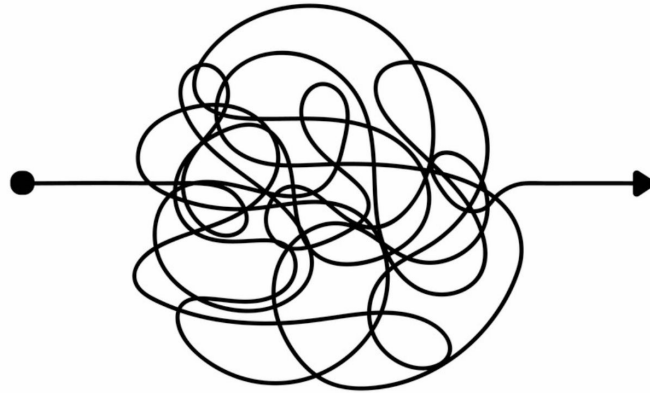# How Data Scientists Use Computational Notebooks for Real-Time Collaboration

**April Yi Wang** | Anant Mittal | Chris Brooks | Steve Oney
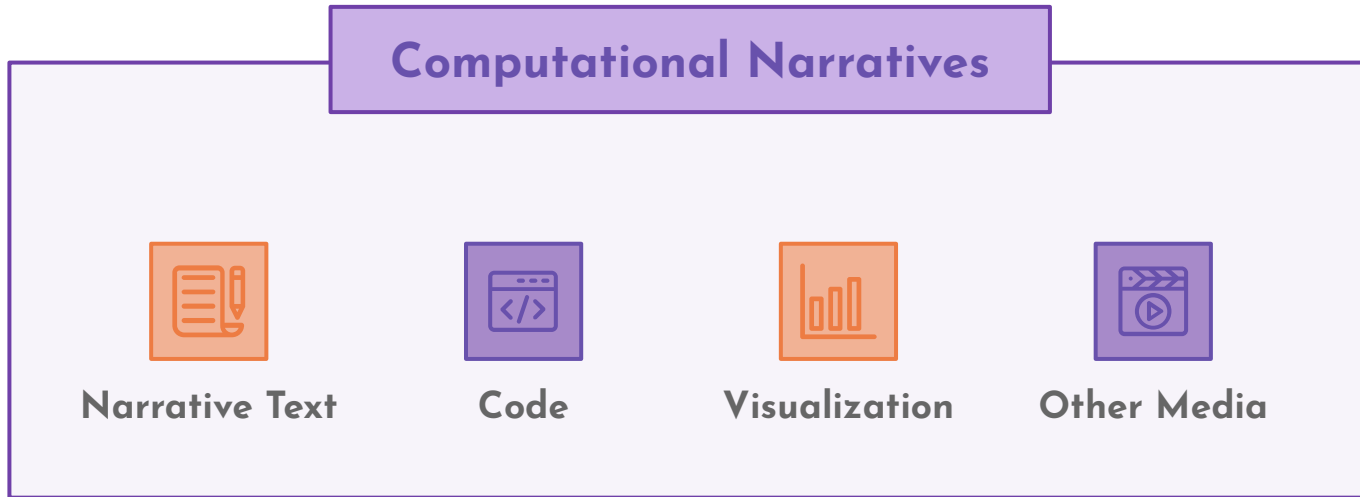
University of Michigan

# The Story Behind Data Analysis

# The Story Behind Data Analysis

## Computational Narratives

Narrative Text    Code    Visualization    Other Media

# Jupyter Notebook

Jupyter notebooks consist of **"cells"** — typically small chunks of code or narrative text in the Markdown format.

Users can execute cells (typically, but not necessarily, from top to bottom) and observe their outputs.

