

MASARYK UNIVERSITY  
FACULTY OF INFORMATICS



# The Categorization of the Darkweb

MASTER'S THESIS

**Bc. Linda Hansliková**

Brno, Spring 2020



## **Declaration**

Hereby I declare that this paper is my original authorial work, which I have worked out by my own. All sources, references and literature used or excerpted during elaboration of this work are properly cited and listed in complete reference to the due source.

Bc. Linda Hansliková

**Advisor:** RNDr. Martin Stehlík, Ph.D.



## **Acknowledgement**

My thanks go to my adviser RNDr. Martin Stehlík, Ph.D. for allowing me to proceed with this topic, for his advice and patience.

# **Abstract**

The future Categorization of the Dark Web abstract

## Keywords

Web Application, Dark-web, Python, RestApi, Django REST, Type-Script, React, Redux, d3 Graph





# Contents

1	<b>Introduction</b>	1
2	<b>Data set analysis</b>	3
2.1	<i>Clear web</i>	3
2.2	<i>Deep web</i>	3
2.2.1	Dark web	3
2.3	<i>Elasticsearch</i>	4
2.4	<i>The data set</i>	5
3	<b>Analysis</b>	7
3.1	<i>Web graph</i>	7
3.1.1	Community structure	9
3.2	<i>Louvain algorithm</i>	9
4	<b>Development</b>	11
4.1	<i>API</i>	11
4.1.1	Technology overview	11
4.2	<i>Front-end</i>	12
4.2.1	User interface	13
4.2.2	Technology overview	17
4.2.3	Implementation	18
5	<b>Conclusion</b>	21
5.1	<i>Evaluation</i>	21
5.2	<i>Future work</i>	21



## List of Figures

- 3.1 Web graph by Citeo made consisting of circa 600 000 domains and 16 billion links. [10]. 8
- 4.1 The view after the application has loaded. 13
- 4.2 The DD for the selection of the grouping-mode. 14
- 4.3 The input field with submit button for filtering. 14
- 4.4 The representation of a community node. The individual sectors of the pie-chart illustrate the category-composition in the community. 15
- 4.5 The representation of a page node. The colour of the square symbolizes the category of the page. 15
- 4.6 The sidebar with details of the selected community. 16
- 4.7 The sidebar with details of the selected page. 16
- 4.8 The pop-up window with detail-options for selection. These options dictate which details will be included in the downloaded text file. 17
- 4.9 The level indicator with the zoom-out button. After the zoom-out button is clicked, the user is shown the communities of the previous level. 17



# 1 Introduction



## 2 Data set analysis

The data set we were working with consisted of pages which had been scraped<sup>1</sup> from the dark web and stored in Elasticsearch (ES) [6]. This chapter will describe the dark web, the Elasticsearch database, and the structure of the stored pages.

The internet consists of networks from all over the world. The networks are connected and together they comprise a global network. The Internet can be divided into the clear web<sup>2.1</sup>, and the deep net<sup>2.2</sup>.

### 2.1 Clear web

The clear web, or Celarnet, represents approximately 4% of the Internet. The content on the clear web is indexed by search engines and is publicly accessible. The users are identified by their IP addresses and are usually not anonymous unless some privacy tools are used. The most used browsers worldwide according to market share are Chrome, Safari and Firefox [29].

### 2.2 Deep web

The deep web, or hidden web, represents approximately 96% of the Internet. The content is not indexed by search engines and is not accessible publicly. Emails or medical information are examples of resources in the deep web. A user needs to be authorized to access this data. A portion of the deep web is the dark web, also called darknet.

#### 2.2.1 Dark web

The dark web constitutes about 6.25% of the deep web. Tor [28] or I2P [1] are two of the networks comprising the darknet.

The Tor network (Tor) makes use of *onion routing* (ORing) [12]. ORing describes routing where each node, except for the originator and the target node, knows only its predecessor and successor. The

---

1. Scraping is a method used for collecting data from web pages with automated software rather than doing so manually.

originator chooses a list of nodes to create a circuit between itself and the target node. Nodes in this circuit are determining their successor leveraging public-key cryptography. For example, *A* is the originator and *D* the target node. *A* creates a circuit consisting of itself and of nodes *B*, *C* and *D*. *A* sends a fixed-sized cell with encrypted information about the circuit nodes using their respective public keys and the information to *B*. *B* removes one layer of the cell using its private key in order to get the address of its successor. *B* cannot remove an additional layer to the successive successor because it lacks the private key of *C*. *B* sends the cell to *C*. *C* removes another layer and sends the cell to *D*. *D* removes the final layer and reads the information. Now *D* sends the response in a similar manner back to *A*. Determining the original source and target node is therefore difficult. ORing was introduced in the 1990s to ensure privacy. Tor is accessible through Tor browser or Tor proxy.

