
d3.js and its potential in data visualization

Creating a diagram showcase using ukrainian refugee data

Luis Rothenhäusler

20202459



Bachelorarbeit - Exposé

Fachbereich Informatik
und Medien
Technische Hochschule Brandenburg

Betreuer: Prof. Julia Schnitzer
2. Betreuer: Prof. Alexander Peterhänsel

Brandenburg, den 24.05.2022
Bearbeitungszeit: dd.mm.yyyy - dd.mm.yyyy

Abstract - German

Englische Arbeiten brauchen eine Zusammenfassung auf Deutsch. Mal abgesehen davon, dass wenn die Zusammenfassung interessant ist man ohne Englisch eh nicht weiter kommt...

Contents

1	Introduction	5
2	Basics	6
2.1	Data	6
2.1.1	Categorical	7
2.1.2	Numeric	7
2.2	Diagrams	8
2.3	D3.js	9
2.3.1	Selections	10
2.3.2	Data Joins	11
2.3.3	General Update Pattern	11
2.3.4	Modules	13
3	Implementation	14
3.1	Datasets	14
3.1.1	Preprocessing	14
3.1.2	Data Types	15
3.2	Diagrams	15
3.2.1	Common features	15
3.2.2	Bar Chart	19
3.2.3	Pie chart	19
3.2.4	Tree map	19
3.2.5	Sankey	19
3.2.6	Area graph	20
3.2.7	Circle graph	20
3.3	Showcase	20
3.3.1	Integration of each diagram	20
3.3.2	Data Updates	20
4	Discussion	22
4.1	Evaluation Criteria	22

5	Conclusion	23
6	Appendix	25

1. Introduction

Describe the background of the thesis, why it is important, what do we want to achieve. Why is this thesis (Motivation)? What do we want to do? What is the status quo? What benefits will result from this thesis? Something about the importance of infographics and comprehensible data.

The postmodern world produces huge amounts of data every second. Analyzing this data leads to better-informed decision-making in every sector. Yet the vast amounts of created data is often hard to comprehend with the human mind. Data visualization is about finding ways to represent this data in visually appealing and easily understandable visual representations. Doing this quickly and always up to date can be crucial. There are many tools available to help with the creation of infographics. This thesis focusses on computer based solutions. Some of the data visualization tools have a graphical-user-interface, others are code based. It is often not easy to decide which tool best suits ones needs. Therefore this thesis will be a deep dive into the possibilities of one of these tools, the 'd3.js'(D3) library for JavaScript. This will help evaluate if D3 is the right tool for the job, as well as provide a starting-point for using the technology. To achieve this there are two main questions that will be answered in this thesis. What is the potential of D3 in data visualization? What are the advantages and disadvantages of using D3? To be able to evaluate these questions a showcase of a several different diagrams is created.

2. Basics

Everything necessary to understand the implementation as well as anything which is done beforehand, will come up here

In the following, all concepts, technologies and required backgrounds for understanding this thesis are explained. Firstly data and data types are described, followed by the need and make of diagrams and finally D3 as the tool to create diagrams.

2.1 Data

We'll talk about data a bit. Where does it come from? How is it structured? What kind of attributes? What even are attributes?

Since ancient times, humans have recorded data. Recording the ins and outs of available resources was one of the driving factors behind the conceptualization of writing. (TODO: Check sources of the beer brewing video series) With the introduction of computers the amounts of gathered data have grown drastically. Vast amounts of data are gathered across all aspects of life.

The vast amounts of data gathered in databases are often hard to comprehend and evaluate with the human mind. They are also unwieldy to present them in the often limited space of articles, dashboards or other informative purposes. Therefore data visualization (Figure 2.1) is used to turn these datasets, collections of data-points, into diagrams. We constantly come across the results of data visualization in everyday life. They can be commonly found across all kinds of news sources, but also in reports, information campaigns or as part of user-interfaces in machinery or control systems.

