Predicting Licenses in Changed Source Code

Xiaoyu Liu¹, LiGuo Huang¹, Jidong Ge², Vincent Ng³

Dept. of Computer Science, Southern Methodist University, Dallas, TX
 State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China
 Human Language Technology Research Institute, University of Texas and Dallas, TX

Emails: {xiaoyul, lghuang}@smu.edu gjd@nju.edu.cn vince@hlt.utdallas.edu









After software changes are made ...









Penalties for Software License Violation

Whenever piracy of responsible for pa with others, they o Just like illegally d may go beyond ju

The legal penaltie



ilty, they will be ed the software

onsequences be charged.

They may be

required to pay fines as high as \$250,000 and may face up to five years in prison. In addition, a permanent felony will be on their record.

Predict license to resolve violation as soon as possible!

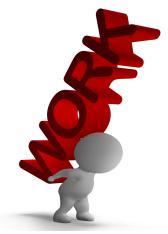






Challenges of Software Licensing

• Time and labor effort consuming



• Difficulty: requires lots of experience and expertise









Software License Prediction

Changed imports in XMLPacker.java XMLPacker.java was licensed with MPL v1.1



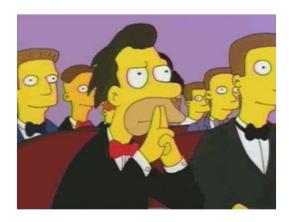






Goal

Automatically predict software license for a changed source code file







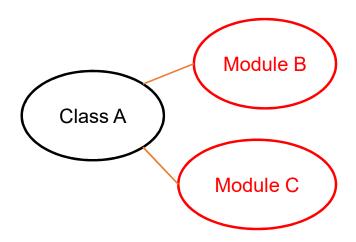


Related Work on File-level License Prediction

- Ninka: [German et.al., 2010]
 - state-of-the-art, leveraging a set of predefined regular expressions built upon these common terms to detect the presence of the license copyrights and terms
- This is the image of the image

- Caller-Callee (CC):
 - license based on code imports' licenses

Not taking into account the license compatibility issues of the different license restrictions









Automatic License Predictor (ALP)

- Novel Contributions
 - Propose ALP, the first learning-based approach for license prediction on changed source code file
 - Conflict resolution leverage additional knowledge resources
- Framework

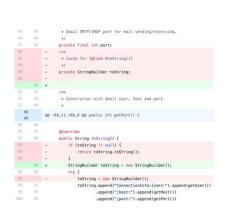




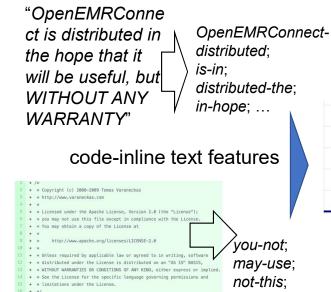




Basic ALP System



changed source code file



Logistic

Regression





diff features

use-file; ...



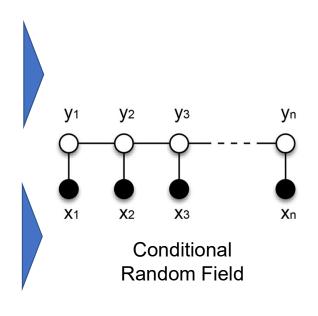




Modeling the Previous License (ALP2)

Previous version













Current version







Adding New Knowledge Resources

New resource added!



Software documents

MIT License

Copyright (c) 2019 Avery O'Banion

LICENSE file

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.







Adding New Knowledge Resources

New resource added!