The I2P network (I2P) is another network of the dark web which anonymizes its traffic.

Both exemplified browsers provide anonymous access to the clearnet as well as to the darknet.

The dark web is used for activities like accessing or publishing illegal or censored material, e.g. the Bible, whistle-blower secrets, or child porn. It is also used for trading with illegal goods, such as drugs or guns. Another way to use the dark web is to offer or order illegal services, for instance money laundering, hacking or murder.

## 2.3 Elasticsearch

ES is a distributed, open source search engine [6] and offers a fast full-text search. Another benefit of ES are documentation and supported tools, such as Logstash [5] for processing of data or Kibana [7] for the visualization of data. ES can be used as a NoSQL database. Such a database consists of indexes, documents and fields as opposed to tables, rows and columns of SQL databases. One of the advantages of NoSQL databases is scalability, therefore they are suited for big amounts of data.



## 2.4 The data set

The dark web pages were acquired via web scraping. The scraped pages belonged to two different networks - I2P and Tor. Both networks provide anonymity for the user. The total number of pages was 221,844. Of those pages 212,851 were Tor pages and 8,993 I2P pages. The total number of unique domains was 5,178 of which 4,912 were from the Tor network and 266 were from the I2P network.

The fields of a page entry will be described in the following list 2.4. Each description will include a concrete example from the database. Fields beginning with an underscore are assigned to every document implicitly.

**\_id** is a unique identifier. For example

*2d622b6fba6f203d790fedbb4f47963e2366c7fd.*

**\_index** informs about the collection the document belongs to. For example *tor*.

**\_type** determines the type of the document. For example *\_doc*.

**content** is the actual content of the page. For example *Purple Kush – 10g – WackyWeed Menu \* Home \* Contact us \* About us .*

**content\_type** describes the type of the content. For example *text/html; charset=UTF-8*.

**domain** of the url address. For example *wacky2yx73r2bjys.onion*.

**h1** is the text with the h1 style. For example *Purple Kush – 10g*.

**links** to other pages this page links to. For example {

*"link": "http://wacky2yx73r2bjys.onion/",*  
*"link\_name": "Home"*  
*}.*

**raw\_text** is similar to *content*. It additionally may contain formatting elements such as *\n*. For example *Purple Kush – 10g – WackyWeed Menu\n\n \* Home\n \* Contact us\n \* About us\n .*

**raw\_title** is similar to title. It additionally may contain formatting elements. For example *Purple Kush – 10g – WackyWeed*

**raw\_url** is the same as *url*.

**title** of the page. For example *Purple Kush – 10g – WackyWeed*.

**updated\_on** depicts the time when the document in the database was last updated. For example *2019-10-22T19:41:09*.

**url** address of the page. For example  
*<http://wacky2yx73r2bjys.onion/?product=purple-kush-10g>*.

### 3 Analysis

The purpose of this thesis was to categorize the dark web and visualize the result. In order to do that, the dark web had been scraped and the acquired pages had been stored in an ElasticSearch database. This was done as part of a bachelor's thesis which was being completed at the same time as this thesis. In order to retrieve the data from the database and perform various operations on them, a back-end (BE) needed to be created. One of the operations was the categorization of the acquired pages. For that, an appropriate topic modeling approach was required. We chose to adopt the Latent Dirichlet Allocation (LDA). To visualize the output a graph was used. However, the data set is rather sizable and to be able to provide the user with useful information, not all pages can be displayed at once. Therefore a proper way to divide the graph into several subgraphs had to be found. We decided to use the well known Louvain algorithm (LA). To display the graph in a comprehensible manner a front-end (FE) was created.

In this chapter we will describe LDA, which is the method used to categorize the scraped pages. Next we will talk about why it is necessary to divide immense numbers of nodes into clusters in order to display them as a graph. And lastly, we will detail communities and LA, the algorithm for dividing pages into communities.

#### 3.1 Web graph

We display our data as a graph. More precisely a web graph [18]. A web graph is a graph representation of the web where nodes are portrayals of the pages and edges depict links between the pages. Web graphs tend to be built from an enormous amount of data. As such, they can be advertised in various ways. One of the visualizations is shown in figure 3.1. The depicted web graph displays all its data at once without any labels or details. The result may be useful for viewing the internet as a whole. However, for our purposes this view was not sufficient. One of the requirements of this thesis was the possibility of the inspection of the relationships between nodes in more detail. The big amount of data described in subsection 2.4