Preprocessing is commonly done in data visualization. Depending on the dataset and the desired result, this can mean different things. One might want to remove excessive information from the dataset, which is not necessary for the representation. On the other hand, additional data can be added by evaluating the existing data-points. These could for example be the median

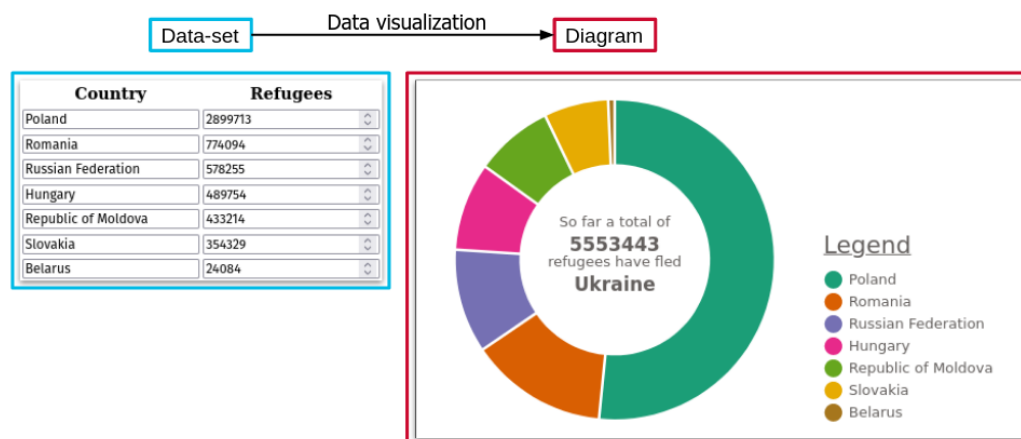


Figure 2.1: Where data visualization comes into play.

of values. It is important to not abuse this step to remove or add arbitrary data-points, to make the dataset align with the desired results.

Even though data comes from a huge variety of sources and can express a plethora of things, there are only four different types of data. They are split into two categories. Categorical and numerical data. Each category has two subtypes. In the following each of the types of data will be explained.

2.1.1 Categorical

What is categorical data? Nominal and ordinal data

Categorical or qualitative data is information collected in groups. It is often of descriptive nature. Whilst the values can be represented in numbers, they do not allow for arithmetic operations. There are two types of categorical data. Nominal and ordinal data.

Nominal data is mostly descriptive in nature. They are independent and have no inherit order. Examples are 'Country', 'Color', 'Brand'.

Ordinal data is mostly similar to nominal data. Yet the data does have some sort of internal order. For example different dates each describe a day, but one day also comes after another.

2.1.2 Numeric

What is numeric data? Continuous and Discrete

Numeric or quantitative data is all data expressed in numbers, where numbers do not represent categories. It allows for arithmetical operations and can be split into Discrete and Continuous data.

Discrete data can only take certain defined values. This usually means whole numbers to represent things that can not be split up further. Like the 'Number of Refugees' or 'Tickets sold'. Discrete data is countable.

Continuous data can be measured. It can have any real number as value. Therefore fractions are possible as well. For example when measuring the temperature, or the length or weight of an object.

2.2 Diagrams

What diagrams exist? Which are the most common? What possibilities do they offer for encoding data? Which considerations for readability? Why do some diagrams not make as much sense? Which considerations were made for fulfilling the showcase requirements?

The selection of which diagram should be used to visualize which dataset is not trivial. There are a plethora of diagrams already in use and anyone can create new diagrams to suit their needs. Yet there are diagrams which are used more commonly. These include bar and pie-charts, scatter plots and heat-maps (TODO: Find source (1)). Mostly there are several possible diagram choices for the given data as well.

All diagrams use a combination of marks and channels to encode data. Marks are used for entries in the diagram. The three possible marks are points, lines and areas. Each mark uses at least one channel to encode data. The most commonly used channels are position, size, color and texture. The position in 2D can be split into the x and y positions. The color is split into hue and luminescence. For example looking at fig. 2.2 we can see lines being used as marks for each entry. It might seem like we are using areas, but the thickness of the line only serves visual understanding. The lines also use three channels to encode data. The y-position is used to represent the categorical data of which country. The hue of the bar encodes the same data. The size, in this case length, of the bar encodes the discrete data of how many refugees have crossed into the country. In fig. 2.2 we see areas used as marks. Just like in the previous example the hue encodes the country and the size encodes the refugee count.

