

Process Discovery using Big Data stack - Implementing the Alpha Algorithm with Map-Reduce

User Manual

Xiangnan Chen, Martin Hashem

July 16, 2019
RWTH Aachen

1 Introduction

Process mining is an approach to extract process models from event logs. Since the distributed nature of modern information systems, event logs are likely to be distributed among different physical machines. Map-Reduce is a scalable approach for efficient computations on distributed data. In this Python application we will present the main idea of a Map-Reduce implementation of the Alpha process mining algorithm, to take advantage of the scalability of the Map-Reduce approach.

This project is a python-based webapp, that integrates the big data capabilities of the Hadoop system into the process mining framework pm4py.

1.1 WebApp Requirements

Our WebApp was designed mainly on Unix-based operating system (MacOS and Debian of Linux), but can be accessed from every major browser.

We recommend to use the latest version of your browser. Both for security reasons and to guarantee the best experience.

1.2 Access WebApp

To use our WebApp, perform the following:

1. Navigate to the link.
2. Enter your username and password.
3. Click on the button Sign in.

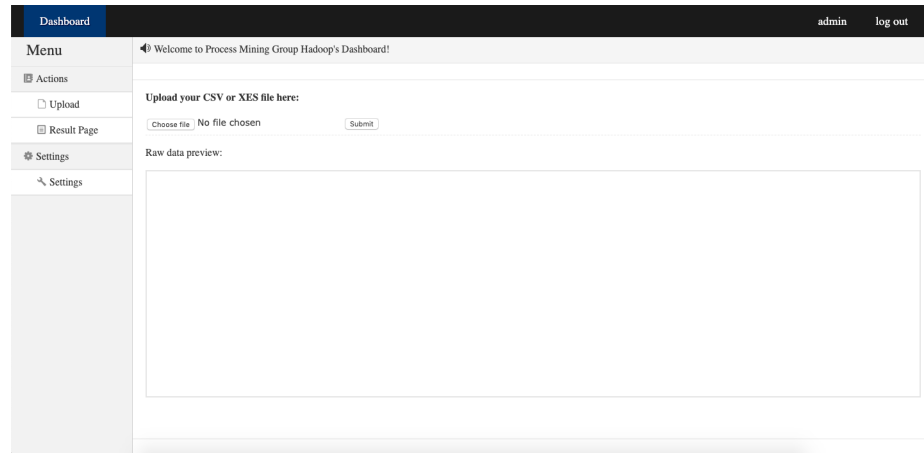


Fig. 1.

1.3 WebApp Overview

The WebApp consists of the following areas, which are showed in Fig. 1.

The main areas are:

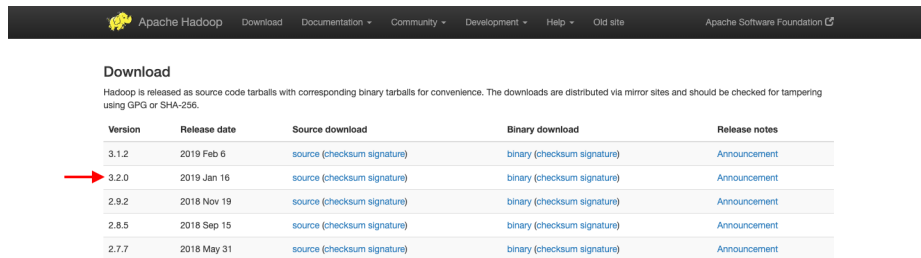
1. **Shortcut Bar** contains shortcuts to the dashboard site itself on the left side. The user settings and logout button can be found on the right side of the shortcut bar.
2. **Main Toolbar** contains buttons with the most important functions for each applications like uploading the required files and the result page.
3. **Main Window** displays the main content of the application. On the upload page there are upload bar and raw data preview box. After selected the file of right format, the content of the uploading file will be showed in the box.

2 Installation

The version control of this project depends on Git, and our team chose to use GitHub to manage our source code and releases.

To deploy our WebApp on a server, you must firstly deploy Hadoop of edition 3.2 (<https://hadoop.apache.org/releases.html> download page like in Fig. 2) on that server. Hadoop 3.2 supports GNU/Linux and Windows as platform. And softwares development environment like Java and ssh will also be required on the server.

