

A clustering-based approach for mining dockerfile evolutionary trajectories

Yang ZHANG^{1,2*}, Huaimin WANG^{1,2} & Vladimir FILKOV^{3,4}

¹Key Laboratory of Parallel and Distributed Computing, National University of Defense Technology, Changsha 410073, China;

²College of Computer, National University of Defense Technology, Changsha 410073, China;

³DECAL Lab, University of California, Davis 95616, USA;

⁴Computer Science Department, University of California, Davis 95616, USA

Received 2 January 2018/Revised 7 February 2018/Accepted 9 April 2018/Published online 15 October 2018

Citation Zhang Y, Wang H M, Filkov V. A clustering-based approach for mining dockerfile evolutionary trajectories. *Sci China Inf Sci*, 2019, 62(1): 019101, <https://doi.org/10.1007/s11432-017-9415-3>

Dear editor,

Docker¹⁾, as a de-facto industry standard [1], enables the packaging of an application with all its dependencies and execution environment in a light-weight, self-contained unit, i.e., containers. By launching the container from Docker image, developers can easily share the same operating system, libraries, and binaries [2]. As the configuration file, the dockerfile plays an important role, because it defines the specific Docker image architecture and building orders [3]. As a project progresses through its development stages, the content of the dockerfile may be revised many times. This dockerfile evolution is indicative of how the project infrastructure varies over time, and different projects can exhibit different evolutionary trajectories. However, projects with similar goals and needs may converge to more similar trajectories than more disparate projects. Identifying software projects that have undergone similar changes can be very important for the discovery and implementation of the best practices when adopting new tools and pipelines, especially in the DevOps software development paradigm. The potential to implement the best practices through the analysis of the dockerfile evolutionary trajectories motivated this work.

tivated this work.

This research studied dockerfile longitudinal changes at large scale and presented a clustering-based approach for mining convergent evolutionary trajectories. An empirical study of 2840 projects was conducted, and six distinct clusters of dockerfile evolutionary trajectories were found. Furthermore, each cluster was summarized, accompanied by case studies, and the differences between different clusters were discussed. The proposed approach quantifies distinct dockerfile evolution modes and reflects the learning curves of project maintainers, and this benefits future project maintenance. Also, the proposed approach is generic and can be used for the study of general infrastructure configuration file evolution.

Dockerfile overview. A dockerfile²⁾ is a text document that contains all the commands a user could call on the command line to assemble a Docker image [3]. Users can create an automated build that executes several command-line instructions in succession by using a docker build. Docker has provided multiple types of instructions for use in dockerfiles (Appendix A). Docker runs the instructions in a dockerfile in order and treats lines that begin with “#” as a comment. A dockerfile must

* Corresponding author (email: yangzhang15@nudt.edu.cn)

1) <https://www.docker.com/>.

2) <https://docs.docker.com/engine/reference/builder/>.

start with a FROM instruction, specifying the base build. Other parts are then added on top of the base one [4], and each instruction represents one layer in a docker image. Thus, the size of a dockerfile (number of lines) can reveal the size and complexity of the corresponding docker image.

To meet the requirements of project development, the content of a dockerfile may be modified at different stages by project maintainers. These changes over time are called the dockerfile evolution. E.g., during a previous stage, the owner inutano of project inutano/wpgsa-docker added a USER instruction and new python scripts to the initial dockerfile. However, during the later stage, he just updated the plugin's wPGSA version, e.g., "0.2.0" → "0.3.0". Thus, these changes reflect the dockerfile evolution, and the changes depend on the practices in individual projects. During the preliminary data gathering of more than 57000 projects from Docker Hub, nearly 70% of them had changed their dockerfile at least once. In addition, over 5000 projects changed their dockerfile more than 10 times.

Our approach. This research proposed a clustering-based approach that defined the dockerfile scale as the number of valid command lines of the dockerfile without blank lines and comments. The proposed approach consists of four steps.

