# Searches and schemas

Hal Varian
Google, Inc.
FESAC presentation
June 12, 2020

# Outline

- Schema.org
- Job search
- Dataset search
- Data commons
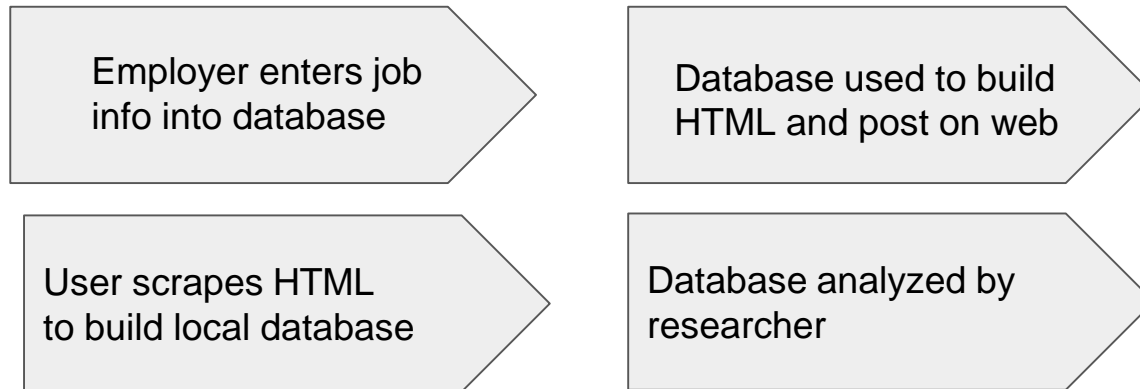- Questions

# Outline

- Schema.org
- Job search
- Dataset search
- Data commons
- Questions

# Scraping data (perhaps using **burningglass®** )

HTML is a markup language to format documents for human readers

It is now common to scrape data from the web, e.g., Computerization of White Collar Jobs. from UpJohn Institute, which scrapes online job postings.
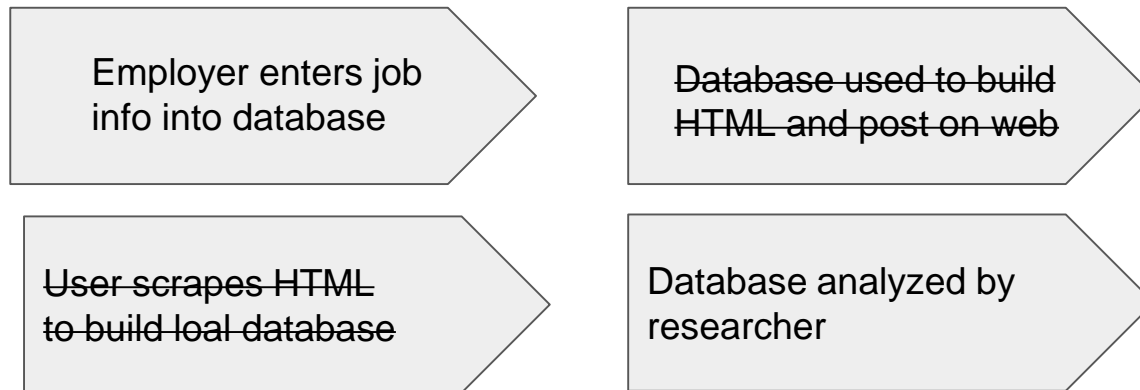
| Employer enters job info into database | Database used to build HTML and post on web |

| User scrapes HTML to build local database | Database analyzed by researcher |

# Scraping data (perhaps using burningglass )

HTML is a markup language to format documents for human readers.

It is now common to scrape data from the web, e.g., Computerization of White Collar Jobs. from UpJohn Institute, which scrapes online job postings.

| Employer enters job info into database | ~~Database used to build HTML and post on web~~ |
| --- | --- |
| ~~User scrapes HTML to build loal database~~ | Database analyzed by researcher |

- **Q:** How make human readable document also machine readable?
- **A:** Embed metadata into the HTML.
- For that we need standards for the metadata.

# schema.org

"Schema.org is a collaborative, community activity with a mission to create, maintain, and promote schemas for structured data on the Internet, on web pages, in email messages, and beyond."