Co-changed files with ConnectionInfo.java







Adding New Knowledge Resources

ALP2+Doc+Co



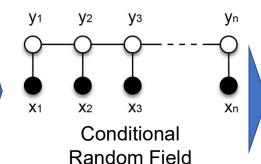
"retain the above copyright notice"

retain, retain-above; the-copyright; above-copyright; copyright-notice; ...

document-text features



code-inline text features
diff features

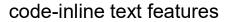




Changed source code file



co-changed source code files

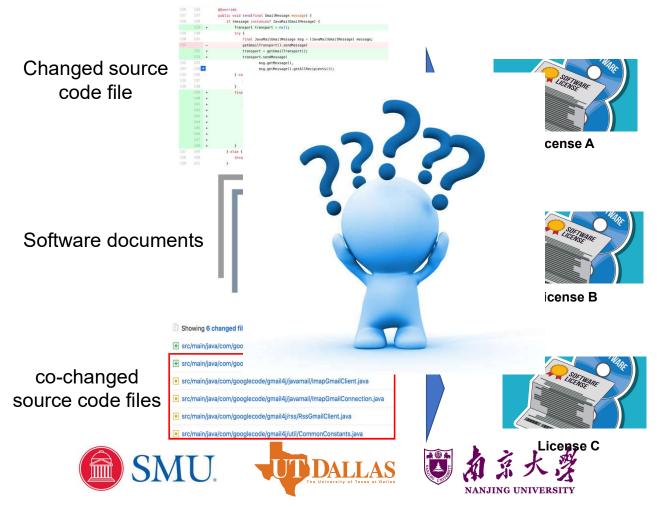


diff features



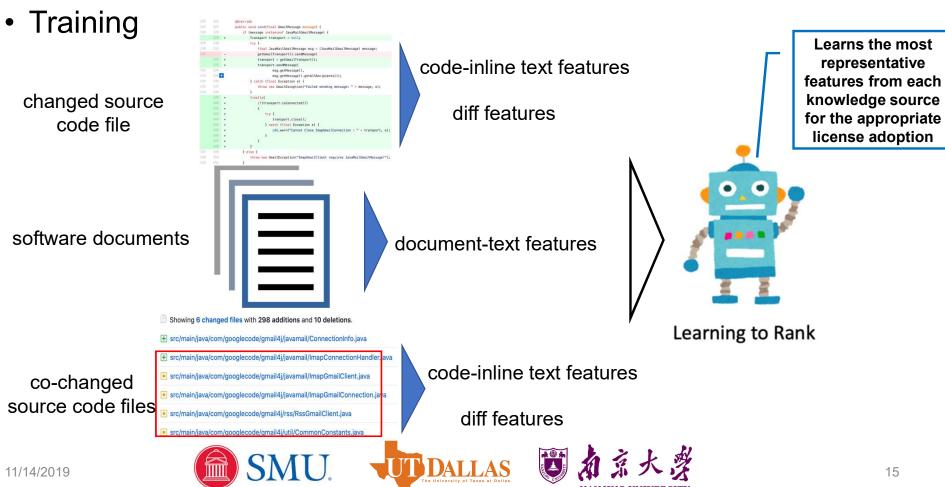


An Example of License Prediction Problem

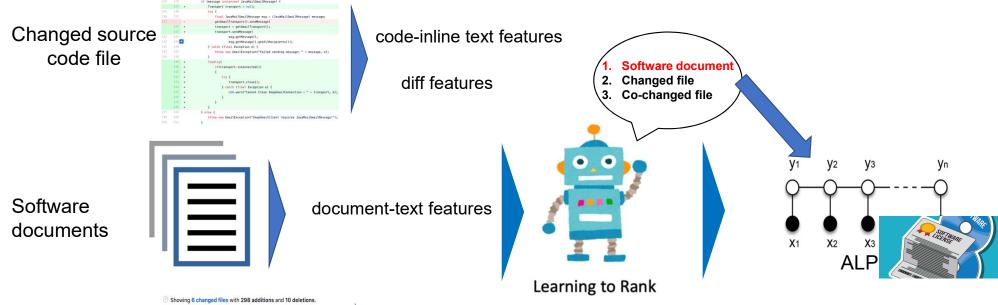


11/14/2019

Modeling Conflicts (ALP2-Ranker)



