

Requirement Specification Document

Label Refinement by Behavioral Similarity

Document owners:

Bianka Bakullari

Christopher Beine

Nicole Ventsch

Juan Garza

Last edited: April 29, 2019

Contents

1	Project Drivers	2
1.1	The Purpose of the Project	2
1.2	The Client, the Customer and other Stakeholders	2
2	Project Constraints	2
2.1	Mandated Constraints	2
2.2	Naming Conventions and Definitions	2
2.3	Relevant Facts and Assumptions	2
3	Functional Requirements	2
3.1	The Scope of the Work	2
3.2	The Scope of the Product	3
3.3	Functional and Data Requirements	5
4	Nonfunctional Requirements	6
4.1	Look and Feel Requirements and Usability	6
4.2	Operational and Environmental Requirements	6
4.3	Maintainability and Support Requirements	6
5	Project Issues	7
5.1	Open Issues	7
5.2	Off-the-Shelf Solutions	7
5.3	New Problems	7
5.4	Risks	8
5.5	User Documentation	8
5.6	Waiting Room	8

1 Project Drivers

1.1 The Purpose of the Project

The purpose of the project is providing a platform that allows refinements of event logs, where activities that originally had the same label are assigned distinct labels if they appear in different patterns in the structure of the event log. In order to do so, an interactive web service should be implemented which allows users to upload event logs in the standard XES or CSV format and set thresholds for the refinements. This event log should contain activities that are carried out multiple times and are named identically. The web service should then apply the refinement algorithm and finally, the user should be able to download the refined event log in XES format. By using the refinement algorithm, more accurate and precise BPMs (Business Process Models) are generated which provide better insights into the processes.

1.2 The Client, the Customer and other Stakeholders

The client of this project is the Chair of Process and Data Science at the RWTH Aachen, which provides the supervision and support for this project. Furthermore, there are many different potential users. For example, students or researchers could use the website as a pre-processing step in order to continue research with the refined log. Additionally, business analysts from companies could use it in order to get a better insight from the refined log than from the original one. Since we will provide a web service, it will be easily accessed and not bound to any specific type of user.

2 Project Constraints

2.1 Mandated Constraints

- The algorithm has to be implemented in Python.
- Each one of the team members must take part in the coding process.
- There are clear deadlines for all deliverables. The deliverables must be submitted before the deadline.

2.2 Naming Conventions and Definitions

We refer to the "Label Refinement algorithm" also as the "Event Label Refining algorithm" or the "Refinement algorithm".

We sometimes use the name "trace variants" to refer to the the different clusters of traces.

2.3 Relevant Facts and Assumptions

The set of candidate labels which are prone to be refined is chosen by the user. The user also picks the thresholds needed as input for the different steps of the algorithm. The quality of the process models resulting from process discovery algorithms applied on the refined log heavily depends on the choices the user makes for the candidate labels and the thresholds.

3 Functional Requirements

3.1 The Scope of the Work

As explained in the Project Initiation document, we aim to provide a web application where the end user can apply the Label Refinement automatically as a pre-processing step to his event log. The new and refined log could later be used as alternative input for process discovery algorithms in other tools, yielding a new model and thus adding an alternative to the solution space of possible process models.

3.2 The Scope of the Product

The final product consists of a web interface where the user can upload event logs, parameterize thresholds and then the Label Refinement algorithm is executed. In order to get an overview of the product, a use case analysis was carried out, where we identify the system actors and their functions. Since the project's focus is the Label Refinement implementation, the project was divided into two systems. This allows a better visualisation of the individual algorithm steps and system interactions. The following use cases were identified within the analysis.

Figure 1 displays the Use Cases for the final web interface and thus the system interactions with the end-user. Figure 2 indicates the Use Cases which are required for the Label Refinement algorithm. So the front-end diagram refers to UI functions and the back-end diagram to a more technical view.

