

Rapid Prototyping: Accelerating Development of Your Organization's Reports Using SAS® Visual Analytics

Elliot Inman and Michael Drutar, SAS Institute Inc., Cary NC

ABSTRACT

One of the most important factors driving the success of requirements-gathering can be easily overlooked. Your user community needs to have a clear understanding of what is possible: from different ways to represent a hierarchy to how visualizations can drive an analysis. Discussions about desktop access versus mobile deployment and/or which users might need more advanced statistical reporting can lead to a serious case of option overload. One of the best cures for option overload is to provide your user community with access to template reports they can explore themselves. Instead of beginning requirements-gathering with a blank slate, your users can begin the conversation with, "I would like something like Template #4," greatly reducing the time and effort required to meet their needs.

In this paper, we describe how you can take a single rich data set and build a suite of template reports that demonstrate the full functionality of SAS® Visual Analytics. We present a suite of the most common, most useful SAS Visual Analytics report structures, from high-level dashboards to statistically deep dynamic visualizations. We show exactly how to build a dozen template reports from a single data source, simultaneously representing options for color schemes and other choices to consider. Although this template suite approach can apply to any industry, our example data set is publicly available data from the Home Mortgage Disclosure Act, de-identified data on mortgage loan determinations.

INTRODUCTION

This paper describes the philosophy behind our approach to Rapid Prototyping (RP) and practical examples for building a suite of prototype reports using this approach. To that end, this paper includes sections on the following topics:

- Rapid Prototyping (RP) and Agile development
- The Use Case Framework (UCF)
- General principles of user-centered design
- Selecting data for the prototype
- Wireframes and examples

Readers who want to skip the development methodology discussion can jump ahead to the Wireframes section to see the end result.

RAPID PROTOTYPING (RP) AND AGILE DEVELOPMENT

This paper is a practical guide to reducing the amount of time that is required to develop operational reporting using SAS Visual Analytics. Our focus here is on practical advice for data scientists and software developers to get started. But it is helpful to place the prototyping process in the context of Agile

software development. In many ways, this approach to Rapid Prototyping supports Agile development. Rapid Prototyping:

- begins the development process with functioning software, not just documentation of requirements
- accelerates opportunities for end-user involvement at the start of the first iteration, enabling users to provide domain-driven feedback to guide modifications and improvements
- enables users to begin testing general functionality at the start of the project, minimizing the risk of late-stage changes to requirements due to misunderstandings about basic functionality of the software
- improves the estimate of effort required for early sprints as work can be defined as the delta of “new” against the benchmark of “existing” prototype reporting
- simplifies the job of a Scrum Product Owner in communicating the features and functionality of the developing software to external stakeholders by providing them with a demonstration – not an explanation – of what the software can and will do
- clarifies the definition of “done” by providing the Scrum Master with a clear starting point, the already developed prototype reports, against which to estimate time and effort to develop the next iteration (not “what will be ‘done’ starting from nothing,” but “what will be ‘done’ when we have done something *more*”)
- provides downstream efficiencies; for example, reusing the same prototype reporting for multiple projects reduces the initial development time for each.

Rapid Prototyping facilitates successful implementation of the Agile process.

THE USE CASE FRAMEWORK (UCF)

While Rapid Prototyping does accelerate development, the key to success is not actually speed. The goal is not simply to brainstorm a large number of use cases and build as many example reports as possible in a short period of time. In fact, prototyping that way might actually exacerbate option overload in the development process, distracting the user community with esoteric features and functionality.

The key to successful rapid prototyping is to start from a coherent framework of general requirements -- the Use Case Framework (UCF). In our experience, there are three main parameters to the framework:

- Data Status (Observed, Estimated)
- Report Focus (One, Many, Combinations)
- Interactivity (None, Traditional, Advanced)

Each is discussed below.

DATA STATUS (OBSERVED, ESTIMATED)

In terms of the nature of the data, statisticians and report developers have a different way of thinking about data and data structures. Statisticians think of data and analytics in terms of the *unit of analysis*. The unit of analysis is the essential thing being measured in the data. A unit of analysis could be a country’s total exports in a particular year or it could be a body temperature reading from a particular patient at a particular time. Statisticians think of the data collected for that unit of analysis as belonging to one of four levels of measurement: nominal, ordinal, interval, and ratio. The unit of analysis and level of measurement are the primary characteristics of data that drive a particular statistical analysis.

For reporting, it is more helpful to think of unit of analysis and level of measurement in terms of a simpler dichotomy: observations versus estimates. For many purposes, reporting of *observed* data might be all that is required: this is the count of instances of X and this is the count of instances of Y; the mean of Z is 4.242; and the standard deviation of a group of measures might be 1.7. Those are counts and measures that reflect the data loaded to the system. While there might be measurement error in terms of capturing those data accurately (reporting oversight, poor data collection, buggy transactional system, and so on), the data that are the focus of this type of report are intended to reflect only what is available in the data set(s) that are loaded.

In contrast, statistical modeling produces *estimates*, forecasted values, groupings based on latent variables, and other measures that were not in the original data. From a mechanical perspective, these estimates are just other numbers or categories in the data. But their meaning sets them apart. These data reflect analytic insights generated by a statistical model. Unlike observed values, these data include associated data regarding the accuracy of the estimates, not just point estimates, but upper and lower confidence intervals and other measures of how well the model fits the data. These additional data must be included in reports that represent estimates and modeled metrics.

Reporting must take the nature of the data into account and clearly delineate values that are observed from those that are estimated. In any situation in which both observed and estimates values will be used, the suite of template reports developed should include examples of both.

REPORT FOCUS (ONE, MANY, COMBINATIONS)

In classic database development, relationships between tables are categorized with terms like *one-to-many* or *many-to-many*. Combinations of these relationships can be categorized into database structures as *first normal form*, *second normal form*, and so on. Borrowing that type of language from relational database development, we can think of an individual report as having a particular focus.

For example, a report might take as its focus a single observation – all of the data associated with *one* case of the unit of analysis. Thus, we might look at one patient, one country, one sale, one financial transaction. At the other extreme, our focus might be *many* or even all cases. We might look at a graph that represents every patient or every country, every sale, every transaction. Finally, there is an opportunity to begin the focus with one or many and move through the data to the other end of that dichotomy. We might focus on one case but provide benchmark comparison data on all such cases. We might start with all cases and allow users to drill down into the data to find one particular case. We can also allow users to search for terms or conditions and filter the data to show *some* of the many. A suite of prototype template reports should include reports demonstrating differences in focus.

INTERACTIVITY (NONE, TRADITIONAL, ADVANCED)

SAS Visual Analytics lends itself to highly interactive reporting, allowing users to cut through data from many different directions. The most striking difference between SAS Visual Analytics reports and previous SAS methodologies like ODS graphics is the degree to which you can use a graph, itself, to drive the analysis. A bar chart is not just a visualization of the data; it is the responsive entry point for exploring the data. Instead of clicking through highlighted links in a hierarchical data table, you can explore the data directly through a visualization.

From the report developer perspective, this completely blurs the line between a traditional business intelligence report and a custom software app for that data. In practice, the only difference between building a SAS Visual Analytics visual data interface and coding a custom app is that you don't need to write and compile the code. Indeed, the line between report writer and software developer becomes very blurry using SAS Visual Analytics.

That said, not every report needs to be developed using the most advanced interactivity. There is and will probably always be a need for more basic reporting (traditional data drill-down, traditional dial and gauge KPIs, and so on). There will also be a need for a representation of the data that is completely static, locked down, and unchangeable. Suffice to say, predictions of a completely paperless office were

forward-looking statements subject to uncertainties in local and global markets...Many organizations still have a significant need for static reports and printable PDF files. SAS Visual Analytics can do that, too.

