

I³ Open Innovation Dataset Index

iiindex.org

Agnes Cameron, agnescam@mit.edu

December 4, 2021

Building an index for innovation data

- A place to curate lists of innovation datasets, tools and projects
- Different projects fit into different repositories: Zenodo, Dataverse, BigQuery all have different merits
- Likewise, multiple curated lists exist (NBER, LensLabs, Google Patents Public Datasets), for different needs
- The goal is not to replicate these, but create a lightweight, editable index to share + update annotated resources
 - > you can also add your own curated lists



Google BigQuery

First steps: Google Sheet

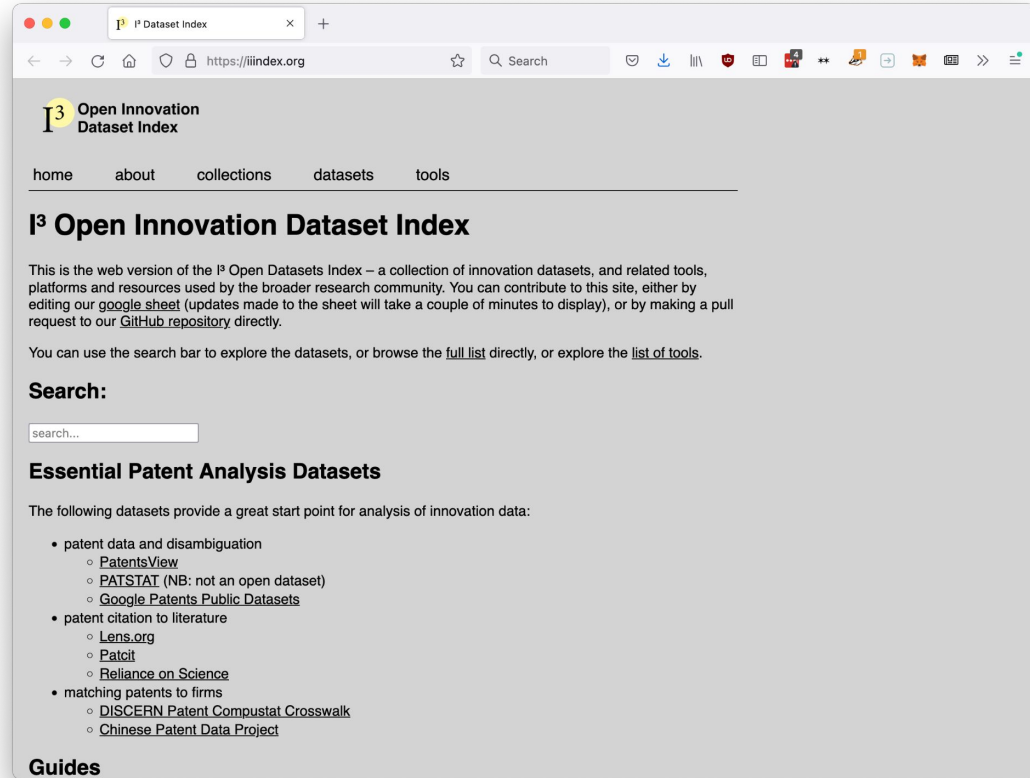
- Created Jan 2020
- Curated by I3 community
- Expanded to include tools and platforms
- Community editable
- Importantly: *still active*