All marks can be used with all channels. But not all data types should be represented by all channels. For example nominal data should not be

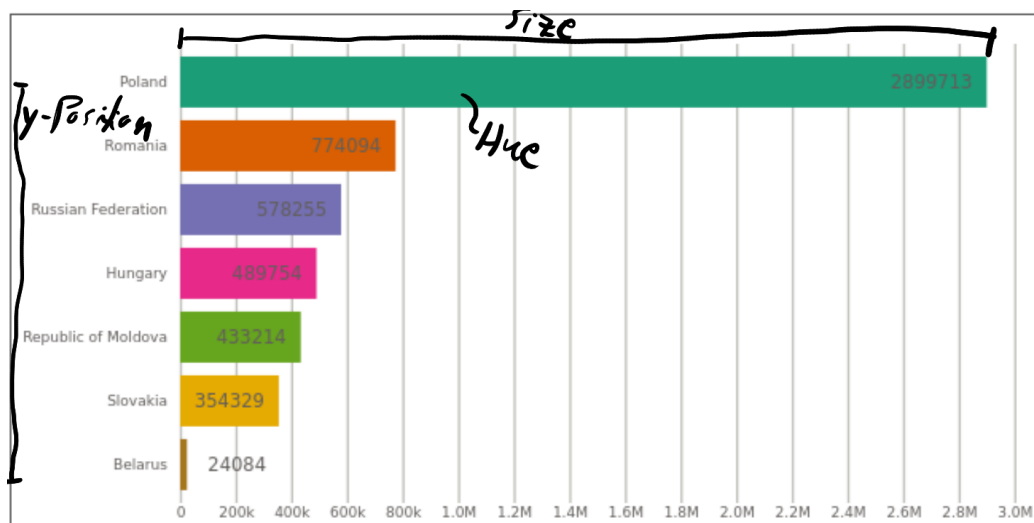


Figure 2.2: This is a bar-chart with used channels shown.(TODO: which needs a frame..? Also draw in marks and channels)

encoded using the size channel. The different sizes would lead to a perceived order, which does not exist in nominal data. As the channels all differ in their appearance they are also not equally good in adequately representing the data types. Therefore it is important to consider which channels are chosen to represent the given data types. Generally the position channels can always be considered the strongest channels [1].

2.3 D3.js

This is all about d3. What is it? Where does it come from? What is it used for? Who uses it? Why should it be used? How does it work? Enter, update and exit pattern. Something about the modular structure of D3 as well. Might be worth mentioning "observables" as well.

"D3.js is a JavaScript library for manipulating documents based on data. D3 helps you bring data to life using HTML, SVG, and CSS." [2]. The name D3 is short for data-driven documents. The D3 library was originally created by Mike Bostock and is published under the BSD-3-Clause open-source license. It is about 500kb in size. It does not require a specific framework and can therefore be easily integrated into all kinds of web based projects. Whilst D3 is not limited to using svg, the visualization created using D3 mostly rely on svg elements for their implementation.

