Strategies for Reuse and Sharing among Data Scientists in Software Teams

Will Epperson willepp@cmu.edu Carnegie Mellon University Pittsburgh, Pennsylvania, USA

Robert DeLine rdeline@microsoft.com Microsoft Research Redmond, Washington, USA

on, USA

Hallmark: a distinguishing

ABSTRACT

Effective sharing and reuse practices have long been hallmarks of proficient software engineering. Yet the exploratory nature of data science presents new challenges and opportunities to support sharing and reuse of analysis code. To better understand current practices, we conducted interviews (N=17) and a survey (N=132) with data scientists at Microsoft, and extract five commonly used strategies for sharing and reuse of past work: personal analysis reuse, personal utility libraries, team shared analysis code, team shared template notebooks, and team shared libraries. We also identify factors that encourage or discourage data scientists from sharing and reusing. Our participants described obstacles to reuse and sharing including a lack of incentives to create shared code, difficulties in making data science code modular, and a lack of tool interoperability. We discuss how future tools might help meet these needs.

KEYWORDS

Data Science, Code Reuse, Code Sharing, Survey

ACM Reference Format:

Will Epperson, April Yi Wang, Robert DeLine, and Steven M. Drucker. 2022. Strategies for Reuse and Sharing among Data Scientists in Software Teams. In 44nd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP '22), May 21–29, 2022, Pittsburgh, PA, USA. ACM, New York, NY, USA, 10 pages. https://doi.org/10.1145/3510457.3513042

1 INTRODUCTION

"da maneira que achar melhor" -

As software engineering developed into a mature discipline, the ability to effectively share and reuse code has become a critical factor for success [9, 11, 16, 17]. Particularly as organizations grow, information management becomes both more difficult and more important. This information takes the form of actual source code but also the documentation for this code, specifications around the problem the code was initially created to solve, and who to talk to in an organization to learn more about the code. Well executed

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

For all other uses, contact the owner/author(s). ICSE-SEIP '22, May 21–29, 2022, Pittsburgh, PA, USA © 2022 Copyright held by the owner/author(s). ACM ISBN 978-1-4503-9226-6/22/05. https://doi.org/10.1145/3510457.3513042 April Yi Wang aprilww@umich.edu The University of Michigan Ann Arbor, Michigan, USA

Steven M. Drucker sdrucker@microsoft.com Microsoft Research Redmond, Washington, USA

software reuse leads to fewer problems in code, less effort spent correcting problems, and higher developer productivity [21]. Technologies such as version control [10] have become commonplace to help keep track of versions between files and principles like "DRY" (Don't Repeat Yourself) are baked into software developers' minds.

However, the field of data science presents new and unique challenges in terms of sharing and reuse. The *people* writing the code come from different backgrounds, the *code* itself lives in a variety of formats including raw text files, computational notebooks, and ad hoc queries, and *data* permeates the entire analysis process.

Data scientists are a relatively recent role on software development teams, and they work alongside established roles like software developers, operations, and program managers [4, 14, 15, 18, 19]. In the industrial setting, data scientists are often not directly responsible for the data they analyze. Their partners, the data engineers, collect, store, and maintain datasets that data scientists access for their analysis [14, 15]. For teams whose services use machine learning (ML), data engineers also deploy, scale out, and maintain ML models that data scientists create [2]. Further, some data scientists work on their own team's data, some act as a centralized service working with several product teams, and some act as consultants working with third-party companies.

A data scientist's work is often exploratory or ad hoc in nature and involves using data to craft analyses, models, and visualizations that are reported outside of the coding environment [15]. At Microsoft, each data science team is free to go about this process however they see fit which leads to a wide variety of approaches. To scale model and inference pipelines into production, this process is handed off to adjacent data engineers.

This unique role for data scientists has led to new approaches and problems for code reuse and sharing. The exploratory, open-ended nature of data science coding impacts the incentives for investing time into reuse. For instance, if code is only used to answer a one-off analysis question, there is less incentive to invest the time into making this a reusable function. Furthermore, reusable code for data analysis must support customization and adaptation since each data analysis is slightly unique.

In this work, we specifically focus on the reuse of analysis *code* as it relates to data science work. This code is almost always coupled with data assets for the analysis, however we consider the versioning and reuse of data itself outside the scope of our study. Decisions about the team's data management are often made at the team level

See/Think fit: to choose to do something because one thinks it is appropriate

and are subject to corporate policies, customer agreements, laws and regulations. Data scientists typically enjoy more agency over the data analysis code than over the data itself. Hence, we focus on their work practices and pain points around the data analysis code as a topic where the research community can usefully intervene. Additionally, we focus on the reuse of code developed within teams or organizations rather than external library use.

By "reuse" we refer specifically to the consumption of code or other artifacts from previous work. This might be the author of that work reusing their own past analysis or that of someone else. On the other hand, by "sharing" we refer to the production of code or other artifacts for oneself and then providing it to someone else for another task. The acts of sharing and reuse might be viewed as two sides of the same coin; unique practices exist for both facets.

To provide a better understanding of how data scientists go about both reusing and sharing past work, we conducted interviews with professional data scientists to understand their current practices related to sharing and reuse of analysis. From these interviews we synthesized five different strategies for sharing and reuse that we developed into a survey to understand how these approaches generalize across a larger population. In the remainder of this paper, we first present related work regarding sharing and reuse in data science, followed by a discussion of our study methodology. We then discuss the strategies for reuse and sharing generated from our interviews and survey, along with determinants of reuse that encourage or discourage sharing and reuse among data scientists. Lastly, we discuss opportunities for future work to address unmet needs for sharing and reuse in the practice of data science. In summary, our contributions are:

- Interviews and then survey
- (1) Based on an interview study and survey with 149 professional data scientists, we characterize five primary strategies for sharing and reuse of analysis code.
- (2) We report our participants' determinants and obstacles for sharing and reuse practices and discuss implications for future tools.

2 BACKGROUND AND RELATED WORK

There are primarily two relevant areas of interest to this work. The first explores the practices of sharing and reuse in traditional software engineering. The second investigates common themes among exploratory data programming.

