

Supporting Exploratory Data Analysis with Live Programming

Robert DeLine and Danyel Fisher
Microsoft Research
Redmond, WA, USA
{rdeline, danyelf}@microsoft.com

Abstract—Data scientists often conduct exploratory data analysis in scripting environments with a read-eval-print loop (REPL), like R, IPython or MATLAB. This user experience requires diligent management of execution and generates lengthy histories of unwanted command responses. This paper explores the alternative of live programming, a user experience in which the user’s edits immediately and automatically update the script results—a “ripple” effect familiar from spreadsheets. Which user experience provides better support for exploratory data analysis, REPL or ripple? We conducted a controlled lab study with 15 data-experienced professionals. Each participant explored four datasets, two in each experience. The REPL sessions left histories with both significantly more data results and significantly more errors than the live sessions. However, both experiences produced comparable numbers of data results that participants self-rated as insightful. Participants largely preferred the live experience for its responsiveness and ability to keep the script content clean, but missed the visible history that a REPL provides.

Keywords—Programming environments; read-eval-print loop (REPL); command loop; live programming; data analysis; data mining; data science.

I. INTRODUCTION

The practice of data science is becoming increasingly important as businesses, governments and other institutes shift towards data-driven decision making. The Harvard Business Review calls data scientist the “sexiest job of the 21st century” [1], yet the tools that many data scientists use are seemingly from the mainframe era [2]. In this paper, we focus on what is often a key stage for generating insights from data, namely, the use of scripting environments to perform exploratory data analysis. During exploratory data analysis, a data scientist investigates, visualizes and summarizes her data in

order to understand them, before moving on to deeper statistical techniques. While some data exploration can be done in specialized direct-manipulation tools, such as Tableau or Microsoft Excel, more powerful data cleaning and manipulation still requires interactive scripting. Today’s most popular environments for this activity—Python, R, and MATLAB—are all based on the read-eval-print loop (REPL), a user experience that dates to the 1960s.

Using a REPL, a data analyst iteratively enters a scripting statement, waits for it to evaluate, and sees the resulting value or error message. During a REPL session, an analyst might try many alternatives, building up a lengthy history of statements and their values. These statements often chain together to form longer computations. For example, in Figure 2, a data scientist enters two statements in a REPL to look at the survival of male Titanic passengers. Switching to female passengers requires a series of clerical steps: she uses the up arrow key to bring back the first statement in her history; she changes “male” to “female” and enters the edited statement; finally, she uses the up arrow key to re-enter the histogram statement, since her previous edit requires the histogram to be recomputed. In short, while data scientists find REPL environments useful for exploratory data analysis, they have downsides like noisy session histories and the diligence needed to manage re-execution.

Researchers have explored the concept of *live programming* to improve programming education and programmer productivity. In a live programming environment, editing a program has an immediate and automatic effect on the program’s execution. This level of responsiveness is a good fit for exploratory data analysis. To replay the earlier Titanic data analysis in a live environment, the data scientist enters

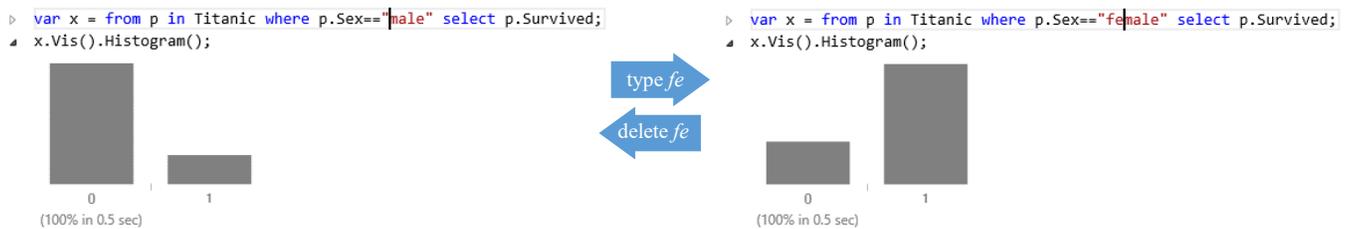
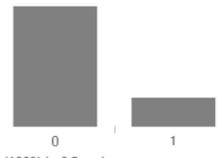


Figure 1. In a live environment, a data scientist enters a two-line script to see the male survival rate (left). She then edits the script to change “male” to “female” (right), which automatically triggers the histogram to update. She can switch between the two results by alternately typing and deleting “fe”.

```

> var x = from p in Titanic where p.Sex=="male" select p.Survived;
> x.Vis().Histogram()

```



```

> var x = from p in Titanic where p.Sex=="female" select p.Survived;
> x.Vis().Histogram()

```

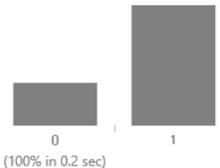


Figure 2. In a REPL, the user looks at male and female survival rates in several steps. She enters the first two script lines to see the results for males. Then she types up arrow twice and changes “male” to “female”. Last, she types up arrow twice to re-compute the histogram.

two statements to investigate male survival, as before (Figure 1 left). This time, to switch to female survival, she edits the document in place to change “male” to “female”, causing the two results to update immediately and automatically (Figure 1, right).