D3 is not a high-level API for creating out of the box visualizations. In-

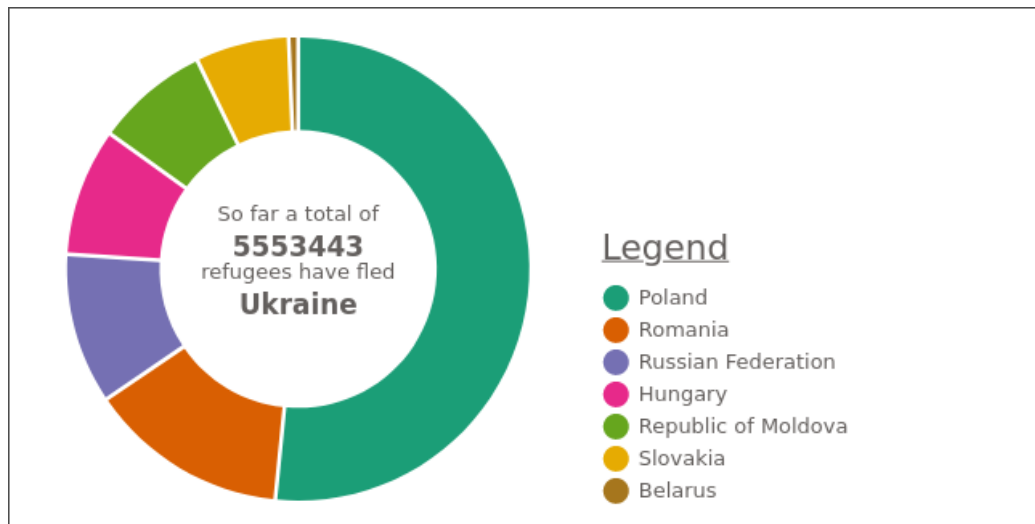


Figure 2.3: This is a donut-chart (TODO: which needs a frame..? Also draw in marks and channels)

stead it is a library which aims at making the tedious parts of DOM manipulation easier. It also provides some helper functions to decrease the amount of mathematical equations needed to convert from the data extends to the necessary coordinates in the desired visualization.

TODO: merge the section above and below to one longer introduction

General functioning of D3.

"D3 allows you to bind arbitrary data to a Document Object Model (DOM), and then apply data-driven transformations to the document." [2]. There are three main concepts that make up the core of D3. Selections, data joins and the general update pattern. All three of these concepts are working closely together. Whilst selections can be used without data joins and the general update pattern, these two aspects both rely on selections. Data joins can also be used without explicitly using the general update pattern. But the general update pattern can not function without data joins. (Importing the d3 library into a project allows access to the d3 namespace.) Usually all three of these concepts are used consecutively. First a selection is created. This collection is provided with a data join. Finally the behaviors for the general update pattern are defined for this data join. In the following, all three of the core concepts of D3 are explained.

2.3.1 Selections

What are they? Why are they useful?

All operations in D3 run on an arbitrary collection of nodes. These collections of nodes are called selections. There are two functions in D3 to create a new selection: `d3.select("selector")` and `d3.selectAll("selector")`. Both functions require a selector for identifying the appropriate elements. The selectors are defined in the W3C Selectors API[3]. Whilst `select` only selects a single element, `selectAll` selects all elements which match the selector. It is important to note, that `select` also propagates the existing information of this node, whilst `selectAll` does not. It is possible to directly access the DOM elements through the selections. The respective DOM elements are linked in the nodes which make up the selection. But usually this is not required, as there are predefined functions for modifying the nodes properties. This includes the modification of attributes and styles, as well as event handling. Selections can also be extended or shrunken by adding or removing nodes, or by combining multiple selections.

2.3.2 Data Joins

What are they and why are they important?

Data joins are a key feature of D3. They link up a specific data-point to a specific DOM element. To create a data-join, one has to call the `.data(dataset)` function on a selection. It takes a dataset, an array of data-points, as parameter. This will bind the data elements to the nodes in the selection. When we want to create diagrams which can respond to data changes over time, we need to be able to correctly identify the data-points. The default identifier function corresponds to the index of the data-point. This can lead to unexpected behavior when the dataset is modified. Especially when data-points are removed or added in arbitrary positions. Therefore we can specify a custom identifier function. This is passed as the second parameter of the data function, will be called for each data-point and has to return some value which will be used as an id.

When creating a data join, it is not important for the number of data-points to match up with the number of nodes in the selection. This is taken care of through the general update pattern.

2.3.3 General Update Pattern

What is it? What can it do? Describe data joins and dom element links.