Not every report is equally well-suited to a static representation. Highly interactive reports can be saved as a PDF with all filtering and selections documented. But it is better to understand the degree to which access to those data needs to be responsive before developing a report. The degree to which the report must be actively responsive to user interaction is an important distinction in our Use Case Framework.

In practice, it would be easy to create 18 different reports simply by crossing each of these parameters: 3 levels of interactivity by 2 types of data by 3 types of focus = 18 reports. Not every cell of that three-way table necessarily needs to be represented in your prototype reporting, but example reports developed for the prototype can be presented within this framework. Having the framework handy can simplify requirements gathering for the final production reporting.

GENERAL PRINCIPLES OF USER-CENTERED DESIGN

Within the context of the Use Case Framework, general principles of sound report development still apply:

- Reports should function left-to-right and top-down or, if not, must be clearly marked as different.
- Color should be used consistently and clearly so that one distinctive color has one meaning.
- Where possible, colors should reflect common expectations (money = green, bad = red).
- A filtered selection should be evident in any affected aspect of the report.
- The same type of multivariate graph or visualization should not be used more than once in a single view, especially not to represent two different types of data.
- If a hierarchy is used, no report should use more than one hierarchy.
- Reports should not exceed the cognitive capacity of the typical user. For example, most people are able to consider the effects of several independent variables on an outcome at the same time. Most people are not able to simultaneously consider a dozen such effects. Polluting a report with options to select or filter using more than a few variables will confuse most users.
- Dynamic aggregation, summarization, and/or other calculations must reflect appropriate rules for that analysis at every level reported. This can be particularly challenging if the data are statistical estimates or other modeled metrics. In those cases, if the report developer is not an expert in that calculation, an expert should be consulted.
- Accuracy of reporting is always more important than interactivity. For example, if the combination of complex business rules at different levels of a hierarchy is such that a metric cannot be presented accurately at all of the levels of a hierarchy, it would be better to sacrifice interactivity and present a more static version of those data than risk a misrepresentation of the data.
- Meaningfulness of the representation is always more important than aesthetic interestingness. With real data, beauty is not always truth and the truth is not always beautiful.
- Exports of data (as data or through exported reporting) should include all parameters that define the export.
- Reports should function in a similar manner to reinforce user expectations or, where one report is built using a different approach, should make clear how that report functions differently.
- Titling should be meaningful, including intended use if possible. "Forecast for Capacity Planning" is a better title than "Capacity Forecast" or, simply, "Forecast."
- Where data were collected over time, the time period of available and displayed data should be made clear.
- Online help should be clearly marked and consistent.

SELECTING DATA FOR THE PROTOTYPE

The Use Case Framework and general guidelines for development help to standardize the rapid prototyping process, but one of the biggest challenges in using template reports is beyond your control. Will your user community be open to starting development using your prototype? Some user groups might find it challenging to review a suite of generic reports without seeing their own specific data in that format. Efforts to explain, “This is a hierarchy like your hierarchy so we could have your three levels right here in the report” are met with blank stares. That happens. In some cases, a flexible suite of reports could be populated with the new user group’s own data before the first review. In other cases, the timeline might not allow for that.

The best strategy is to build the template reports using a robust enough data source to support many different types of reports and visualizations. Ideally, a single rich data set can be used. Having developed template SAS Visual Analytics reports for a wide variety of industries, we know what an ideal data set would include. We recognize that such a data set is not always available, but if we could source our data from an infinite database of records, we would choose a single data source with the following characteristics.

The ideal data:

- are recognizable and meaningful to the organization
- *can* be organized according to well-known business structures (geo, division, department) to demonstrate traditional BI reporting (KPIs, drill-down, dynamic aggregation)
- include both observations (transactions, records) and statistical estimates (predicted values, groupings based on latent trait/segmentation/decision-tree models, and so on)
- are wide enough and deep enough to provide ample opportunity to demonstrate interactivity
- can be reduced to a unique case of interest to provide an opportunity to demonstrate search and needle-in-a-haystack reporting
- include geolocation information for mapping
- include dates/datetimes for time series analysis and reporting

For the examples below, we are using publicly available data collected according to laws and regulations under the Home Mortgage Disclosure Act. “The Federal Reserve Board created these statements under Regulation C of the Home Mortgage Disclosure Act (HMDA) to help determine whether lending institutions are serving the needs of the community, assist public officials in distributing public-sector investments, and identify any possible discriminatory lending practices” (<https://catalog.archives.gov/id/4699299>). The original data were a flat file of records from the “Home Mortgage Disclosure Act (HMDA) / Community Reinvestment Act (CRA) Combined Census Data Files” set for 2009. The data are available online here: <https://www.ffiec.gov/hmda/hmdaflat.htm> and here: <https://catalog.archives.gov/id/4700101> .

These data fulfill almost all of the characteristics of an ideal data set for prototyping. These data are relevant to business processes at many lending institutions and the data include categorical variables, numeric metrics, geolocations, and multiple ways of classifying relationships between agencies and purchasers and applicants and co-applicants. The source data have already been de-identified, but the data we are using here have been scrubbed and modified such that specific conclusions about home mortgage determinations in 2009 are not valid. This reporting is provided here only to demonstrate the application of the wireframes to real-world data.

Note that these data do not include forecasting estimates or latent trait groupings from statistical models, but they do include estimates of population parameters for census tracts (minority percent ownership, income, and so on). We have marked those report examples as including some estimated data, although

confidence intervals were not available for reporting. These data do not include datetime values that would allow us to generate time series reporting, but users can gather additional HMDA data files and create a longitudinal structure.

WIREFRAMES AND EXAMPLES

Where possible, you should use standard layouts and repeatable designs that can be used for many different types of data. Below is a set of common wireframes for SAS Visual Analytics reporting. Some of these structures are better than others in terms of the three parameters of our framework: data (observed, estimated), focus (one, many, combinations), and interactivity (none, traditional, advanced), but many can be modified to serve more than one purpose. For example, the majority of visualizations below allow for a medium-to-high degree of interactivity, but the same check boxes, sliders, and other controls in the table could be fixed for one version of the report that is then generated automatically as a fixed-format PDF.

The key (Figure 1) below explains the wireframe diagrams. It provides examples of some specific SAS Visual Analytics report design elements (Table, Graph, Control, and Other) that can be used to build actual reports. Following that is a gallery of wireframes (Figure 2) that shows a dozen generic design examples. There are many different types of software for generating wireframes, but in practice, a box of colored pencils and a pad of paper are good enough. When available, a whiteboard and markers work just fine.