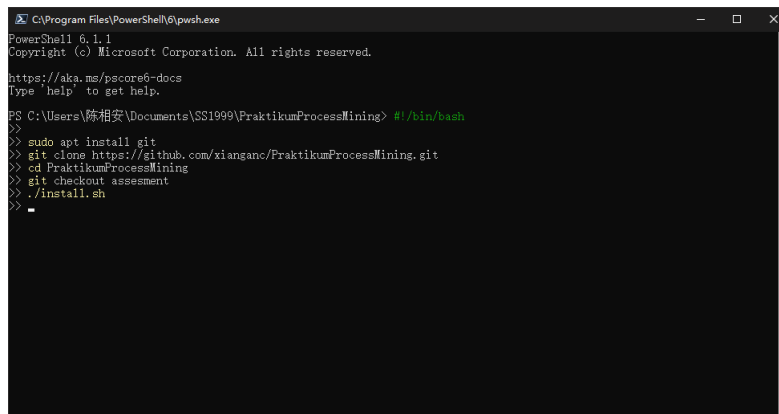
There are three modes of deploying Hadoop: Local(Standalone) mode, Pseudo-Distributed Mode and Fully-Distributed Mode. In our project we will use the Pseudo-Distributed mode (<https://hadoop.apache.org/docs/r3.2.0/hadoop-project-dist/hadoop-common/SingleCluster.html#Pseudo-Distributed.Operation>), since the limitation of our server, however a real deployment of a App with Hadoop should be in Fully-distributed mode (<https://hadoop.apache.org/docs/r3.2.0/hadoop-project-dist/hadoop-common/ClusterSetup.html>).



Version	Release date	Source download	Binary download	Release notes
3.1.2	2019 Feb 6	source (checksum signature)	binary (checksum signature)	Announcement
3.2.0	2019 Jan 16	source (checksum signature)	binary (checksum signature)	Announcement
2.9.2	2018 Nov 19	source (checksum signature)	binary (checksum signature)	Announcement
2.8.5	2018 Sep 15	source (checksum signature)	binary (checksum signature)	Announcement
2.7.7	2018 May 31	source (checksum signature)	binary (checksum signature)	Announcement

Fig. 2.

After finishing deploying Hadoop software, we can just go to our GitHub releases page (<https://github.com/xianganc/PraktikumProcessMining/releases>), then download the *setup.sh* to local, run the contents inside the *.sh* file like in Fig. 3.



```

C:\Program Files\PowerShell\pwsh.exe
PowerShell 6.1.1
Copyright (c) Microsoft Corporation. All rights reserved.

https://aka.ms/powershell-docs
Type 'help' to get help.

PS C:\Users\陈相安\Documents\SS1999\PraktikumProcessMining> #!/bin/bash
>
> sudo apt install git
> git clone https://github.com/xianganc/PraktikumProcessMining.git
> cd PraktikumProcessMining
> git checkout assesment
> ./install.sh
>

```

Fig. 3.

3 Upload

In this section we explain how to use the WebApp. After reading this section, you will be able to preview the to be uploaded data and upload a CSV or XES file to the server.

3.1 Uploading a CSV or XES File

1. Click the *Choose file* button on the top part of the main window.
2. To upload a CSV or XES file, choose the file from local disk.
3. If the selected file's format matches the server requirements, you will see the data inside the file in the text box below.
4. Press *Submit* button to upload the file to the server for further calculation.

3.2 Replace a wrong selected file

If you unexpectedly select a wrong file, you can do the followings:

2.2.1 Wrong File Format

If the selected file were in the wrong format (e.g. a pdf file), then a pop-up window will show up, to notify the user for choosing a wrong file like Fig. 3.

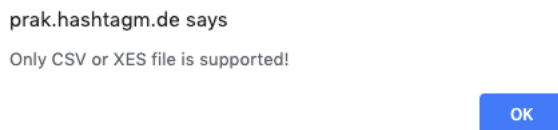


Fig. 4.

2.2.2 Wrong File Content

If the selected file were with the wrong content, maybe the wrong selected has a similar file name as the should be selected one. When the user selected a wrong file, you can re-select the file by clicking the *Choose file* button or just refresh the website to re-select the file.

4 Result

After successfully submitted the file, it should be uploaded to our server and being sent directly to calculation.

4.1 Check out the Output

Then you can simply click the *Result Page* on the left side in the main toolbar to check out the output.

5 Map Reduce Configuration

Since we are using the Map Reduce's feature of Hadoop, for the further configuration like *job* and *pipes* for Map Reduce, the user can look up the commands guide on the documentation's website of Hadoop (<https://hadoop.apache.org/docs/r3.2.0/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapredCommands.html>).

6 Log out

If anytime the logged user want to log out, he can just click the *logout* button on the top right corner to log out.