

VizRec 2.0: Towards a holistic workflow for visual data exploration

Doris Jung-Lin Lee, Aditya Parameswaran
{jlee782,adityagp}@illinois.edu
University of Illinois, Urbana-Champaign

Abstract

test abstract

1 Introduction

2 Precise Visual Querying

Visual analysis often reveal important anomalies or trends in the data[6]. However, it is often challenging to find the appropriate piece of information to realize these insights.

2.1 Motivating Example

Astronomers from the The Dark Energy Survey (DES)[2] are interested in finding anomalous time series to discover astrophysical transients (objects whose brightness changes dramatically as a function of time), such as supernova explosions or quasars. When trying to find celestial objects corresponding to supernovae, which have a specific pattern of brightness over time, scientists need to individually inspect the visualizations of each object until they find ones that match the pattern. With more than 400 million objects in their catalog, each having their own set of time series brightness measurement, the process of manually exploring a large number of visualizations is not only error-prone, but also overwhelming for scientists who do not have extensive knowledge about their dataset.

The astronomy use case highlights a common challenge in exploratory data analysis (EDA). There is often a large space of possible visualizations that could be generated from a given dataset and manual search through this large collection is inefficient. Visualization authoring tools such as Tableau and Excel focusses on presenting one visualization at a time. There is no systematic way to create, compare, and query large collections of visualizations.

2.2 Effortless Data Exploration with *zenvisage*

The challenges presented earlier points to a need for tools that enables users to create and search through large collections of visualizations. Therefore, we developed *zenvisage* a visual query system that allowed users to

Copyright 0000 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

Bulletin of the IEEE Computer Society Technical Committee on Data Engineering

search through large collections of visualizations. *zenvisage* is built on top of a querying language called ZQL, which provides a mechanism for managing collections of visualizations[12]. Contrary prior work on visualization languages for specifying visual encodings of individual visualizations [13, 16], ZQL supports high-level queries over visualization collections, such as composing, sorting, filtering a collection of visualization. ZQL functionals and primitives can be constructed into rich and expressive query semantics, with functionalities including :

- Finding top-k visualizations whose y values are most or least similar from a queried visualization (e.g. Find other cities with sold price over time similar to Manhattan. Varying along CITY while keeping $X=TIME, Y=AVG(PRICE)$ fixed.)
- Comparing across a collection of visualizations by iterating over one or more x, y, z attributes while fixing other attributes (e.g. Find a y attribute that varies with time similarly how average price changes over time)
- Finding a pair of X and Y axes where two specific visualization instances differ the most. (e.g. For what pair of attributes does the products ‘stapler’ and ‘chair’ differ the most?)

Given a ZQL query, *zenvisage* parses the query into a graph of visual component nodes(containing visualization information, such as X, Y columns) and task nodes (common and user-defined primitives for processing visual components, such as sort-filter). *zenvisage* then performs query optimization to merge together multiple nodes, as well as reducing the processing time required for individual visualization components. Using the optimized query plan, the executor compiles visual nodes into SQL queries for retrieving the visualization data and postprocesses the result via the defined operations.

While ZQL provides powerful mechanism for expressively specifying queries on large collections visualizations, writing ZQL queries can be daunting for novice users. Therefore, we extracted a typical workflow of visualization querying (finding top-k most similar visualization from a collection with fixed X,Y while varying Z) to allow users to formulate ZQL queries through interactions. The user can either directly input ZQL queries through a frontend table input or their frontend interactions is mapped into ZQL queries. The query results are rendered as a ranked list of visualizations in the Results panel in the frontend. *zenvisage* is a full-fledged visual querying system that supports a variety of querying interactions as illustrated in Figure 2.2. In the following section, we will discuss the design process of how we developed this visual query system and the lessons that we have learned for designing future visual data exploration systems.