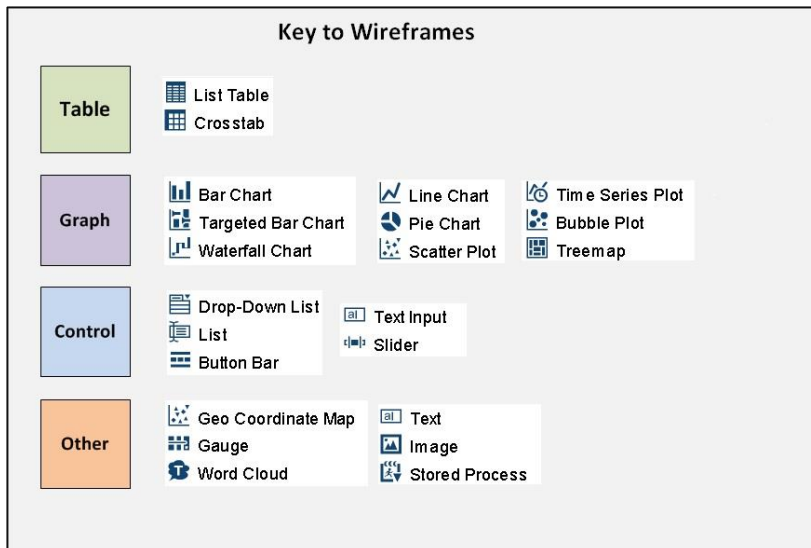


Figure 1. Key to the Wireframes

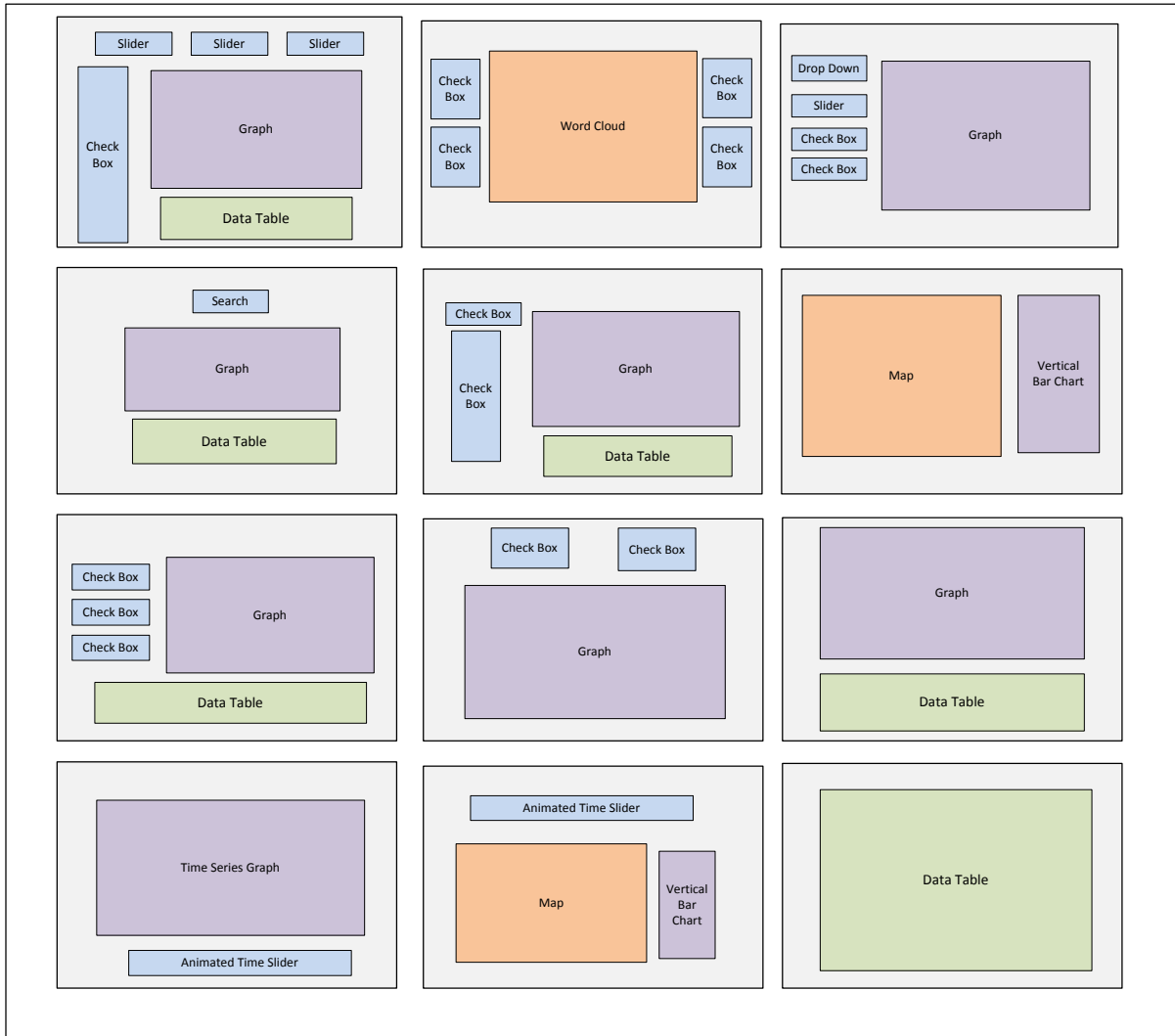


Figure 2. Gallery of Wireframes

In the examples that follow, we demonstrate a variety of color and theme options. For prototype development, we recommend a common theme that reflects corporate standards and, where appropriate, whatever visual branding is typically implemented. The color and theme options here are intended only to represent some possibilities. The examples categorize each report within the Use Case Framework and describe how the example report functions. There is only one example of each wireframe report. There are, in fact, many ways to use each wireframe structure, substituting different graphics and controls as appropriate for that report.

Example 1

Data: Observed and/or Estimated
 Focus: Many-to-one
 Interactivity: Advanced

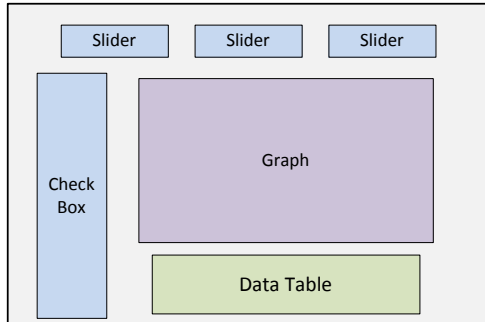


Figure 3. Example 1 Wireframe

This report allows analysts to select one or more states and then explore the effect of three independent interval/ratio variables on the total number of loans and the percent of those loans approved disaggregated by percent minority home ownership in that census tract. This provides a great deal of power to the analyst to see the interaction of these variables. The data table at the bottom includes data down to the unique loan level and can be used as an export tool to download data for further analysis.