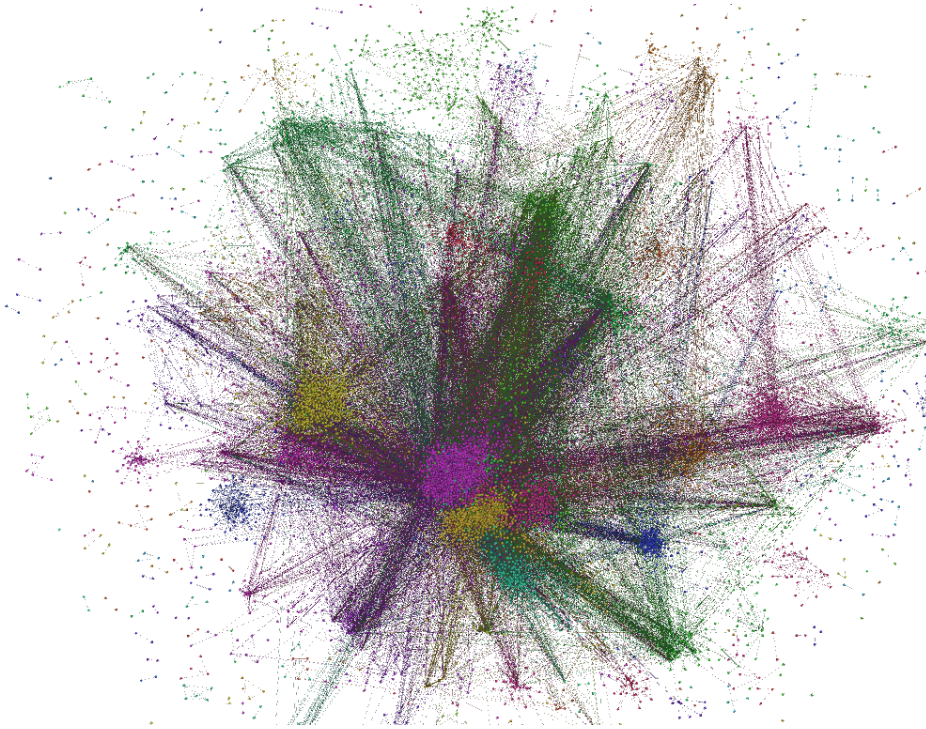


Figure 3.1: Web graph by Citeo made consisting of circa 600 000 domains and 16 billion links. [10].

prevented the displaying of all pages at once. The information the user would read from such a graph would be either incompatible with the requirements or too cluttered. Too much visual data worsens the quality of the user's information processing. The web graph therefore needed to be composed of a significantly smaller amount of nodes. This needed to be achieved in order for the view port not to be cluttered. This would result into obtaining the sought knowledge without hindrances. The majority of the nodes in our graph were not isolated <sup>1</sup>. The nodes were in fact part of a single connected component. As a result it was possible to partition the graph based on the density of its nodes into so called communities.

---

1. An isolated node is a node with zero incoming and outgoing edges.

### 3.1.1 Community structure

If a graph can be partitioned into several subgraphs so that nodes from one subgraph are internally connected densely and are connected scarcely to nodes from other subgraphs, we can claim it has a community structure. Each subgraph of such a graph is a community [13]. Each community can be portrayed as a meta node of the graph. This way the number of nodes in the graph can be reduced. The quality of such a partition can be measured using modularity [26]. L. Wenye and D. Schuurmans describe modularity in their work with the following words:

For a candidate partition of the vertices into clusters, the modularity is defined to be the portion of the edge connections within the same cluster minus the expected portion if the connections were distributed randomly [22]

[26].

## 3.2 Louvain algorithm

A widely used algorithm for finding communities in graphs is LA [3]. It is a greedy algorithm which maximizes modularity locally. Each node is assigned a community. Afterwards the node is taken out of its community and randomly appointed to the communities of its neighbours. After the node visited communities of all its neighbours it is left in the one with the maximum modularity value, which can also result in it remaining in its original community. Next, the algorithm runs on the newly gathered communities and tries to assign each whole community to its neighbouring communities in a similar manner. This is repeated until the modularity cannot be improved further.

LA is favoured for its simplicity, speed and accuracy. Since its introduction, in 2008, it was possible to detect communities in graphs with billions of nodes in a relatively timely manner. LA was compared to other algorithms for community detection [3]. Namely the algorithm of Wakita and Tsurumi [32], of Pons and Latapy [27], and of Clauset, Newman and Moore [9]. The used graphs were of sizes varying between 34 nodes and 77 edges to as much as 118 million nodes

and 1 billion edges. The difference between the computing times of the previously stated algorithms grows with the size of the graphs and favours LA. In fact, it took 152 minutes for LA to detect the communities of the greatest graph whereas the computation time of the other algorithms was more than 24 hours. In terms of precision, LA was also the most precise one with slightly better modularities.

## 4 Development

The aim of the development was to categorize the scraped pages into categories. This was to be done depending on the content of the pages. The second goal was to provide a way to observe the structure of the pages, links between them, and categories. One of the partner requirements was for the application to function on a UNIX system. Another requirement was a user friendly UI. Also, the option of retrieving all the available information about the pages was demanded.