In a live scripting environment, the user’s edits have a “ripple” effect on the displayed values, like in a spreadsheet. This reduces clerical work and clutter, at the cost of hiding history. Given these tradeoffs, we wonder: which user experience is better for exploratory data analysis, REPL or ripple? In particular, does the choice of user experience affect how data scientists handle errors or the number or quality of the data results they produce?

To investigate these questions, we created a prototype live programming environment for exploratory data analysis. Next, we conducted a controlled lab study of 15 participants at a large, data-driven software company. The contributions of this paper are:

- an algorithm for implementing a live user experience, given a scripting language with a parser, interpreter and introspection;
- a controlled lab study comparing REPL and live user experiences in a scripting environment; and
- data scientists’ subjective assessments of REPL and live user experiences for exploratory data analysis.

II. BACKGROUND

We chose to explore live scripting rather than direct manipulation because of its ubiquity within data science. Kandel *et al.* [3] interviewed 35 data analysts and categorized them into three “archetypes.” Two of the three archetypes make extensive use of scripting languages in their process of data analysis and visualization.

Live programming was first created for visual programming languages, focused on specific problem domains like physics simulation [4] and image processing [5]. This domain focus continues today. Much of the current live programming community works on *live coding*, in which a programmer-musician continuously updates a MIDI program to create a live musical performance [6] [7].

Other live programming systems help students learn programming, using multiple variations of the live model. For example, Alvis Live provides live algorithm animation side-by-side with program text, to help novice programmers learn data structures and control flow [8]. Flogo II, an environment for end-user robot programming, provides *live text*, in which the appearance of the program text reflects the current execution behavior, for example, by graying out untaken branches [9]. AgentSheets provides a live spreadsheet representation of rewrite rules that encode the logic of a simulation [10]. YinYang supplements live programming with a debugging experience that allows the user to probe the values of expressions [11]. In a variation of live programming, TheSeus supplements a standard program environment with always-on visualizations of execution behavior [12].

According to Tanimoto’s liveness taxonomy [5], a programming environment can offer four levels of liveness: (level 1) a user’s edits have no effect on the computation (e.g. a typical text editor); (level 2) a user explicitly submits edits to cause updates to the computation (e.g. a REPL or the IPython Notebook [13]); (level 3) a user’s edits automatically trigger any necessary re-computation (e.g. a spreadsheet); and (level 4) a user’s edits trigger updates to ongoing computations (e.g. changing a game while it is played [14]). In this paper, we focus on level-3 liveness because data analysis scripts are often functional programs, rather than interactive behaviors. Our environment exhibits level-4 liveness when a data scientist uses it on live streaming data, which we cover in another paper [15].

III. LIVE PROGRAMMING FOR DATA SCIENTISTS

We begin our explorations in a live programming environment with the Tempe analytics system [15] [16]. Here, we briefly discuss several design decisions in Tempe’s implementation of live coding, to clarify the user study.

Tempe’s overall look and feel is a research notebook, in the style of Mathematica [17] or the IPython Notebook [13]. Whereas the notebook pages in both Mathematica and the IPython Notebook host REPL sessions, a Tempe notebook page instead contains a live script, which we believe is a better fit for the notebook model. Hosting a REPL session in a notebook makes it possible to create notebook pages that look incorrect. For example, consider the IPython notebook in Figure 3. If the user deletes the second command box, labeled “In [2]”, from the notebook page, the remaining statements look nonsensical. The statement numbering providing the only hint about the root cause. IPython Notebook users

```
In [1]: x = 3
In [2]: x = x + 1
In [3]: y = x
In [4]: y
Out[4]: 4
```

We cut this command.

Figure 3. Hosting a REPL session in a notebook is awkward. If the user erases the second command from the page, the sequence looks strange.

avoid this problem by grouping dependent statements together a single command box, which forces them to execute together. Instead, Tempe tracks dependencies and maintains the invariant that the evaluation results always reflect the script’s current content.