I3 Open Innovation Dataset Index

</

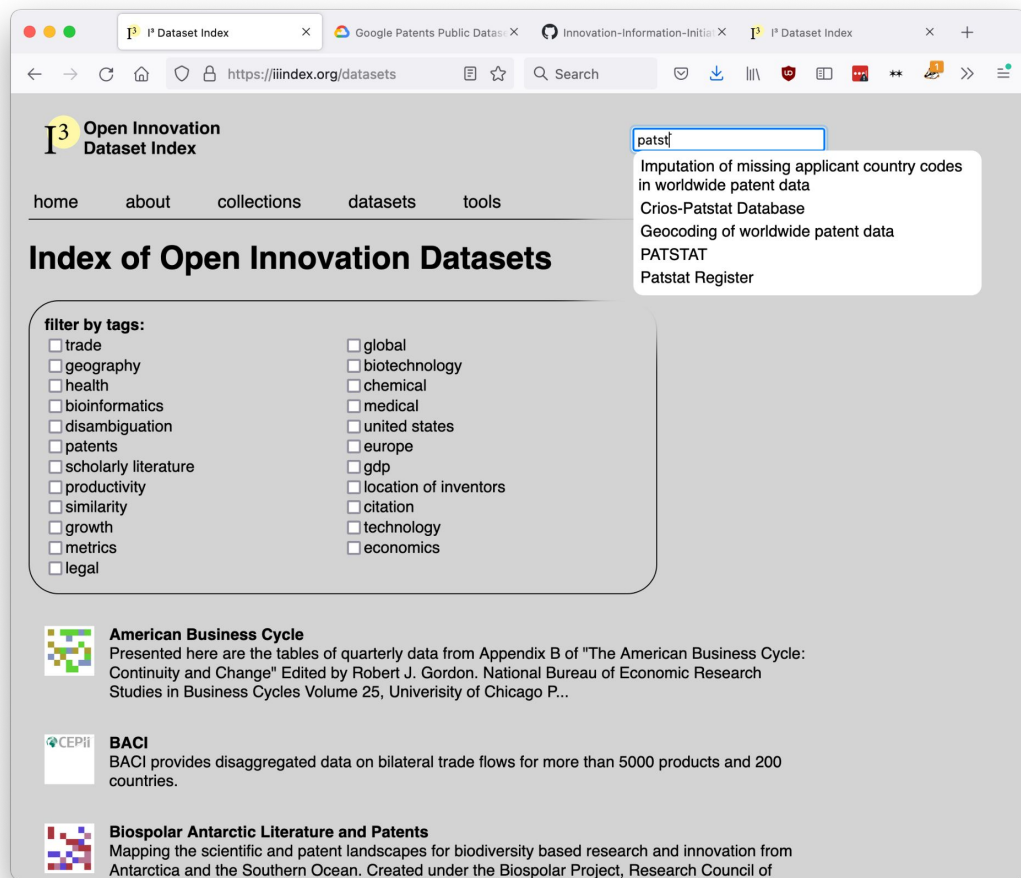
I³ Open Innovation Index Website

- *Still* editable through google sheet
- Also editable through github
- Open edits, full version history
- Searchable and linkable
- Allows for community-curated dataset indexes and guides
- Makes space for notes, annotations, links and examples