4

# Writing and Sharing Computational Notebooks in Various Contexts

**Data Science Education**
Kross and Guo, CHI 19

**Open Science**
Randles et al., JCDL 17

**Professional Data Analytics**
Kery et al., CHI 18

# From Sharing to Synchronous Editing



Deepnote



Google colab

# Issues with Synchronous Editing



- Reluctant to write together when collaboratively constructing a document

- Social embarrassment to be watched by others when typing

~ Wang et al. CSCW'17

# Issues with Synchronous Editing

**Collaborative Writing**

Wang et al. CSCW'17

D'Angelo et al. CSCW'18

**Collaborative Programming**

Goldman et al. UIST'11

Oney et al. CSCW'18

# What about collaborative data science?

data science ≠ writing + coding

**RQ1** What **tools and strategies** do data scientists currently use for collaboration?

**RQ2** Compared to **working on individual notebooks** in a collaborative setting, how does **synchronous notebook editing** change the way data scientists collaborate in computational notebooks?

**RQ3** What **challenges**, if any, do data scientists perceive in synchronous notebook editing?

**Study 1 Formative Survey**

**RQ1** What **tools and strategies** do data scientists currently use for collaboration?

**RQ2** Compared to **working on individual notebooks** in a collaborative setting, how does **synchronous notebook editing** change the way data scientists collaborate in computational notebooks?

**RQ3** What **challenges**, if any, do data scientists perceive in synchronous notebook editing?

**Study 2 Observational Study**

# Demographic

## Data Source

**195**
Valid Responses

**35**
Umich

**160**
Coursera

## Experience with Data Science



## Job Roles



- Students
- Data Scientists
- Software Engineers
- Researchers
- Managers
- Business Analysts
- Others

## Choices of Tools

| | |
|---|---|
| **Programming** | Jupyter Notebooks (88.72%), IDEs (51.79%), **Google Colab (12.31%)** |
| **Communication** | Emails (79.49%), Face-to-face Communication (68.72%) |
| **Project Management** | Version Control Tools (49.74%), Task Tracking Tools (21.03%) |

# Strategies for Keeping a Shared Understanding

| | |
|---|---|
| **Discussions and Meetings** | Weekly meeting among team members; |
| **Frequently Check-in** | Communicate actively and frequently; |
| **Documentation** | Keep notes in Google Docs; ... comments in code; |
| **Organization** | Divide up the work into definable parts; |
| **Shared Assets** | Common repository for files; |
| **Others** | Code review to ensure code matched intent |

**RQ1**   What **tools and strategies** do data scientists currently use for collaboration?

### Traditional Collaboration Setting

Working on individual Jupyter notebooks

### Emerging Collaboration Setting

Working on notebooks with synchronous editing

15

**RQ2** Compared to **working on individual notebooks** in a collaborative setting, how does **synchronous notebook editing** change the way data scientists collaborate in computational notebooks?

**RQ3** What **challenges**, if any, do data scientists perceive in synchronous notebook editing?

**Study 2 Observational Study**

16

# Participants

- 24 participants (12 from the survey)
- Randomly assigned to pairs
- Work collaboratively on a predictive modeling problem remotely

CA (5)

PK (1)

US (9)

CN (2)

IN (6)

BR (1)

17

# Study Setup

## Non-Shared Condition

Participants worked on individual notebooks

- ✓ Exchange the notebook file
- ✓ Set up a git repository
- ✓ Send code snippets through other tools if necessary

## Shared Condition

Synchronous editing was supported.

- ✓ Share notebook edits and actions (e.g., moving cursor, adding cells) in real-time
- ✓ Execute code on a single interpreter
- ✓ Update output and runtime variables among collaborators

18

## Task

- Predict house sale prices using 80 features (e.g., lot size, year built)

- Additional incentives for the group with the lowest error score

- Submit prediction results as well as one Jupyter notebook report

- Choose from text-messaging (Slack) or video-conferencing (Google Hangouts) for communication

19

## Procedure

The study consisted of four sessions, each of which lasted an hour.