The general update pattern is another core concept of D3. Every time a data join is created or updated, it comes into play. The general update pattern differentiates between three different cases. For each of these cases a selection is created. For each of these three selections the behavior can be

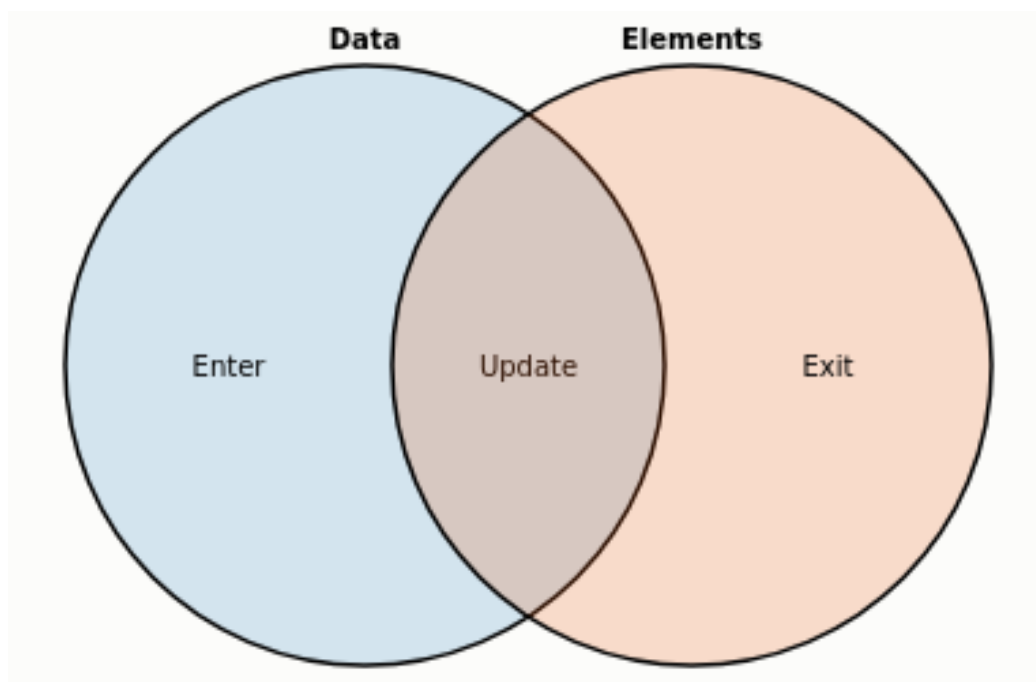


Figure 2.4: A visual representation of the make up of the general update pattern[4]

defined. The first selection is the enter selection. All data-points which do not have a matching element are provided with a placeholder element while creating the data join. All these placeholder elements make up the enter selection. In the behavior for the enter selection, a corresponding element is usually created as the first step. All the elements which are already linked to a data-point make up the update selection. The last selection, the exit selection, is made up of all the elements for which the corresponding data-point has been removed. The behavior of the exit selection is defined by default, to remove these elements.

When the goal is to create only static diagrams, which are only initially created from data, it is enough to define the behavior for the enter selection. If diagrams should be able to react to data changes and update their appearance, like in this thesis, it is also important to define the update behavior. The exit behavior can be defined, if a more visually pleasing removal of elements is desired.

2.3.4 Modules

The way D3 is split up into modules, the core package and what kind of extensions are there.

D3 provides the most used, general functionalities in the core library. Yet there are many more modules which can be added, which add functionalities for more specific use-cases.

3. Implementation

How are the chosen diagrams implemented? Which D3 modules have been used? How was the implementation done?

In the following sections the process of creating the showcase and the diagram are described. There are several parts to this. At first the data-sets, which should be represented, are chosen. In most real world usages, this is already given. Afterwards the possible diagrams are considered and chosen. Their implementation and usages of D3 are described. Finally the showcase bringing all the diagrams together is described.

3.1 Datasets

What are our datasets about? Where do they come from?

In this thesis, two data-sets are used. They are both from UNHCR[5]. The first data-set contains information about the total number of refugees per country[6]. The second dataset is about the total cumulative total amount of refugees per day[7]. Both data-sets are in the JSON format. As different data-types allow for different representations and require varying parts of D3, the data-sets have been specifically chosen to cover all data-types.

3.1.1 Preprocessing

What is done in preprocessing? Python script which removes all excess / maybe do that in JS as well..?

