

FilterTree: a Repeatable Branching XES Editor (Extended Abstract)

Sander J.J. Leemans¹

¹RWTH, Aachen, Germany

Abstract

A large fraction of process mining efforts is spent on event data preparation: the step between data extraction and the subsequent analysis using process mining tools. Data preparation may be repetitive and is typically performed in a trial-and-error way. In this paper, we introduce the FilterTree XES and CSV editing tool, which allows for the programmatic chaining of XES and CSV filters, allowing for repeatable event data preparation. The FilterTree tool is platform-independent and open source.

Keywords

process mining, event log filtering, XES editor

1. Introduction

Common folklore in process mining tells that of the time spent on process mining projects, 80% is spent on preparing event data, and only 20% is spent on analysis. Massaging event data into an event log is a trial-and-error process, which may involve selecting activity columns, selecting case columns, altering data types, combining columns, parsing timestamps, filtering, addressing data quality issues, selecting activities, computing aggregate columns, etc. The repetitiveness of this process is captured by several process mining methodologies, such as [1]: even in the final phases of analysis, the data preparation may have to change, for instance after discovery of data quality issues or after analysis questions have been adjusted, thus requiring a slightly different view on the event data.

In this paper, we propose a tool to import CSV or XES files, and edit XES event logs by means of filters. A *filter* reads an XES (or CSV) file from disk and writes an adjusted XES file to disk. Filters are organised by the user into a *filter tree*, which specifies the filters with their parameters. In a filter tree, most filters are sequential, that is, one is applied to the result of its predecessor. It is also possible branch the filter chain, where a single filter may have more than one subsequent filters.

2. Significance, Innovations & Main Features

Every process mining project starts with extracting event data. Except in straightforward “projects” where public data is used in a standard process mining technique, event data preparation is a necessary next step before the actual process analysis can commence. Many existing

ICPM 2022 Doctoral Consortium and Tool Demonstration Track



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).



CEUR Workshop Proceedings (CEUR-WS.org)

tools can perform event data preparation and apply filters to an event log, and a few tools can perform edit operations on XES logs [2, 3, 4, 5]. To the best of our knowledge, there is no tool that is a combination of:

- XES-based, with support for CSV. Importing CSV files is critical, but advanced process mining operations require XES concepts, such as trace attributes, log attributes, summing trace outcomes, etc.
- Repeatable. The same set of filters can be applied to a new event log, thereby repeating the analysis without manual filtering steps, or to slightly change how an event log is prepared, requires easy repeatability.
- Branchable. In process mining projects, several different perspectives may be necessary to answer the analysis questions. These perspectives may require different event logs. Branching allows to re-use parts of chains of filters.
- Disk based and file-manager friendly. A common problem encountered is CSV or XES files that are too big to visualise or load into process mining tools. While filtering, logs should be handled disk-to-disk as to support any log that fits on disk.
- Offline. An offline tool provides privacy and confidentiality, and does not need to upload datasets.
- Extensible. There is always another, unsupported, filter, thus the tools needs to be easily extensible.

The FilterTree tool aims to satisfy all of these properties.

2.1. Notable Plug-ins

CSV files can be either row-based or column-based; for both, FilterTree has a plug-in. In a row-based CSV file, each row represents an event. Using the filter `CSV to XES`, the only parameter necessary is the name(s) of the column(s) of the CSV file that contain the trace identifier, that is, the column that tells us which case the event belongs to. This column or combination of columns is copied to the trace level as its concept `:name`.

In a column-based CSV file, each row represents a trace, and the columns contain timestamps, indicating when the activity belonging to that column was executed. This structure of data is often encountered in healthcare settings, where standard forms that are used to log treatment steps use this structure. The plug-in `CSV to XES - trace per row` converts such a file into an XES event log, where each row becomes a trace, and every cell that has a timestamp is converted to an event (with the concept `:name` being the name of the column); every cell that does not parse as a timestamp becomes a trace attribute. This plug-in optionally takes a list of Java-based timestamp formats, such as `yyyy-M-d H:mm:ss.SSS`. Both of these CSV-plug-ins set some default log attributes, and attempt to guess the data type of each cell as accurately as possible.

Using the `map events` and `map traces` plug-ins, a particular event/trace attribute can be transformed using a provided map (in a separate CSV file), and written as another event/trace attribute.

The `add start events` filter copies each event, copies a chosen trace timestamp attribute to `time:timestamp` of the new event, sets `lifecycle:transition` to `start`, and adds a

corresponding concept : instance to both events. The plug-in `sort_events` sorts the events based on `time:timestamp`.

3. Usage

3.1. File Format

A filter tree is represented in a simple text file (with the `.ftree` extension). Comment lines start with `%`. The first line in this file contains the `import event log`, and each line thereafter contains one filter. On such a line, the name of the filter comes first, followed by the bar `|` symbol, followed by the parameters necessary for the filter, separated by spaces. If a parameter contains a space, it must be enclosed in double quotes. Indenting a filter line starts a new branch.

3.2. User Interface

The user interface shows a filter tree file, with syntax highlighting and auto-completion; a screenshot is shown in Figure 1. When the user changes something by typing, after a small timeout the filters are automatically (re-)computed as necessary. The resulting logs, including all intermediate steps, are kept in a managed sub-folder; NB: the tool removes irrelevant files from this sub-folder. In the sub-folder, the final results of all branches are kept with a consistent and predictable filename, for compatibility with file management systems.

The user can enable a visualisation of the last log of the last branch, however, please note that this will attempt to load the log in memory; henceforth, this option is not enabled by default.

3.3. Download

FilterTree is platform-independent (Java) and available with a GPL license from <https://leemans.ch/filtertree>. A full list of supported filters is included on this website, as well as a screencast demoing the tool. The source code is available at <https://svn.win.tue.nl/repos/prom/Packages/SanderLeemans/FilterTree/>. An empty text file can be used to start the editor.