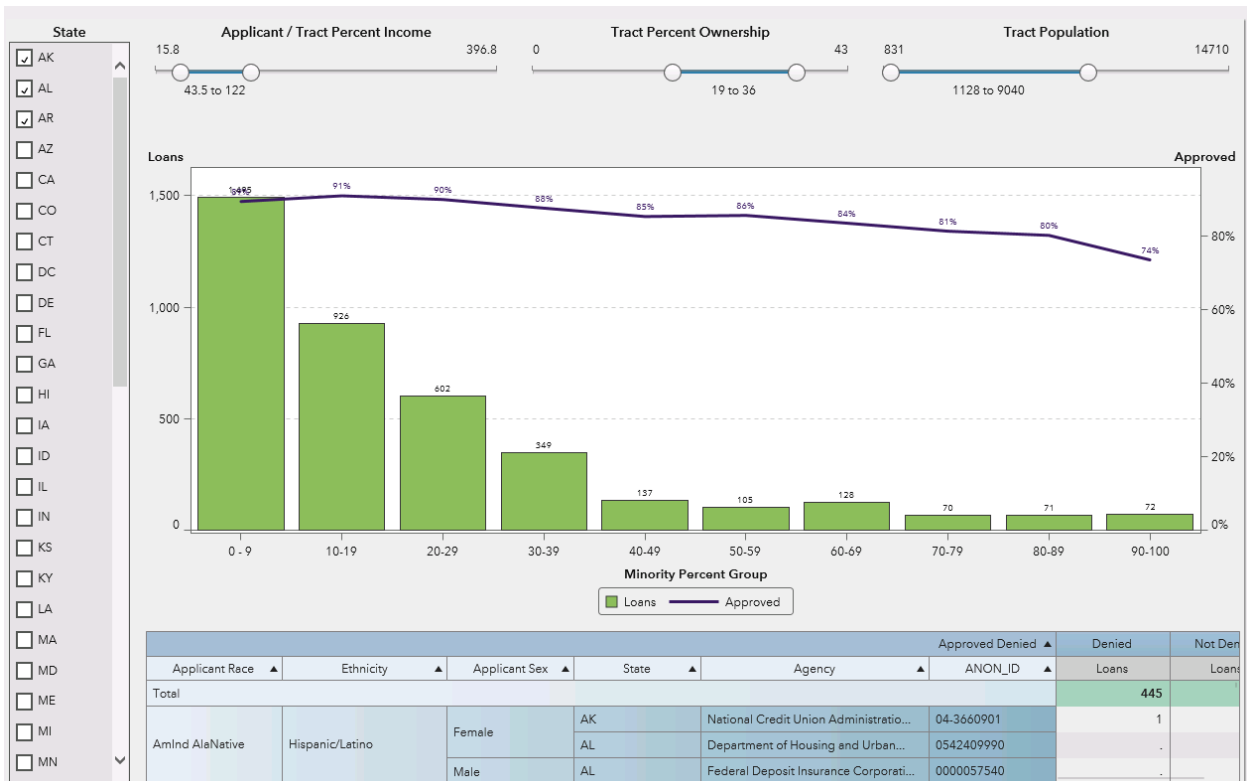


Figure 4. Example 1 Report

Example 2

Data: Observed

Focus: Many

Interactivity: Advanced

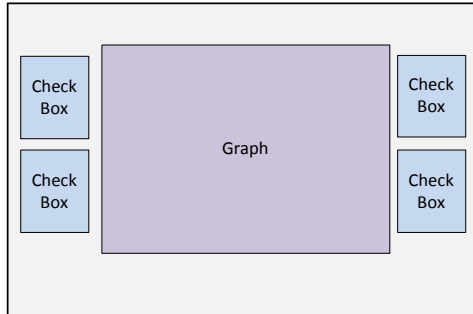


Figure 5. Example 2 Wireframe

This report uses a word cloud to visualize denial reasons for a loan determination. The size of the word represents the number of loans denied for that reason and green coloring moves from a lighter to darker hue representing lower to higher loan amounts. While the four check boxes are easily understood, we categorize this as advanced interactivity because the check boxes operate largely independently. A user could select options for all four or just one, for example, co-applicant race. In that way, this actually captures the effect of four separate and/or interacting variables on the text analysis of denial reason descriptions.

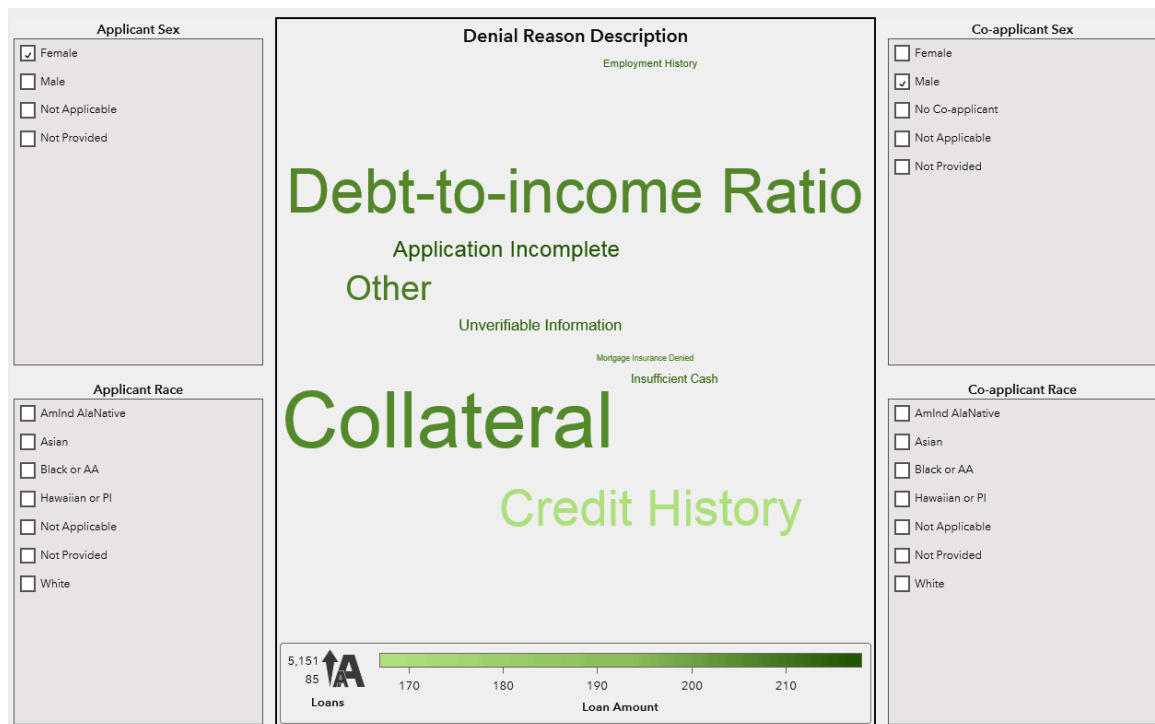


Figure 6. Example 2 Report

Example 3

Data: Observed and/or Estimated
Focus: Many-to-one
Interactivity: Traditional to Advanced

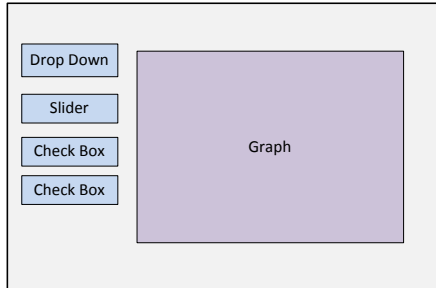


Figure 7. Example 3 Wireframe

This report allows users to select levels of four variables and see differences in the visualization of individual loans (loan amount by income). In some ways, that would seem to be a more advanced mode of interactivity, but the prompts at the left actually operate in a more traditional manner. The top pull-down menu selects one county, immediately drilling into the data and restricting the view to just that county. The prompts cascade so that moving top to bottom, the only selections available are based on the selections made above. The percent minority slider affects the loan purpose options, which affects the owner occupancy selections. Each choice is immediately reflected on the graph which, via a tooltip, can show the unique loan for each bar. Multiple loans are displayed, but users can filter data down to a single unique loan.