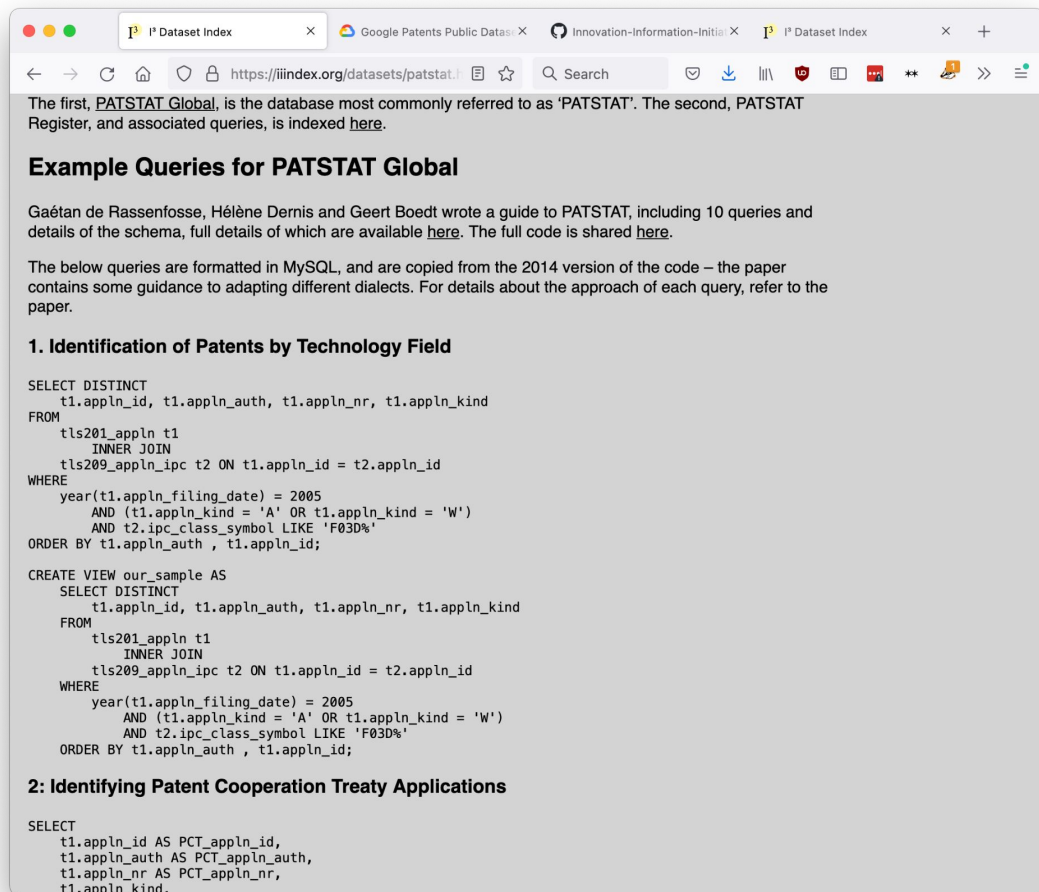
Searchable index

- Lists of tools and datasets
- Filter by common tags
- In the future: advanced search queries
 - By author
 - By timeline
 - By related datasets



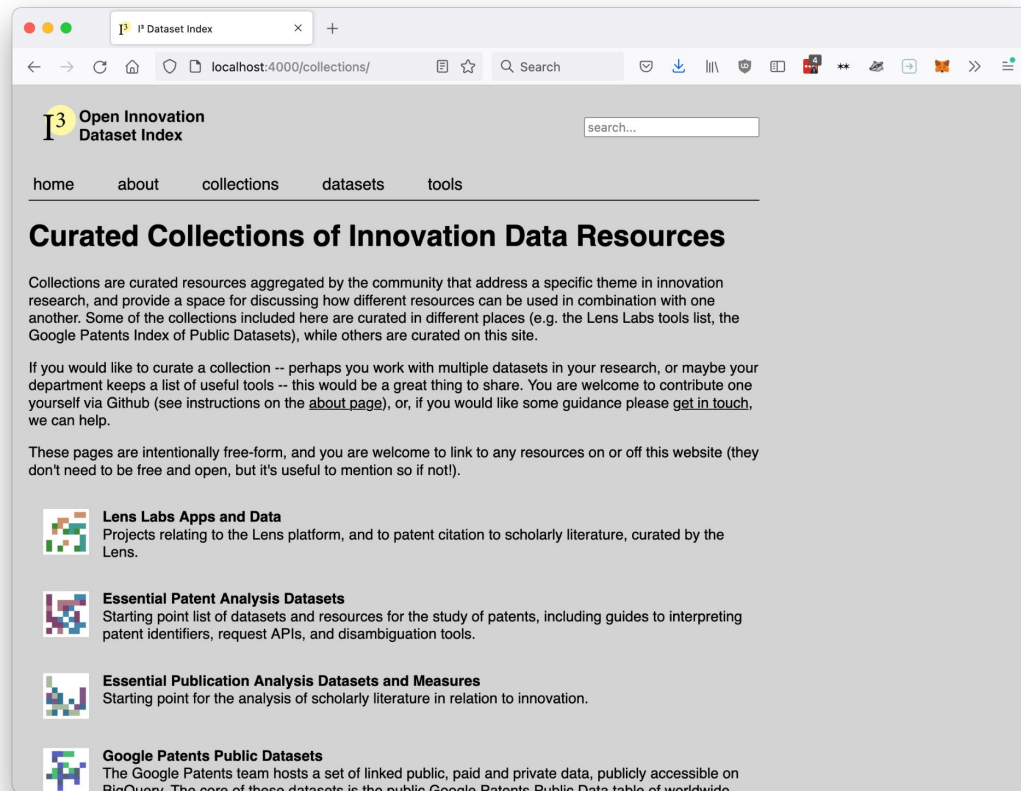
Annotatable files

- Each page in the index is an editable markdown file
- Headers include metadata from the google sheet
- Markdown body provides open-ended space to add notes and queries