**Understand the Data**
Session 1

**Pre-Processing**
Session 2

**Basic Predictive Model**
Session 3

**Advanced Model**
Session 4

20

## Collaboration Style

| Collaboration Style | GID | Definition |
|---|---|---|
| **Single Authoring** | | One team member contributed the majority of ideas and did the majority of the implementation, while the others did not contribute much. |
| **Pair Authoring** | | |
| **Divide and Conquer** | | |
| **Competitive Authoring** | | |

## Collaboration Style

| Collaboration Style | GID | Definition |
|---|---|---|
| **Single Authoring** | | One team member contributed the majority of ideas and did the majority of the implementation, while the others did not contribute much. |
| **Pair Authoring** | | One team member did the majority of implementation while the others contributed ideas, engaged in discussions and reviewed the results. |
| **Divide and Conquer** | | |
| **Competitive Authoring** | | |

22

# Collaboration Style

| Collaboration Style | GID | Definition |
|---|---|---|
| **Single Authoring** | | One team member contributed the majority of ideas and did the majority of the implementation, while the others did not contribute much. |
| **Pair Authoring** | | One team member did the majority of implementation while the others contributed ideas, engaged in discussions and reviewed the results. |
| **Divide and Conquer** | | Members divided the task into subgoals and explored the subgoals independently. |
| **Competitive Authoring** | | |

23

# Collaboration Style

| Collaboration Style | GID | Definition |
| --- | --- | --- |
| **Single Authoring** | | One team member contributed the majority of ideas and did the majority of the implementation, while the others did not contribute much. |
| **Pair Authoring** | | One team member did the majority of implementation while the others contributed ideas, engaged in discussions and reviewed the results. |
| **Divide and Conquer** | | Members divided the task into subgoals and explored the subgoals independently. |
| **Competitive Authoring** | | Team members wrote the code for the same purpose and reached the consensus to use the code by whomever finished first. |

# Collaboration Style

| Collaboration Style | GID | Definition |
|---|---|---|
| **Single Authoring** | S2, S5 | One team member contributed the majority of ideas and did the majority of the implementation, while the others did not contribute much. |
| **Pair Authoring** | S6 | One team member did the majority of implementation while the others contributed ideas, engaged in discussions and reviewed the results. |
| **Divide and Conquer** | N2, N5, S1, S3, S4 | Members divided the task into subgoals and explored the subgoals independently. |
| **Competitive Authoring** | N1, N3, N4, N6 | Team members wrote the code for the same purpose and reached the consensus to use the code by whomever finished first. |

# Communication Channels

| | Non-Shared Condition | Shared Condition |
|---|---|---|
| Choices of Tools | Text Messaging (6/6) | Text Messaging (3/6) Video Conferencing (3/6) |

26

Participants in the non-shared condition send files, code snippets, and output more often.

➔ Working in the shared notebook may reduce the communication costs by establishing a shared context.

# Final Submissions

→ Groups in the shared condition achieved a better prediction result.

■ Non-Shared Condition
■ Shared Condition

Error Score

0.27
0.17

→ Groups in the shared condition explored more alternative models.

**27**

Number of Alternative Models*

3.00
6.17

p= 0.05

Lines in the Notebook*

90.33
186.67

p= 0.04

# Work Across Phases

| Preparing | Cleaning | Feature Engineering | Modeling | Submission |



Session 1   Session 2   Session 3   Session 4

Person 1
Person 2

N6

S5

switch

Participants in the shared condition switched more frequently (p<0.001).

➔   Working on the same notebook provides collaborators with convenience to branch through tasks

28

## Work Across Phases



Shared Condition     Non-shared Condition

switch

Participants in the shared condition switched more frequently (p<0.001).

➔ Working on the same notebook provides collaborators with convenience to branch through tasks

29

# Benefits of Synchronous Editing in Notebook

➔ Reducing communication costs

➔ Flexibility to branch through tasks

➔ Enabling explorations of more alternative models

➔ Leading to a better prediction result

30

## Challenges of Synchronous Editing

# Challenges of Synchronous Editing

1. Interference with each other

| train_df | → | df_train |
|---|---|---|

*"...* ***When using Jupyter Notebook together, it's hard to keep track of variable names.*** *Everyone might use a different name and may cause issues. For example, my teammate used train_df as name, and later changed it to something else, but I wanted him to keep using the original name..." (P2 from S1)*

32

# Challenges of Synchronous Editing

## 2. Lack of Strategic Coordination

Why competitive authoring happens in the non-shared condition?

Alice: 80%    Bob: 60%

Why single authoring happens in the shared condition?