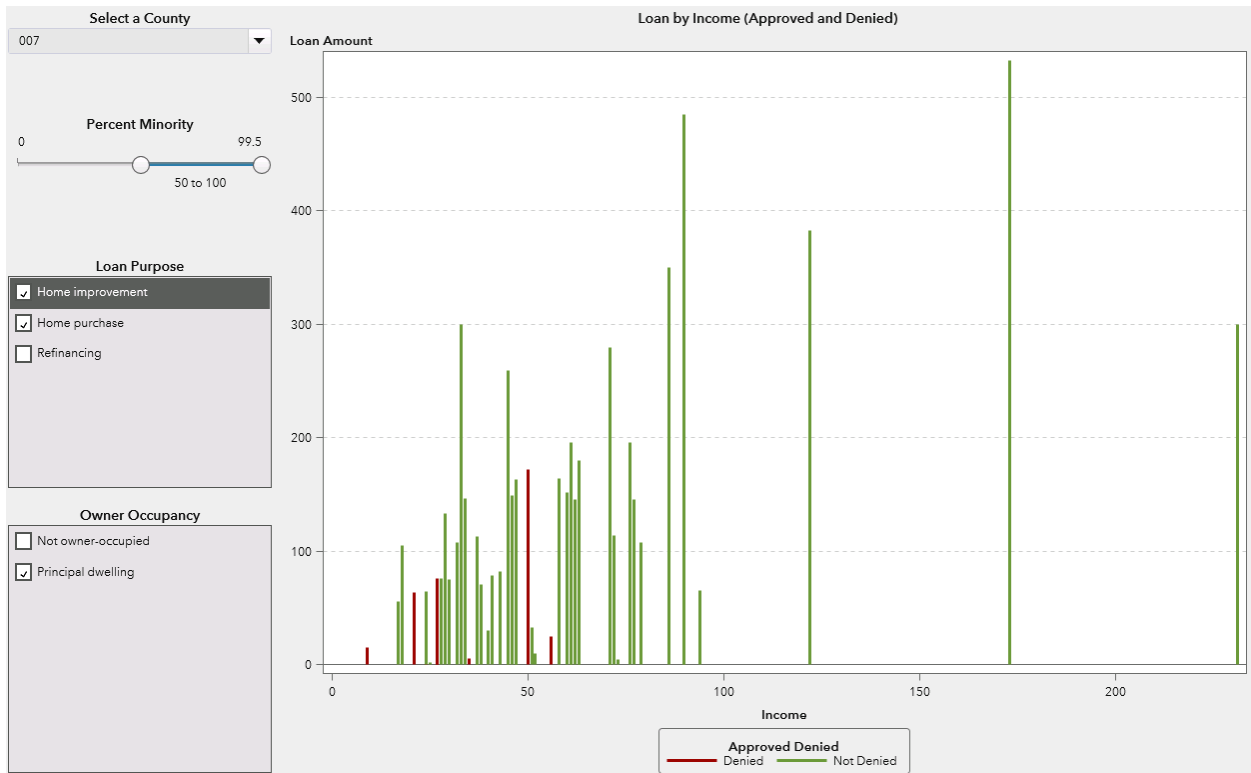


Figure 8. Example 3 Report

Example 4

Data: Observed

Focus: One or Many-to-one

Interactivity: Traditional Search

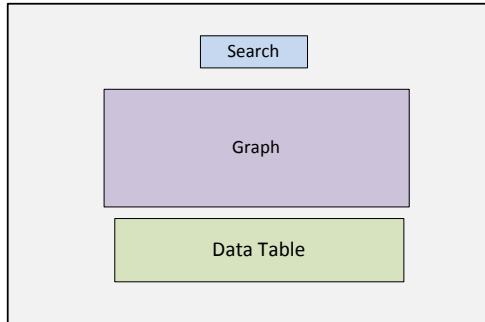


Figure 9. Example 4 Wireframe

This report demonstrates a classic search that can be used to find one census tract or even one loan within a tract. There are thousands of census tracts in the data set. A user can search for one using the search field at the top of the report. That selection cascades down to the graph and data table. The graph shows all loans, approved and denied, using the applicant's income on the X-axis, loan amount on the Y-axis, and approval on a lattice row. By selecting one of the bars in the graph, the user will get a table of the data for that one unique loan. Similarly, users could select multiple bars and the data table would reflect that group.

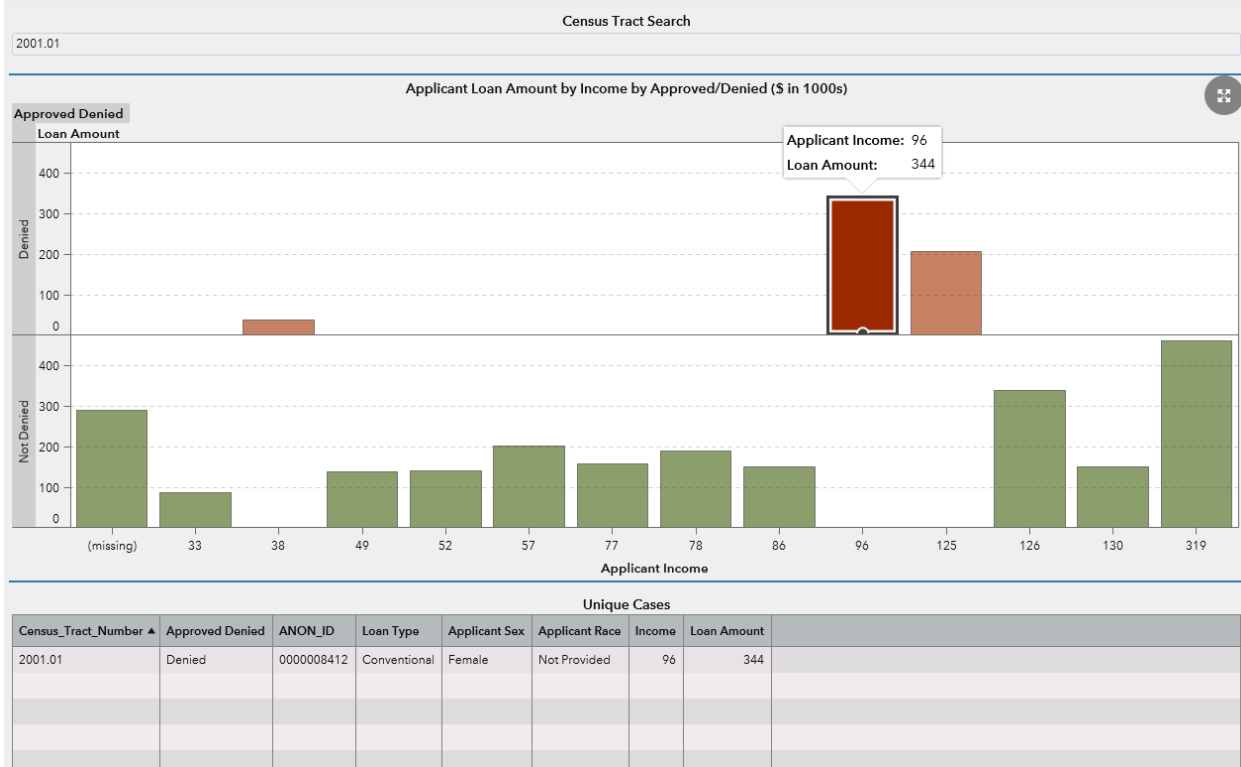


Figure 10. Example 4 Report

Example 5

Data: Observed
 Focus: Many
 Interactivity: Traditional

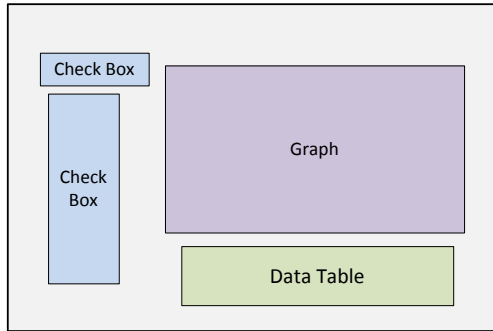


Figure 11. Example 5 Wireframe

This report takes a very traditional approach to interactivity. Users select one (but not both) types of loans at the top left, choose one or more denial reasons in the check box below it, and see a dual-axis bar chart of the number of loans and loan amounts. A data table at the bottom reflects all selections, including the use of the graph, itself. Data are presented only at the aggregated level and do not provide access to a unique loan.

