少量**代码样例**。同时自然语言的描述有可能也包含**代码元素**。

Zhong H, Su Z. Detecting API documentation errors, International Conference on Object Oriented Programming Systems Languages & Applications (OOPSLA). ACM, 2013:803-816.

# API文档的质量

研究者对API文档的质量做了实证研究，调研**API学习的障碍**和**API文档失效**的原因。

| Subcategories/descriptions | | Associated respondents |
|---|---|---|
| Obstacles caused by inadequate or absent resources for learning the API (for example, documentation) | | 50 |
| Examples | Insufficient or inadequate examples | 20 |
| General | Unspecified issues with the documentation | 14 |
| Content | A specific piece of content is missing or inadequately presented in the documentation (for example, information about all exceptions raised) | 12 |
| Task | No reference on how to use the API to accomplish a specific task | 9 |
| Format | Resources aren't available in the desired format | 8 |
| Design | Insufficient or inadequate documentation on the high-level aspects of the API such as design or rationale | 8 |

| Category | Problem | Description | E* | D* |
|---|---|---|---|---|
| Content | Incompleteness | The description of an API element or topic wasn't where it was expected to be. | 20 | 20 |
| | Ambiguity | The description of an API element was mostly complete but unclear. | 16 | 15 |
| | Unexplained examples | A code example was insufficiently explained. | 10 | 8 |
| | Obsoleteness | The documentation on a topic referred to a previous version of the API. | 6 | 6 |
| | Inconsistency | The documentation of elements meant to be combined didn't agree. | 5 | 4 |
| | Incorrectness | Some information was incorrect. | 4 | 4 |
| Total | | | 61 | 57 |
| Presentation | Bloat | The description of an API element or topic was verbose or excessively extensive. | 12 | 11 |
| | Fragmentation | The information related to an element or topic was fragmented or scattered over too many pages or sections. | 5 | 5 |
| | Excess structural information | The description of an element contained redundant information about the element's syntax or structure, which could be easily obtained through modern IDEs. | 4 | 3 |
| | Tangled information | The description of an API element or topic was tangled with information the respondent didn't need. | 4 | 3 |

**质量问题**

| | |
|---|---|
| 重要知识不凸显 | ➡ 自动识别并高亮重要知识 |
| 重要信息（代码）缺失 | ➡ 自动填充缺失的重要信息 |
| 信息冗长且碎片化 | ➡ 自动切分并推荐相关内容 |

Uddin G, Robillard M P. How API Documentation Fails[J]. IEEE Software, 2015, 32(4):68-75.
Robillard M P. What Makes APIs Hard to Learn? Answers from Developers[J]. IEEE Software, 2009, 26(6):27-34.

南京航空航天大学
Nanjing University of Aeronautics and Astronautics

## API文档内容解析



**解决重要知识不凸显的问题**

## API文档信息增强



**解决重要信息（代码）缺失的问题**

## API文档参考推荐



**解决信息冗长且碎片化问题**

# 汇报提纲

01 背景介绍

02 内容解析

03 信息增强

04 参考推荐

05 汇报总结

# API文档内容解析

**已有研究者对API文档包含的知识进行了分析，发现了12种知识类型以及他们的出现模式。**

| Knowledge Type | Description (Excerpt) |
|---|---|
| **Functionality** and Behavior | Describes what the API does (or does not do) in terms of functionality or features. Describes what happens when the API is used (a field value is set, or a method is called). |
| **Concepts** | Explains the meaning of terms used to name or describe an API element, or describes design or domain concepts used or implemented by the API. |
| **Directives** | Specifies what users are allowed / not allowed to do with the API element. Directives are clear contracts. |
| Purpose and Rationale | Explains the purpose of providing an element or the rationale of a certain design decision |

**Maalej W, Robillard M P. Patterns of knowledge in API reference documentation, IEEE Transactions on Software Engineering, 2013, 39(9): 1264-1282.**

# API文档内容解析

**定义：** API directive是正确调用API所必须遵守的约束和限制

**目标：** 自动的识别和高亮API文档中的directive，提醒开发者

南京航空航天大学
Nanjing University of Aeronautics and Astronautics

**解决方案:** 带采样的深度学习（LSTM）模型

# API文档内容解析

## 1. 我们提出的方法是否超过对比方法?

| API Specification | Precision | | | Recall | | | F-Measure | | |
|---|---|---|---|---|---|---|---|---|---|
| | Baseline | DeepDir | Improv. | Baseline | DeepDir | Improv. | Baseline | DeepDir | Improv. |
| Java | 15.18% | 34.43% | +19.25% | 73.81% | 76.17% | +2.36% | 25.18% | 47.39% | +22.21% |
| JFace | 28.75% | 54.61% | +25.86% | 78.40% | 81.09% | +2.69% | 42.07% | 64.90% | +22.83% |
| commons.collections | 49.37% | 65.82% | +16.45% | 67.41% | 85.44% | +18.03% | 57.00% | 74.30% | +17.30% |
| Average | 31.10% | 51.62% | +20.52% | 73.21% | 80.90% | +7.69% | 41.42% | 62.20% | +20.78% |