Tempe shows a scripting error inside an orange box, with the relevant text underlined in orange:

```
select p from Titanic |
; expected
```

This is similar to the “red squiggle” feedback in modern programming environments, except that Tempe errors encompass both compile-time errors and run-time exceptions.

We also made several design choices to strike a balance between responsiveness and distraction. To keep the backend from trying to execute incomplete thoughts, the frontend sends updates only when the user is idle for a threshold amount of time (currently, 500 ms). Tempe also attempts to maintain spatial stability as results change. In particular, when a user’s edits produce a new, non-error result, Tempe replaces the old result. However, if an edit produces an error or no result, Tempe grays out the obsolete result, but leaves it in place, in anticipation that the user will shortly produce a new result.

Statements execute asynchronously. If the user enters a non-terminating or unwanted statement, Tempe provides a global Stop button that cancels all ongoing computations.

A. Implementation

Tempe uses C# as its scripting language, chosen in part for its LINQ feature. LINQ is an API for querying and transforming data [18] and provides a syntax extension for embedding SQL-like queries in C# code (Figure 2, first statement).

Tempe invokes our live programming algorithm whenever the user changes the script (Figure 4). The first half of the algorithm computes the top-level statements that need to be re-executed. The initial set consists of those statements that have syntactically changed since last execution, based on Wagner and Fischer’s dynamic programming algorithm for shortest edit distance [19]. (Most descriptions of this algorithm have it return the number of edits; our version *Diff*

```
interface ILanguage
{
    SyntaxTree Parse(string text);
    List<object> Interpret(List<Statement> stmts);
    List<Result> Inspect(List<object> obj);
}

class LiveProgram
{
    ILanguage language;
    SyntaxTree oldTree;

    List<Result> UpdateResults(string newContent)
    {
        var newTree = language.Parse(newContent);

        var diffs = Diff(oldTree.Statements,
            newTree.Statements);

        var changed =
            from diff in diffs
            where diff.Kind == Change ||
                diff.Kind == Insert
            select diff.Statement;

        var du = new DefUseAnalysis(newTree);
        changed.AddRange(
            du.GetUses(du.GetDefs(changed)));

        var values = language.Interpret(changed);

        oldTree = newTree;
        return language.Inspect(values);
    }
}
```

Figure 4. Our live programming algorithm, with its dependencies on an interpreter.

returns the edits themselves.) Finding the minimal set of changes between consecutive versions reduces the re-computation that the backend performs, which is particularly important for large datasets. The algorithm then extends the set of statements to re-compute, based on the results of a definition-use chain analysis [20]. Namely, the algorithm adds a statement to the set if the statement (transitively) uses a variable that an edited statement defines. The second half of the algorithm uses an interpreter to execute all the changed statements and summarizes the resulting objects.

Data analyses written in LINQ tend to be functional, i.e. free of side effects. Our implementation of the *DefUseAnalysis* provides a conservative over-approximation, but will nonetheless miss any side effects that occur in library code outside the script’s text.

This algorithm requires three capabilities of a scripting language: a parser to produce a list of top-level statements; an interpreter to execute statements to produce values; and the ability to inspect values. Tempe’s implementation uses [Microsoft Roslyn](#) for the first two and standard .NET reflection for the third. This algorithm could also be implemented

for R and Python since their environments have these capabilities.

IV. CONTROLLED STUDY

To evaluate the effect of the user experience on how data scientists do exploratory data analysis, we conducted a controlled lab study. We chose a within-subjects design for two reasons. First, participants could directly compare the two user experiences. Second, this design mitigates individual variability in skill level and work pace.

A. Software

In addition to the live scripting environment described above, we added a REPL experience to Tempe, which directly submits the user’s input to the underlying language interpreter. In designing and implementing the REPL user experience, we carefully balanced two goals: providing a representative implementation of each class of user experience, while minimizing the differences between them.

Each participant completed each task in its own notebook page. Notebook pages in each condition had the same features, except for the choice of scripting experience. Within the scripting experience, both versions included familiar programming features, like code completions, method tooltips, and undo/redo. In the REPL version, these features were available for the current expression at the command prompt; in the live version, these features were available in the entire script. Both versions presented the same kinds of error messages and visualizations in response to the user’s input. In the REPL, the response appeared after hitting Enter; in the live version, a response appeared whenever the user’s editing is idle.

We modeled the REPL user experience after the one in [RStudio](#), a popular scripting environment for data analysis. In particular, we copied its key bindings, its command history semantics, and its command prompts (the `>` sign for the first line of an expression, and the `+` sign for continued lines of the same expression).