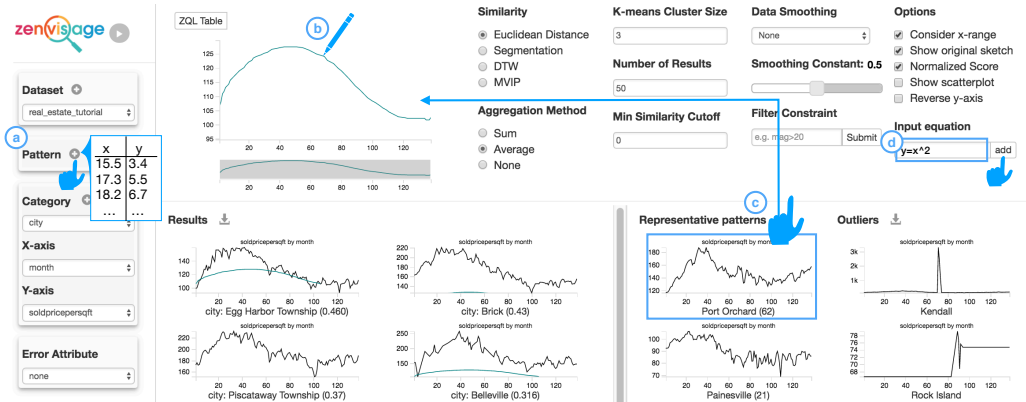


Figure 1: *zenvisage* offers a variety of querying modalities, including: a) uploading a sample pattern from an external dataset as a query, b) sketching a query pattern, c) dragging-and-dropping an existing pattern from the dataset, and d) inputting an equation as a query.

3 Hypothesis Formation and the cycle of visual analysis

3.1 Visual Querying in the Sensemaking Framework

In developing *zenvisage*, we collaborated with scientists from astronomy, genetics, and material science in a year-long participatory design process [5]. In particular, we study how various features impacts analyst’s ability to rapidly generate new hypothesis and insights and perform visual querying and analysis. In addition, we were interested in how a visual query system (VQS) like *zenvisage* fits into the participant’s analysis workflow. Our findings offered design guidelines for improving the usability and adoption of next-generation VQSs. More importantly, in this section, we focus on applying our findings on visual querying to the context of supporting the full cycle of visual data exploration.

Supporting Complex Vague Expressive Querying

zenvisage is an example of a *precise visual querying system (PVQS)*, which accepts precise queries as an input, expressed through interactions or directly specified via ZQL. When users are querying with a PVQS, they often need to translate their ambiguous, high-level questions into an plan that consists of several interactions in conjunction to address their desired query incrementally.

We found that participants in our *zenvisage* design study often creates unexpected workflows that chain together multiple analysis steps, including interactions, controls, and queries in order to address a higher-level research question. For example, the geneticists in our study repeatedly explored representative trends to gain an overall sense of what typical profiles exist in the dataset and queried mainly through drag-and-drop of these representative trends. Variations to their main workflow also include changing cluster sizes and display settings to offer them different perspectives on the dataset and filtering on data attributes.

The expressiveness of PVQSs comes from the multiplicative effect of stringing together combinations of interaction sequences into a customized workflow. Designing features that diversifies potential variations and helps the creation of multiple analysis workflows expands the space of actions that could be performed during the analysis.

However, even with many supported interaction, there were still complex multi-step queries that could not be expressed through this framework. One solution is to divert to ZQL. There is an inevitable design trade-off between the query expressivity and interface usability[6, 4]. Balance between language/querying complexity versus expressiveness. The extensibility of these systems or querying language also comes with the cost of potentially overloading the users with too many potential options to chose from. In Section 4, we survey related work in this area and advocate for —vague querying. More importantly pointed towards a need for vague querying. Give some examples of vague querying.

Top-down and Bottom-up Querying Modalities We employed Pirolli and Card’s [8] information foraging framework for domain-experts to contextualize our study results. Pirolli and Card’s notional model distinguishes between information processing tasks that are *top-down* (from theory to data) and *bottom-up* (from data to theory). In the context of visual querying, users employ top-down approaches by starting with a hypothesis on what patterns to look for and express it through sketching or inputting an equation (Figure 2.2b,d). On the other hand, bottom-up approaches originate from the data (or equivalently, the visualization). For example, the user may drag and drop a visualization of interest in the dataset as the input query or upload a visualization from an external dataset (Figure 2.2a,c).

Our interactions with the scientists showed that *bottom-up querying via drag-and-drop was more intuitive and more commonly used than top-down querying methods when the users have no desired patterns in mind*, which is commonly the case for exploratory data analysis. One of the main reason why participants did not find sketching useful was that they often do not start their analysis with a pattern in mind. Later, their intuition about what to query is derived from other visualizations that they see in the VQS, in which case it made more sense to query using those visualizations as examples directly. Similarly, while functional fitting is a common operation in scientific data analysis, querying by equation is also unpopular, since it is challenging to formulate functional

forms in an prescriptive, ad-hoc manner without seeing what the common patterns in the dataset are.