4. Maturity

The FilterTree tool has been used in several of our own projects, including on private healthcare data (Figure 1), historical bureaucrat career data, road traffic fine collection data [6], etc. These settings ranged from simple (e.g. lifting a few event attributes to trace level) to complex (see Figure 1). In this latter case, the initial log was too complex (260MB CSV) to be of use directly in ProM or any other process mining tool we were allowed to try, due data confidentiality and semi-commercial nature of the data. The FilterTree tool allowed us to transform the log from CSV into XES and to filter it down to a manageable sub-view, which could be analysed in standard process mining tools.

```

FilterTree editor
Input event log COPY_Eventlog_2020Census_210705_07.csv Show result
1 CSV to XES | journey_id
2 remove empty event attributes |
3 %data types
4 make event attribute numeric | total_cost_abf,
5 make event attribute a timestamp | end_date_2 "d/M/yyyy H:mm"
6 make event attribute a timestamp | start_date_2 "d/M/yyyy H:mm"
7 make event attribute literal | activity
8 make trace attributes boolean | 0 death_14d death_30d death_6m q1_1920 q2_1920
9 %concept:name
10 copy trace attribute | journey_id concept:name
11 replace event attribute value | activity "service type - -Angio" "service type
12 copy event attribute | activity concept:name
13 %cost
14 copy event attributes to trace level | diag_code_trunc_d total_cost_abf cost:total
15 annotate trace with sum of trace attribute values | cost:total allieddirect_w
16 remove empty and zero-valued trace attributes | cost:total
17 %compute start events
18 add literal event attribute if not exists | lifecycle:transition complete
19 copy event attribute | end_date_2 time:timestamp
20 copy event attribute if not exists | start_date_2 time:timestamp
21 %add start events - conditional | start_date_2 end_date_2
22 sort events |
23 %lift trace attributes and clean up
24 copy event attributes to trace level |
25 remove constant trace attributes |
26 %select attributes
27 keep events with attribute prefix value | concept:name "service type"
28 remove empty traces |
29 %select traces of diagnosis (two branches)
30 set log attribute | concept:name L03
31 *keep traces with values of attribute | diag_code_trunc_d "L03 Diseases of
32 set log attribute | concept:name R55
33 *keep traces with values of attribute | diag_code_trunc_d "R55 Symptoms, s
Press ctrl+space for a list of available filters; use % for line comment

```

Figure 1: Screenshot of the FilterTree editor interface.

5. Conclusion

Preparing event data for analysis remains a rather ill-supported task, especially in settings with repeated small changes, large and complex event logs, data quality issues, or changing analysis questions. In this paper, we proposed a tool to edit XES and CSV files by means of filters. The FilterTree tool is repeatable as it keeps a full filter chain specification; it supports branching in the chain to allow multiple chains to share the same initial filters.

References

- [1] M. L. van Eck, X. Lu, S. J. J. Leemans, W. M. P. van der Aalst, PM²: A process mining project methodology, in: *Advanced Information Systems Engineering - 27th International Conference, CAiSE 2015, Stockholm, Sweden, June 8-12, 2015, Proceedings*, volume 9097 of *Lecture Notes in Computer Science*, Springer, 2015, pp. 297–313. URL: https://doi.org/10.1007/978-3-319-19069-3_19. doi:10.1007/978-3-319-19069-3_19.
- [2] A. Berti, S. J. Van Zelst, W. van der Aalst, Process mining for python (pm4py): bridging the gap between process-and data science, arXiv preprint arXiv:1905.06169 (2019).
- [3] A. Polyvyanyy, Process query language, in: A. Polyvyanyy (Ed.), *Process Querying Methods*, Springer, 2022, pp. 313–341. URL: https://doi.org/10.1007/978-3-030-92875-9_11. doi:10.1007/978-3-030-92875-9_11.
- [4] B. F. van Dongen, A. K. A. de Medeiros, H. M. W. Verbeek, A. J. M. M. Weijters, W. M. P. van der Aalst, The prom framework: A new era in process mining tool support, in: *Applications and Theory of Petri Nets 2005, 26th International Conference, ICATPN 2005, Miami, USA, June 20-25, 2005, Proceedings*, volume 3536 of *Lecture Notes in Computer Science*, Springer, 2005, pp. 444–454. URL: https://doi.org/10.1007/11494744_25. doi:10.1007/11494744_25.
- [5] H. M. W. Verbeek, J. C. A. M. Buijs, B. F. van Dongen, W. M. P. van der Aalst, Xes, xesame, and prom 6, in: *Information Systems Evolution - CAiSE Forum 2010, Hammamet, Tunisia, June 7-9, 2010, Selected Extended Papers*, volume 72 of *Lecture Notes in Business Information Processing*, Springer, 2010, pp. 60–75. URL: https://doi.org/10.1007/978-3-642-17722-4_5. doi:10.1007/978-3-642-17722-4_5.
- [6] S. J. J. Leemans, S. Shabaninejad, K. Goel, H. Khosravi, S. W. Sadiq, M. T. Wynn, Identifying cohorts: Recommending drill-downs based on differences in behaviour for process mining, in: G. Dobbie, U. Frank, G. Kappel, S. W. Liddle, H. C. Mayr (Eds.), *Conceptual Modeling - 39th International Conference, ER 2020, Vienna, Austria, November 3-6, 2020, Proceedings*, volume 12400 of *Lecture Notes in Computer Science*, Springer, 2020, pp. 92–102. URL: https://doi.org/10.1007/978-3-030-62522-1_7. doi:10.1007/978-3-030-62522-1_7.