- Step-1. Preprocessing. Initially, a number of projects with different dockerfile versions were imported. For each dockerfile, its scale was computed after removing blank lines and comments. Then, for project p with n dockerfile versions, the scale of the i -th version was defined as s_i , and the evolutionary trajectory of project p 's dockerfile was represented by the vector $S_p = \{s_1, s_2, \dots, s_n\}$.

- Step-2. Smoothing. The kernel regression smoothing [5] was used to reduce the noise in the original dockerfile evolutionary vectors. Traditional parametric estimation methods may not capture all evolutionary trajectories with a limited mathematical formula, but kernel regression is a non-parametric linear smoother that can reduce noise by estimating the conditional expectation of a random variable. For each dockerfile scale s_i in project p , a tuple vector was constructed, (x_i, s_i) , $i = 1, \dots, n$, $x_i = i$. For the smoothing, our goal was to estimate r with some function \hat{r} , so that $s_i = r(x_i) + \epsilon_i$, $i = 1, \dots, n$. The Gaussian kernel was used to define the kernel function as $K(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$.

Then, the function \hat{r} of kernel regression [6] was

defined as follows:

$$\hat{r}(x) = \frac{\sum_{i=1}^n K(\frac{x_i - x}{h})}{\sum_{i=1}^n K(\frac{x_i - x}{h})} \times s_i, \quad (1)$$

where h was the bandwidth, computed by the cross-validation least-squares function. After smoothing, a new evolutionary vector was obtained, $S'_p = \{s'_1, s'_2, \dots, s'_n\}$. On average, the R-squared coefficient for each kernel regression model is 0.85 (median is 0.95) in our study.

- Step-3. Resampling. Next, one-dimensional linear interpolation [7] was used to resample our data, which returns one-dimensional piecewise linear interpolants to a function with given values at discrete data-points. For a given x_r , the value of \hat{s}_r is computed by the following:

$$\hat{s}_r = \frac{s'_i \times (x_j - x_r) + s'_j \times (x_r - x_i)}{x_j - x_i}, \quad (2)$$

where (x_i, s'_i) and (x_j, s'_j) represent two known points in our smoothed data, and the x_r is represented in the interval (x_i, x_j) . For each project p , its resampling interval is set as $\frac{n}{20}$. Thus, after resampling, each project has a 20-dimensional vector $(\{\hat{s}_1, \hat{s}_2, \dots, \hat{s}_{20}\})$ that characterizes its evolutionary trajectory. Since different projects may have different dockerfile scales, each project's vector value was normalized for comparison.

- Step-4. Clustering. Finally, the well known K-means algorithm [8] was used to cluster our dockerfile scale evolutionary trajectory vectors. K-means stores k centroids used to define clusters. A point is in a cluster if it is closer to that cluster's centroid than any other centroid. Based on the initial set $\hat{S} = \{\hat{S}_1, \dots, \hat{S}_i, \dots, \hat{S}_N\}$, $\hat{S}_i \in \mathbb{R}^{20}$ (N is the total number of projects), and a fixed positive integer, k , K-means can converge to a set of k cluster centroids, $C = \{C_1, C_2, \dots, C_k\}$ in \mathbb{R}^{20} , by minimizing the "K-means cost":

$$\Phi_S(C) = \sum_{S_i \in S} \min_{C_j \in C} \|S_i - C_j\|^2. \quad (3)$$

Afterwards, each S_i ($i \in N$) was marked with a specific cluster label, i.e., Cluster- j ($j \in k$).

Our results. Data from two communities Docker Hub³⁾ and GitHub⁴⁾ were obtained (Appendix B). Before performing our approach, the heatmap of the dockerfile scale evolutionary trajectories of 2840 projects was drawn (Figure B1). By observing the color change trends, six clusters with significant differences were found. Thus, the number of clusters was set as $k = 6$ in our K-means clustering process. After clustering all the projects'

3) <https://store.docker.com>.

4) <https://github.com>.

evolutionary vectors, we marked each project as belonging to one of the six categories (C1–C6).

For each cluster, the dockerfile scale variation curves of samples were drawn and case studies on randomly selected projects were conducted to validate the clusters. After our manual analysis, the specific paradigms describing the evolution trajectories of the six clusters were summarized, as follows (Appendix C).