Reusing and Sharing in Software Engineering. Reusing and sharing code benefits software developers by saving their time and resources to build and maintain applications, while maintaining code simplicity [9, 11, 16, 17]. Effective reuse relies on concise and expressive abstractions [17]. Using various form of abstractions, common strategies for software reuse include high-level languages, ad hoc code scavenging, and source code components such as libraries [17]. In particular, with the growth of open source software, library reuse has become a prevalent practice in software engineering [1, 12, 23]. Library repositories like NPM and PyPI have facilitated the sharing, discovering, and management of third-party libraries, which lead to the increasing usage of third-party libraries among software developers [1, 23]. This huge demand further incentivizes the creation and implementation of third-party libraries. However, library reuse has its limitations. Xu et al. found that

developers would replace an external library with their own implementation if the library is over complicated or not flexible to satisfy their needs [26]. In data science, code is less formal compared to traditional software engineering [18]. Although data scientists generally benefit from libraries for performing common data operations and computations, these libraries tend to be low level. The practice of sharing and reusing entire analyses or workflows in data science remains unexplored and worth investigating. Thus, our work aims to reveal the reuse and sharing practice in data science programming, understand the different reuse decisions between traditional software engineering and data science, and identify design opportunities for facilitating reuse to improve work efficiency.

The Process of Data Science. Several researchers have investigated what steps data scientists go through in their work and how they seek to coordinate their efforts. Machine learning product development involves iterations of a process that begins with gathering model requirements and ends with model deployment and monitoring [2]. Throughout this process, data scientists, domain experts, team leaders and software engineers collaborate in unique roles as indicated by tool use (technical vs non-technical users) [27]. In this process, Jupyter notebook users tend to think less about future use of their code even though notebooks are touted as a more readable platform [27].

Several large scale reviews of computational notebooks have documented that despite the benefits of computational notebooks, they can encourage bad coding practices because of unexpected execution order and a lack of modular code [22, 24]. For example, Rule et al discovered that nearly half of the 1 million Jupyter notebooks they scraped from Github were uploaded with non-linear execution orders [24]. These notebooks serve a variety of purposes like data exploration, places to store code, or for ad hoc versioning of analyses. However, analysts must invest time and effort to clean up their notebook to make it reproducible or reusable by others. Our investigation is distinct in that we take a broader look at how data scientists specifically reuse and share their work across all platforms.

Tools for Data Science Workflow Management. Several platforms have been developed to help data scientists share code while programming. Git is a widely used version control system that lets users manage version of raw files for software projects [10]. However git does not work well for comparing versions of rich text files like computational notebooks. Systems like Verdant augment notebooks to better support versioning by tracking a user's analysis history [13]. Our work discusses unmet needs for sharing and reuse in data science that future tools might address.

3 METHODS

To better understand the state of the art in sharing and reuse practices in data science today, we conduced semi-structured interviews with 17 data scientists at Microsoft. To ensure the generality of our findings, we then developed a survey that was completed by an additional 132 data scientists. Both studies were approved by an Internal Review Board, and all participants signed consent forms. Interview participants were each compensated with \$25 USD gift card. Survey participants were entered into a raffle for three \$100 USD gift cards.

Interview Study 3.1

For our interviews we recruited 17 participants from a pool of data scientists chosen at random from the employee database, based on their job titles, levels, and business units. In particular, we recruited data scientists from software product or service teams and excluded those from non-product units like Research and Legal. We conducted the interviews in one-hour sessions where we asked about the participants' sharing and reuse practices as an individual and as a member of their team. These interviews were transcribed and then analyzed for themes. Specifically, two of the authors reread the transcripts, coded them, and card sorted the codes to come up with themes. The authors then discussed the themes until a consensus was reached. The interviews revealed both common strategies for reuse that are discussed more in depth in Section 6 as well as factors that encourage and discourage reuse discussed in Section 7. Throughout this paper, interview participants are called informants, and their quotes are designated with "IP".

3.2 Survey

To assess the generalizability of our identified themes, we ran a survey with data scientists drawn at random from the same pool as the interviews. In total, 132 participants filled out our survey, out of 563 invited (23% response rate). The survey asked for the following information:

- Background about the respondent's years of experience, tool usage, and team size;
- For each of the five strategies:
 - Frequency of reuse;
 - How the reuse happens (e.g. copy-paste, function call, etc.);
 - What functionality is reused;
 - How reused work is found;
 - Frequency of sharing;
 - Additional work required for sharing
- Influences on their willingness to do reuse and sharing
- Best practices and pain points for reuse and sharing

The survey took 10-15 minutes to complete. Throughout this paper, survey participants are called respondents, and their quotes designated with "SP".

PARTICIPANT BACKGROUNDS

Our informants (9 male, 8 female) had experience ranging from 2-14 years of professional data analysis work. All informants were from different teams. They had a range of educational backgrounds, including fields like statistics, math, or computer science. Our respondents (85 male, 46 female, 1 did not say) reported a range of professional experience in data science from 3 months to 35 years, with a mean of 8.6 years.

Context. Our study is conducted at Microsoft, a large software corporation. At Microsoft, data scientists work both on dedicated teams and also within software engineering teams. There exists little standardization across teams mandating how data scientists must go about their work. This leads to a wide variety of tool usage and reuse strategies.

Team Composition. The informants are either members of dedicated data science teams, situated within larger product teams or work on product teams alongside software engineers. The survey

respondents report team sizes ranging from 1-22 people, with a mean of 8 people (after removing outlier responses). However even on teams with many data scientists, most projects involve only one or two of them: 74% of respondents work either alone or in pairs on projects. Analysts communicate with team members throughout their project to communicate results, seek help, and to find code for reuse (see Section 6 for details), however most projects only involve one or two data scientists actually touching the code or data. This is starkly different from large software projects within Microsoft that may have dozens of developers iterating on a single, interdependent, code base.

As noted in prior literature, much of the work done by data scientists is exploratory in nature and their coding practices reflect this. Data scientists will often begin an analysis, but if it turns into something recurring or the model they built must be scaled up for production deployment, the code will be handed over to an adjacent software engineering team. Several of our informants expressed admiration for the high quality code produced on production teams and claimed it was far superior to their own in terms of organization, commenting, and style. One of our interview participants (IP6) noted that the software engineers they work with "write pretty fabulous code in terms of readability and actual usability".