To be true to each user experience, there was a small difference in the allowed scripting input. In both experiences, the user could enter C# statements, classes and methods. In the REPL experience only, the user could also enter C# expressions. This difference is because the live experience is about producing a legal C# program; whereas, the REPL experience accepts a series of “commands”. As a simple illustration, the lines

```
2+3
4*5
```

are valid as a pair of commands, but are not syntactically valid as a C# program.

In total, then, there are three interrelated differences between the REPL and live conditions:

- The live experience updates automatically, while the REPL updates only when the user presses enter.
- Live code can be edited anywhere in the script, while the REPL is append-only.
- The live code tracks dependencies for selective re-evaluation, while REPL results did not update after initial evaluation.

B. Hardware

We conducted the study on an HP Z420 PC, with four 64-bit processors and 16 GB of memory. The PC had two 22-inch screens. Each participant conducted the tasks on the left screen. The right screen displayed reference material about the LINQ and visualization APIs.

C. Participants

The role of data scientist is still emerging at the company where we recruited participants. As a result, we could not use job title as a reliable marker for potential participants. Instead, we chose a random sample from three internal distributions lists. We sent personalized email invitations that screened for experience in C#, LINQ, and data analysis.

From this pool, 16 participants took part in the study, with one dismissed for lacking the required skills. The remaining 15 (1 female) reported an average of

- 9.0 years of professional experience;
- 5.7 years of C# experience;
- 2.8 years of LINQ experience;
- 3.9 years using scripting environments;
- 4.7 years analyzing data with plots and charts;
- 3.0 years analyzing data to help their team make better business or engineering decisions.

In terms of job roles, eight described themselves as developers, six as testers, and one as a data scientist.

D. Tasks

Each study session lasted roughly 1½ hours, in which the participant:

- signed a consent form and filled in a short questionnaire about work experience (5–10 min)
- read a short tutorial about both versions of the user experience and practiced using each one (15–20 min)
- analyzed four datasets, two using the REPL user experience and two using the live user experience (4 × 15 min per dataset = 60 min)
- filled in a questionnaire rating and comparing the user experiences (5–10 min)

The analysis tasks used five [datasets](#) downloaded from public web sites, which we selected for their size and intuitive domains:

- A table about passengers on the Titanic (practice session)
- A table about recent visitors to the White House
- A table about recent house sales in a major US metropolis

- A table about loan applications from a lending club
- Three tables with Yelp reviews, users, and rated businesses

The four non-warmup datasets had 10^5 – 10^6 rows and 13–28 columns. For each dataset, we instructed the analysts to explore and analyze the data. Our instructions read, in part, “During each play session (15 min), you’ll be given some data that we hope you’ll find interesting. Your task is to think up questions about the data and to use Tempe to try to answer those questions. You are completely free to choose whatever questions you like and however many you like.” Pilot testing confirmed that while these instructions were open-ended, they were meaningful to the analysts. Each analyst proceeded immediately on the task; we frequently had to interrupt them at the end of the fifteen-minute period.

We used counterbalancing to mitigate order effects. Each participant completed the four tasks in either the order REPL-live-REPL-live or live-REPL-live-REPL, in counter-balanced fashion. We further counterbalanced the order in which the participants used the four datasets in the tasks.

At the end of each of the four exploratory data analyses, we asked the participant to review his/her work and to mark those data results that he/she personally found insightful, interesting, or worth sharing with a friend. To do the marking, the user selected text in the query generating the result and clicked an Annotate button.

The final questionnaire contained three sections asking the participant to rate the REPL experience, to rate the live experience, and to compare the two. The REPL and the live sections each contained six standard usability questions on a 5-point Likert scale (strongly disagree to strong agree), plus two free-response questions asking for their favorite and least favorite things about that experience (Figure 11). The comparison section contained 11 questions on a 5-point Likert scale (1 expressed strong preference for REPL and 5 expressed strong preference for live), plus two free-response questions asking when they prefer REPL to live and vice versa (Figure 12).

V. RESULTS

Because Tempe is a web application, the web server keeps a database of the user’s content, including notebook pages, datasets, edits, error results, data results, and annotations. We mined this database to analyze how the participants used the REPL and live user experiences.

A. REPL sessions make frequent use of the history

We wanted to understand how often participants needed to edit and resubmit a statement before getting the desired result (or giving up and moving on). We mined the participants’ REPL edits for sequences of consecutive resubmissions. A histogram of the lengths of these sequences find that many commands were submitted once. However, it was also not uncommon for participants to re-submit a statement many times, up to a maximum of 18.