While the usage of each querying feature may vary from one participant to the next, a key design principle that came from this finding was the need for visual query systems to provide visualization recommendations that can help analysts jumpstart their exploration. We found that many users made use of the representative trends and outliers visualizations provided by *zenvisage* as contextual information to better understand their data (e.g. after a filter is applied) or to query based on these recommended visualizations (e.g. find visualizations that are similar to the one in the largest representative clusters).

Recommendation facillitate smoother flow of analysis, ensures that user is never stuck or out of ideas. it does this by going towards better data understanding, accurate understanding of the context of analysis and scope of data. should not only close the loop between the two modalities of querying and exploration, but also contribute towards — data understanding. In Section 5, we advocate the importance of building recommenders that contributes towards data understanding but also —.

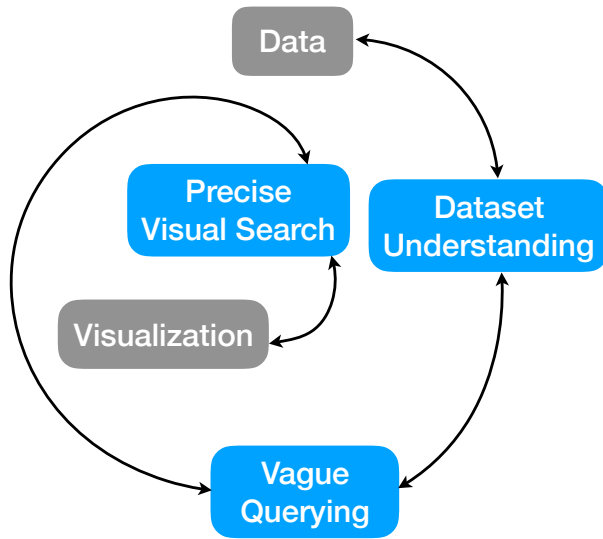


Figure 2: Cycle of visual data exploration.

3.2 Challenges Ahead

The goal here is to help novice submit precise queries without SQL background, easy to use interface. Our study found that VQS does more than just this, but still not enough.

- Precise Search Fail to understand intricacies of user need/intent, need more expressivity/flexibility for querying.
- No perfect training workload, real-world data + task is noisy and complex.
- towards more holistic model for insight discovery

4 Vague Intelligent Search

Accounting for user interaction, mental models. More global objective taking into account user with the goal of dataset understanding rather than task completion.

4.1 Challenges

- Inferring user intent in querying and context is important (both in terms of user input and what is recommended)
- tools can not assume user has querying intention. exploration without intention, user don't know what they are searching for – Recommendation.
- The important thing here is identifying what should be done by the system v.s. requested from user. Inappropriate choice of these will result in lack of expressibility and user feeling lack of control of analysis, limiting exploration.
- Need for a unified framework of inference to take all of these into account (e.g. natural language, etc)

4.2 Related Works

Most systems design exhibits a trade-off between how expressive can the query be and how usable the interface is. For example, while querying language (such as SQL) are highly expressive, formulating SQL queries that maps user's high-level intentions to specific query statements is challenging. Therefore, query builders have been developed to address this issue. Form-based query builders often consist of highly-usable interfaces that ask users for a specific set of information mapped onto a pre-defined query. However, form-based query builders are often based on templated queries with limited expressiveness in their linguistic and conceptual coverage, which makes it difficult for expert users to express complex queries. The extensibility of these systems or querying language also comes with the high engineering cost, as well as potentially overloading the users with too many potential options to choose from.

Given that there is no one size fit all interface for query specification for users of different expertise levels and workload, future visual data explorations systems needs to take into account a wide spectrum of queries of different input types and varying degrees of specificity that could be potentially generated from different interfaces. There is a need for vague or ambiguous specification.

As illustrated in Fig.??, these can range from cold-start (no supervision) to input examples, input relations to complete specification.

natural language

4.3 Towards 3'I's of rapid hypothesis generation support

Given our observations from the participatory design study, we distill several desiderata for the next generation VQSs. Towards 3'I's Interactive, Iterative, Informative (Give examples from the ZV-TVCG paper)

Integrated: should always be aware of the context of data and user

Interactive flow: (how natural is it to move between analysis steps, facilitate fluid analysis and not get “stuck”) : interactivity, feedback (latter is quite unexplored), and recommendation, expressivity (how easy is it to express what to do via interactions) and diversity of actions that could be performed.