Resolve Software License Conflicts (ALP2-Ranker)



Co-changed source code file



code-inline text features

diff features







Empirical Evaluation

- Datasets: Open Source Java Projects from GitHub
 - ❖700 Java projects from GitHub

| # of systems | 700 |
|--------------------|-------|
| # of commits | 8128 |
| # of changed files | 57450 |

(a) Overall statistics







Procedure of Annotating Licenses for Each File

- •We looked at the terms and conditions for copying, distribution and modifications of each license of each resource
- Coders: two PhD students who have extensive experience in industry as developers
- Agreement ratio: 73.7%
- Disagreements are resolved by open discussion







Results of Data Annotation

• 25 unique licenses

| Licenses | # of changed files |
|--------------|--------------------|
| Apache v2 | 18770 (32.7%) |
| GPL v2 | 9458 (16.5%) |
| GPL v3+ | 5943 (10.3%) |
| MIT | 3125 (5.4%) |
| LGPL v3+ | 2609 (4.5%) |
| LGPL v2.1+ | 1542 (2.7%) |
| BSD | 1404 (2.4%) |
| EPL v1 | 1276 (2.2%) |
| Other | 4960 (8.6%) |
| Non-licensed | 8363 (14.6%) |

(b) Per-license frequencies







Changed file

(Apache v2) ... Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, nocharge, royalty-free, irrevocable copyright license to reproduce, prepare Derivative Works of, publicly display, publicly perform, sublicense, and distribute the Work and such Derivative Works in Source or Object form...

Software document

(GPL v2) ... Therefore, by modifying or distributing the Program (or any work based on the Program), you indicate your acceptance of this License to do so, and all its terms and conditions for copying, distributing or modifying the Program or works based on it ...

Co-changed file

(LGPL v2.1+) ...linking a "work that uses the Library" with the Library creates an executable that is a derivative of the Library (because it contains portions of the Library), rather than a "wo k that uses the library". The executable is therefore covered by this License.....







Empirical Evaluation

- Datasets: Open Source Java Projects from GitHub
 - ❖700 Java projects from GitHub, 25 unique licenses
- Baseline Systems
 - * Ninka
 - **❖** Caller-Callee (CC)
 - * Previous Version (Prev): predicts license to be the same as previous version







Empirical Evaluation

- Datasets: Open Source Java Projects from GitHub
 - ❖700 Java projects from GitHub, 25 unique licenses
- Baseline Systems
 - * Ninka
 - **❖** Caller-Callee (CC)
 - * Previous Version (Prev): predicts license to be the same as previous version
- Metrics
 - ❖ Macro F1: treat each license as equally important
 - ❖ Micro F1: treat most frequent licenses as more important
 - ❖ Five-Fold Cross Validation

