CENTER ON RURAL INNOVATION

# BEYOND RAW DATA

## Streamlining Insights with Curated Data Packages
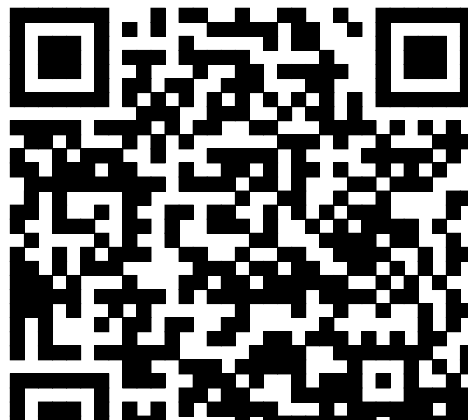
Olivier Leroy

Center On Rural
Innovation

2024-10-20

CORI — Center on Rural Innovation

## SLIDES ARE ONLINE:

- Created in Quarto (RevealJS)

- Running R and cori.data.fcc

- Hosted and deployed on GitHub
  https://ruralinnovation.github.io/
  prez_auber_2024/



CORI Center on Rural Innovation

## GOALS:

1. Primary data is great! (but can be a pain...)

2. Use code (other pain?) to work with primary data

3. Packaging it will make live easier!

CORI Center on Rural Innovation

# PLAN: RETRACE OUR JOURNEY

[TODO: rework at the end]

# 👋 HI, I AM OLIVIER SENIOR DATA ENGINEER

- Working with primary data: "messy, unstructured data" is upon us

**Analytics And Data Science**

## Data Scientist: The Sexiest Job of the 21st Century

Meet the people who can coax treasure out of messy, unstructured data. by Thomas H. Davenport and DJ Patil

From the Magazine (October 2012)

Source: Harvard Business Review

- Socials: LinkedIn | Mastodon | Personal Website
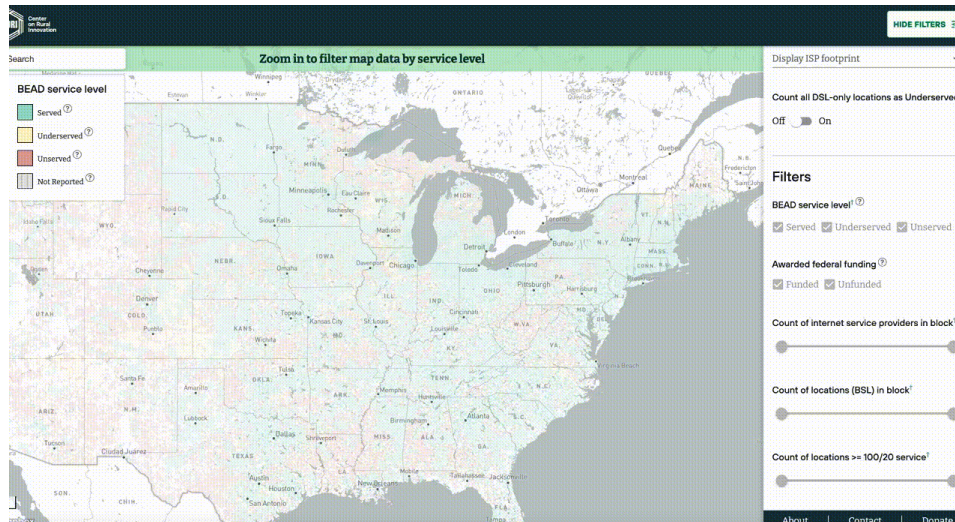
# DATA AT CENTER ON RURAL INNOVATION (CORI)

1. Rural data: we can't work with aggregate and need to go deeper

I plan to delete this unless good points against

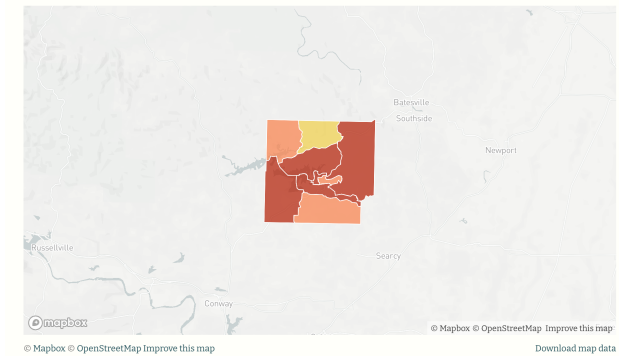# BROADBAND: AN IMPORTANT PART OF OUR WORKS

