

Final Report

Michael Gottlieb mikemgottlieb@gmail.com

Emily Hindalong ehindalong@gmail.com

Dmitry Tebaykin dmitry.tebaykin@gmail.com

December 16, 2015

Abstract

concise summary of your project. do not include citations.

Introduction

give the big picture. establish the scope of what you did, some background material may be appropriate here.

Neuroscientists have conducted extensive research on the electrophysiological (ephys) properties of different neuron types, but there are barriers to comparing and aggregating results across different studies. This can be attributed to a lack of standard definitions and procedures as well as paywalls maintained by closed-access journals. To alleviate this, Tripathy et al.¹ have developed NeuroElectro - a freely available web-tool that allows users to directly compare data from different neuroscience articles. The primary goal is to “facilitate the discovery of neuron-to-neuron relationships and better understand the role of functional diversity across neuron types.”

NeuroElectro is a Django text mining and curation application [link](#) developed mainly in C-Python with a javascript-based front end and an SQL database on the back end. It currently hosts experimental data from over 500 articles and is expected to grow to host the experimental data of thousands of articles. The data for each article can be accessed by the type of neuron, its electrophysiological properties, or via a table of articles.

However, the current visualizations provided by NeuroElectro are lacking: the user is only able to view data for one neuron type or one ephys property at a time using scatter plots. Our goal was to develop a new interface that supports seamless browsing and analysis of select subsets of the data.

Neuroelectro is a rare breed amongst text-mining projects in the fact that it allows end users to interact with curated data directly. Most text-mining tools in the biomedical domain assume that the end user will want an association matrix for terms in a controlled vocabulary, such as MEDLINE or MeSH terms [26,27]. These tools automatically generate and output an association matrix without providing the user with a way to interface with the original data. This limits the analyses a user can perform.

Taking that flaw into account, we designed our NeuroElectro visualization mainly to provide the user with the overview of the gathered data and with the ability to explore it as they see fit. To accomplish said goals in the most efficient manner we performed an extensive literature search into online biomedical databases visualization efforts and the most efficient general ways to display as well as interact with the data similar to ours (on the scope of hundreds to thousands data rows with a few dozen variables).

Related Work

include both work aimed at similar problems and work that employs similar solutions to yours structure into subsections based on your own synthesis of themes in the related work although there is no requirement to establish research novelty since it's a course project, you should still discuss how the previous work is similar to or different from your own work (either individually or with respect to an entire group) definitely cover academic work; it's often good to cover non-academic

work as well (commercial software, thoughtful blog posts) you may choose to reorder sections to put this one after Abstractions, if it will help you write this section more concisely/clearly.

Solutions to Similar Problems

Neuroelectro provides some crude data visualizations in the form of static scatterplots and a single PCA analysis plot. Most text-mining tools in biomedicine do not use visualization at all, and those that do are restricted to analyses on the derived association matrix. For example, VOSviewer [28] uses colour and spatial position to visualize the semantic clustering and strength of association across text mined terms. The Trading Consequence project [29] focuses on mined trading documents supported by controlled vocabularies to generate maps of commodity trading over time.

Exploring relationships Exploration of relationships between some of the properties in Neuroelectro’s dataset is supported by the current version of Neuroelectro (Brain region vs Electrophysiology only, metadata is only accessible via database). However, it is a static set of strip plots that do not account for all combinations of variables and do not allow any interaction.

Exploration of relationships is a common task in many analytics platforms such as Tableau [17], SAP BOBJ ALOP [32] and Microsoft Excel [33]. Generally, these platforms provide a tabular view of the data in addition to customizable visualizations to enhance users’ exploration of the data. Our goal differs from these platforms as we are not including a tabular view, we are limiting the users’ choices to provide a simpler experience, and we are using plots that are not easily achieved with these platforms (e.g. interactive plots, hive plot).