Iterative: query refinement, dialogue (not a one-shot query) Joining the flow: Section 4 focuses on the first two items .

Informative: not just task-based interestingness but more explanation-based (causality, introduce distribution awareness notion in viz-sum), focussed on data understanding, which we will discuss in Section 5

5 Towards Dataset Understanding

One of the key goals of visual data exploration is to promote a better understanding of the dataset that enables users to make actionable decisions. While our focus in the previous sections have focussed on intention-driven

queries, where users have some knowledge of what types of questions he may be interested in. This section discusses general query-free recommendations and continual provenance that helps users become more aware of the dataset with respect where they are in their analysis workflow.

Situations where there is an absence of explicit signals from the user can happen in two scenarios: 1) user is at the beginning of their analysis (commonly known as the ‘cold-start’ problem) and 2) user doesn’t know what to query for, which is the situation derived from our *zenvisage* finding in Section 3. In this section, we will describe STORYBOARD, a system that provides data summaries and guides users through informative subsets of data, as an example of a system that promotes —. Then, we will discuss two other types of data understanding during dynamic visual data exploration to highlight the challenges and opportunities ahead in this space.

5.1 STORYBOARD: Promoting Distribution Awareness of Data Subsets with Summary of Visualizations

Common analytics tasks, such as causal inference, feature selection, and outlier detection requires studying the distributions or patterns at different levels of data granularity [1, 17, 3]. However, it is often hard to know *what* subset of data contains an insightful distribution to examine. In order to explore different data subsets, a user would first have to construct a large number of visualizations corresponding to all possible data subsets, and then navigate through this large space of visualizations to draw meaningful insights. The lack of a systematic way to perform these exercises makes the process of manually exploring distributions from all possible data subsets tedious and inefficient [9, 10].

To this end, we present STORYBOARD, an interactive visualization summarization system that automatically selects a set of visualizations to summarize the distributions within a dataset in an informative manner. Figure ?? illustrates an example dashboard generated by STORYBOARD from the Police Stop Dataset [7], which contains records of police stops that resulted in a warning, ticket, or an arrest. The attributes in the dataset include driver gender, age, race, and the stop time of day, whether a search was conducted, and whether contraband was found. We requested STORYBOARD to generate a dashboard of 9 bar chart visualizations with x-axis as the stop outcome (whether the police stop resulted in a ticket, warning, or arrest/summons) and y-axis as the percentage of police stops that led to this outcome. First, at the top of our dashboard, STORYBOARD highlights three key data subsets that results in a high arrest rate, which looks very different trend than the overall (where the majority of stops results in tickets). Following along the leftmost branch, we learn that even though in general when a search is conducted, the arrest rate is almost as high as ticketing rate, when we look at the asian population, whether a search is conducted had less influence on the arrest rate and the trend resembles more like the overall distribution.

While such summary dashboards are useful for making sense of relationships between data subsets, finding effective visualizations to summarize a dataset is not as trivial as picking individual visualizations that maximizes some statistical measure, such as deviation [14], coverage [11], or significance testing [1], which can often result in misleading summarizations. The key insights of our work is

For example, our insight regarding

The above example demonstrates a scenario where the selection of an improper reference (female) for comparing the visualization (black female) against results in misleading insights. In STORYBOARD, we formulate an objective where a visualization is *actually* interesting when it deviates from and can not be explained by *even* its most informative reference.

The effectiveness of STORYBOARD largely comes from how it helps analysts become more *distributionally aware* of the dataset. We define *distribution awareness* as the aspect of data understanding in which analysts make sense of the key distributions across different data subsets and their relationship in the context of the dataset. So that even though it may be infeasible to examine all possible data subsets, with distribution awareness, the analyst will still be able to draw meaningful insights and establish correlations about related visualizations by generalizing their understanding to make predictions regarding the unseen visualizations. Our user