This chapter describes the design and implementation of the application which is composed of a representational state transfer (REST) application program interface (API) and a front-end (FE) web application.

### 4.1 API

The scraped pages of the dark-web were stored in *ElasticSearch*. We created a back-end application (BE) in order to perform various operations on the data-set before sending it to the FE. Such operations involve resource intensive processing of large volumes of data, and caching. We decided to create a Python BE. The reason behind this decision was the requirement for the application to function on a UNIX system. Other reasons are described in subsection 4.1.1. The requirements for the BE was the categorization of the nodes along with their division into groups. The groups were either represented by nodes of the same category or communities. The BE also needed functionality to return details for specific pages or groups if requested.

#### 4.1.1 Technology overview

*Python* is a widely used interpreted programming language known for readability and portability [17]. Python is open-source and is considered to have an extensive documentation and community available.

Python is popular in the science community because it is easy to learn and has simple syntax. There is therefore a great amount of use-

ful libraries for research purposes such as *NetworkX*<sup>1</sup> [11] or *cylouvain*<sup>2</sup> [15] because of this. As we wanted to follow the REST architecture we decided to make use of the *Django* framework [14]. It is responsible for tasks such as running the server or managing web requests. Another advantage of Django is its *Django REST framework* (DRF). DRF offers a convenient way for creating restful endpoints and responses [23]. Both frameworks are open-source with helpful documentation and community.

We had to solve performance issues. It took approximately 30 minutes to retrieve and categorize about 220,000 pages from the database and circa xxx seconds to divide such a response into communities. We considered this wait time to be too long. Therefore we decided to cache the data of the first response. For that purpose *Redis* [21] was used. It is an open-source solution which we used as a key-value store. It supports basic data structures<sup>3</sup> as values, but not custom objects. Since the API uses custom objects for *communities*, *pages* and *links*, an object serializer was leveraged along with Redis. We decided not to write our own but to utilise the Python *pickle* module<sup>4</sup> [16] (pickle). The reason behind this decision was the simplicity of pickle. Pickle also fulfilled all our needs for serializing. Specifically the serialization or deserialization of data in the form of the previously mentioned models for Redis to store in an acceptable time. Pickle took xxx seconds to serialize xxx pages. The deserialization of the same number of pages took xxx seconds.

## 4.2 Front-end

For users to be able to see the data acquired from the BE in a reasonable way a FE application was created. The goal of this FE was to visualize the scraped pages in a graph. The pages or communities, and links from the BE were to depict nodes, and links of the graph respectively. The category of the pages or communities was to be readable

- 
1. A library used for creating and working with graphs.
  2. A library with a fast implementation of LA.
  3. Simple structures, e.g. strings, numbers or sets.
  4. A module used for converting Python objects to streams of bytes and vice versa.



from the graph. Additional information about the nodes needed to be displayed or retrieved on demand.

#### 4.2.1 User interface

We designed the UI the following way.

After the application is loaded the UI is composed of a *header* with the name of the application *Dark web categorization*, a *loader* and a *sidebar* on the right hand side. The application resembles at this point the image in Figure 4.1.

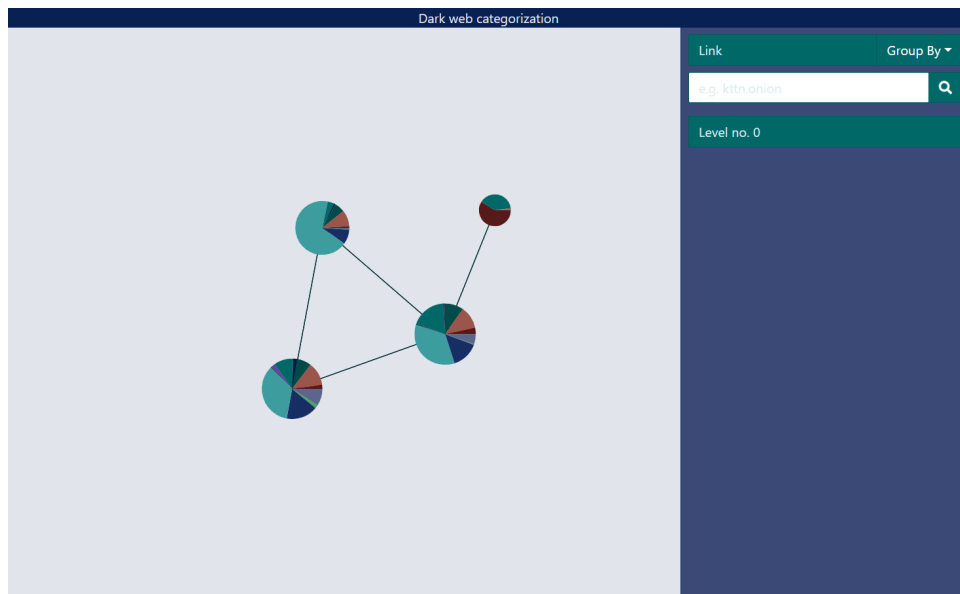


Figure 4.1: The view after the application has loaded.