Providing an overview of data Essentially, we are facing a problem of visualizing a network when we are trying to give an overview of our data. Over the years many solutions have been proposed for this type of task: hairball [30], matrix [31], arc diagram [31], call network [3], hive plot [29] are among the most common. Simply visualizing the network as a collection of nodes connected with edges (the hairball approach) seems impractical due to a large number of nodes and connections (currently: 150 nodes and ~10k edges), scaling is also a problem since the hairball only gets bigger with time. The matrix approach deserves some credit in terms of data visibility and it is a familiar visualization style to biologists, but there are 2 issues with utilizing matrices for this task: 1. Our data is 3 dimensional (neuron type, ephys. property, metadata) and 3D matrices are usually very hard to interpret, we could provide a faceted view of 1 matrix per metadata as a possible solution, but the amount of screen space that would require is enormous. 2. Matrices do not scale well, the labels get too small to be legible at some point. That said, a count matrix could provide a similar overview information to a hive plot as long as the number of parameters is small enough.

In the implementation section we will discuss further our decision to use both a hive plot and a matrix approach.

A call network visualization would end up looking very similar to a hairball in our case, as a result we had to discard this possibility due to scalability issues. Arc diagrams came in as a close second as our visualization of choice - they are easy to interpret, pleasant to look at and they can scale reasonably well with the amount of evidence in the database. The problem with arc diagrams is that all nodes would end up being on one line and that does not represent the 3 distinct groups of nodes (neuron types, ephys. properties, metadata) in our data.

As a result, we decided to use hive plots for providing an overview of our data. Krzywinski, Birol, Jones and Marra [4] describe the advantages of hive plots in terms of gaining quantitative understanding when visualizing networks. They also support: multiple axes, information encoding in the nodes and edges, scaling. The one issue with hive plots is that they are a fairly new visualization style and researchers may have trouble understanding what they are looking at. However, we plan to provide a guide to interpreting the hive plot as well as provide links to the supporting literature. This feature would be on-demand and can be disabled in the application preferences.

Applications of Similar Solutions

Filter panels The filter panel paradigm, where one panel is used to control what data appears in the main panel, is well established in visualization domain 5. An alternative solution is the filter bar, which uses less screen real estate [22]. However, we have opted to stick with filter panels because they will never interfere with the main view and will make it easier for the user to track which filters are applied at any given time. Furthermore, the number of filtering options that we offer will require a larger section of the screen.

There are two basic attribute-based filtering paradigms: drill-down and parallel selection [21]. As the referenced blog post describes, Amazon uses drill-down filtering and Kayak uses parallel filtering. Our solution uses a hybrid of these, allowing the user to drill-down categories and apply parallel selection within. We intend to refer to this blog post when designing the specific details of our filter panel.

Connected scatterplots Since Neuroelectro data is rather diverse (dozens of electrophysiology properties for each of over one hundred neuron types), we plan to utilize scatterplots 8 and connected scatterplots ??? for answering research-oriented questions. Haroz et al. showed the effectiveness of the latter in representing time-series data: even though connected scatterplots are novel to many users, they are excellent at being intuitive to understand and capturing and holding the viewer’s attention.

Linked highlighting There are a number of interaction approaches to linked highlighting in scatter plots [23]. Through our consultation with the stakeholder, linked highlighting on hover was emphasized as a critical element. However, this is not the only means of linked highlighting available. For example, linked brushing, where the user selects a subset of points to be highlighted, is a popular choice for multi-selection [23, 36].

Hive plot Overviews of the data will make use of the work done by Krzywinski et al. 4 and Hanson ??? for our proposed hive plot. Various good examples of hive plot visualizations of networks are shown on the hive plot website [29]. We have developed our visualization based on the features that are most meaningful for our data and the questions we are trying to answer with the view (sparsity of data, node degree, edge weights, outliers and general trends in the data). We also use a matrix and a table approach to supplement the hive plot, since alone the hive plot does not provide enough details about connections between data types.

Colour We have decided to use colour as a channel for linked-highlighting. Previous research suggests that colour is one of the most powerful visualization channels 9 and, when used correctly, it can provide insight into the data so intuitively that the user wouldn’t even need a legend to understand what kind of data the channel encodes for. On the overview panel colour is used as a measure of how many articles exist for each connection between nodes.