## 2. 自动的识别高亮directive是否对开发者有利?

| Part | Q | Question | Answer Option | Percentage |
|---|---|---|---|---|
| Part 1 | Q1 | How often do you refer to API documentation when facing unfamiliar APIs? | Always | 48% |
| | | | Often | 46% |
| | | | Seldom | 6% |
| | | | Never | 0% |
| Part 2 | Q2 | Have you ever pay attention to API directives in API documentation? | Yes | 91% |
| | | | No | 9% |
| | Q3 | Is it difficult to find out API directives from API documentation? | Very difficult | 2% |
| | | | Difficult | 19% |
| | | | Normal | 64% |
| | | | Easy | 13% |
| | | | Very easy | 2% |
| | Q4 | What factors prevent you to find out API directives from API documentation? | Too long explanation | 52% |
| | | | Unstructured text without uniform style | 68% |
| | | | Useless information embedded | 40% |
| | | | Lack of focus in long text | 61% |
| | | | Other (please specify) | 16% |
| | Q5 | Do you think automatically reminding and highlighting API directives is helpful to avoid bugs? | Very Helpful | 46% |
| | | | A little helpful | 48% |
| | | | helpless | 6% |
| Part 3 | Q6 | Do you think API directive detection tools are helpful? | Yes | 84% |
| | | | No | 16% |
| | Q7 | Which is important when detecting API directives? | Higher precision | 33% |
| | | | Higher recall | 18% |
| | | | Both higher precision and recall | 49% |
| | Q8 | What suggestions do you have for API directive detection tools? | Blank option | -- |

# 汇报提纲

# API文档信息增强

Kim J, Lee S, Hwang S W, et al. Enriching Documents with Examples: A Corpus Mining Approach, ACM Transactions on Information Systems, 2013, 31(1):157-160.

# API文档信息增强

## 解决方案: 基于群智的代码样例填充方法

**我们方法增强的API文档能否提高开发者效率?**

**方法步骤:**

1. 邀请21个参与者，设计问卷调查调研他们的编程背景。根据他们的学习编

| Id | Programming task | Related topic | Potentially Useful API | Test case | | Level of difficulty |
|----|------------------|---------------|------------------------|-----------|-----------|---------------------|
| | | | | Test input | Expected output | |
| 1 | Split a text based on a specified character | String processing | java.util.String. Tokenizer | class1=1\nclass2=2\nclass3=3\nclass4=4\nclass1=5\nclass2=6\nclass3=7\nclass4=8\nclass1=9\nclass2=10\nclass3=11\nclass4=12 | class1:1, 5, 9 class2:2, 6, 10 class3:3, 7, 11 class4:4, 8, 12 | easy |
| 2 | Read and print the source code of a webpage | Network connection and interaction | java.net.URL. Connection | http://global.bing.com/?FORM=HPCNEN&setmkt=en-us&setlang=en-us | The source code of bing search | moderate |
| 3 | List the first menu when press the ctrl key | GUI design and implementation | javax.swing. JMenuBar | Run the program and press the ctrl key | The menus in the first menu bar are shown | hard |

1. 完成编程任务的数量
2. 完成编程任务的时间

# API文档信息增强

## 我们方法增强的API文档能否提高开发者效率?



Fig. 8. Comparison between different tasks

Fig. 9. Comparison between different groups

Fig. 10. Comparison between different API specifications

1. 编程问题设置的难易程度符合预期。
2. 我们的人员分组策略有效。
3. 利用我们方法增强的API文档，可以提高开发者的编程效率。

# 汇报提纲

01 背景介绍

02 内容解析

03 信息增强

04 参考推荐

05 汇报总结

# API文档参考推荐

## 问题：
- ✓ **单个API文档冗长，不能快速定位想要的内容。**
- ✓ **解释某个API的文档片段分布在文档的各个部分。**
- ✓ **并不是某个API出现，该文档就一定是解释这个API的。**

> (1)JodaTime is like an iceberg, 9/10ths of it is invisible to user-code. (2)Many, perhaps most, applications will never need to see what's below the surface. (3)This document provides an introduction to the JodaTime API for the average user, not for the would-be API developer. (4)The bulk of the text is devoted to code snippets that display the most common usage scenarios in which the library classes are used. (5)In particular, we cover the usage of the key DateTime, Interval, Duration and Period classes. (6)We finish with a look at the important topic of formatting and parsing and a few more advanced topics.

|  | IR | GMR | FITSEA |
|---|---|---|---|
| Input | Unfamiliar APIs | | |
| Output | Relevant API tutorial fragments explaining unfamiliar APIs | | |
| Method | Information Retrieval | Text Classification | Text Classification |
| Drawbacks | Precision is low | Different corpora require their corpus-specific annotated data. The effectiveness of supervised approaches depends on the features. | |

He Jiang, Jingxuan Zhang, Xiaochen Li, Zhilei Ren, and David Lo. A More Accurate Model for Finding Tutorial Segments Explaining APIs, *In Proc. of 23th IEEE International Conference on Software Analysis, Evolution, and Reengineering (SANER'2016)*, Osaka, Japan. pp. 157-167.