The *sidebar* contains several *input fields* and *buttons* in a column. At the very top of the *sidebar* a *drop-down button* (DD) is present.

The DD sets the mode according to which the pages are grouped into communities. There are two modes available, *link-mode* and *category-mode* as can be seen in Figure 4.2. If *link-mode* is selected, the pages are divided by the LA3.2 depending on the connections between them only. If *category-mode* is selected, the pages are divided into groups by categories. The pages in the groups are further divided into communities according to links between them, as if in *link-mode*.



Figure 4.2: The DD for the selection of the grouping-mode.

An *input field* with a *submit button* was placed underneath the DD for the purpose of filtering nodes according to a search phrase as can be seen in Figure4.3.



Figure 4.3: The input field with submit button for filtering.

The last element shown is an indicator of the current level - how many times the user zoomed into a community. The indicator has a *zoom-out button* placed next to it.

The moment the data is retrieved from the BE the *loader* gets replaced for a graph. The graph-nodes represent either the communities the pages are partitioned into or the pages themselves. If the graph contains any isolated nodes (isolates) a *mock-community* is displayed containing all isolates. This community cannot be zoomed into. A community node is depicted as a *pie chart*. The individual sectors of the *pie chart* represent the categories of the pages belonging into the community. This is portrayed in Figure4.4. It was difficult to distinct between a small community of only one category and page nodes. A page node is therefore differentiated from a community node by depicting the node as a square shaped symbol. The colour of the square corresponds to the category of the page. An example of a page is illustrated in Figure4.5. It is possible for a level to depict communities and pages at once. The links between the nodes visualize the links between pages or communities.

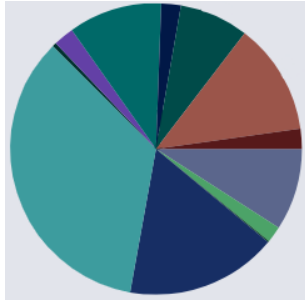


Figure 4.4: The representation of a community node. The individual sectors of the pie-chart illustrate the category-composition in the community.

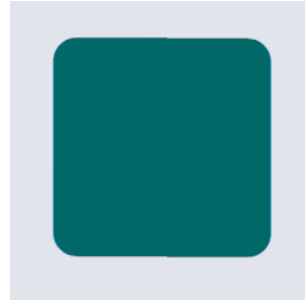


Figure 4.5: The representation of a page node. The colour of the square symbolizes the category of the page.

After single-clicking a node additional information is exemplified in the *sidebar*. The details vary depending on whether the node is representing a community or a single page. The details of an individual page are as follows:

**Url** which also serves as a unique identifier of the page.

**Category** of the page. Each page belongs to exactly one category.

**Links** to other pages displayed as a list of url addresses. There are up to ten links visible in the sidebar. The remaining links, if any, are downloadable in a text file. Figure 4.7 contains a sidebar view with page details.

The details of a community consist of grouped information of its members and include the following:

**Category composition** which is aggregated from the categories of all the pages of the community. Each category is represented by its name and the percentage of its relevance in the community.

**Page url** addresses (urls) of members belonging into the community. There are up to ten urls visible in the sidebar. The remaining urls, if any, are downloadable in a text file.

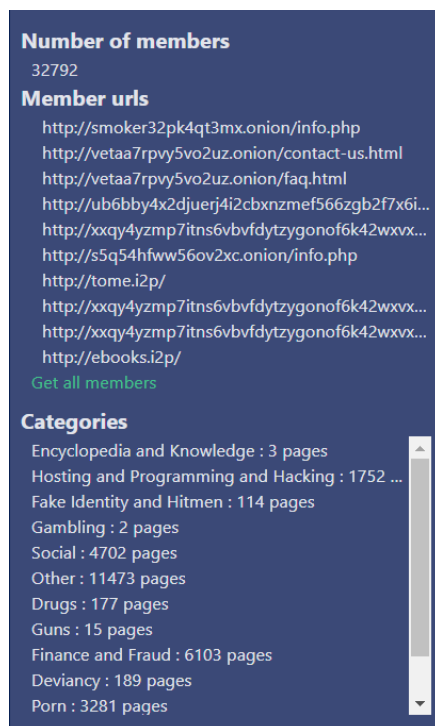


Figure 4.6: The sidebar with details of the selected community.