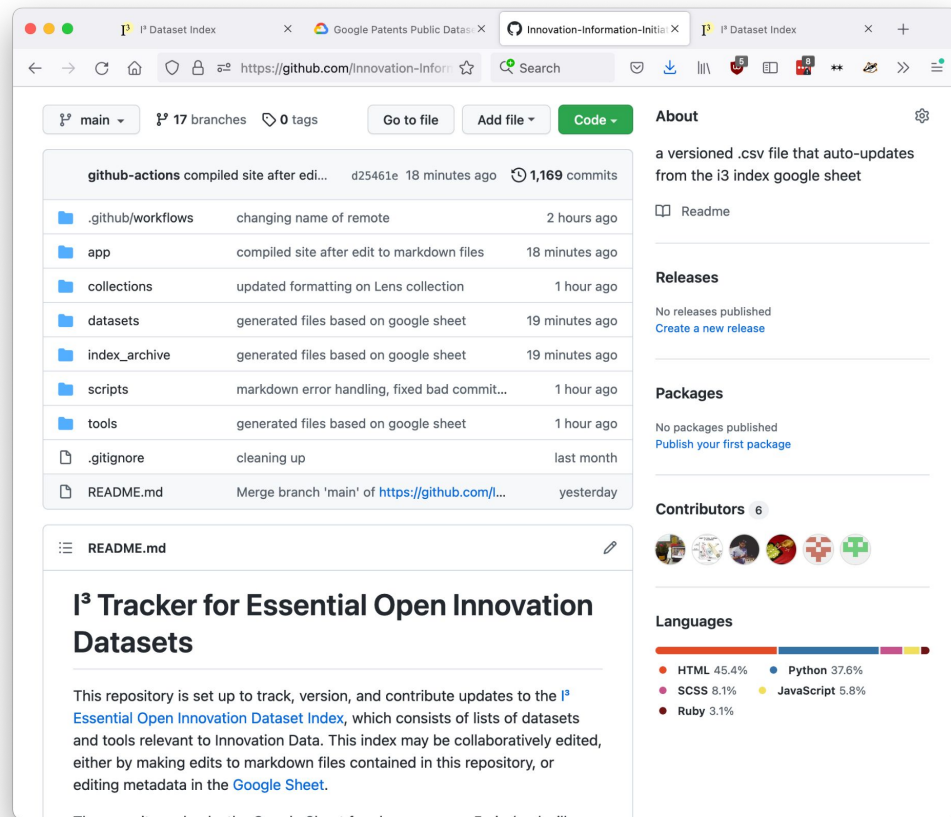
Curated collections

- An open way to list thematic data
- Some collections track external repositories (e.g. Lens Labs, NBER)
- **This is a great place to contribute!**



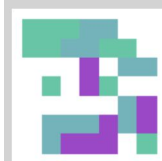
How to contribute

- Did you publish a dataset?
 - add it to the google sheet
- Do you use code to process a dataset?
 - add a sample query on Github
- Have you made a list of datasets for a project, or to give to your students?
 - create a collection
- **Guidelines for contributing** are on the [about page](#) and on the [github](#)



Future extensions

- Automated updates of metadata and versions (cron job that fetches metadata)
- Use APIs to get greater range of citation metadata, rather than just MediaWiki citations
- Expand use of 'superceded by' tag



Matching SIPO patents to Chinese listed firms ("Main Board")

page superceded by: [Chinese Patent Data Project](#)

location: <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/CF1IXQ> 

timeframe: through 2016?

tags: China, SIPO, disambiguation, patents, firms

documentation: <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/QUH8KT>

description: Matching SIPO patents to Chinese listed firms ("Main Board"). Please refer to the user documentation "Chinese Patent Database User Documentation: Matching SIPO Patents to Chinese Publicly-Listed Companies and Subsidiaries" for more details about this dataset.

Technical details

- Updates to the site are automated via Github actions
- Pull requests trigger scripts that update the information
- A versioned copy of each table in the google sheet is published in .csv form
- Citation metadata for datasets & tools is harvested automatically when submissions are made