# PUBLIC-FACING APPS:



Rural Broadband Mapping Tool



Broadband Climate
Risk Mitigation Tool

**Beyond Connectivity:**
The Role of Broadband in Rural Economic Growth and Resilience

**Dr. Amanda Weinstein**
Director of Research, Knowledge, and Evaluation
Center on Rural Innovation

**Dr. Adam Dewbury**
Researcher
Center on Rural Innovation

# BROADBAND RESEARCH IS EVOLVING!

# FEDERAL COMMUNICATIONS COMMISSION (FCC)

- Agency in charge of regulating "everything" communications

- Their data landscape is evolving fast (now they have a platform!)

**CHALLENGE: HOW TO WORK WITH CONSTANT CHANGE UPSTREAM?**

CORI Center on Rural Innovation

# FEDERAL COMMUNICATIONS COMMISSION (FCC)

|  | Form 477 | National Broadband Map |
|---|---|---|
| US Census Boundaries | 2010 | 2020 |
| Type of recording | self declarative | self declarative |
| Granularity | Census blocks | Locations |
| Services | Mobile/ Fixed | Mobile/ Fixed |
| Timeframe | 2014 - 2021 | 2022 - *Ongoing* |
| Releases | **twice a year** | **twice a year**[1] |

1. One release can have multiple versions

**BROADBAND ACCESS SINCE 2014 BUT:**

- 2 **voluminous** datasets, hundreds of files

- No "**road safety barriers**": multiple encodings, erronous values (and more..)

CORI | Center on Rural Innovation

# BEHIND THE SCENE:

# HOW IT STARTED: MANUEL EXTRACTION



- 🔁Repeat for every State (56)

- 🔁Repeat for every version (??)

- 🔴Error prone (500 hundred clicks)

# ... AUTOMATION WITH CODE ...

```
 1  # [...]
 2  # Download all FCC data needed
 3  all : dirs dl_file list_file
 4
 5  include config.mk
 6
 7  ## dirs:    Create specific dirs to
      download data
 8  dirs : config.mk
 9      @pwd
10      @mkdir -p $(DATA_PATH)
11      echo "create $(DATA_PATH)"
12
13  ## dl_file: Create a list of files
      availiable, download them
14  dl_file : $(CSV_SRC) get_csv_dl.R
15      Rscript get_csv_dl.R $(CSV_SRC)
16      Rscript list_file_dl.R
      $(DATA_PATH) $(CSV_SRC)
17  .PHONY : dl_file
18
```

- New projects -> New specs

**CHALLENGE: COPY/PASTING AND MAINTAINING X SIMILAR CODE SOURCES**

- Spoiler: Not everyone like SQL

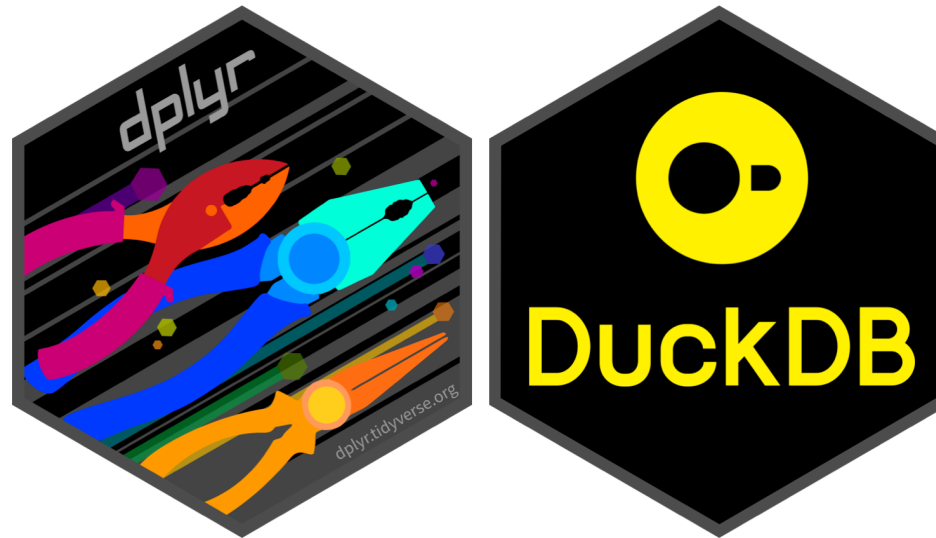**CHALLENGE: MAKE IT MORE ACCESSIBLE**

- Out of memory -> Process in DB

**CHALLENGE: SLOW, LACK OF**

💡 **LIGHT FROM** POSIT::CONF

Workshop on DuckDB (Parquet / Arrow / OLAP)

# SLACK IS BUZZING!

> **Warning**
>
> Contains real slack messages with typos!