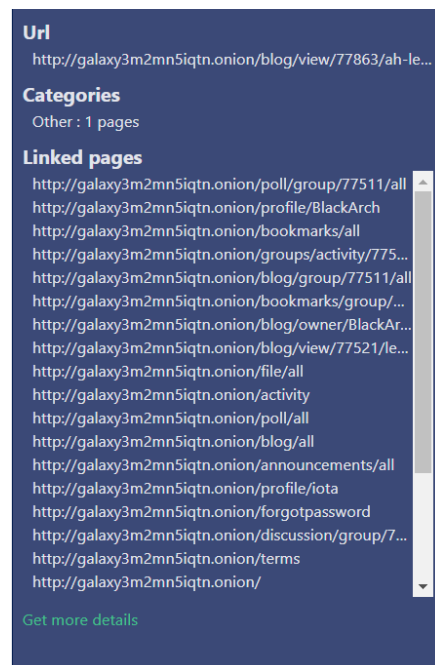


Figure 4.7: The sidebar with details of the selected page.

**Urls count** represents the number of all the pages belonging into the community.

Figure4.7 contains a sidebar view with community details.

In case of the need of further information a link, also present in the *sidebar*, can be clicked. If clicked, a *pop-up window* with detail-options is displayed. Details may include the *title*, *category*, *links* and *page-content*, depending on the user's selection. *Urls* of the pages are present by default. The window can be seen in Figure4.8. After selecting the needed options and clicking the *download button*, a *text file* with the desired information is downloaded. The page or community-members with their details are depicted in JSON format. This functionality was required for the scenario of the user needing to download all available data about members of a specific community.

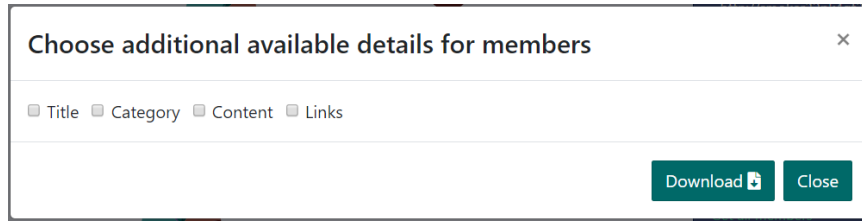


Figure 4.8: The pop-up window with detail-options for selection. These options dictate which details will be included in the downloaded text file.

A graph node representing a community can be double clicked. After doing so, a new graph is shown. The data of this graph consists of the members of the clicked community. We call this process *zooming*. In case the parent community contains too many sub-communities, it is displayed as a single community-node. The zooming-in or -out of communities adjusts the *zoom-level*. This level is represented by a number and is visible below the *node-filter*. If the current zoom level is more than zero, i.e. at least one community-node was double clicked, a button for zooming out appears next to the *level indicator*. It is shown in Figure 4.9.



Figure 4.9: The level indicator with the zoom-out button. After the zoom-out button is clicked, the user is shown the communities of the previous level.

If further zooming is not possible, the user reached the maximum level. Each community may have a different maximum level, depending on the number of its pages and its structure.

#### 4.2.2 Technology overview

*JavaScript* [25] (JS) is the most favoured programming language used for creating web applications [33]. It is an interpreted language supported by all modern browsers. It is open-source and as such dis-

poses of a considerable community with convenient documentation. JS is not strongly typed and the code might therefore be complicated to read or navigate. For this reason the FE was written in *TypeScript* (TS) [24] which is a superset of JS with the advantage of being typed. There is a significant number of tools and libraries for implementing user interfaces (UI) in a clean and timely manner created for both JS and TS. One of such tools is the framework *React.js* [19] (React). React is one of the most favoured JS frameworks [20]. The advantages of using React are readable code and improved performance by managing the re-rendering of page elements.

To achieve a satisfying UX the app needed to be interactive and obtain new or modified data frequently. Repeated requests to the BE would mean longer wait time for the user. However, a proper mechanism for data-storing would present a convenient solution to this issue and make the requests unnecessary. A JS library which handles the app state and works well with TS and React is called *Redux.js* [2] (Redux). Redux is a single store approach. This ensures easy hydration<sup>5</sup>. Redux also provides a custom set of TS typings and provides the developers with easy to use debugging tools. Another advantage of this library is the documentation.

The visualization of the web graph alone was realized using the *react-d3-graph* library [8] (RD3). RD3 is an implementation of the library D3.js [4] made more convenient for the use with React.js.

### 4.2.3 Implementation

The FE project consists of three folders and several configuration files. The folder `node_modules` contains imported libraries including *React.js*, *Redux* or *d3*. The next folder named `public` encloses a `.ico` file<sup>6</sup> and a `html` file which is the default entry point when the application is started. The last folder `src` contains the source code itself.