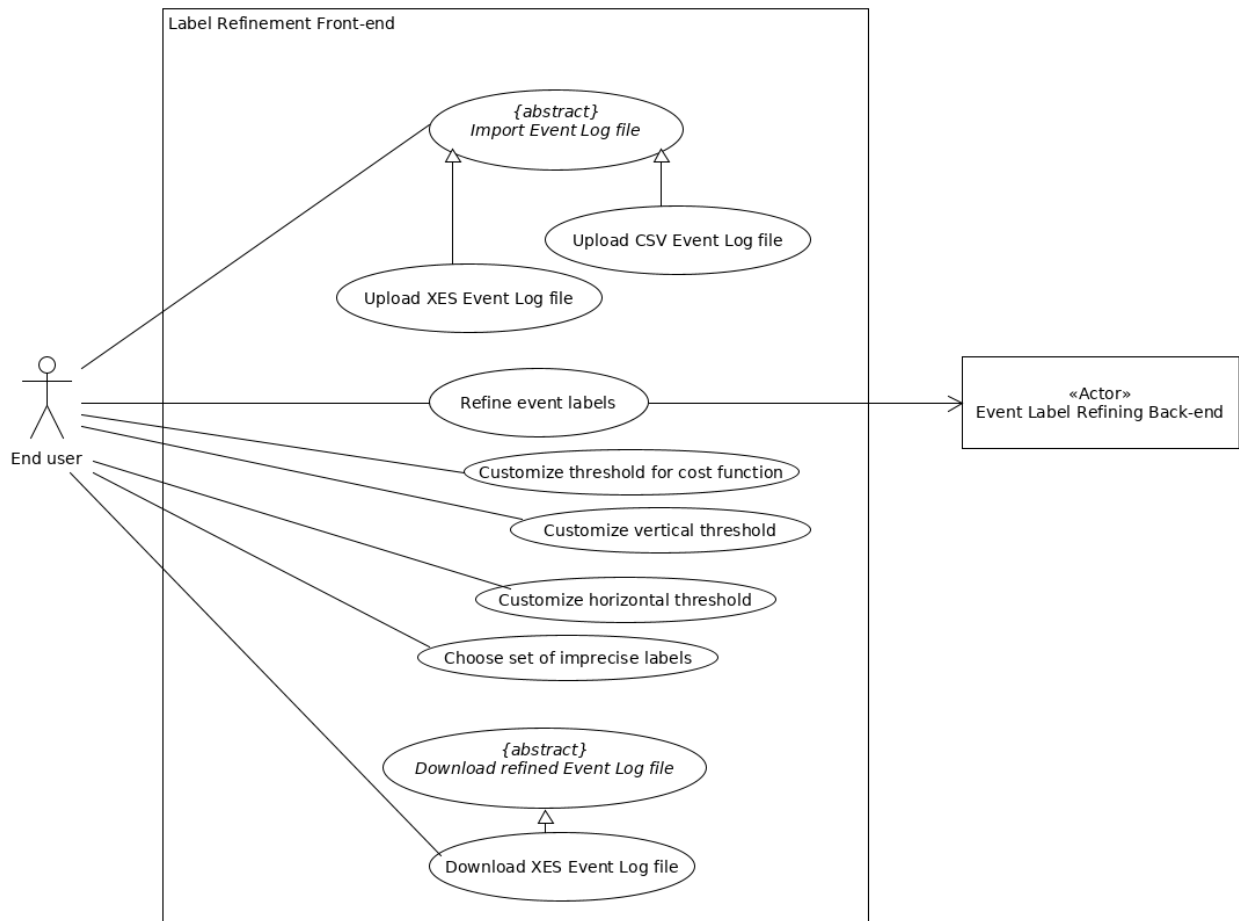


Figure 1: Use cases front-end

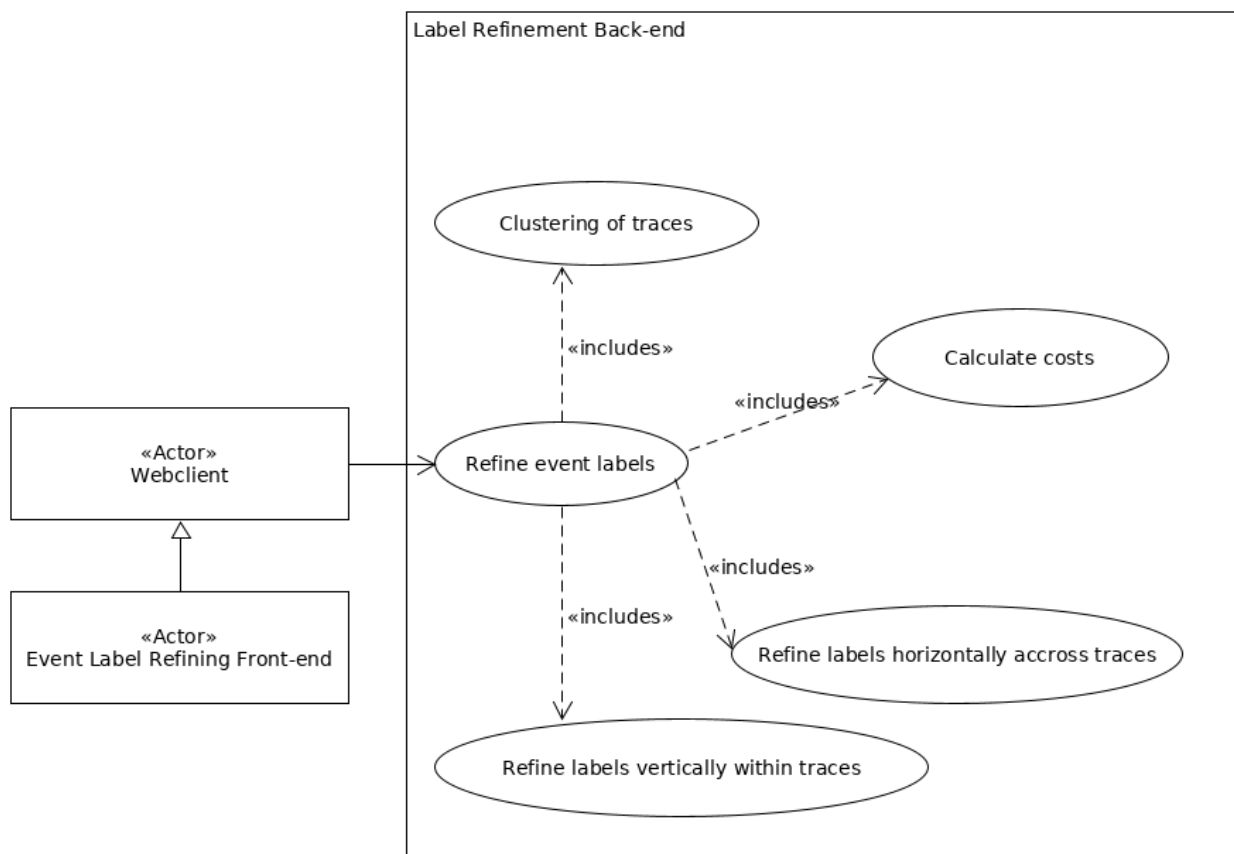
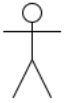





Figure 2: Use cases backen-end

The following actors exist in the final system:

Actor	Symbol	Description
End user	 End user	The end user interacts with the web interface of the system and does not need any information about the technical implementation.
Webclient	 «Actor» Webclient	The Webclient provides the API access on simple abstraction layer to provide an easy algorithm execution.
Label Refinement Front-end	 «Actor» Event Label Refining Front-end	The Label Refinement Front-end is the Webclient which is delivered with the project.
Label Refinement Back-end	 «Actor» Event Label Refining Back-end	The Label Refinement Back-end is responsible for the label refining algorithm. It provides access via an API.