**Olivier Leroy** 6:16 PM
I have the csv

TY 1

I am stuck with my touch pad copy pasting in slack

😂 1

# DUCKDB SOLVED THE LACK OF REACTIVITY!

# ... CODE PACKAGED

| Visual | Code |
|--------|------|

```
library(cori.data.fcc)

dir <- "data_swamp/nbm/"

get_nbm_release()

system(sprintf("mkdir -p %s", dir))

dl_nbm(
  path_to_dl = "data_swamp/nbm",
  release_date = "June 30, 2023",
  data_type = "Fixed Broadband",
  data_category = "Nationwide",
)
```

0:00 / 0:17

- DuckDB: in-process

- Reduce errors

- Complexity is in **one** place and **abstracted** so we can focus on what bring **value**!

CORI Center on Rural Innovation

# CORI.DATA.FCC

# WHAT IS PACKAGED CODE?

**ONE WHITE LIE: CODE THAT WORK OUTSIDE OF "MY" COMPUTER**

**BENEFITS FROM THE LANGUAGE ECOSYSTEM**

- Tested against GitHub runner: it does not just work on my computer!

- Website to share it/distribute it

# OPEN SOURCE: WHY SHARING IT?

- Way to **publish**: we get more value that way
- Way to **share**: code is used!
- A place were we can **exchange** in a constructive way

CORI Center on Rural Innovation

# HOW TO GET IT:

The source code is hosted on GitHub (Version control)

You need the R package {remotes}.

```r
1  install.packages("remotes")
2  remotes::install_github("ruralinnovation/cori.data.fcc")
```

🚧Check the version you have and see what new versions are offering 🚧

```r
1  packageVersion("cori.data.fcc")
```
```
[1] '0.0.1'
```

Center on Rural Innovation

## WHAT CAN IT DO: CHOOSE YOUR OWN ADVENTURE!

- I need to go back to the source: Download data from FCC

- I need primary data but working with hundreds of CSV is not for me: Download raw data from CORI (NBM / Form 477)

- Census block is perfect for me: Download tranformed data for NBM from CORI (ISP / County)

- Guide you in those steps!

# PACKAGING CODE = EXTRA TEAM MEMBER 👷

> we've talked about how packages can act like team members such as the **IT Guy**, **Analyst**, Tech Lead, or Project Manager.[1]

> [...] developed packages are [...] extra expert team members.[2]

1. https://www.emilyriederer.com/post/team-of-packages/#collaboration

2. https://milesmcbain.xyz/posts/data-analysis-reuse/

CORI Center on Rural Innovation

# PACKAGING CODE = INCREASE SPEED TO INOVATION / INSIGHT

- Functions designed to get data and load it in your environment

- **Capitalize on it**!

  - New project that is requiring FCC staff block estimates -> add it!

Center on Rural Innovation

# EXAMPLES OF USES CASES

# FIRST EXAMPLE: BLOCK COVERED BEAD LIKE

Place holder for Camdem blog post

# WHAT ARE THE ISP IN OHIO?

Visual | Code

Show 5 ∨ entries

Search: 

| Provider Name | Cnt Block | Cnt Rel. |
|---|---|---|
| Cox Communications, Inc | 2163 | 9414 |
| Ottoville Mutual Telephone Company | 270 | 992 |
| Telephone Service Company | 570 | 2516 |

CORI Center on Rural Innovation

| Provider Name | Cnt Block | Cnt Rel. |
|---|---|---|
| SYCAMORE TELEPHONE COMPANY | 404 | 1437 |
| The Champaign Telephone Company | 531 | 2287 |

Showing 1 to 5 of 121 entries

Center on Rural Innovation

# SUMMARY

- Primary data can give more power (especialy in new field)!

- (Primary) data is only useful if you can **effectively** use it.

- Use code that are abstracting the pain

- Easier to **capitalize/build** on it: generating less tech debt

- we need those extra "team members"

# CONTACTS

Website: https://ruralinnovation.us

LinkedIn | Twitter | Facebook | Instagram | YouTube

CORI | Center on Rural Innovation