Tools shape sharing and reuse. The interview participants use a variety of tools to do their work. The most popular languages used were Python, R, a query language for large-scale relational queries (Cosmos/SCOPE), and a query language for large-scale telemetry data (Kusto Query Language/KQL [20]). Analysis was done in a combination of computational notebooks (Jupyter, Visual Studio Code Notebooks, and Databricks) and integrated development environments. However, even data scientists on the same team rarely used the exact same tool stack. The tool chosen to write code is often shaped by how and where the data is stored, which in turn shapes the sharing practices afforded by the tool.

For example, several teams used Databricks for writing their analyses [7]. Databricks allows code to be written in a variety of languages including Python and R in a computational notebook interface. Teams using these tools developed reuse strategies around notebooks such as template notebooks and notebook libraries (Section 6.4). Yet other teams that do most of their analysis using KQL developed strategies to share their work through the tool where they write their queries.

Lastly, many participants noted that the final output of their work is often not the code itself. This in turn shapes both the practice of reuse, as well as the incentives to invest in creating reusable code. Results are presented as PowerPoint presentations, text documents, or dashboards. This finding is consistent with prior literature about how artifacts of data analysis are shared outside of the computing environment [27]. The consumers of these results seldom ask for rarely the code that generated the results they are seeing. This is distinct from traditional software engineering where the outcome is the code. For data science, the outcome is the *result* of the analysis.

WHAT FUNCTIONALITY IS REUSED?

Informants mentioned a variety of tasks for which they share and reuse code. In this section, we discuss these tasks in aggregate across all reuse strategies, from most to least common task.

Data preprocessing, transformation, cleaning. The most commonly reused data science code is for data preprocessing. This includes cleaning data, transforming data between formats, and generally processing data into a usable state before analysis can begin. Since data analysis pipelines often start with raw, uncleaned sources, reusing code that cleans and formats data for analysis helps speed up this process and ensures that the cleaning is done consistently. Reusable cleaning code can also encapsulate idiosyncratic details of how data is represented in different data sources.

Reading and writing data. The second most common reuse task was reading and writing data from sources. Since many data sets are too large to fit locally on an analyst's machine, data scientists frequently read and write data (samples) from the same data stores or submit jobs to run on cloud servers. Working in the cloud often involves configuration details like access keys and resource IDs, which are hard to remember and therefore often copy-pasted from previous work. Informants also reuse queries that are shared within a team for commonly accessed data or to make sure different analyses are looking at the same slice of the data. For example, IP5 mentioned that engineering teams often share a query to pull anomalous data that needs further analysis: "we interact with the engineers on the team who can help us understand the telemetry better and identify what is the right query to use for certain things."

Modelling and evaluation. Code for creating and evaluating models is also frequently reused. Data scientists are often creating very similar models since their data is semantically similar over time. For example, some teams only do time series forecasting and thus use models appropriate for time series data. Other teams deal with natural language text data and so use state of the art deep neural nets for language processing tasks. This reusable modelling code can be individual snippets or entire modelling pipelines.

ate similar visualization and reporting. Our informants often create similar visualizations during the course of a project and reuse these across projects. They reuse visualization code when the data is similar across projects, when the task is similar (for example, showing model performance), or to reuse the work of finding the desired visualization parameters. Furthermore, several informants mentioned that they like to maintain a similar style across their charts so reuse common code to do this styling, depending on which visualization library they are using. Informants also create dashboard templates for tools like Power BI to save time and maintain consistency.

Are the last two highlights connected?

Miscellaneous tasks. Various other tasks were mentioned with less frequency, but still offer some breadth as to what data scientists view as worth sharing and reusing. For instance, some teams focusing on ML model development maintain library functions for postprocessing, model selection, and model ensembling. Others reuse code that runs pipelines for analysis, composing many of the aforementioned steps together. These pipelines allow data scientists to iterate faster on models by automating tedious steps such as hyperparameter selection or data formatting.

6 APPROACHES TO REUSE AND SHARING

Since each informant works on a different team, each reported a slightly different approach to sharing and reuse in their work, shaped by the kinds of data they used and their team's sharing structure. We thematically grouped their sharing and reuse strategies into five distinct strategies. Some strategies are personal, namely copying previous work (6.1) and keeping a personal utility library (6.2). Other strategies are team-wide efforts: sharing notebooks (6.3), creating template notebooks (6.4), and developing shared libraries (6.5). We describe these strategies in increasing order of effort required to produce a shared artifact. These strategies are not mutually exclusive. For instance, some participants used both a personal utility library and template notebooks. The modal survey respondent used three distinct strategies at least yearly, and two different strategies at least monthly. Every survey respondent used at least one of these strategies to share and reuse their code.

6.1 Personal analysis reuse

The most basic form of reuse is to look at one's own past analyses as a reference for current work. This code is only maintained by a single individual and is not shared. Nearly every respondent participated in this form of reuse at least sometimes (97%). Local analysis code is reused for all of the tasks mentioned in Section 5, from data processing to model evaluation. As one respondent (SP63) described, "Anything I know I did before, I reuse it."

6.1.1 Access and Search. Data scientists find prior work almost exclusively by memory since they are the primary author of these files. Furthermore, there is a range in which local analysis files are tracked through version control. Some informants put all of their analysis in a (personal) Git repository. Though, it is of note that this is primarily so that the code files can be accessed on another machine rather than to track versions of files. Others leave their analyses on the local file system. For storage and easy access (personal and team)

The majority of analysis code stored locally is reused via copy and paste. A common case is cloning code from one computational notebook to another, then adjusting the code to suit the current context. When the reused code is a series of top-level statements rather than class or function definitions, it cannot be imported and called as an API. (Turning the reused code into a library is discussed in the next section.) In our survey, the next most common way of reusing this code was simply looking at it for reference and the least common was importing the code.