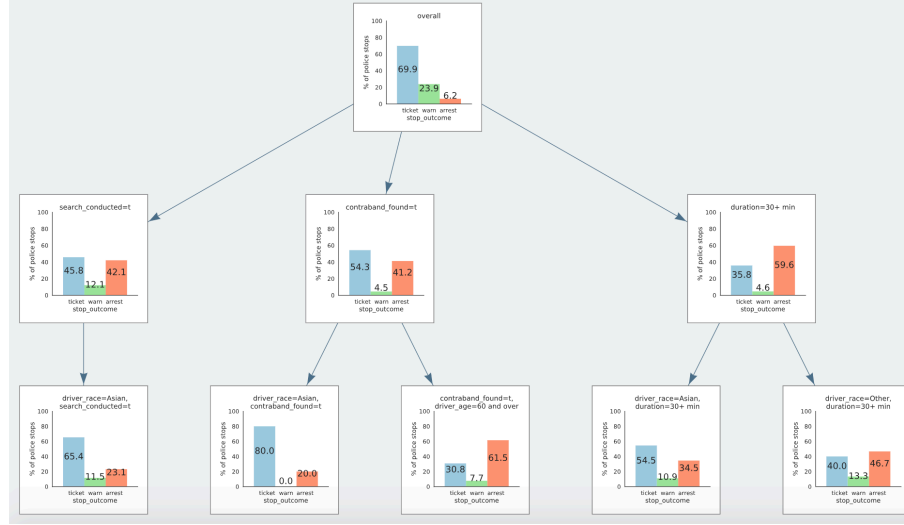


Figure 3: Example dashboard generated by STORYBOARD summarizing the key insights in the Police dataset.

study evaluations show that facilitating distribution awareness through STORYBOARD guides users to make better predictions regarding unseen visualizations, ranking attribute importance, and retrieval of interesting visualizations compared to the baselines.

5.2 Contextual and Situational Awareness: Challenges and Opportunities

The notion of distribution awareness is useful when we are looking at user understanding at one static point in time of the analysis (e.g. during cold start). As a — future research, in this section, we introduce two complementary notions of data understanding (contextual and situational awareness) when considering dynamic visual exploration in the context of an analytic workflow.

Contextual awareness is the aspect of data understanding that relates to understanding what is current —, including making sense of the data attributes and schema, keeping track of which filters or transformations have been applied to the displayed data.

serves to — in a dynamic exploration situation, keep track of what filter is in play, what dataset/ schema am I looking at, which operations have been applied to the data that I’m looking at? Related to provenance for the current time. Within a dataset, structure and provenance is essential to help users navigate and provide users with sense of coverage and completion. This is an important but underexplored area. (viz-sum, Sarvghad et al 2017)

Mechanisms that facilitate distribution awareness for users can effectively couple with contextual awareness in a dynamic exploration situations to help update the user’s mental model on the current data context. For example, the representative and outlier patterns in *zenvisage* provides summaries of data in context. When a dataset is filtered, the representative trends are updated accordingly. By understanding both the context (i.e. I’m only looking at data filtered with ...), the users becomes distributionally aware of the data in the particular context. an overview of typical trends for the data to be queried.

Situational awareness: related to provenance, as a function of time. - provenance of schema and attribute understanding (coverage, etc) . Similar to situational awareness (cite Tory)

The key difference between contextual and situational awareness is that contextual awareness is focussed on how the current context of the data came about, whereas situational awareness accounts for all prior user actions to — intent. For example, a user may be interested in how a measure variable changes when groupby upon different dimensions variables. The analyst may examine and clear various x-axis in sequence: contextual

awareness informs the user that the context is the last attribute that has not yet been cleared, whereas situational awareness provides a notion of coverage informing the analyst that they have already explored this attribute.

Note that while the discussion above have been focussed on how to design systems that can help facilitate these aspects of user's awareness in dataset understanding, these ideas can be generalized to principles in deising the types of intelligent querying systems discussed in Section 4. An intelligent visual exploration system needs to be distributionally, contextually and situationally aware, by make use of information about the data (distribution awareness), the analytic context, and situation jointly in making timely recommendations. For example, contextual awareness can inform the system that the user's current — x,y, main visualization, while a distributionally aware system may recommend a highly-skewed data subset as interesting, a sitational aware system may realize a variable have been explored extensively in the past and recommends it accordingly. In other words, these intelligent visual query system not only needs to facillatate these aspects of data understanding, but also need to make use of this information to make inference and recommendations in an interpretable manner that can guide analysts towards meaningful stories and insights for further investigation.

inference and descisions intepretable.