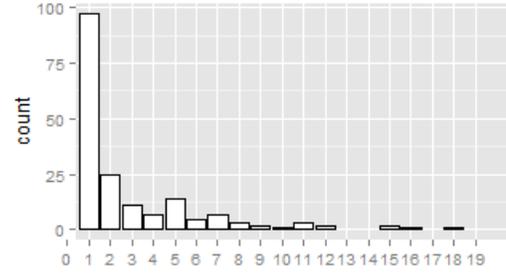


Figure 5. Number of times that REPL users re-executed statements, which was often two or more times.

One challenge in a REPL environment is the phenomenon of “up, up, up, enter,” when a user updates a variable early in a sequence, then resubmits those statements that use that variable in order to see the effect of the new value. Nine of the 15 participants re-executed previous statements, sometimes as far back as eight steps in the history.

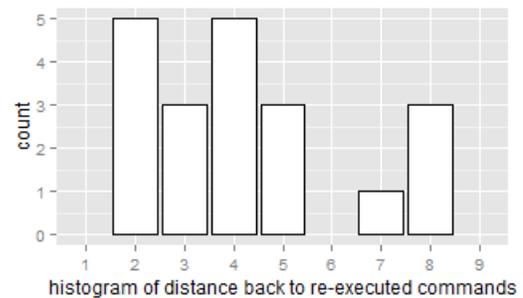


Figure 6: Distance to re-executed commands in REPL condition. On three occasions, users reached back eight steps to re-execute a command.

There were also 21 instances where participants re-executed their previous command (distance=1), which do not represent re-executions to deal with changed variables. After inspecting several of these cases, these seem to be occasions where the participant did not believe an error message until seeing it twice.

B. Edits in live sessions were often spread across the script

All participants had previous experience using REPLs ($M=3.9$ years), and none had experience using a live environment. This creates the potential for them to treat the live experience as though it were a REPL, by consistently adding new statements to the bottom of the script.

To test this, we looked at the line numbers of consecutive edits. In a REPL session, the first edit is on line 1 of the script, the second is on line 2, and so on, which means that the difference in line numbers between consecutive edits is always 1. In the live experience, the difference in line numbers between consecutive edits is typically zero (many edits on the same line), but ranged from 10 lines above the previous edit to 17 line below. All participants except P11 edited the script two or more lines away from the previous edit. Most participants did it several times:

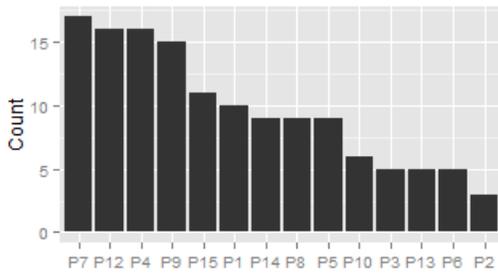


Figure 7. Number of times each participant edited at least one statement two or more lines away from their previous edit in the live condition. All participants edited in a variety of places around the script.

This demonstrates that all but one participant took advantage of the live programming environment to edit the live script as a document, rather than treating it like a REPL.

[I like] how easy it was to make small changes to existing queries [P1]

Quick updates to the visualizations as I tweak my ideas. [P9]

[REPL makes it] hard to edit a chain of commands. I needed to retype everything [P5]

Many participants particularly appreciated the ripple effect of re-evaluation.

[I like] the waterfall effect of editing a small value on the top of a chain of queries. [P5]

It is easy to build on data and evolve queries as you progress through your thinking. [P14]

However, some participants were concerned about the performance cost of re-evaluation.

I really like how the live version automatically reevaluates data, but this might become a bottleneck if computation takes a lot of time or happens on server [P13]

C. REPL sessions leave a history of errors

In the REPL experience, participants often corrected errors iteratively over many commands, leaving an unwanted history of errors. In the live experience, participants similarly iterated to correct errors, but the live experience removes obsolete errors. For many participants, this was a key difference between the experiences.

Queries that are genuine errors (syntax errors e.g.) don't need to stay in history forever. [P12]

You can fix mistakes quickly. [P1]

It hides all the mistakes I made along the way, showing only the final results. [P3]

Much easier to go back and tweak when I made a mistake or changed my mind. [P10]

To measure the difference, we can inspect each session's final notebook page and measure the fraction of responses that were errors versus non-error outputs.

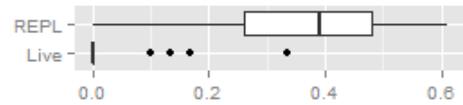


Figure 8. Fraction of error responses per session in the final session output. In the live condition, users corrected almost all errors.