# API文档参考推荐

**现有工作的不足**：监督模型需要大量人工标注构建训练集，预测效果依赖于提取的特征，在应用中并不准确。

**目标**：提出一种无监督模型。



**相关性识别阶段:** 发现API与文档片段的相关性。

**片段推荐阶段:** 为用户输入的不熟悉的API推荐文档解释片段。

He Jiang, Jingxuan Zhang, Zhilei Ren, and Tao Zhang. An Unsupervised Approach for Discovering Relevant Tutorial Fragments for APIs, *In Proc. of 39th IEEE International Conference on Software Engineering (ICSE' 2017)*, Buenos Aires, Argentina. pp. 38-48.

# API文档参考推荐

## 我们提出的无监督方法效果如何?

| Corpus | Tutorial | Precision (%) | | | | Recall (%) | | | | F-Measure (%) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | FRAPT | FITSEA | GMR | IR | FRAPT | FITSEA | GMR | IR | FRAPT | FITSEA | GMR | IR |
| McGill Corpus | JodaTime | 85.19 | 69.00 | 58.82 | 73.00 | 76.67 | 74.17 | 50.00 | 73.00 | 80.70 | 70.24 | 54.05 | 73.00 |
| | Math Library | 84.78 | 67.89 | 52.00 | 67.00 | 73.58 | 72.70 | 49.06 | 65.00 | 78.79 | 61.53 | 50.49 | 66.00 |
| | Col. Official | 62.03 | 55.74 | 62.69 | 30.00 | 87.50 | 48.62 | 31.79 | 94.00 | 72.59 | 48.10 | 42.18 | 45.00 |
| | Col. Jenkov | 61.19 | 90.44 | 82.14 | 33.00 | 97.62 | 85.17 | 58.97 | 88.00 | 75.23 | 85.17 | 68.66 | 48.00 |
| | Smack | 77.94 | 83.38 | 70.00 | 74.00 | 94.64 | 88.33 | 93.33 | 52.00 | 85.48 | 83.90 | 80.00 | 61.00 |
| Android Corpus | Graphics | 49.21 | 50.42 | 45.60 | 35.80 | 75.61 | 42.52 | 44.50 | 67.44 | 59.62 | 43.73 | 45.04 | 46.77 |
| | Resources | 65.22 | 75.83 | 55.00 | 40.32 | 66.67 | 66.17 | 21.11 | 55.56 | 65.93 | 66.80 | 30.51 | 46.73 |
| | Data | 71.43 | 56.52 | 19.29 | 33.33 | 55.56 | 52.00 | 14.76 | 44.00 | 62.50 | 54.17 | 16.72 | 37.93 |
| | Text | 57.58 | 36.19 | 66.67 | 37.21 | 76.00 | 48.33 | 22.22 | 57.14 | 65.52 | 39.56 | 33.33 | 45.07 |



### Approach Comparison

| | Precision | Recall | F-Measure | Precision | Recall | F-Measure |
|---|---|---|---|---|---|---|
| | McGill Corpus | | | Android Corpus | | |
| FRAPT | 74.23% | 86.00% | 78.56% | 60.86% | 68.46% | 63.39% |
| FITSEA | 73.29% | 73.80% | 69.79% | 54.74% | 52.26% | 51.07% |
| GMR | 65.13% | 56.63% | 59.08% | 46.64% | 25.65% | 31.40% |
| IR | 55.40% | 74.40% | 58.60% | 36.67% | 56.04% | 44.13% |

■FRAPT ■FITSEA ■GMR ■IR

1. **作为一种无监督方法，我们提出的方法超过现有监督方法。**

2. **考虑到无监督方法的优点，用我们提出的方法来发现解释API的文档片段是个很好的选择。**

# 汇报提纲

# 汇报总结

- ✓ 作为学习API的最重要的资源，API文档的质量并**不理想**。

- ✓ 对API文档进行分析，进而提高API文档的质量是当前研究**热点问题**。

- ✓ 我们分别介绍了API文档**内容解析**，**信息增强**和**参考推荐**方面的工作。

- ✓ 但这些研究还无法完全消除API文档的全部问题，**还有很多研究问题**。

# 谢谢各位专家

# 敬请批评指正

---

张静宣
南京航空航天大学
2019年11月24日