As previously mentioned, the FE is written in TS which has the advantage of readability and easy navigation. There are, however, also disadvantages. One of them is the need of a TS file with the types

---

5. The process of an object being provided with information

6. A picture with the dimensions 16x16 pixels used by the browser to represent the web page or application. It is usually displayed in the tab in which the application is opened.

(typing) for every used library. Typings for popular libraries are often downloadable as modules. If a library has no ready-to-download typings own ones need to be written. In our case the typings for the library *react-d3-graph* were custom made. They can be found in the folder `@types/react-d3-graph`. The file `common.d.ts` holds types used heavily across the application, e.g. *Action*. Types in this file are available without importing them to all files in the project.

Objects passed between functions also need to be typed. Those models are stored in the folder `models`. Each file contains one server-model and one client-model. The conversion between these models is conducted in specific helper functions. The advantage of this approach is the independence of client-models from the BE models.

The visual aspect is implemented using Less [30] which is a language extending CSS with improvements such as the possibility of using variables. The Less classes are divided into files depending on the element they are meant to modify. These files were placed into the folder `styles`.

The remaining folders each represent a different part of the UI. The structure of their sub-folders is similar. Therefore it is sufficient to describe them as a whole. Folders named `utils` contain files with helper functions such as converters between server and client models. `Constants` contains folders with string constants or simple functions which return a string depending on the input. The rest of the folders represent some part of the Redux framework.

The most basic files which only include string constants are situated in the folders named `actionTypes`. These are utilised as action types in actions. An action is a simple objects containing a type and an optional payload. Actions themselves are returned by action creators (AC). AC are functions returning an action and can be found in folders called `actions`. They can be as simple as those present in the file `nodesActionCreators.ts`. But they can be more complicated such as the AC `fetchNodes.ts` and dispatch multiple simple ACs. The purpose of an AC is to be injected into reducers, i.e. dispatch them.

`fetchNodes` and the folder it is placed in share the same name. For easier testing purposes the main logic of this AC is put into a function which receives the simple ACs as dependencies. When this AC is called it first dispatches a simple AC to indicate the fetching has begun. After that an identifier (id) is created. This id is later used to

create an error object in case of failure. Next, the fetching itself begins. The fetching in `fetchNodes` is realized with the library *isomorphic-fetch*. The `fetch` function of this library expects the first argument to be the url address of the resource. The second argument is an object describing further details of the request and is optional. Such an object may contain the request method, headers or the payload. If the request does not result in error the response status is checked. After the fetching is complete a success AC with the acquired response is dispatched. If an error is caught during the fetching a failure AC is dispatched. The payload of this AC is an error-object with the id and error message if any.

The dispatching of actions enables the changing of the state via reducers situated in the reducers folders. The state is a single immutable<sup>7</sup> object and is used in the whole application. A reducer is a pure function<sup>8</sup> receiving the current state and the dispatched action as its arguments. It then returns the newly computed state. A reducer creates a new state only if the type of the given action is recognized. If not, the previous state is returned unmodified. The state object received or returned by the reducer does not need to be the entire state (app-state). A reducer may be responsible for just a part of the app-state. However, the root reducer is responsible for the whole app-state.

The folders components contain files with React components. They define the skeleton of the UI with the specified behaviour. The folders containers hold files with React containers. A container is a file with access to the app-state. It is responsible for passing data to components.

---

7. The object cannot be adjusted directly. Instead, a new modified object is returned and the original one stays unmodified.

8. The return value of a pure function is only dependent on its input values. A pure function has no side effects.



## **5 Conclusion**

### **5.1 Evaluation**

### **5.2 Future work**

The Leiden algorithm (LeA) might be used in the future instead of LA. LeA is another algorithm for community detection on large graphs. The paper in which LeA is described claims LA has a major flaw which is eliminated in LeA [31]. It needs to be taken into consideration if it becomes adopted widely.



## Bibliography

- [1] Introduction to the i2p network. <https://geti2p.net/en/about/intro>. [cit. 2020-23-01].
- [2] Dan Abramov and the Redux documentation authors. Redux - a predictable state container for javascript apps. <https://redux.js.org/introduction/getting-started/>, 2019. [cit. 2019-10-09].
- [3] Vincent Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics Theory and Experiment*, 2008, 04 2008. [cit. 2019-08-16].
- [4] Mike Bostock. D3.js - for manipulating documents based on data. <https://d3js.org/>, 2019. [cit. 2019-10-09].
- [5] Elasticsearch B.V. Centralize, transform & stash your data. <https://www.elastic.co/logstash>, 2020. [cit. 2020-03-03].
- [6] Elasticsearch B.V. What is elasticsearch? <https://www.elastic.co/what-is/elasticsearch>, 2020. [cit. 2020-01-03].
- [7] Elasticsearch B.V. Your window into the elastic stack. <https://www.elastic.co/kibana>, 2020. [cit. 2020-03-03].
- [8] Daniel Caldas. React-d3-graph - interactive and configurable graphs with react and d3 effortlessly. <https://goodguydaniel.com/react-d3-graph/docs/>, 2019. [cit. 2019-10-09].
- [9] Aaron Clauset, M. E. J. Newman, and Cristopher Moore. Finding community structure in very large networks. *Phys. Rev. E*, 70:066111, Dec 2004. [cit. 2019-08-28].
- [10] criteo engineering. Web graph by criteo. <http://engineering.criteolabs.com/2014/05/the-web-graph-as-seen-by-criteo.html/>, 2014. [cit. 2019-08-16].