It would be possible to use both these data-sets as is. Yet the vast amounts of filler data, which does not contain any valuable information, makes data accesses unnecessarily complicated. Therefore both data-sets are preprocessed. In both cases the data-sets are read and the valuable information extracted and saved in the csv format. Both of the newly created csv files have two columns and one header row. The resulting csv of the refugees per country dataset, contains the two columns of country and refugees. The

other resulting csv contains a column for the date and one for the cumulative refugees. The data-set about the refugees per country can also easily be converted into using percentages. After adding up the total amount of refugees from each data-point, one can convert the absolute number of refugees into percent.

3.1.2 Data Types

Which data types can be found in our data-sets? Where?

The two chosen data-sets already cover most of the data-types. The refugees per country dataset contains two attributes per data-point. The country is a categorical attribute. The number of refugees is discrete. When converting this data-set into using percentages, the percentage of refugees becomes a continuous attribute. In the refugees per date data-set, the amount of refugees is still a discrete attribute. The date itself is an ordinal attribute though. As one day clearly comes before and after another day.

Choosing data-sets which cover all types of data-types was an important consideration. Different data-types can have different ways of representation, as well as different ways of implementation on the programming side of things.

3.2 Diagrams

describe all the diagrams and why they are special and what makes them tick. Why have they been chosen?

The following section is all about the selection and implementation of each diagram. As the selection of diagrams was heavily influenced by differences in functions of D3 they use, the common features are explained first, before diving into the more detailed explanation of each diagram and their specifics.

3.2.1 Common features

What does every diagram use?

All diagrams consist of three parts. A HTML, a CSS and a JavaScript file. The HTML loads the D3 library in the header. The body of the HTML only consists of a `svg` tag where the diagram will be drawn, and a script tag which loads the JavaScript file. The CSS defines the general styling of the diagram which is not dependent on the input data. The main part of the implementation is done in the Javascript section. The JavaScript part follows a specific pattern in its implementation. At first there is a general setup section, which is run once as the website is loaded. It is followed by

a render function, which is responsible for drawing the diagram. The render function requires the data-set as input. It is also split into two parts. At first the requires scales are initiated. It is followed by the section responsible for creating and manipulating the diagram. This is done by firstly creating a data join, before specifying the behavior of the general update pattern. The render function is called once initially and every time the data changed. Therefore all diagrams are implemented in a way which allows for them to react to data updates.

Initialization

What happens here? Generally all things which are data independent are done during the initialization. It starts with setting some core variables. A reference to the svg tag which will be used as the container is made. It is followed by a margin definition, where the margin of our diagram content in relation to the container size is defined. The resulting values for `ourHeight` and `ourWidth`, which we will use as space to draw the diagram, are saved.

Following there are a few group elements which are added to the svg tag. These group elements provide a general hierarchy and spacing for different aspects of the diagram. For example the diagram and the legend are usually split in two different groups on the highest level. The diagram can further be separated into groups for the axes and the actual content of the diagram. This general hierarchy is only created on a level which is independent of the provided data and is therefore sufficient to be only created at the beginning.

Finally if there data independent scales, they are defined here. A common example here is a color scale for discrete values. It is not important to already know the specific input values, to be able to create a list of colors which is used by the scale. When queried, it will then return a new color from the list, for each new query value. It is important to note that when the color list runs out of new colors, it restarts at the beginning of the list.

Scales

What are scales? Which kinds of scales exist? How do they work? The first part of the render function is setting up all the required scales.

Scales convert from a domain to a range. The domain is filled with the inputs values. The range is consist of the available outputs. Depending on the type of scale and which data types are used, the domain and range can look quite different. In this thesis the scales are mostly used to find the appropriate coordinates to be able to accurately draw svg elements in the available space. Either by converting the number of refugees to coordinate

space, or to find the appropriate spacing and coordinate position of different categories. Specific scales are described in the implementation of the specific diagrams.

As scales depend on the input data they need to be recreated with each render call, to make sure that they are still up to date. Furthermore they are required while describing the behavior in the general update pattern, and therefore need to be defined in the beginning of the render call.

Data joins

How are they created?

After the scales are defined, the data join is created. When there are several parts to the diagram, like the content and a legend, additional data joins might be required.