Data and Task Abstractions

you should analyze your domain problem according to the framework of the book, translating from domain language into abstract descriptions of both tasks and data typically data will need to come first, since you will need to refer to that data in your task descriptions. it is very likely that you will need to first have domain-specific descriptions, followed by the abstracted versions.it is often helpful to split these sections into two pieces: first have subsection where data is described in domain-specific terms, and then follow up with a next subsection where you’ve abstracted into domain-independent language. you may choose to have a separate Domain Background section before this one if there’s a lot of specialized vocabulary/material to explain in order to have a comprehensible and concise data description. similarly, it can be easier to write the tasks section by first providing a domain-specific list of tasks, and then present the abstracted version. you should decide whether to split by domain vs abstract (first have both data and task domain-specific sections, followed by the abstractions for data and task) or to split by data/task

(first have data domain-specific then data abstract, then have task domain-specific and finally task abstract).

Domain-specific data and tasks

NeuroElectro is a database of neuroscience articles and it applies text-mining as well as curation approaches to extract electrophysiology measurements, neuron type information and experimental setup conditions (metadata) from the html-encoded articles. At this point, text-mining alone is not reliable enough since neuroscientists authoring the articles in questions were not writing them using guidelines. As a result, each article is a snowflake of sorts - even with very well written algorithms automated text-mining is not at human text interpretation level yet. Hence the need for training undergraduate curators to verify the text-mining results and correct its errors. We focus on visualizing only the curated data, meaning that we have ~1000 articles worth of data. Note that not all articles contain all data types that NeuroElectro is able to store.

Domain-specific data types

1. Electrophysiology measurements

- These are intrinsic neuron properties: membrane potential at rest, spike threshold (minimum membrane potential that causes a spike), input resistance, rheobase (minimum amount of current one needs to inject to cause an spike), etc.
- Neurons communicate with the help of action potentials (voltage spikes) which are caused by cell's membrane voltage rising above the spike threshold causing a cascade of Sodium ions to flood into the neuron, propagating the action potential signal down its axon and to other neurons. The cell then closes Sodium channels and opens Potassium channels in order to return to its original state.
- Neuron signalling is an electrochemical process and electrophysiology aims to record all meaningful characteristics that describe this process.

2. Neuron types

- It is no secret that brain contains many different kinds of neurons. Neuroscientists have not decided on exactly how many neuron types there are and the debate has been ongoing for over a hunder years. Nevertheless, there are resources on the Internet that attempt to offer a classification for neuron types. NeuroElectro utilizes enhanced NeuroLex neuron classifications, eventually NeuroElectro may be offering its own neuron type hierarchy as we gather more data. For the purposes of our visualization, we distinguish 2 levels in the neuron type classification hierarchy - all neurons are assigned a brain region and a neuron type within that region. Each neuron type is assigned exactly one brain region (Neuron types that are present in many places or if their location is unknown comprise the "Other Region" brain region). These assignments were performed by a Dr. S.J. Tripathy - a neuroscience postdoc and the original developer of NeuroElectro.

3. Experimental conditions (metadata)

- This data type stores information about the electrophysiological experiment itself, such as: species, strain, age and weight of the animal used, electrode type with which the measurements were taken, chemical solutions used to keep the brain slice moist and semi-alive, recording temperature, etc.
- This data is important in order to compare ephys measurements from different experiments and labs.

NeuroElectro also stores data about each article: title, publication year, authors, etc.

Domain-specific tasks

1. Explore relationships between neuron types, ephys properties, and experimental conditions.

- Find the neuron type, ephys measurement and metadata of interest.
 - View specific electrophysiology measurement for a specific neuron type (e.g. view rheobase values for Hippocampal CA1 pyramidal cells).
 - Compare ephys measurements across neuron types (e.g. resting membrane potential across all or a selected set of neuron types).
 - Explore the effect of metadata on an ephys measurement (e.g. action potential amplitude change with animal age).
2. Identify regions of the brain where particular neuron types are found.
 3. Identify how many data points exist for different combinations of neuron types, ephys properties, and experimental conditions.
 4. Find out how many articles support a specific analysis.
 5. Summarize the data in the current analysis scope.
 6. Lookup details for individual evidence lines extracted from articles.