3.3 Functional and Data Requirements

The following requirements result from the use case analysis:

<p>Requirement #1: System: <i>Back-End</i></p> <p>Use Case: <i>Calculate costs</i></p> <p>Description: <i>Implement distance function between all pairs of traces.</i></p> <p>Details: <i>Given a pair of traces, calculate a cost function with three weighted components.</i></p> <p>Fit Criterion: <i>$n \cdot (n-1)$ values in range $[0,1]$.</i></p>
<p>Requirement #2: System: <i>Back-end</i></p> <p>Use Case: <i>Clustering of traces</i></p> <p>Description: <i>Cluster the traces.</i></p> <p>Details: <i>A threshold for the cost function decides which traces are clustered together.</i></p> <p>Fit Criterion: <i>A partitioning of the traces into sets.</i></p>

Requirement #3: System: *Back-end*

Use Case: *Refine labels horizontally across traces*

Description: *Refining labels horizontally across traces*

Details: *Each imprecise label has the same identifier in traces of the same cluster and different identifiers in traces of different clusters.*

Fit Criterion: *Precise labels remain unchanged, relabeling equivalent to variants of traces.*

Requirement #4: System: *Back-end*

Use Case: *Refine labels vertically within traces*

Description: *Refining labels vertically within traces*

Details: *An unfolding threshold is used to decide if and where relabeling takes place within a trace.*

Fit Criterion: *Identical traces must get identical relabeling, dissimilar traces remain dissimilar.*

Requirement #5: System: *Front-end*

Use Case: *Upload CSV event log file*

Description: *Upload event log in CSV format.*

Details: *Upload function for extracting event log data from a CSV file, a clear sign showing where the user should click for upload.*

Fit Criterion: *Files with the right format are loaded successfully, otherwise return error message.*

Requirement #6: System: *Front-end*

Use Case: *Upload XES Event Log file*

Description: *Upload event log in XES format.*

Details: *Upload function for extracting event log data from a XES file, a clear sign showing where the user should click for upload.*

Fit Criterion: *Files with the right format are loaded successfully, otherwise return error message.*

Requirement #7: System: *Front-end*

Use case: *Customize threshold for cost function*

Description: *Set thresholds for the cost function.*

Details: *The selected thresholds appear on the screen.*

Fit Criterion: *Thresholds within right range are accepted successfully. Otherwise, return error message.*

Requirement #8: System: *Front-end*

Use case: *Customize horizontal threshold*

Description: *Set thresholds for the horizontal clustering.*

Details: *The selected threshold appears on the screen.*

Fit Criterion: *Thresholds within right range are accepted successfully. Otherwise, return error message.*

Requirement #9: System: *Front-end*

Use case: *Customize vertical threshold*

Description: *Set thresholds for the vertical clustering.*

Details: *The selected threshold appears on the screen.*

Fit Criterion: *Thresholds within right range are accepted successfully. Otherwise, return error message.*

Requirement #10: System: *Front-end*

Use case: *Choose set of imprecise labels*

Description: *Select imprecise labels from the uploaded file.*

Details: *Possibility to select an arbitrary set of labels to be refined from the algorithm.*

Fit Criterion: *Only selected labels will be considered.*

Requirement #10: System: *Front End*

Use case: *Download XES Event Log file*

Description: *Download refined log in XES format.*

Details: *The user chooses to download the refined log through a click.*

Fit Criterion: *Download is successful and the refined log has XES format.*

4 Nonfunctional Requirements

4.1 Look and Feel Requirements and Usability

The user interface must be self-explanatory and easy to use. It should be clear where the event log can be uploaded and/or downloaded. For each threshold, the user should be able to type in his choice and see it on the screen. In case that the event log has the wrong format or the chosen threshold is out of range, an error message displayed on the screen should imply what has gone wrong.

Our product is aimed at users who have some knowledge in Process Mining, are familiar with the algorithm or

are previously informed about the concept and effects of relabeling.

4.2 Operational and Environmental Requirements

The project must meet the following technical requirements:

- Back-end must be deployable on a production-grade server.
- Front-end must be deployable on production-grade server.
- Front-end must support Google Chrome (Version ≥ 65) and Mozilla Firefox (Version ≥ 60).

To validate the Operational and Environmental requirements, the front-end requirements are tested in Mozilla Firefox and Google Chrome.