6.1.2 Additions. As data scientists work on a variety of projects over time, they store their analysis code. Often, this takes the form of messy analysis notebooks or scratch files to produce a certain result for an analysis. These notebooks and files are not cleaned before storage by adding extra comments or documentation as participants do not necessarily anticipate reusing them in the future. In this sense, reusing local analysis code is opportunistic reuse – when the code is initially written the data scientist is not sure if they will use it again in the future and does so when the time arises. This also means reusing in this way requires little effort since minimal prior planning is required.

6.1.3 Benefits and Drawbacks. The benefit, and drawback, of this strategy is its ease – all it takes is putting files in some form of file structure and then finding relevant files later. Projects that begin in local file storage might later make it into a different form of reuse if they need to be shared with others. This low effort

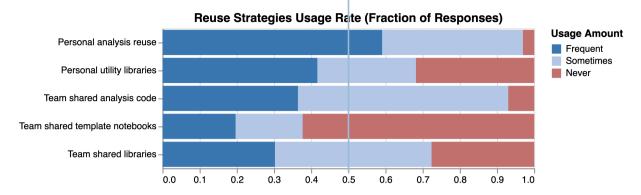


Figure 1: We bin usage rate into three categories: daily or weekly usage as *Frequent*, monthly, seasonal, or yearly usage as *Sometimes*, and *Never* for respondents who never use a strategy. Using personal past work is the most frequently used strategy, and notebook templates the least frequent.

setup can lead to issues like multiple saved versions of the same file or the same analysis with incremented suffixes ("analysis-v1", "analysis-v2", "analysis-v3", …). However, since data scientists are searching over code that they wrote they generally find it easy to find past analysis stored in this form. Several interview informants intentionally organized previous analyses by project or data store to aid in future search.

6.2 Personal utility libraries

Similar to local analysis code, some participants have developed a personal "utility library" of common functions or code snippets they find themselves repeating. This utility library is only used and maintained by a single individual. However, this strategy is distinct from the local folders of analysis code in that code is intentionally cleaned up before being stored in the library. The code stored in personal utility libraries is often short (no more than a dozen lines), restructured into a callable API (class and function definitions), and parameterized so it can be used in new contexts. This code might be used for all of the different tasks mentioned in Section 5 from data cleaning to visualization.

6.2.1 Access and Search. Since this code has been restructured into a callable API, the most common way respondents use this code is by importing or calling it directly. However, copy and pasting from the personal utility library is still very common. For such small snippets of code, it can be just as easy to copy and paste into the development environment such as the current analysis notebook. In a similar fashion to personal analysis code, most functions in a personal utility library are found by memory since the author has an intimate knowledge of the code contained therein.

6.2.2 Additions. Data scientists add new things to their personal utility library when they notice they repeat a task often enough it is worth cleaning up into a reusable function. Over time, this evolves into a utility library of assorted functions and templates for various tasks:

These files are just calls which I have been using for three or five years now. I just constantly go to them again and again and again. So, I extracted these very, very generic common

things in single repo called utils that each file just does a single thing. - IP12

6.2.3 Benefits and Drawbacks. Personal utility libraries benefit data scientists by increasing productivity for common tasks. In particular, the search cost for finding code is very low since data scientists know what is in the library by memory. The only real cost of maintaining a personal utility library is the time investment to create reusable components for common tasks. However, since these libraries are only made for personal use there is little risk that this work will go to waste.

What about maintenance and keeping dependencies up to date?

6.3 Team shared analysis code

The previous two strategies only involved reuse at the personal level. However data scientists also share previous analyses among their team to find how others have analyzed similar data or to get help with tricky bits of code. This shared code often starts in a local analysis folder and then particularly useful parts of an analysis are extracted out and shared to the central, team-wide store. Team sharing also encompasses more ad hoc sharing between team mates over email or messages. Using and sharing code with teammates is very common – over 90% of our survey respondents use this strategy at least yearly and 63% at least monthly.

6.3.1 Access and Search. Team-wide sharing can take many forms, however the most typical is a shared code repository on Github or a similar version-control platform. Other storage strategies discussed included copying code to a file share (without version control) or attaching code to a team's shared documentation, like a team wiki. The code on shared Git repositories includes both computational notebooks and source files in analysis languages like Python or R. Our informants described using Git as a global file store rather than using it for version control. To this end, different teams adopted different ways of organizing the shared files to make it more conducive to reuse. One team had a separate folder for each person on their shared repository where they pushed their past analyses. Another organized the files by business area and data source so that future data scientists looking at the same or similar data can find past analyses. Sometimes this structure developed organically; other teams were very intentional about how they stored this shared code

Different file structures

to maximize the opportunities for reuse. For example, one team hired a dedicated "information architect" to develop a structure for them to organize and share their past analyses after the team size grew and experienced some turnover and it became increasingly difficult to find relevant analyses.

Ad-hoc sharing of analysis code between individual teammates has less structure than when code is stored on a central team-wide store. This type of sharing involves simply sending one of the files stored locally for *personal analysis reuse* to a teammate. Whereas code posted to git is a complete analysis, ad hoc sharing among team mates might also be used for debugging purposes to get help on a certain part of an analysis or query.

The most commonly mentioned way that our survey respondents reused this shared code was by looking at it as reference material. However, they also copy and pasted the code and even imported parts of these previous analyses to reuse. To find relevant past code, data scientists most commonly ask their teammates if something exists (70% of responses). They might post a question to a team chat channel or message someone directly who they think has an analysis they are searching for. Social searching in this manner is far more frequent than the next ranked options of reading team-wide documentation (13%) or searching over the shared code directly by keyword (10%). Data scientists may not know if the code on the shared repository is up to date or if teammates have other relevant analyses that have not been uploaded to the shared repository.

6.3.2 Additions. Code shared to a central analysis store is almost always cleaned before being shared. This includes adding more comments and documentation, restructuring the code to make it easier to understand, and generally making the code better *styled*. These cleaning steps take a time investment, and thus data scientists want to make sure this investment is worthwhile. If someone has explicitly asked for a past analysis, they will put it on the shared Git or share directly over email. Sometimes data scientists will proactively post an analysis if they think someone might benefit from it in the future. Some teams have a review process for code shared to the central analysis repository. This is often to make sure the code is clean and error free for others to use it in the future and to make sure it is a sound analysis.