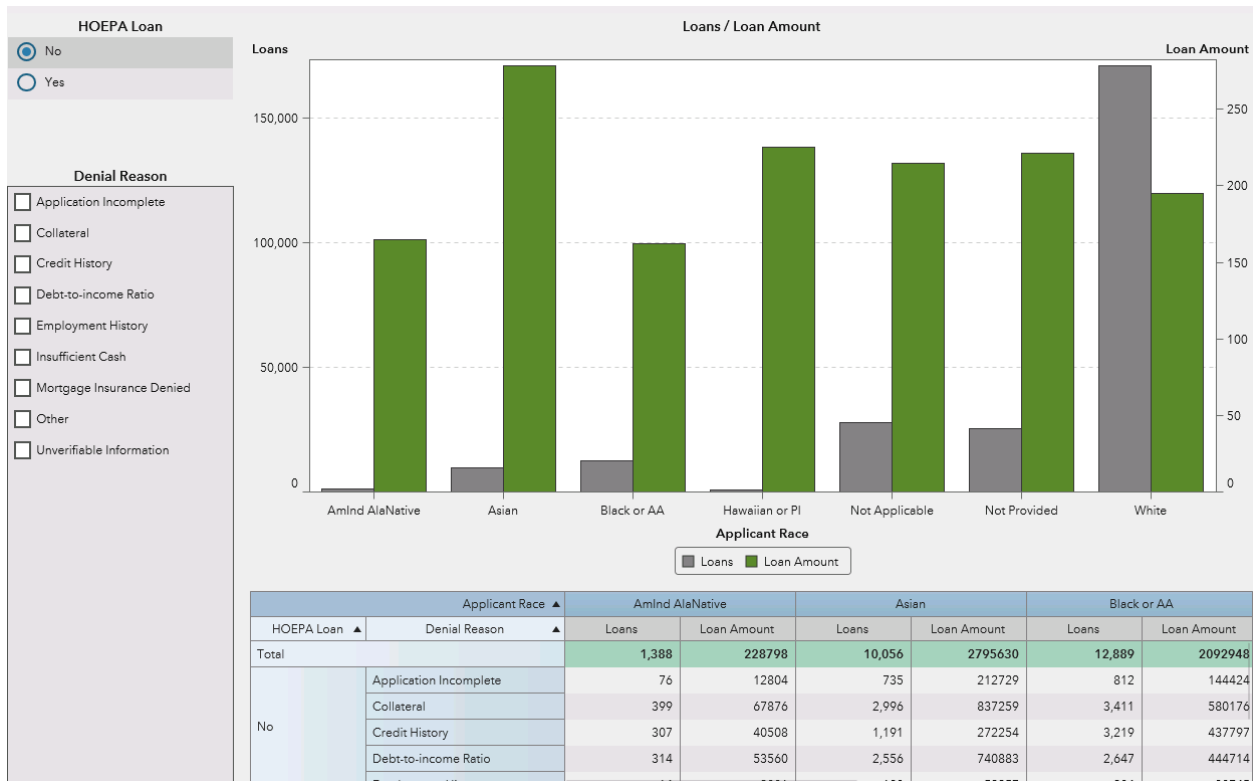


Figure 12. Example 5 Report

Example 6

Data: Observed
Focus: Many
Interactivity: Traditional

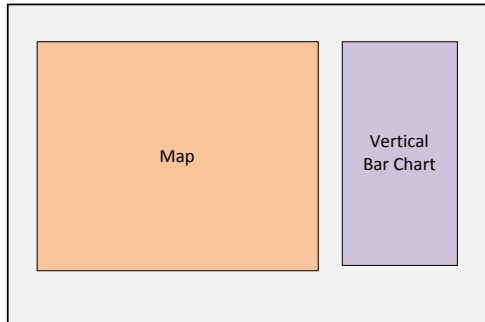


Figure 13. Example 6 Wireframe

This report is typical of public reporting of secure data. It uses observed values only, no statistical estimates. The map and corresponding vertical bar chart provide some interactivity. Users can select a state or group of states (including geographically discontinuous states) and the bar chart reflects those selections. Not shown in the screenshot, the bar chart does include a three-level hierarchy of the data, enabling users to drill into the bars. This type of report is ideal for introducing a suite of reports, especially for users who have a geo-centric view of their business, but it will not provide the type of analytic power that more statistically savvy users demand.

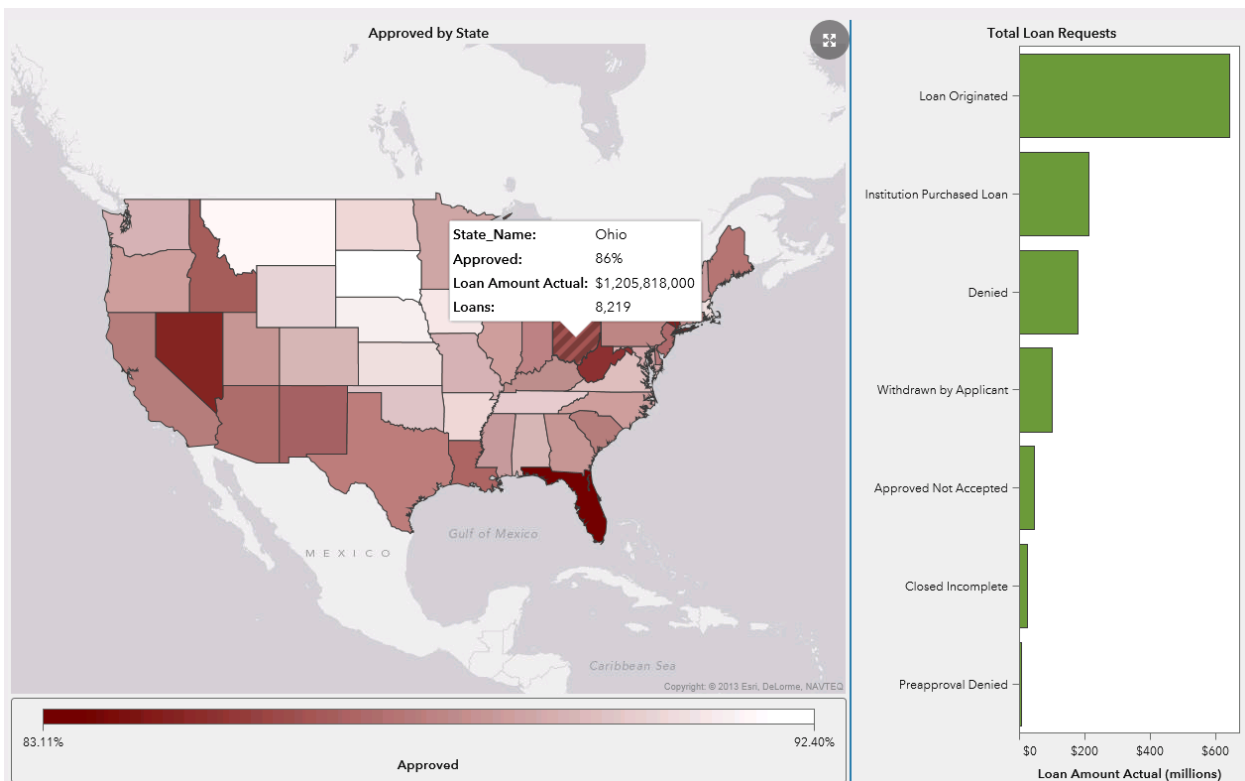


Figure 14. Example 6 Report

Example 7

Data: Observed
 Focus: Many
 Interactivity: Traditional

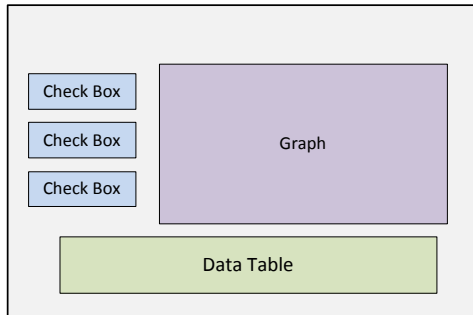


Figure 15. Example 7 Wireframe

This is another example of a report that presents aggregated data on many loans. All values in the graph represent money, so we have used two shades of green. The interactivity is more traditional with cascading check boxes on the left, although the data presented in the graph are built on a hierarchy that allows for additional exploration of the data. The data table at the bottom can be locked to the report or provided for export. The report has been designed to provide significant insights into loan trends and patterns, but it does not provide any access to a unique loan.