Alice: 80%    Bob: 20%

33

*"... I feel I am not splitting work well enough. **I was thinking about how to get the work done and just tried the ideas on myself.**..." (P11 from S2)*

# Challenges of Synchronous Editing

S3 wrote down subtasks in the notebook.

**Pre-processing and cleaning**

**Steps**

1. Replace discrete values with indices
2. Remove data samples with too many missing features
3. Normalize continuous variables
4. Compute correlation, or use other techniques to select features

```python
In [24]:    1  def count_nans(data):
            2      for name in data:
            3          count_nan = len(data[name]) - data[name].count()
            4          print(name, 'num of nans:', count_nan)
```

```python
In [74]:    1  def label_encoding(data):
            2      for name in data:
            3          if data[name].dtype == 'object':
```

# Challenges of Synchronous Editing

## 3.  Contextual Chatting

P14 and P15 were looking at the scatterplots of
independent variables together.



2:18 PM
In my opinion there are outliers in all of our features
there are 1 or 2 points that outlies

2:19 PM
which ones?

P14 downloaded the graph, opened MS Paint,
annotated the graph and sent it back to P15.

# Challenges of Synchronous Editing

1. **Interference with each other**
2. **Lack of Strategic Coordination**
3. **Contextual Chatting**
4. Lack of Awareness
5. Problems with the Linear Structure
6. Privacy Concerns

36

➔ Working on the same notebook results in different collaboration styles compared to working on individual notebooks.
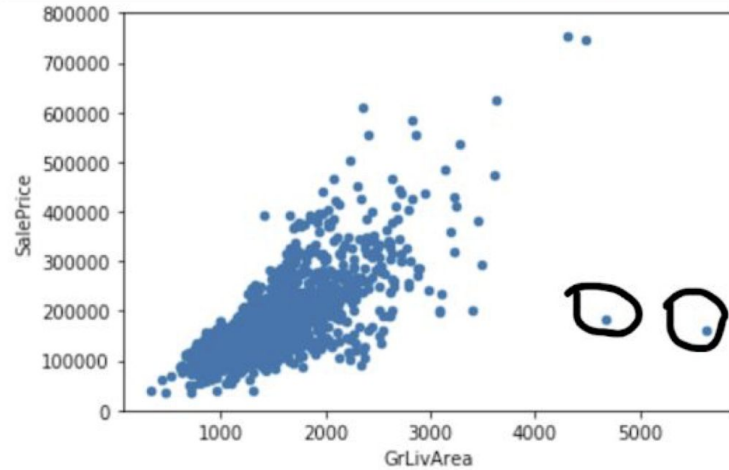
➔ Synchronous editing tools improve collaboration by helping data scientists maintain a shared context and improve work efficiency.

➔ However, the current real-time collaborative editing features may lead to several problems (e.g., interference with each others' work, unbalanced contributions).

**Extending Our Understanding of Collaborative Editing Across Contexts**

- Collaborators may hold different programming backgrounds and domain knowledge

- Different roles in collaborative data science

# Example: How to deal with the two dots?

```
In [8]:  # scatter plot grlivarea/saleprice
         label = 'GrLivArea'
         data = pd.concat([df_train['SalePrice'], df_train[label]], axis=1)
         data.plot.scatter(x=label, y='SalePrice', ylim=(0,800000));
```

## Design Implications

- Improve Awareness of Collaborators' Activity

- Provide Access Control

- Enable Discussions within Notebooks

40

## Limitations

- Generalizability
  - the type of data science problems
  - the expertise of collaborators
  - the size of the team
  - the synchronicity of the collaboration

41

# How Data Scientists Use Computational Notebooks for Real-Time Collaboration

What tools and strategies do data scientists currently use for collaboration?

Study 1 - Formative Survey on Collaborative Data Science

Traditional Collaboration Setting + Emerging Collaboration Setting

How does synchronous notebook editing change the way data scientists collaborate?

What challenges do data scientists perceive in synchronous notebook editing?

Study 2 - Observational Study on Collaborative Data Science

Having synchronous editing is great for collaborative data science, but not perfect!

Presenter: April Yi Wang | aprilww@umich.edu

Co-authors: Anant Mittal, Chris Brooks, Steve Oney