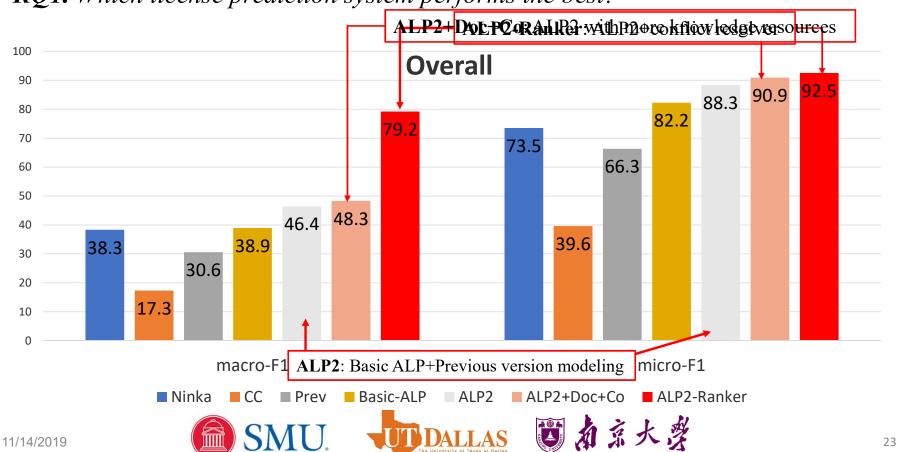




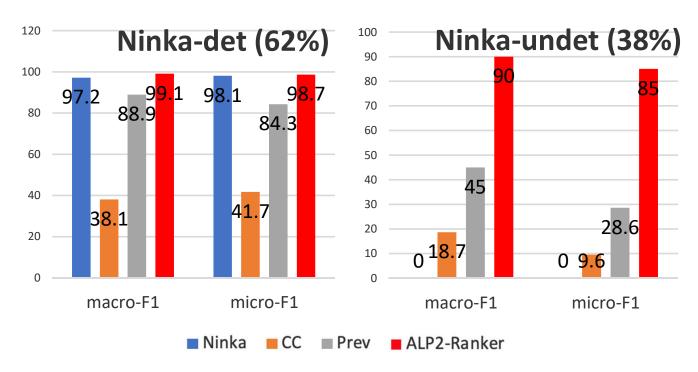


Overall Performance

RQ1. Which license prediction system performs the best?



RQ2. How do the systems perform on the easy and difficult tasks?

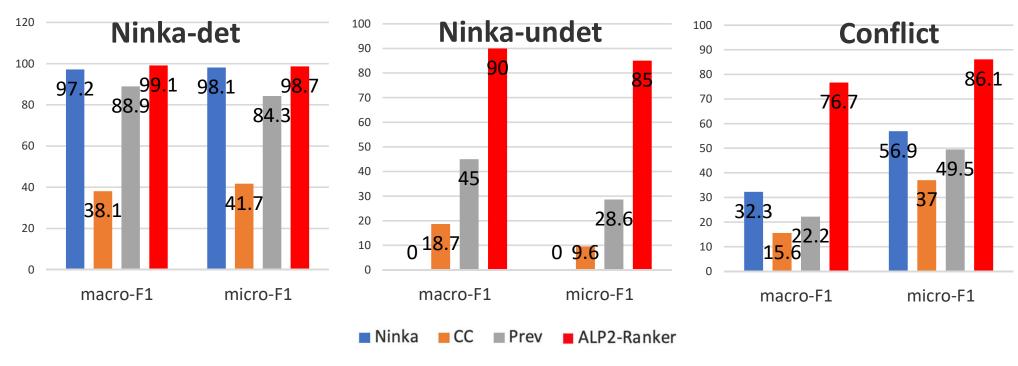








RQ2. How do the systems perform on the easy, difficult, and conflict instances?











Goal

Automatically predict software license for a changed source code file



11/14/2019

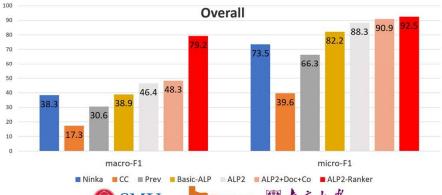






Overall Performance

RO1. Which license prediction system performs the best? ALP2: Basic ALP+Previous version modeling; ALP2+Doc+Co: ALP2+Implicit conflict resolver; ALP2-Ranker: ALP2



Automatic License Predictor (ALP)

- Novel Contributions
 - · Propose ALP, The first learning-based approach for license prediction on changed source code file
 - · Conflict resolution leverage additional knowledge resources
- Framework













Xiaoyu Liu1, LiGuo Huang1, Jidong Ge2, Vincent Ng3

1 Dept. of Computer Science, Southern Methodist University, Dallas, TX ² State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China 3 Human Language Technology Research Institute, University of Texas and Dallas, TX

Emails: {xiaoyul, lghuang}@smu.edu gjd@nju.edu.cn vince@hlt.utdallas.edu



11/14/2019



11/14/2019