A data join is created when binding data to a selection. This is achieved by first calling the `.data(DATA)` function of a selection. The data function creates pairs of elements and data entries. By default, these are matched through their index in the selection and data arrays. This can lead to unexpected behavior when entries are removed or inserted at the not last position. Therefore the default identifier function can be overwritten by passing a custom identifier function as the second optional parameter to the data function. A custom identifier function should return a value and is called for each element in the data array. For the refugees per country data-set in this thesis, the identifier is usually `d => {return d.country}`.

When we initially create the data join, or when data-points are added, we do not have sufficient elements in the selection to pair them with data entries. D3 will therefore create empty placeholders for these elements. To make these placeholders become a part of the DOM, we add the `.join()` after the `data()` call. There are two ways to use the join function. We can either pass a string which will result in adding a matching tag to the DOM. The attributes and style for each new element can then be defined by method chaining. This approach is reasonable for diagrams that do not need to react to data changes. In this thesis we want all diagrams to implement the full extend of the general update pattern, to be able to react to changing data and use the full possibilities of D3.

General Update Pattern

How is this implemented? Where does it come into play?

When the join function is called, instead of passing a single string as parameter, three functions can be passed as parameters. These three functions

correspond to the three cases of the general update pattern and describe their respective behavior. Each of the three function has one input parameter, corresponding to the respective sub-selection. In the enter function we usually want to add some element to the DOM. In the exit selection we remove elements again. The update function is optional, but always used in this thesis, as this is the place to update existing elements to accommodate for data changes and therefore possibly removed or newly added elements as well.

The enter function should add the applicable placeholder element to the svg as actual content. Therefore the first part of the enter function is usually an `.append(string)` call. The string describes the tag which will be added to the DOM. Following this the applicable styles, attributes and sub-elements are added. It is important to provide enough information, that the provided selector which was used to create the selection for the data join whose behavior we are defining, can also match the newly created element. Part of creating the new element is using the scales to position and size the elements accordingly.

The update function is necessary when we want to react to data changes. It is usually similar to the enter function, in that it adjusts the positioning and sizing of the elements according to the possibly changes scales. The exit function is defined by default to simply remove the applicable elements.

All three functions can make use of animations and transitions to improve their feel.

Animations

What are they? How do they work? Why are there two kinds? What makes them tick?

Animations can improve the feel, appeal and readability of diagrams. Especially when reacting to data changes, it is easier to understand and see the changes when for example bars in a bar-chart shift to their new positions, instead of a seemingly entirely new diagram popping up out of thin air. The animations allow the viewer to keep track of the existing entries and visually follow any changes. It is also possible to see the changes of existing values, by following, for example, the growth or shrinking of the length of a bar in a bar-chart. Animations can also be used when initially drawing the diagram, to make it seemingly build itself from nothing, instead of popping into existence.

Animating elements in D3 is achieved by using transitions. Transitions are called from a selection and run on all the elements of the selection. A transition requires a duration and can also be provided with a delay and an easing function. The duration and delay are both in milliseconds. Animating

numerical, color or string values is very easy with transition. It is only required to call the attribute or style with the target value and the transition will take care of the rest. This makes it very fast and easy to animate for example positioning or sizing.

```
1      enter.call(enter => enter.transition(t)
2      .attrTween('d', (d, index, nodes) => {
3          const i = d3.interpolate(0, d.startAngle);
4          const j = d3.interpolate(0, d.endAngle);
5
6          nodes[index].previousStartAngle =
              d.startAngle;
7          nodes[index].previousEndAngle = d.endAngle;
8
9          return time => {
10             d.startAngle = i(time);
11             d.endAngle = j(time);
12             return arc(d);
13         })})
```

Instead of using the default behaviors for numbers, string and colors or when trying to animate other values like svg paths, a tween function can be defined using `attrTween` or `styleTween`. Both tweens need to return a function which will be invoked for each frame of the animation, with a time value between 0 and 1, depending on the frame. The returned function must itself return a value, which is applied to the desired style or attribute every frame. In this thesis tweens are only specifically defined to animate svg path tags.