Abstracted data and tasks

Abstracted data Ignoring all domain-specific complications, we are dealing with 1 csv spreadsheet that contains ~1000 lines of ~150 variables. We can split these variables up into 4 types: quantitative measurement data, qualitative location data, mixed type explanatory variables and qualitative information about each line. The line information can be used for tooltips or tables to provide more context, but in itself, it is not interesting for the analysis.

Abstracted tasks

1. Explore relationships between three different data types.
 - Browse data available for analysis (data overview).
 - Identify distributions of quantitative attribute values for a particular categorical attribute values.
 - Compare distributions of quantitative attribute values for different categorical attribute values.
 - Identify correlations between quantitative attributes.
2. Navigate data type hierarchy to identify data of interest.
3. Identify how many data points exist for different combinations of data types.
4. Identify how many data points support a specific analysis.
5. Summarize the data in the current analysis scope.
6. Lookup details for individual data points (items).

Solution

describe your solution idiom, analyze it according to the framework of the book, and justify your design choices with respect to alternative possibilities if you have done any significant algorithmic work, discuss the algorithm and data structures. You might choose to split out Interface into its own section

Implementation

medium-level implementation description. you must include specifics of what you did yourself versus what other components/libraries/toolkits you built upon. this section is one major divergence from standard research paper format, you need to provide much more detail than would normally be appropriate in a research context.

R-Shiny app (general)

Our solution was built using the Shiny web application framework[10] for the R language[11]. Shiny server and ui components handle all transactions between the front and back end of our application. Each major component took advantage of a number of existing libraries, which is explained in more detail in the following subsections.

The first step of our applications performs data loading, cleaning and wrangling. Our app takes a csv dump of Neuroelectro’s database as input. We load, clean and manipulate data using base R and dplyr[12] functions. At this point we also generate and cache or load the cached version of a modified dataset to speed up matrix view generation as the filters are changed. Once the data is loaded and prepped, it is passed to the ui and server components that house the core functionality of our app.

Navigation/filtering panel

The filtering panel uses collapse panels from the shinyBS package[13]. The Neuron Type and Organism trees use the shinyTree package[14]. The shinyTree source code was modified to improve appearance and introduce text wrap to long labels that encroached on the space of other components. The shinyjs and V8 packages[15 & 16] were used to add JavaScript commands on startup to modify the shinyTree component’s unruly behaviour. The filter options for continuous features use slider bars from the base Shiny framework. As with the shinyTree component, they were not perfectly suited for our needs. We modified the sliders to use log scales via JavaScript commands called via shinyjs and V8. We implemented how the filter states were applied to the data set and observers to update the plots only when the selected data had been changed. Default behaviour resulted in all plots being redrawn whenever a filter element was touched, even if its value was not changed (e.g. expanding a node on a filter tree or moving a slider without deselecting any data points).

Explore panel

The four plots of the explore panel share data that is filtered based on the state of the filter panel. The data displayed on each plot is determined by the two axis selectors above each plot. The axis selectors are standard Shiny ui components that did not require modification. The plots in the explore panel are generated dynamically depending on the type of data that the user selects to view. We used the ggvis package[17] to generate the plots in the explore panel as they promised easy interactivity. Our dataset changed based on filtering rules and axes selected and was passed to a function we implemented to determine what type of plot to show and to reduce repeated code. This lead to problems, as interactive ggvis features, such as hover and brush handlers, do not work with well or at all with dynamic datasets and inside functions. While we consider other plotting options, we have implemented an on click handler that highlights points in all plots in which they appear. We also implemented action buttons to clear highlighting, remove highlighted points from all plots and restore removed points.

Overview panel

The overview panel has three distinct components: The hiveplot, the heat map and the table. The hiveplot uses the HiveR package[18] as well as the RColorBrewer package[19] for color selection. It does some other stuff based on whatever Dmitri did.

The heatmap uses a derived dataset that contains whether or not each datum contains information regarding features of particular interest to the stakeholder. It computes an association matrix on demand which is passed to the pheatmap package[20] to display the associations and annotations. Like the hive plot, it uses the RColorBrewer package for its color palette.