These error fractions are significantly less in the live experience than the REPL (Wilcoxon signed rank, $p=4.4 \times 10^{-6}$). For REPL sessions, up to half the responses on the final page were errors; whereas for live sessions, almost all the errors were corrected and removed.

On the other hand, some participants felt that the live experience obliged them to correct mistakes, even when they felt it was unnecessary.

[REPL had] less state associated with each query, so one is not maintaining code after it has been written. [P14]

Although participants could remove unwanted erroneous code in the live experience, this requires additional effort; whereas, in the REPL experience, “what’s done is done.”

D. REPL sessions retain more results than live sessions

For a data scientist, the goal of exploratory data analysis is to gain insights into the data. We operationalize this by looking at the number of non-error responses participants produced per session and the number that they marked as insightful.

A direct comparison between the numbers of responses produced in each user experience is not meaningful, since the REPL produces a response only when the user types Enter and the live experience potentially produces a response for each burst of user input. Indeed, the mean number of responses for REPL sessions is 13.2 (SD=5.9) and for live sessions is 33.5 (SD=16.8). Instead, we can look at the *final* number of responses per session, which eliminates all the intermediate responses that the live experience produces.

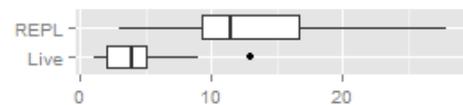


Figure 9: Count of final, non-error responses per session by condition

The number of final non-error responses per session is significantly less in live experience than the REPL experience (Wilcoxon signed rank, $p=9.9 \times 10^{-7}$). However, this comparison offers REPL an advantage, since REPL sessions retain all responses, even intermediate or unwanted ones.

E. REPL and live sessions produce as many responses marked insightful

If we look instead at the number of responses that participants marked as insightful in each session, the distributions between the REPL and live conditions are not significantly different (Wilcoxon signed rank, $p=.17$).

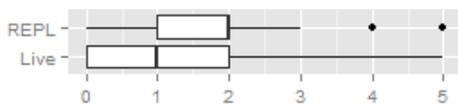


Figure 10: Number of responses marked "insightful" per session

This measure might undercount the insightful responses in live sessions, since participants could not mark responses that had been overwritten during live sessions. However, no participants mentioned this issue, and many chose to mark throughout the session.

F. Subjective Assessments

Based on participants' subjective assessments of both the REPL and live experiences, they enjoyed both (Figure 11). Given the participants' considerable experience with REPLs, it is not surprising that they gave the REPL experience high

marks for ease of learning and ease of use. More surprisingly, however, less than half agreed that the REPL experience makes them effective or productive at analyzing data. The live experience received high marks across the board. Almost everyone agreed that the live experience is easy and productive to use and should be recommended to a friend.

The second half of the subjective assessment asked participants to compare the two experiences on a custom Likert scale:

- 1 REPL is definitely better.
- 2 REPL is a bit better.
- 3 Both versions are about the same; or, sometimes one is better, sometimes the other.
- 4 Live is a bit better.
- 5 Live is definitely better.

The live experience was clearly preferred on almost all items (Figure 12). The two items where the REPL experience had the upper hand were ease of learning and ease of use. However, this preference may simply reflect the participant's prior experience with REPL environments, given that they

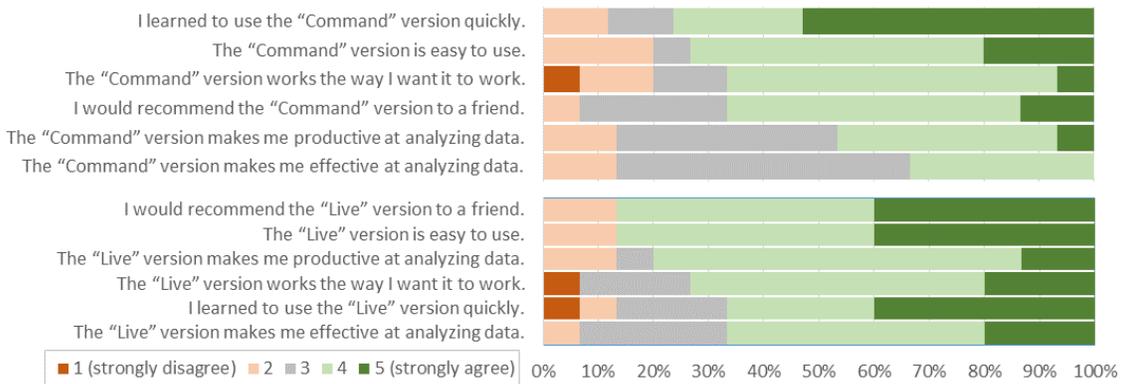


Figure 11. The participants' subjective assessments of the two user experiences, shown as the fraction of participants who chose each of five Likert-scale responses. The items are sorted by percent agreement. (The REPL experience was called "Command" in the participant materials, to avoid jargon.)