, rather than the system's awareness of the user's context, situation ,etc. Ideally, an intelligent system should related works have focussed on making specification easier, but not really trying to understnad user intent or what might the user want to see.

6 Concluding Remarks

Data is agnostic to the user, intention —, by building tools—, Section 2 to 4 have focussed on extracting what user want from data. bridging together what user want from data, what data has to offer, supporting interactive discourse between the two. Either using one-size-fits-all statistics, templates, heuristics as a solution or problem only applicable to a subset of analytic tasks[14, 15].

References

- [1] Anushka Anand and Justin Talbot. Automatic Selection of Partitioning Variables for Small Multiple Displays. 2626(c), 2015.
- [2] Drlica Wagner et al. Dark Energy Survey Year 1 Results: Photometric Data Set for Cosmology. 2017.
- [3] Jeffrey Heer and Ben Shneiderman. Interactive Dynamics for Visual Analysis. *ACM Queue*, 10(2):30, 2012.
- [4] H. V. Jagadish, Adriane Chapman, Aaron Elkiss, Magesh Jayapandian, Yunyao Li, Arnab Nandi, and Cong Yu. Making database systems usable. In *Proceedings of the 2007 ACM SIGMOD International Conference on Management of Data*, SIGMOD '07, pages 13–24, New York, NY, USA, 2007. ACM.
- [5] Doris Jung-Lin Lee, John Lee, Tarique Siddiqui, Jaewoo Kim, Karrie Karahalios, and Aditya G. Parameswaran. Accelerating scientific data exploration via visual query systems. *CoRR*, abs/1710.00763, 2017.
- [6] Kristi Morton, Magdalena Balazinska, Dan Grossman, and Jock Mackinlay. Support the Data Enthusiast: Challenges for Next-Generation Data-Analysis Systems. *Proceedings of the VLDB Endowment*, Volume 7, pp. 453–456, 2014, 7:453–456, 2014.
- [7] E. Pierson, C. Simoiu, J. Overgoor, S. Corbett-Davies, V. Ramachandran, C. Phillips, and S. Goel. A large-scale analysis of racial disparities in police stops across the united states, 2017.

- [8] Peter Pirolli and Stuart Card. The Sensemaking Process and Leverage Points for Analyst Technology as Identified Through Cognitive Task Analysis. A Notional Model of Analyst Sensemaking.
- [9] S Sarawagi, R Agrawal, N Megiddo, Generalitat Valenciana Ajuntament Valencia Univ Politecn Valencia, and E T H Zurich Oracle Sybase Softlab Iberia Edbt Fdn. Discovery-driven exploration of OLAP data cubes. *6th International Conference on Extending Database Technology (EDBT 98)*, pages 168–182, 1998.
- [10] S. Sarawagi. User-adaptive exploration of multidimensional data. *Proc of the 26th Intl Conference on Very Large*, pages 307–316, 2000.
- [11] Ali Sarvghad, Melanie Tory, and Narges Mahyar. Visualizing Dimension Coverage to Support Exploratory Analysis. *IEEE Transactions on Visualization and Computer Graphics*, 23(1):21–30, 2017.
- [12] Tarique Siddiqui, Albert Kim, John Lee, Karrie Karahalios, and Aditya Parameswaran. Effortless Data Exploration with zenvisage: An Expressive and Interactive Visual Analytics System.
- [13] C. Stolte, D. Tang, and P. Hanrahan. Polaris: a system for query, analysis, and visualization of multidimensional relational databases. *IEEE Transactions on Visualization and Computer Graphics*, 8(1):1–14, 2002.
- [14] Manasi Vartak, Samuel Madden, and Aditya N Parmeswaran. SEEDB : Supporting Visual Analytics with Data-Driven Recommendations. 2015.
- [15] Manasi Vartak, Silu Huang, Tarique Siddiqui, Samuel Madden, and Aditya Parameswaran. Towards Visualization Recommendation Systems. *ACM SIGMOD Record*, 45(4):34–39, 2017.
- [16] Leland Wilkinson. *The Grammar of Graphics (Statistics and Computing)*. Springer-Verlag, Berlin, Heidelberg, 2005.
- [17] Eugene Wu and Samuel Madden. Scorpion: Explaining Away Outliers in Aggregate Queries. *Proceedings of the VLDB Endowment*, 6(8):553–564, 2013.