When code is only shared with a single individual in a more ad hoc manner, our informants described less pressure to clean the analysis code before sharing. This also allows sharing to happen more quickly.

Store as a more generic term than repository

6.3.3 Benefits and Drawbacks. Although the shared analysis store may not always be up to date, it offers a far better alternative than no central place for storing analysis. Without a central, shared, location work is very difficult to find and often ends up getting repeated, leading to frustration. When describing a previous work experience where there was no structured way of sharing analysis code, IP16 commented:

Oh yes, tons of work was repeated in many places. Because if you weren't a part of the team that had access to some shared drive somewhere, then you didn't know that [work] existed. If you didn't have access to this one instance that was still running on the desktop of some guy who left three years ago, you didn't know that [work] existed.

There exists a division in how often data scientists *reuse* previously shared code versus *contribute* their own shared code. The most typical survey respondent reused shared team code weekly, with some even doing so daily. However, the most typical rate at which code was shared to the central analysis store was monthly. This speaks to both the utility of referencing past work, as well as the time investment required to clean and share one's own code that is unlikely to be done daily. Additionally, over time an existing analysis might already exist and so data scientists feel the benefit of sharing their new, slightly different, analysis is marginal. This trend of reusing shared code more frequently than contributing new shared code holds for all of the remaining strategies.

Sharing cleaned analysis code allows team members to benefit from one another's prior work. This shared knowledge extends beyond code but can also include documentation shared about previous experiments and what did *not* work so that future analyses do not make the same mistake. For this sharing to be effective however, more effort is required prior to sharing. By nature of sharing work with other people, it takes more effort and data scientists want to share high quality code with their peers. Answering this question of "when is my work good enough to share with other people?" as IP17 called it, can be difficult to navigate.

6.4 Team shared template notebooks

Typically, computational notebooks have issues of non-reproducible code that limit their ability to be shared and reused [24]. However, as they have grown into the de facto tool for data science, users have developed strategies for making notebooks more reusable, namely through *template notebooks*. To a certain extent, template notebook sharing is a subset of team-wide sharing mentioned in Section 6.3, however the methods presented here are unique enough to merit their own strategy.

Template notebooks are normal computational notebooks that have been cleaned up and generalized for a certain task. This cleaning can take many forms. Some teams put parameters at the top of the notebook and leave "TODOs" in the comments to fill in these parameters before running the notebook like a function. Over time, each of these task-specific notebooks come together as a sort of "library" for common code on the team. These notebooks are intentionally cleaned so that they will run top to bottom without issue. In this way, template notebooks are run more like a traditional python script rather than an interactive notebook. Several informants also discussed using the %run command in Databricks, which allows one notebook to invoke another like a function, including parameter passing.

Alternatively, a single notebook may contain numerous common functions in a *notebook as a library* (NAL). For example, IP6's team has a centrally shared Databricks notebook that has numerous commonly used functions. These functions are cleaned and abstracted before addition to the NAL. Yet when these functions are reused, they are copy and pasted out of the shared notebook rather than imported. This notebook is essentially a traditional software engineering library – similar functions grouped together under one umbrella. However it lives in *same place* analysts do their work.

Some tools have extra functionality that supports the use of template notebooks and NALs. For example, Databricks allows notebooks to be parameterized so that they can be run as a function (top to bottom) with new parameters or a default value (see [8] for more details on this syntax). Some of our informants have developed whole analysis pipelines using this functionality where one notebook does data cleaning and then calls another notebook to do modelling, and so on.

Template notebooks were used least frequently among the surveyed responses. This can be attributed to two main factors. The first is that not every notebook shared with teammates is a *template* notebook. Past analyses that are tightly coupled with a particular data source are shared on team-wide stores but unless they are cleaned to the extent they can be run independently, they are not a template. Secondly, there is limited tool support for template notebooks outside of the Databricks environment, which was used on a limited number of teams. Adding *TODO* comments to Jupyter notebooks was one alternative strategy, however executing Jupyter notebooks as a function is uncommon.

6.4.1 Access and Search. Template notebooks are most commonly accessed from Git repositories for Jupyter based templates or within an analysis ecosystem such as Databricks, where tool features best support reuse. Even for notebooks stored in a version-controlled Git repository, participants seldom use common Git operations like branching and merging. Rather, they use Git as a remote store for their work that others have easy access to, more similar to any cloud storage location. Several participants even mentioned they are not very comfortable with Git workflows and it is something they would like to work on.

Code is most commonly reused by copy and paste or (if the tool supports it) by calling the code directly. As is the case with all forms of shared code artifacts, the most common way among survey respondents to find relevant notebook templates is by communicating with teammates. However, the next most popular response was based on memory indicating that most shared notebooks are for small, modular tasks that data scientists easily remember.

- 6.4.2 Additions. Our interview informants mentioned that most notebook template additions are done by a select few members of their team. However, everyone has the ability to share new templates. The division between reuse versus sharing frequency also holds for notebook templates: the modal response for adding new template notebooks is seasonally, whereas they are most commonly used weekly.
- 6.4.3 Benefits and Drawbacks. Why do data scientists use shared template notebooks rather than a library? One reason is that template notebooks are more compatible with where they do most of their work. If a reused notebook contains lots of visualization or table output, this most easily re-run as a notebook rather than a separate library. Another reason for template notebooks is that they support tweaking the code to the analysis at hand. As SP131 says, they use template notebooks "when the code requires customization when applied to different scenarios." Another motivation, from SP46, is that template notebooks are useful for "Just about any situation involving analytics. Sharing notebooks is far easier and portable than making a code library. For Data Scientists at least."

Creating template notebooks is often faster than creating a traditional library. For example, survey respondents mentioned they chose to create template notebooks because of "Extreme time constraint which limits the time I can spend on doing things the right way" (SP43) or "Most of the time because the flow is easier" (SP24).