The table is generated based on filter rules and is made entirely using the DT package[21] with some non-default parameters. As its title implies, it is simply a wrapper for the DataTables JavaScript library.

Results

should include scenarios of use and multiple screenshots of your software in action. walk the reader through how your interface succeeds (or acknowledge how it falls short) in solving the intended problem. if you did any evaluation (deployment to target users, computational benchmarks), do report on that here.

Discussion and Future Work

Strengths, weaknesses, limitations (reflect on your approach) Lessons learned (what do you know now that you didn't when you started?) Future work (what would you do if you had more time?)

Conclusions

summarize what you've done in a way that's different from the abstract because you can count on the reader having now seen all of the content of the paper in between

Bibliography

make sure to use real references for any work that's been published academically, not just URLs do pay particular attention to my instructions for checking reference consistency

References

- [1]S. J. Tripathy, J. Savitskaya, S. D. Burton, N. N. Urban, and R. C. Gerkin, "NeuroElectro: A window to the world's neuron electrophysiology data," *Frontiers in neuroinformatics*, vol. 8, 2014.
- [2]N. Henry and J.-D. Fekete, "MatrixExplorer: A dual-representation system to explore social networks," *Visualization and Computer Graphics, IEEE Transactions on*, vol. 12, no. 5, pp. 677–684, 2006.
- [3]K.-K. Yan, G. Fang, N. Bhardwaj, R. P. Alexander, and M. Gerstein, "Comparing genomes to computer operating systems in terms of the topology and evolution of their regulatory control networks," *Proceedings of the National Academy of Sciences*, vol. 107, no. 20, pp. 9186–9191, 2010.
- [4]M. Krzywinski, I. Birol, S. J. Jones, and M. A. Marra, "Hive plots—rational approach to visualizing networks," *Briefings in bioinformatics*, vol. 13, no. 5, pp. 627–644, 2012.
- [5]R. Bearavolu, K. Lakkaraju, W. Yurcik, and H. Raje, "A visualization tool for situational awareness of tactical and strategic security events on large and complex computer networks," in *Military communications conference, 2003. mILCOM'03. 2003 IEEE*, 2003, vol. 2, pp. 850–855.
- [6]M. Dumas, J.-M. Robert, and M. J. McGuffin, "Alertwheel: Radial bipartite graph visualization applied to intrusion detection system alerts," *Network, IEEE*, vol. 26, no. 6, pp. 12–18, 2012.
- [7]A. Sopan, A. S.-I. Noh, S. Karol, P. Rosenfeld, G. Lee, and B. Shneiderman, "Community health map: A geospatial and multivariate data visualization tool for public health datasets," *Government Information Quarterly*, vol. 29, no. 2, pp. 223–234, 2012.
- [8]M. Friendly and D. Denis, "The early origins and development of the scatterplot," *Journal of the History of the Behavioral Sciences*, vol. 41, no. 2, pp. 103–130, 2005.
- [9]T. Munzner, *Visualization analysis and design*. CRC Press, 2014.
- [10]W. Chang, J. Cheng, J. Allaire, Y. Xie, and J. McPherson, *Shiny: Web application framework for r*. 2015.

- [11]R Core Team, *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing, 2015.
- [12]H. Wickham and R. Francois, *Dplyr: A grammar of data manipulation*. 2015.
- [13]E. Bailey, *ShinyBS: Twitter bootstrap components for shiny*. 2015.
- [14]Trestle Technology, LLC, *ShinyTree: JsTree bindings for shiny*. 2015.
- [15]D. Attali, *Shinyjs: Perform common javaScript operations in shiny apps using plain r code*. 2015.
- [16]J. Ooms, *V8: Embedded javaScript engine*. 2015.
- [17]W. Chang and H. Wickham, *Ggvis: Interactive grammar of graphics*. 2015.
- [18]B. A. Hanson, *HiveR: 2D and 3D hive plots for r*. 2015.
- [19]E. Neuwirth, *RColorBrewer: ColorBrewer palettes*. 2014.
- [20]R. Kolde, *Pheatmap: Pretty heatmaps*. 2015.
- [21]Y. Xie, *DT: A wrapper of the javaScript library 'dataTables'*. 2015.