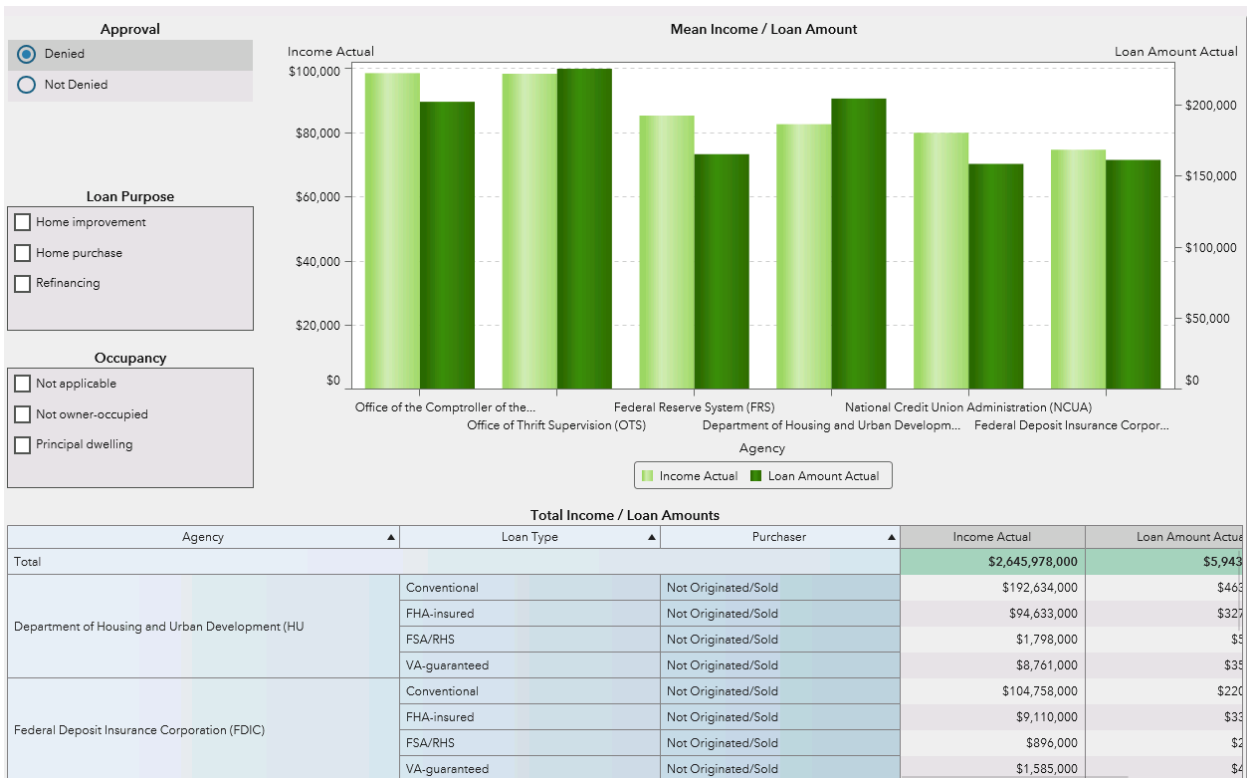


Figure 16. Example 7 Report

Example 8

Data: Observed
Focus: Many
Interactivity: Traditional

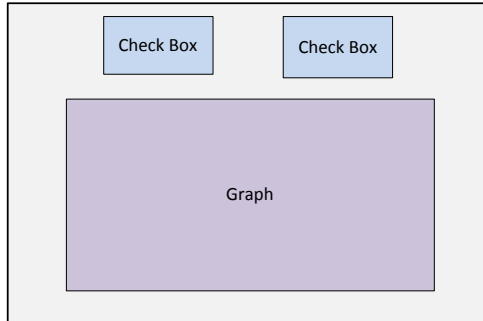


Figure 17. Example 8 Wireframe

This type of structure is ideal for exploring the relationship between two entities, in this case, agency and purchaser. The graph below can be used to explore how agency and purchaser relationships affect the number and size of loans.

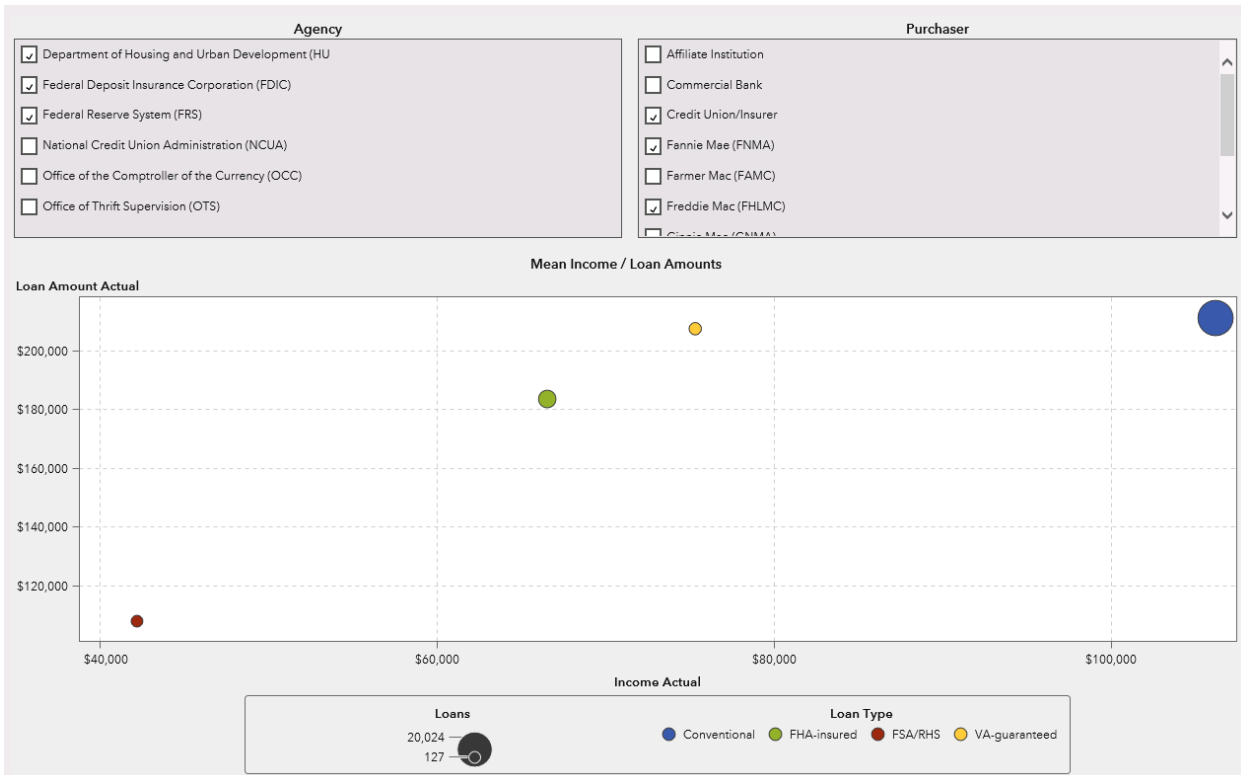


Figure 18. Example 8 Report

Example 9

Data: Observed and/or Estimated
 Focus: Many
 Interactivity: Advanced

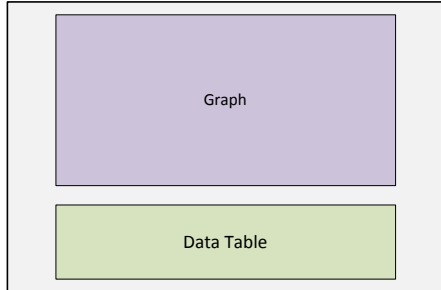


Figure 19. Example 9 Wireframe

While this report might appear simple, it provides a unique way to explore these data. Traditional interactivity starts with a filter and proceeds to a visualization; advanced interactivity starts with the visualization. The graph slows the total number of loans by the minority percent home ownership in that census tract. Users can use the graph to capture one agency, one minority percent group, or multiple continuous or discontinuous combinations of those variables in the graph. Data in the table below reflect the details of that selection. This is a good example of a visualization that drives the data analysis. Note that the same wireframe structure can be useful with bar charts, line plots, treemaps, and other visualizations.