6.5 Team shared libraries

The last strategy we observed data scientists using is the traditional software library. These are no different than any other library: similar code is grouped together into one place that can be imported and run from another analysis context. For example, IP3's team does time series analysis and so has developed a library in R of the common cleaning methods, models, and postprocessing methods they use. In fact, their library is so robust, they claimed it might even hinder them trying new models or other methods if they are not in the library because the marginal benefit of writing a new function does not out-weight the cost.

6.5.1 Access and Search. Although these shared libraries are more stable than some of the other forms of reuse, the most common way that data scientists find code in the library is still by talking with team mates. Other common responses were finding code based on memory or by reading team-wide documentation that references library functions. Regardless of the level of structure in sharing, data science workers default to finding reusable code through social interactions rather than searching over the shared code themselves.

Unsurprisingly, the primary way that libraries are used is through importing. They are also stored on Git and typically go through a code review process between the whole team or the library maintainers on the team when new functions are added in. Library functionality once again spans the entire range of data science activities. The most common functions were data processing and cleaning, followed by data reading and writing.

Some of the teams maintained a library for certain functions that do not change much or need to be customized, and then used one of the other team-wide sharing strategies for analysis code that is updated more often:

In terms of the library that we have, that's actually more of a how to interact with our catalogs and our file systems. This [notebook] right here, this is mostly just time series analytics. They could have probably been added to [the library], but this is something that we add to it every day. Whereas the file system interaction stuff only happens... We only change it once a month that we need to add a new dataset or something. - *IP6*

template notebooks rather than a library? One reason is that template notebooks are more compatible with where they do most of their work. If a reused notebook contains lots of visualization or format degree. For example, IP3 described their release cycles in the following way:

Suitable for bootstrapping analyses we have a monthly sync up to let all the teammates know what happens with the new library. And also others will let me and the other data scientists know what new functions they would like to add in.

Among our survey respondents that used libraries, the most common response for library *usage* rate was monthly followed by weekly. However, the most common rate of library contribution

was never then seasonal. This shows that data scientists are using these shared libraries far more often than contributing new code. This has several possible explanations. One is that the library code becomes more stable over time and so less things need to be added, as is the case for IP3 above. Furthermore, the type of code shared to libraries is more stable and so needs less updating than analysis code or templates that necessitate frequent tweaks.

6.5.3 Benefits and Drawbacks. The preference for a shared library versus shared notebooks templates was largely influenced by the tool data scientists typically work in. Some teams that used the R ecosystem, like IP3, find that using libraries is very convenient in that environment. Others who use Python and notebooks primarily like IP9, prefer template notebooks:

Honestly, I think it was just the ease of being able to view it [in a notebook], edit it and at the same time, run it, just using the command...Because abstracting it as a Python [module] would have worked too, but then it hides it a bit too much in my opinion. In this way, people can at least, like if they need to make a change somewhere let's say tomorrow, some of these client secrets are no longer valid. Then, it still provides the ability to easily debug what's going on by just coming to the common notebooks and swapping some of these out. It's just easier.

7 DETERMINANTS OF REUSE

Throughout our interviews, we noticed common factors that encouraged or discouraged data scientists from sharing and reusing past analysis code. We refer to these as the determinants of reuse and split them up among qualities of shared code, and incentives for sharing and reusing analysis code.

7.1 Qualities of shared code

When evaluating whether to reuse past analysis code or write it from scratch, data scientists must evaluate a trade off: do they invest the time to find and adapt past code for their current analysis, or write the code from scratch?

The data scientists surveyed agreed on several factors that helped them navigate this trade off and encouraged them to reuse past code. First, they must be able to quickly understand what the code is doing (85% agreement¹). If code is messy, hard to read, or confusing then it is not worth taking the time to parse the past analysis. Along these same lines, data scientists agreed that well documented and "clean" code encourages reuse (87% agreement). Furthermore, if only a small number of edits are required to adapt the code to their current context this also promotes reuse (83% agreement). Of course, if this code is contained in a library then small tweaks are more difficult than in analysis files or template notebooks. Additionally, if data scientists perceive that the code is doing something complex, they will be more likely to reuse it since it would take them a long time to rewrite themselves (92% agreement).

Trust in the author of code also encourages reuse (82% agreement). This manifests when asking a specific person for their previous analysis, but also if data scientists look at shared repositories

and can see who uploaded an analysis. If they know that someone has a reputation for good work and they trust them, they will be more likely to reuse this person's code.

Two other aspects that encourage reuse were met with more mixed responses. Survey respondents were ambivalent if they would be more likely to reuse code that is a small number of lines (39% thought this was important, 32% unimportant, and 29% were in the middle). We attribute this to the fact that the aforementioned factors are more important than the actual length of the code – long, well documented, but complex code is more likely reused than short, messy, and simple code that is more easily written from scratch. Lastly, survey respondents also had mixed responses on whether recently updated code encouraged them to reuse. Since some shared libraries experience less updates over time this does not mean they have diminished utility. However, some respondents mentioned they would prefer more recently updated analysis code since it is less likely to contain stale references to databases or API keys.

7.2 Incentives & Culture

We also discovered some of the incentives that encourage or discourage sharing and reusing on data science teams. Some of these are also related to the quality of the shared code, but others speak to team culture and how well tools support sharing. The full range of questions and responses are presented in Figure 2.

Data science work takes place in a wide variety of tools. A single data scientist might use several tools just to analyze their data and someone else on their team might use a different set of tools for analysis. This can lead to difficulties finding and reusing code, as demonstrated by the high agreement among survey respondents that analyses in different languages and locations can make reuse difficult (74% agreement). Furthermore, errors in shared code make it harder to reuse other's analyses (43% agree they often experience errors, 30% disagree).

During our interviews, data scientists repeatedly mentioned how their team culture influences their sharing practice. Respondents almost unanimously agree that sharing code takes a time investment (93% agreement). Code must be commented, cleaned, and presented in a way conducive for others to reuse that would not have to be done otherwise. In order for this time investment to feel worthwhile, the team culture must reward sharing. This means when things are shared, people will actually use them, and that data scientists have (or are allotted) the time to invest in creating high quality, shared code. Some participants felt their team supported and encouraged them to invest this time into making analysis code sharable. IP15 says:

Another benefit is that we have a really strong learning culture, a knowledge sharing culture which was really nurtured by the management and saying "Hey you know, yes you own your project, but we're not going to be proprietary about anything within our team". So we want to make sure that everybody has access to everything and that's how we help each other out.