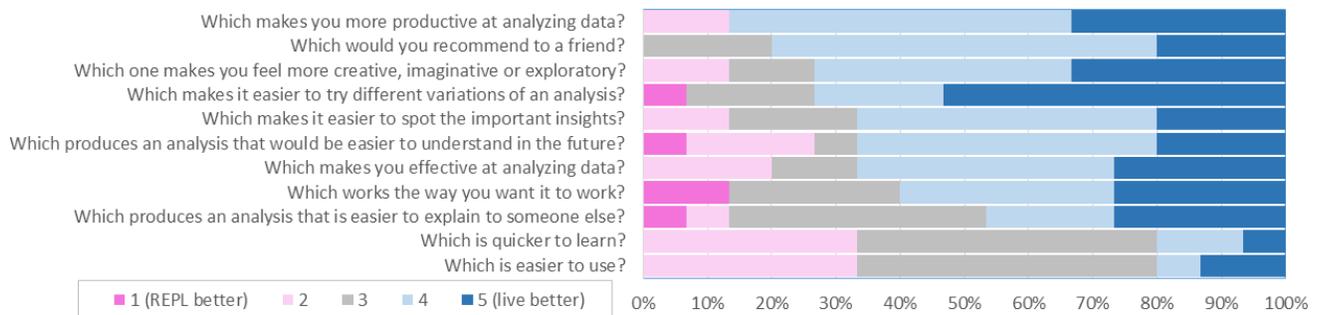


Figure 12. The participants' subjective comparison of the two experiences, shown as the fraction of participants who chose each of five Likert-scale responses. The items are sorted by the percent who favored the live experience. While the REPL version seems quicker to learn, the live version seems to be stronger for exploration, iteration, spotting insights, and productivity.

separately gave the live experience high marks on ease of learning and ease of use.

VI. DISCUSSION

A. Live programming hides history

Several participants remarked that keeping history was their favorite aspect of the REPL experience.

[The REPL keeps] the history/footprint of what I have tried [P5]

[The REPL] keeps all the history. [P11]

The lack of history was a major downside of the live experience.

It doesn't remember the history, so I need to be explicit on when to save something in a variable, just to see it. [P3]

Easy to get carried away with modifying current query and intermediate data is lost. [P12]

Ironically, the web server maintains a complete history of all edits and responses (which we mined for our data analysis), but the live experience lacks a user interface for surfacing this history. In designing a new feature for browsing and resurrecting live programming history, we hope to take advantage of recent work by Yoon *et al.* [21] [22]

B. Live programming helps convey intent

In choosing a notebook model for our environment, we want to encourage users to share and preserve their analyses. This is also the reason for the rich-text annotations that participants used to mark their insightful results. Participants felt that the live experience itself encouraged them to be intentional about the contents of the page.

The ability to modify line by line is much more effective—it allows me better control over what data points to visualize and helps bridge the gap between exploratory & prod[uction] better. [P9]

It lets me design what I want to see. [P11]

No useless query remains. [P12]

For the live experience, this is the flip side of lacking a complete history: the user can curate what computations are worth preserving.

VII. LIMITATIONS

Conducting a controlled lab study required some compromises that affect external validity. In a professional context, exploratory data analyses typically last longer than 15 minutes and use datasets that are larger, messier and more complex than those we chose for the lab. Nonetheless, our datasets had thousands to millions of rows, with dozens of columns. They were sufficiently complex that participants pursued different questions on the same dataset.

Using C# as the scripting language, while a good fit for our participants' skills, is unusual for exploratory data analysis. Most data scientists use languages like R, IPython, and MATLAB. Our live programming algorithm and user experience are equally applicable to these languages, but without further studies we cannot conclude that our results generalize to other languages.

Our participants were recruited from a single company, although from different businesses within the company. The participants may share common training or culture that biases the results.

VIII. CONCLUSIONS

In our controlled lab study, we saw participants exercise the core behaviors associated with the REPL and live user experiences. In the REPL experience, participants repeatedly edited and resubmitted statements to reach desired results, leaving histories with many errors and data results they considered uninteresting. They also manually resubmitted previous statements to bring their results up to date. In the live experience, participants corrected errors in place and made some edits spread across the script, allowing them to experience the ripple of automatic updates.