4.3 Maintainability and Support Requirements

After finishing the project, the web service should be self-explanatory for the users. We do not intend to offer support other than documentation to the users. Moreover, the website should only include the aspects we mentioned above and is not intended to be expanded by us after this project. We do not intend maintaining the web service after finishing this project.

5 Project Issues

5.1 Open Issues

Our team has a rough idea of the steps needed to implement the algorithm and the interface. However, we are unsure about the relative time needed for each of these components. During the implementation of the algorithm, many unexpected bugs might hinder the time we have in our disposal to design a user-friendly web-interface.

While we will strive to write an efficient algorithm, there are no guarantees about its performance when the number of requests increases beyond the server capacity. Drawbacks will probably be found during implementation and decisions will have to be made on the spot.

Since each of the team members has to be fully aware of and actively participate during the coding process, it is still unclear how the functionalities will be divided between the members. Whether or not some components should be implemented together or individually depends on the type and importance of the component, deadlines and other personal issues.

5.2 Off-the-Shelf Solutions

There has already been an implementation of the algorithm in ProM. This might be helpful in the designing steps of our algorithm.

We will use a code written in Java to automatically generate events logs which we will then use as small tests for our code.

5.3 New Problems

The Label Refinement algorithm intends to give the user an alternative event log on which process discovery algorithms can be applied. Whether or not the models resulting from the new refined log have better precision or fitness compared to the original models is beyond our scope. Since the result also depends on the thresholds and set of candidate activities which are chosen by the user, we assume that the user has some background on Process Mining and has clear intentions for trying to work with a new refined log. The algorithm does not filter out events or features in the data.

Optimally, the new models resulting from the refined event log should keep the same old labels even for the refined ones. This requires from the process discovery algorithm itself to compare the original event log with the refined one and then map the refined labels to their old names after having obtained the refined model. We consider this beyond our scope. Since the set of new labels to use for refinement is unspecified, our

implementation will assign new endings to the original label, so that the context is still clear even with a model containing the new labels.

A common mistake the user could make is upload an event log in some other format (e.g. .xlsx format) instead of XES or CSV format.

5.4 Risks

Risks	Description	Category	Mitigation
Inaccurate expectations.	Inaccurate expectations are developed(believe that the project will achieve something not in the requirements, plan, etc.).	Stakeholder	Look up the requirements document and consult with the client.
Process inputs have low quality.	Inputs from stakeholders that are low quality (e.g. business case, requirements, change requests).	Stakeholder	Kindly ask for a more detailed and clearer version of any input they may provide i.e., requirements, business cases.
Misunderstood requirements.	When requirements are misinterpreted by the project team.	Communication	Meet and discuss the requirements again until the team is sure that they have completely understood them.
Learning curves.	Project team needs to acquire new skills for the project.	Team	Motivate the project team, give them the best practices on the IT field and make experts instruct them using their knowledge and own experience.
Integration failure.	Product components will fail to integrate with each other.	Integration	Establish standards for product development and make sure that the individual components passed the tests flawlessly.
Requirements are incomplete.	Requirements are not fully captured or are overlooked.	Requirements	Make a peer-review of the requirement documentation and make sure that nothing is being left out.

5.5 User Documentation

1. Technical documentation:

- Software code documentation.
- Technical specifications.

2. User documentation including:

- How to use the UI.
- Examples of inputs and outputs.
- Explanation of error messages.
- Information to contact the developers (in case of further questions).

5.6 Waiting Room

In the following we gather some ideas on which functionalities can be added or extended in our product:

- Additional feature which enables the user to choose a Business Process Model Discovery (BPMD) technique from another tool to visualize the resulting process model and to pick the one which is considered to be the best one (according to user's expertise).
- Additional feature that allows for the automatic detection of "imprecise labels" by using properties of the Inductive Miner (IM).
- Additional feedback describing the changes being made to the original event log.

References

- [1] Lu, Xixi, et al. "Handling duplicated tasks in process discovery by refining event labels." International Conference on Business Process Management. Springer, Cham, 2016.