- [11] NetworkX developers. Networkx - software for complex networks. <https://networkx.github.io/>, 2019. [cit. 2019-09-09].
- [12] Roger Dingledine, Nick Mathewson, and Paul Syverson. Tor: The second-generation onion router. In *Proceedings of the 13th Conference on USENIX Security Symposium - Volume 13, SSYM'04*, page 21, USA, 2004. USENIX Association. [cit. 2020-04-03].
- [13] Santo Fortunato. Community detection in graphs. *Physics Reports*, 486(3-5):4–8, Feb 2010. [cit. 2019-08-16].
- [14] Django Software Foundation and individual contributors. Django - a high-level python web framework. <https://www.djangoproject.com/>, 2019. [cit. 2019-09-09].
- [15] Python Software Foundation. Cylouvain - a fast implementation of the louvain algorithm. <https://pypi.org/project/cylouvain/>, 2019. [cit. 2019-09-09].
- [16] Python Software Foundation. Pickle - python object serialization. <https://docs.python.org/3/library/pickle.html>, 2019. [cit. 2019-09-09].
- [17] Python Software Foundation. Python - an interpreted, high-level, general-purpose programming language. <https://www.python.org/about/>, 2019. [cit. 2019-09-09].
- [18] Jean-Loup Guillaume and Matthieu Latapy. The Web Graph: an Overview. In *Actes d'ALGOTEL'02 (Quatrièmes Rencontres Francophones sur les aspects Algorithmiques des Télécommunications)*, Mèze, France, 2002. [cit. 2019-08-16].
- [19] Facebook Inc. React.js - a javascript library for building user interfaces. <https://reactjs.org/>, 2019. [cit. 2019-10-09].
- [20] Stack Exchange Inc. Most popular frameworks according to stack overflow. <https://insights.stackoverflow.com/survey/2018#technology-frameworks-libraries-and-tools/>, 2020. [cit. 2020-06-01].

- [21] Redis Labs. Redis - a in-memory data structure store. <https://redis.io/>, 2019. [cit. 2019-09-09].
- [22] Wenye Li and Dale Schuurmans. Modular community detection in networks. In *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence - Volume Volume Two*, IJCAI'11, pages 1366–1371. AAAI Press, 2011. [cit. 2019-08-28].
- [23] Encode OSS Ltd. Django rest framework - a flexible toolkit for building web apis. <https://www.django-rest-framework.org/>, 2019. [cit. 2019-09-09].
- [24] Microsoft. Typescript - a typed superset of javascript. <https://www.typescriptlang.org/>, 2019. [cit. 2019-10-09].
- [25] Mozilla and individual contributors. Javascript - most well-known as the scripting language for web pages. <https://developer.mozilla.org/en-US/docs/Web/JavaScript>, 2019. [cit. 2019-10-09].
- [26] M. E. J. Newman. Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577–8582, 2006. [cit. 2019-08-16].
- [27] Pascal Pons and Matthieu Latapy. Computing communities in large networks using random walks. *Journal of Graph Algorithms and Applications*, 10(2):191–218, 2006. [cit. 2019-08-28].
- [28] The Tor Project. Introduction to the tor network. <https://www.torproject.org/about/history/>. [cit. 2020-23-01].
- [29] StatCounter. Browser market share. <https://gs.statcounter.com/browser-market-share>, 2020. [cit. 2020-04-03].
- [30] the core Less team. Less - a little more than css. <http://lesscss.org/>. [cit. 2019-12-09].
- [31] Vincent A. Traag, Ludo Waltman, and Nees Jan van Eck. From louvain to leiden: guaranteeing well-connected communities. In *Scientific Reports*, 2018. [cit. 2019-08-16].

- [32] Ken Wakita and Toshiyuki Tsurumi. Finding community structure in mega-scale social networks. *CoRR*, abs/cs/0702048, 2007. [cit. 2019-08-28].
- [33] Carlo Zapponi. Active repositories per language on github. <https://github.info/>, 2014. [cit. 2019-10-09].