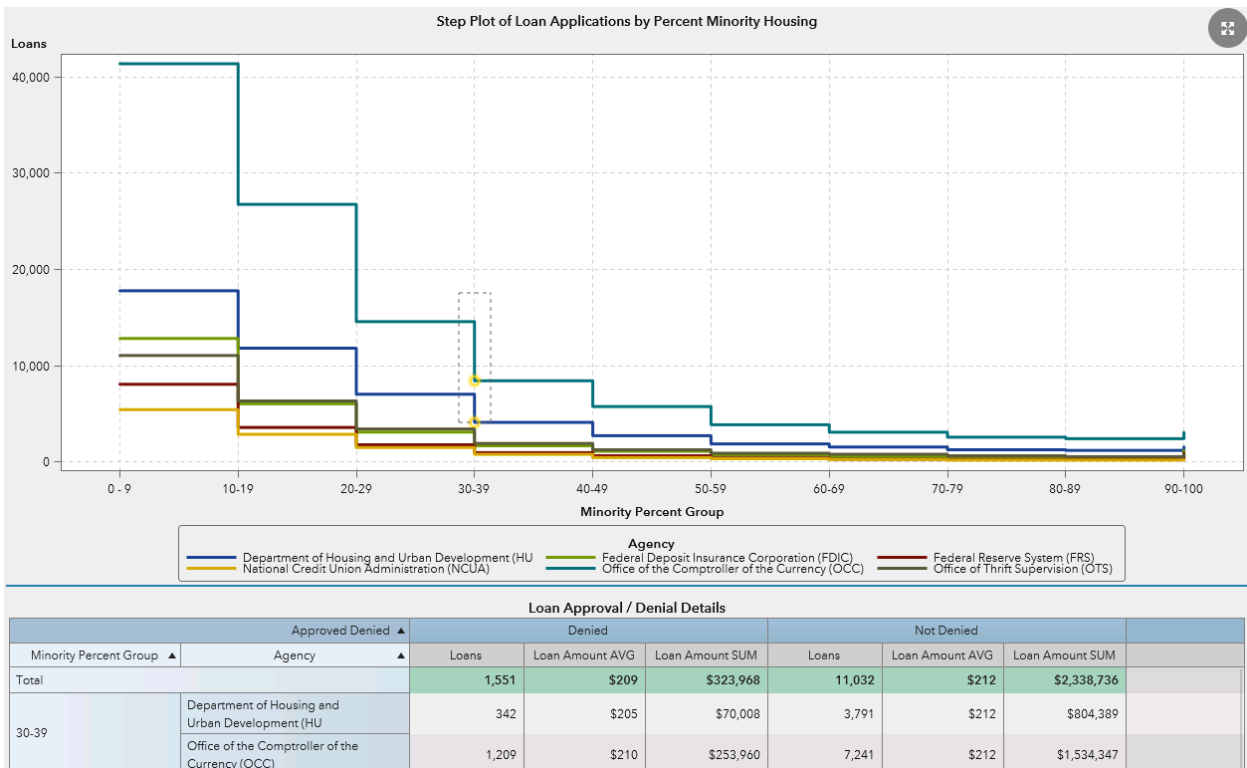


Figure 20. Example 9 Report

Example 10

Data: Observed
 Focus: Many, Many-to-one
 Interactivity: Traditional

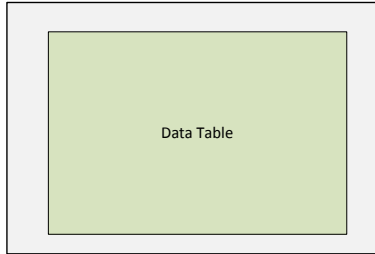


Figure 21. Example 10 Wireframe

This report requires no explanation. This is a typical data table driven by a hierarchy. This is best used for observed data and can be set to enable drill-down to any level within the data from larger groupings to individual cases. This most closely resembles traditional OLAP reporting, which might be required for users expecting this functionality.

Loan Data (\$ in 1000s)						
	Denied			Not Denied		
	Loans	Loan Amount AVG	Loan Amount SUM	Loans	Loan Amount AVG	Loan Amount SUM
Total	29,650	\$200	\$5,943,207	220,348	\$202	\$44,403,433
Alabama	401	\$127	\$50,836	3,503	\$147	\$516,252
001	2	\$105	\$210	38	\$154	\$5,844
Department of Housing and Urban Development (HU)	.	.	.	9	\$148	\$1,332
Federal Deposit Insurance Corporation (FDIC)	.	.	.	5	\$119	\$593
Federal Reserve System (FRS)	.	.	.	2	\$123	\$246
National Credit Union Administration (NCUA)	.	.	.	3	\$155	\$466
Office of the Comptroller of the Currency (OCC)	1	\$117	\$117	12	\$158	\$1,901
Office of Thrift Supervision (OTS)	1	\$93	\$93	7	\$187	\$1,306
003	32	\$217	\$6,954	176	\$176	\$31,020
Department of Housing and Urban Development (HU)	8	\$179	\$1,435	36	\$154	\$5,538
Federal Deposit Insurance Corporation (FDIC)	2	\$240	\$479	20	\$133	\$2,667
Federal Reserve System (FRS)	5	\$310	\$1,551	43	\$203	\$8,713
National Credit Union Administration (NCUA)	.	.	.	5	\$105	\$527
Office of the Comptroller of	13	\$221	\$2,868	62	\$189	\$11,747

Figure 22. Example 10 Report

CONCLUSION

The Use Case Framework presented here provides a useful structure for building a suite of template reports. In our experience, almost every type of report can be classified within the context of these parameters:

- Data Status (Observed, Estimated)
- Report Focus (One, Many, Combinations)
- Interactivity (None, Traditional, Advanced)

The framework is not a theoretical framework. It was developed in consideration of the ways in which other data professionals work with data. It respects the perspective of statisticians, database developers, and software developers. It has been field-tested with numerous prototyping projects. While there are probably some rare exceptions, almost every type of report will fit into one of the 18 cells of this 2x3x3 matrix.

Using a set of standardized wireframes, including the ones presented here and others that you develop yourself, you will have a ready-made set of structures on which to build your prototype. Following the general principles for user-centered design, you have a rulebook for ensuring that individual reports in the prototype suite work as expected for users.

Real-world prototyping and software development is often done under extreme time constraints. One of the criticisms of some software development methodologies is that those methods add more time and managerial overhead to the development process. The implementation of a template suite of reports following this Use Case Framework offers considerable time savings at the start of any reporting project. Once implemented, that suite of prototype reports will provide additional time savings on any future projects.

FURTHER READING

For an example of a large scale implementation of SAS Visual Analytics using this approach, see the SAS Visual Analytics for UN Comtrade project here:

http://www.sas.com/en_us/software/visual-analytics-comtrade.html

For more by the authors of this paper, see:

“Visualizing Clinical Trial Data: Small Data, Big Insights”

<http://support.sas.com/resources/papers/proceedings15/SAS1888-2015.pdf>

“Instant KPI: From Data to Dashboard in Record Time”

<http://support.sas.com/resources/papers/proceedings10/323-2010.pdf>

Full documentation of SAS Visual Analytics is available here:

<http://support.sas.com/software/products/va/>

CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the authors:

Elliot Inman, Ph.D.
100 SAS Campus Drive
Cary, NC 27513
SAS Institute Inc.
Elliot.Inman@sas.com
<http://www.sas.com>

Michael Drutar
100 SAS Campus Drive
Cary, NC 27513
SAS Institute Inc.
Michael.Drutar@sas.com
<http://www.sas.com>

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.