However, other respondents mentioned their team is still working to establish a culture where sharing is rewarded:

¹We calculate agreement as the sum of positive Likert responses of Important, Somewhat Important, and Very Important. Similarly, disagreement is the sum of negative Likert responses.

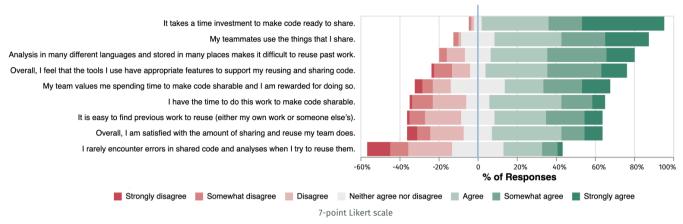


Figure 2: Likert scale responses to questions about the determinants of reuse.

What about code snippets in frontend development?

Sharing code in data science is often not rewarded at the level needed to offset the investment necessary to do it well. - SP104

Many survey participants felt their team valued and rewarded them for investing time to share code (54% agreement), however 18% disagreed and felt their team did not recognize their efforts. Regardless of tool usage or coding style, a team culture that encourages data scientists to invest the time to share their culture is critical for reuse to succeed.

8 DISCUSSION

In preceding sections, we have presented interview and survey results on *how* data scientists go about sharing and reusing past analysis; here we delve into a discussion of *why*. We also distinguish how reuse in data science is distinct from reuse in traditional software engineering and present opportunities for future tools to improve the experience of sharing and reuse in data science.

8.1 Code Is Not a Deliverable

Throughout our interviews and survey, data scientists repeatedly mentioned how the code that they write is separate from their project deliverables. In data science, the outcome of an analysis is not the analysis code itself, but whatever insights, models, or datasets are produced from that analysis. These insights are delivered through slide decks, word documents, and interactive dashboards. However, the consumers of these artifacts rarely care about the analysis code itself or might not have the technical skill set to understand all the code that went into producing the analysis output. Future tools might investigate how to tie presentation artifacts such as charts or tables back to the code and data used to create the artifact to help speed up this iteration cycle and support reuse at the presentation level as well.

Furthermore, many analyses are one-off, unrelated requests. Non-code deliverables and one-off analyses ostensibly combine to create disincentives for sharing and reuse. Despite this, reuse and sharing are still very common. Data scientists realize how much reusing past analyses improves their productivity by avoiding re-work. This different incentive structure leads data scientists to adapt common software engineering reuse strategies to their needs. The reuse of personal code and team wide libraries have strong parallels with

typical software reuse whereas personal utility libraries, shared analysis stores, and notebook templates developed out of the unique needs for code reuse in data science. Future work might attempt to quantify the benefits of reusing past code, for instance in terms of the time required to complete an analysis.

8.2 Modular Data Science

Given the need for customization, it is difficult to create modular data science analysis components. As mentioned in Section 6.5, libraries are typically used for components that change little over time, like data access APIs. However full analyses need to be customized almost every time and so are more likely to be shared using a customizable interface like a computational notebook.

There are few tools that support this kind of interaction. The best example in the tools surveyed was the %run syntax in Databricks for running notebooks as functions discussed in Section 6.4. This feature allows the parameterization of notebooks. However, if notebooks need to be customized beyond the available parameters, data scientists will often just clone and edit the notebook.

Observable offers a notebook style interface for JavaScript programming that lets users import cells from any other Observable notebook [5]. Observable is most often used for visualization creation; future work might explore how cell-level imports can aide sharing and reuse in other data science programming environments. Additionally, future tools might investigate how to combine modularity at the functional level with *modular analyses* that can be customized to the current data by allowing the addition or deletion of entire analysis steps.

8.3 Tool Interoperability

A common pain point among data scientists is the fact that analyses can exist in many different tools and languages. Many of these tools are customized to a certain type of analysis or part of the analysis pipeline or are driven by data privacy and regulatory concerns. As SP68 describes, "We have to use different platforms for different purposes, which makes code reuse difficult." Even those who are proficient in one or more languages can struggle to understand code in unfamiliar languages [25].

Existing integrated development environments allow data scientists to work in multiple languages in the same user interface. Ongoing efforts like Apache Arrow [3] allow in-memory data exchange between language runtimes, which might encourage code reuse across languages. Another approach is to provide a domain-specific language for data science that can be compiled to multiple languages and runtimes. For example, OpenAI's Codex project [6] allows data science workflows to be expressed in natural language that is compiled to Python. Writing workflows at this higher level of abstraction might make them more reusable in other analysis contexts and AI-assisted code authoring can help data scientists search for past code to reuse.

8.4 Data Privacy

In data science, analysis code can be difficult to understand and impossible to run without the data it analyzes. Due to customer privacy or regulatory restrictions, sharing the data that accompanies an analysis may not be possible, even within the same company. Notebooks help since analysis results can be shared without exposing access to the underlying data. However, notebooks need access to the original data to *re*-run. Future work might explore how analysis code can be shared with anonymized, private data so that the code can still execute and thus provide help for future analyses without exposing private data.

9 LIMITATIONS

Our interviews and survey are subject to several limitations. We only interviewed data scientists from a single company with a relatively mature data science practice. However, even within one company not every data science team had uniformly mature reuse practices. We expect these results will generalize to other data scientist populations; some of the issues described may even be felt more acutely by smaller organizations with less organized data science practices. However, future work might explore this explicitly. For our survey, incorrect branching logic caused not all participants to see the Likert rating questions at the end. We report all survey results as percentages of those that responded; skipped questions or no responses are excluded from reported counts.

10 CONCLUSION

In summary, we present the results of 17 interviews and a 132 person survey with professional data scientists about how they share and reuse work. Our investigation revealed data scientists reuse work through five strategies: personal analysis reuse, personal utility libraries, team shared analysis code, team shared template notebooks, and team shared libraries. We discuss factors that encourage and discourage reuse, and how future tools might better support this practice in data science.