[C1] Increasing and holding. This cluster comprises 21.8% of the projects. In this paradigm, developers added new instructions or new settings to the dockerfile in early periods. However, after reaching to a certain size, they just updated the basic environment variables or changed the location of instructions. Thus, the dockerfile scale stayed stable.

[C2] Constantly growing. This cluster comprises 31.6% of the projects. In this paradigm, developers kept adding services, support, or plugins to the dockerfile, causing its scale to continually increase. Compared with other paradigms, the evolution path of this paradigm is easiest to understand.

[C3] Holding and increasing. This cluster comprises 19.2% of the projects. In this paradigm, developers just updated the basic environment variables in early periods, keeping the dockerfile scale stable for a while. Then, new instructions, plugins, or support were added, and the dockerfile scale increased.

[C4] Increasing and decreasing. This cluster comprises 10.2% of the projects. In this paradigm, developers kept adding new instructions in early periods, increasing the dockerfile scale. However, after reaching a large size, maintainers tried to reconstruct the dockerfile by moving useless plugins/services or moving settings to additional script files, reducing the dockerfile scale in later periods.

[C5] Holding and decreasing. This cluster comprises 9.5% of the projects. In this paradigm, developers change little in early periods and the dockerfile scale is stable. In later periods, developers tried to change the dockerfile structure by moving some instructions to additional script files or using a more light-weight base image, so the dockerfile scale dropped.

[C6] Gradually reducing. This cluster comprises 7.7% of the projects. In this paradigm, developers kept removing useless instructions or changed the base image to reduce image layers, continually decreasing the dockerfile scale. However, the decrease slows in later periods.

The differences between projects in the six clusters were further contemplated, and we found that they differ in project age and average dockerfile

scale, indicating they have different project development stages and goals (Appendix D).

Conclusion. This research studied the dockerfile evolution in docker projects and proposed a clustering-based approach for mining dockerfile evolutionary trajectories that reveal projects of similar stages and goals. After performing an empirical study on 2840 projects, six notable paradigms of the dockerfile scale evolutionary trajectories were obtained. Through case studies, each paradigm was summarized in detail. Our approach benefits future and current docker project maintenance by providing a discrete characterization of dockerfile scale evolution over many existing projects, thus allowing co-localization of similar projects. Also, the proposed approach is generic and can be used to study general infrastructure configuration file evolution.

Acknowledgements This work was supported by National Natural Science Foundation of China (Grant Nos. 61502512, 61432020), China Scholarship Council, and National Science Foundation of USA (Grant No. 1717370). Part of this study was performed during the visit in 2017 by the first author at the DECAL lab, UC Davis.

Supporting information Appendixes A–D, and Figure B1. The supporting information is available online at info.scichina.com and link.springer.com. The supporting materials are published as submitted, without typesetting or editing. The responsibility for scientific accuracy and content remains entirely with the authors.

References

- 1 Anderson C. Docker [software engineering]. IEEE Softw, 2015, 32: 102
- 2 Bernstein D. Containers and cloud: from LXC to docker to kubernetes. IEEE Cloud Comput, 2014, 1: 81–84
- 3 Merkel D. Docker: lightweight Linux containers for consistent development and deployment. Linux J, 2014, 2014: 2
- 4 Jaramillo D, Nguyen D V, Smart R. Leveraging microservices architecture by using docker technology. In: Proceedings of SoutheastCon, Norfolk, 2016
- 5 Altman N S. An introduction to kernel and nearest-neighbor nonparametric regression. Am Stat, 1992, 46: 175–185
- 6 Cawley G C, Talbot N L C. Fast exact leave-one-out cross-validation of sparse least-squares support vector machines. Neural Netw, 2004, 17: 1467–1475
- 7 Powell M J. A direct search optimization method that models the objective and constraint functions by linear interpolation. In: Advances in Optimization and Numerical Analysis. Berlin: Springer, 1994. 51–67
- 8 Hartigan J A, Wong M A. Algorithm as 136: a K-means clustering algorithm. J Roy Stat Soc, 1979, 28: 100–108