Comparing the two experiences, REPL sessions produced both more errors and more results than live session. In REPL sessions, the errors amounted to almost half of the total responses on the page. REPL sessions also left a history of more results than live sessions. However, the numbers of results that participants marked as insightful were not significantly different in the two experiences. The difference in the raw number of results may be because the live experience allowed participants to overwrite uninteresting intermediate results.

Participants strongly preferred the live experience, particularly for its error correction and its control over script content. Nonetheless, they liked a REPL's retention of complete history. This suggests that the live experience would benefit from a history browsing mechanism, which might help achieve the best of both worlds.

IX. REFERENCES

- [1] D. Patil and T. Davenport, "Data Scientist: The Sexiest Job of the 21st Century.," *Harvard Business Review*, October 2012.
- [2] D. Fisher, R. DeLine, M. Czerwinski and S. Drucker, "Interactions with big data analytics," *interactions*, vol. 19, no. 3, pp. 50-59, May/June 2012.
- [3] S. Kandel, A. Paepcke, J. M. Hellerstein and J. Heer, "Enterprise Data Analysis and Visualization: An Interview Study," in *IEEE Visual Analytics Science & Technology (VAST)*, 2012.
- [4] R. B. Smith, "The Alternate Reality Kit: An animated environment for creating interactive simulations," in *IEEE Computer Society Workshop on Visual Languages*, 1986.

- [5] S. L. Tanimoto, "VIVA: A visual language for image processing," *Journal of Visual Languages and Computing*, vol. 1, no. 2, p. 127–139, June 1990.
- [6] A. R. Brown and A. Sorensen, "Interacting with Generative Music through Live Coding," *Contemporary Music Review*, vol. 28, no. 1, pp. 17–29, February 2009.
- [7] S. Aaron, A. F. Blackwell, R. Hoadley and T. Regan, "A principled approach to developing new languages for live coding," in *Proceedings of New Interfaces for Musical Expression*, 2011.
- [8] C. D. Hundhausen and J. L. Brown, "What You See Is What You Code: A "Live" Algorithm Development and Visualization Environment for Novice Learners," *Journal of Visual Languages and Computing*, vol. 18, no. 1, pp. 22–47, 2007.
- [9] C. M. Hancock, *Real-Time Programming and the Big Ideas of Computational Literacy*, Massachusetts Institute of Technology, 2003.
- [10] A. Repenning, "AgentSheets: An interactive simulation environment with end-user programmable agents," *Interactions*, 2000.
- [11] S. McDirmid, "Usable Live Programming," in *SPLASH Onward!*, 2013.
- [12] T. Lieber, J. R. Brandt and R. C. Miller, "Addressing Misconceptions About Code with Always-On Programming Visualizations," in *CHI*, Toronto, 2014.
- [13] F. Pérez and B. E. Granger, "IPython: A System for Interactive Scientific Computing," *Computing in Science and Engineering*, vol. 9, no. 3, pp. 21–29, May/June 2007.
- [14] B. Victor, *Inventing on principle*, Invited talk at the Canadian University Software Engineering Conference (CUSEC). <http://vimeo.com/36579366>, 2012.
- [15] R. DeLine, D. Fisher, B. Chandramouli, J. Goldstein, M. Barnett, J. F. Terwilliger and J. Wernsing, "Tempe: Live Scripting for Live Data," *Proc. of IEEE Symp. on Visual Languages and Human-Centric Computing*, 2015.
- [16] M. Barnett, B. Chandramouli, R. DeLine, S. Drucker, D. Fisher, J. Goldstein, J. Platt and P. Morrison, "Stat! - An Interactive Analytics Environment for Big Data," in *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 2013.
- [17] S. Wolfram, *The Mathematica Book*, Wolfram Media, Inc., 1996.
- [18] E. Meijer, B. Beckman and G. Bierman, "LINQ: reconciling object, relations and XML in the .NET framework," in *Proceedings of the 2006 ACM SIGMOD international conference on Management of data*, 2006.
- [19] R. Wagner and M. Fischer, "The String-to-String Correction Problem," *Journal of the ACM*, vol. 21, no. 1, p. 168–173, 1974.
- [20] A. V. Aho, R. Sethi and J. D. Ullman, *Compilers: Principles, Techniques and Tools*, Addison-Wesley Publishing, 1986.
- [21] Y. Yoon and B. A. Myers, "A Longitudinal Study of Programmers' Backtracking," in *IEEE Symposium on Visual Languages and Human-Centric Computing*, 2014.
- [22] Y. Yoon, B. A. Myers and S. Koo, "Visualization of Fine-Grained Code Change History," in *IEEE Symposium on Visual Languages and Human-Centric Computing*, 2013.