ACKNOWLEDGMENTS

We would like to thank our interview and survey participants for their help in this research, and the anonymous reviewers for their valuable feedback.

REFERENCES

[1] Rabe Abdalkareem, Olivier Nourry, Sultan Wehaibi, Suhaib Mujahid, and Emad Shihab. 2017. Why do developers use trivial packages? an empirical case study on npm. In Proceedings of the 2017 11th joint meeting on foundations of software engineering. 385–395.

- [2] Saleema Amershi, Andrew Begel, Christian Bird, Robert DeLine, Harald Gall, Ece Kamar, Nachiappan Nagappan, Besmira Nushi, and Thomas Zimmermann. 2019. Software engineering for machine learning: A case study. In 2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP). IEEE. 291–300.
- [3] Apache. 2021. Apache Arrow. Retrieved October 15, 2021 from https://arrow.apache.org/
- [4] Andrew Begel and Thomas Zimmermann. 2014. Analyze this! 145 questions for data scientists in software engineering. In Proceedings of the 36th International Conference on Software Engineering. 12–23.
- [5] Mike Bostock. 2018. Introduction to Imports. Retrieved August 25, 2021 from https://observablehq.com/@observablehq/introduction-to-imports
- [6] Mark Chen and et al. 2021. Evaluating Large Language Models Trained on Code. CoRR abs/2107.03374 (2021). arXiv:2107.03374 https://arxiv.org/abs/2107.03374
- [7] Databricks. 2021. Databricks. Retrieved August 25, 2021 from https://databricks.com
- [8] Databricks. 2021. Notebook workflows. Retrieved October 11, 2021 from https://docs.databricks.com/notebooks/notebook-workflows.html
- [9] William B Frakes and Kyo Kang. 2005. Software reuse research: Status and future. IEEE transactions on Software Engineering 31, 7 (2005), 529–536.
- [10] Git. 2021. Git. Retrieved August 25, 2021 from https://git-scm.com/
- [11] Martin L Griss. 1993. Software reuse: From library to factory. IBM systems journal 32, 4 (1993), 548–566.
- [12] Lars Heinemann, Florian Deissenboeck, Mario Gleirscher, Benjamin Hummel, and Maximilian Irlbeck. 2011. On the extent and nature of software reuse in open source java projects. In *International Conference on Software Reuse*. Springer, 207–222.
- [13] Mary Beth Kery and Brad A Myers. 2018. Interactions for untangling messy history in a computational notebook. In 2018 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC). IEEE, 147–155.
- [14] Miryung Kim, Thomas Zimmermann, Robert DeLine, and Andrew Begel. 2016. The emerging role of data scientists on software development teams. In 2016 IEEE/ACM 38th International Conference on Software Engineering (ICSE). IEEE, 96–107
- [15] Miryung Kim, Thomas Zimmermann, Robert DeLine, and Andrew Begel. 2018. Data Scientists in Software Teams: State of the Art and Challenges. *IEEE Transactions on Software Engineering* 44, 11 (2018), 1024–1038. https://doi.org/10.1109/TSE.2017.2754374
- [16] Yongbeom Kim and Edward A Stohr. 1998. Software reuse: survey and research directions. Journal of Management Information Systems 14, 4 (1998), 113–147.
- [17] Charles W Krueger. 1992. Software reuse. ACM Computing Surveys (CSUR) 24, 2 (1992), 131–183.
- [18] Tim Menzies. 2016. How Not to Do It: Anti-Patterns for Data Science in Software Engineering. In Proceedings of the 38th International Conference on Software Engineering Companion (Austin, Texas) (ICSE '16). Association for Computing Machinery, New York, NY, USA, 887. https://doi.org/10.1145/2889160.2891047
- [19] Tim Menzies, Ekrem Kocaguneli, Fayola Peters, Burak Turhan, and Leandro L. Minku. 2013. Data Science for Software Engineering. In Proceedings of the 2013 International Conference on Software Engineering (San Francisco, CA, USA) (ICSE '13). IEEE Press, 1484–1486.
- [20] Microsoft. 2021. Kusto query overview. Retrieved October 8, 2021 from https://docs.microsoft.com/en-us/azure/data-explorer/kusto/query/
- [21] Parastoo Mohagheghi and Reidar Conradi. 2007. Quality, productivity and economic benefits of software reuse: a review of industrial studies. *Empirical Software Engineering* 12 (2007), 471–516.
- [22] João Felipe Pimentel, Leonardo Murta, Vanessa Braganholo, and Juliana Freire. 2019. A Large-Scale Study About Quality and Reproducibility of Jupyter Notebooks. In 2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR). 507–517. https://doi.org/10.1109/MSR.2019.00077
- [23] Israel J Mojica Ruiz, Meiyappan Nagappan, Bram Adams, and Ahmed E Hassan. 2012. Understanding reuse in the android market. In 2012 20th IEEE International Conference on Program Comprehension (ICPC). IEEE, 113–122.
- [24] Adam Rule, Aurélien Tabard, and James D. Hollan. 2018. Exploration and Explanation in Computational Notebooks. Association for Computing Machinery, New York, NY, USA, 1–12. https://doi.org/10.1145/3173574.3173606
- [25] Nischal Shrestha, Colton Botta, Titus Barik, and Chris Parnin. 2020. Here We Go Again: Why Is It Difficult for Developers to Learn Another Programming Language?. In 42nd International Conference on Software Engineering (ICSE).
- [26] Bowen Xu, Le An, Ferdian Thung, Foutse Khomh, and David Lo. 2020. Why reinventing the wheels? An empirical study on library reuse and re-implementation. Empirical Software Engineering 25, 1 (Jan. 2020), 755–789. https://doi.org/10.1007/s10664-019-09771-0
- [27] Amy X. Zhang, Michael J. Muller, and Dakuo Wang. 2020. How do Data Science Workers Collaborate? Roles, Workflows, and Tools. Proceedings of the ACM on Human-Computer Interaction 4 (2020), 1 – 23.