3.2.2 Bar Chart

How does it work? Which d3 features does it use? how do they work?

3.2.3 Pie chart

How does it work? Which d3 features does it use? how do they work?

3.2.4 Tree map

How does it work? Which d3 features does it use? how do they work?

3.2.5 Sankey

How does it work? Which d3 features does it use? how do they work?

3.2.6 Area graph

How does it work? Which d3 features does it use? how do they work?

3.2.7 Circle graph

How does it work? Which d3 features does it use? how do they work?

3.3 Showcase

How is the showcase structured? How can you get there? Why does it exist? Who might benefit? How can you reuse a part the interesting parts?

There are actually two showcases. One for each data-set. Each showcase is split into two main parts. Firstly there is a section with all the different diagrams for the applicable data-set. All diagrams in one showcase represent the same data-set. This allows for an easy visual comparison, as well as easier comparison of the code. When each diagram would show different data, it would be harder to distinguish between implementation differences which are due to the different representation and differences which are caused by accommodating different data-sets. Additionally each showcase has a section where the input data can be seen and modified.

3.3.1 Integration of each diagram

How is each diagram integrated? How can you access them? Where can you grab them standalone?

Each diagram is implemented to work on its own and without the showcase. Each diagram is also designed to use all the space available in its container. When loading one of the diagrams htmls directly, it will therefore fill the whole browser window. The showcase loads each of the diagrams into a separate IFrame tag with consistent width and height.

3.3.2 Data Updates

How can you simulate data changes? Why is this useful?

The showcases have a section which allows for data manipulation. Rows of data can be modified, added or removed here. The data changes here are not persistent and therefore do not get written in the original data csv files. When data is changed, the diagrams are provided with the updated data and adapt accordingly. It is important to be able to modify the data, as one of the core features of D3 is reacting to changes in data. This manual style of

changing data is probably not so common in real world applications. Yet it is easy to replace these manual data changes to regular API calls or other automatically updating data-sources. As the source of the data changes does not matter for the functionality of D3, the manual approach chosen here is sufficient in demonstrating the possibilities of D3.

4. Discussion

Discuss pros and cons and all things that came up. Pro and cons.

4.1 Evaluation Criteria

How will we evaluate? What?

The two main questions this thesis answers are both subjective to a certain extent. The potential of D3 in data visualization will be determined and evaluated by looking at the many use-case examples which can be found online, as well as the resulting show-case of this thesis. Besides having a look at how D3 can be used, we will also look at where it is possible to use D3. As there are many scenarios in which the JavaScript environment is not required, desired or possible. This leads us into the second question of the advantages and disadvantages of D3. To evaluate this question, we will evaluate the initial learning curve, the effort to create static and dynamic diagrams, the amount and functionalities of available modules, the completeness of documentation as well as pros and cons which come with the web environment.

5. Conclusion

How well did it work? Was it worth the effort? What could be improved?

Bibliography

- [1] J. Mackinlay, “Automating the design of graphical presentations of relational information,” *Acm Transactions On Graphics (Tog)*, vol. 5, no. 2, pp. 110–141, 1986.
- [2] M. Bostock, “Data-driven documents,” accessed:31.03.2022. [Online]. Available: <https://d3js.org/>
- [3] A. v. Kesteren and L. Hunt. [Online]. Available: <https://www.w3.org/TR/selectors-api/>
- [4] M. Bostock, Feb 2012. [Online]. Available: <https://bost.ocks.org/mike/join/>
- [5] UNHCR, “Operational data portal,” accessed:12.05.2022. [Online]. Available: <https://data2.unhcr.org/en/situations/ukraine>
- [6] —, “Refugees per country,” accessed:12.05.2022. [Online]. Available: https://data.unhcr.org/population/get/sublocation?widget_id=312121&sv_id=54&population_01-01
- [7] —, “Refugees per day,” accessed:12.05.2022. [Online]. Available: https://data.unhcr.org/population/get/timeseries?widget_id=312123&sv_id=54&population_01-01

6. Appendix

I guess this should contain all the source code. Maybe there are ways to import it automatically too?