- Creative works: CreativeWork, Book, Movie, MusicRecording, Recipe, TVSeries ...
- Embedded non-text objects: AudioObject, ImageObject, VideoObject
- Event
- Health and medical types
- Organization
- Person
- Place, LocalBusiness, Restaurant ...
- Product, Offer, AggregateOffer
- Review, AggregateRating
- Action
- **JobPosting**
- **Dataset**

# Schema.org is already widely used

- See CACM article on [Schema.org: Evolution of Structured Data on the Web](#)
- ~45% of pages in the Google search index have schema.org markup
- Most popular mark up use
  - [6 billion product](#) listings across 5m domains
  - [96m job postings](#) (including duplicates)
- ...and many more from [comic books](#) to [recipes](#) to [mortgages,](#) to [events](#) …
- Lesson: use schema.org standards if you want your data to be used

# Outline

- Schema.org
- Job search
- Dataset search
- Data commons
- Questions

# schema.org

## JobPosting

**Thing** > **Intangible** > **JobPosting**

A listing that describes a job opening in a certain organization.

| Property | Expected Type | Description |
|---|---|---|
| **Properties from JobPosting** | | |
| applicantLocationRequirements | AdministrativeArea | The location(s) applicants can apply from. This is usually used for telecommuting jobs where the applicant does not need to be in a physical office. Note: This should not be used for citizenship or work visa requirements. |
| applicationContact | ContactPoint | Contact details for further information relevant to this job posting. |
| baseSalary | MonetaryAmount or Number or PriceSpecification | The base salary of the job or of an employee in an EmployeeRole. |
| datePosted | Date or DateTime | Publication date of an online listing. |
| educationRequirements | EducationalOccupationalCredential or Text | Educational background needed for the position or Occupation. |
| eligibilityToWorkRequirement | Text | The legal requirements such as citizenship, visa and other documentation required for an applicant to this job. |
| employerOverview | Text | A description of the employer, career opportunities and work environment for this position. |

# Outline

- Schema.org
- Job search
- Dataset search
- Data commons
- Questions

# Dataset Search

Search for Datasets

Try coronavirus covid-19 or global temperatures.

Learn more about including your datasets in Dataset Search.

Last updated ▾ | Download format ▾ | Usage rights ▾ | Topic ▾ | Free

**100+ datasets found**

**D** Proportion of Population Aged 65 or Older

catalog.data.gov
datadiscoverystudio.org
+1more

Updated Mar 31, 2016

**F** Population ages 65 and above for the United States

fred.stlouisfed.org

Updated Sep 20, 2019

statista📊 Development of the global population aged 65 and over...

www.statista.com

Updated Mar 29, 2012

# Proportion of Population Aged 65 or Older

| Explore at catalog.data.gov | Explore at datadiscoverystudio.org | Explore at data.wu.ac.at |

*12* scholarly articles cite this dataset ([View in Google Scholar](#))

**Dataset updated** Mar 31, 2016

**Dataset provided by**

(Point of Contact)

**Description**

Variable was created as part of a set of indicators that demonstrate links between the condition of natural areas and human concerns and that quantify dependencies on resources. More information about these resources, including the variables used in this study, may be found here: https://edg.epa.gov/data/Public/ORD/NERL/ReVA/ReVA_Data.zip.

# schema.org

Home    Schemas    Documentation

## Dataset

**Thing** > **CreativeWork** > **Dataset**

A body of structured information describing some topic(s) of interest.

| Property | Expected Type | Description |
|---|---|---|
| **Properties from Dataset** | | |
| **distribution** | DataDownload | A downloadable form of this dataset, at a specific location, in a specific format. |
| **includedInDataCatalog** | DataCatalog | A data catalog which contains this dataset. Supersedes catalog, includedDataCatalog.<br>Inverse property: dataset. |
| **issn** | Text | The International Standard Serial Number (ISSN) that identifies this serial publication. You can repeat this property to identify different formats of, or the linking ISSN (ISSN-L) for, this serial publication. |
| **measurementTechnique** | Text or URL | A technique or technology used in a Dataset (or DataDownload, DataCatalog), corresponding to the method used for measuring the corresponding variable(s) (described using variableMeasured). This is oriented towards scientific and scholarly dataset publication but may have broader applicability; it is not intended as a full representation of measurement, but rather as a high level summary for dataset discovery.<br><br>For example, if variableMeasured is: molecule concentration, measurementTechnique could be: "mass spectrometry" or "nmr spectroscopy" or "colorimetry" or "immunofluorescence".<br><br>If the variableMeasured is "depression rating", the measurementTechnique |

# Metadata about datasets

Table 3: Percentage of datasets with specific properties. Column 2 lists the source predicates for each property. Properties not listed in the table have values in fewer than 1% of the datasets.

| Property | Source predicates | Percentage |
|---|---|---|
| description | so#description, purl#description | 100.00% |
| title | so#name, purl#title | 100.00% |
| provider | so#publisher, so#provider, purl#publisher | 84.59% |
| keywords | so#keywords, dct#keyword, purl#keyword | 80.08% |
| URL | so#url, dct#accessurl, dct#landigpage | 68.30% |
| temporal coverage | so#temporalCoverage, so#temporal, purl#temporal | 45.41% |
| data download | so#distribution, dct#distribution | 44.34% |
| spatial coverage | so#spatialCoverage, so#spatial, purl#spatial | 38.69% |
| date modified | so#dateModified, purl#modified | 37.46% |
| license | so#license and so#license on so#distribution | 34.80% |
| date published | so#datePublished, purl#published | 30.83% |
| catalog | so#includedInCatalog | 29.74% |
| variable | so#variableMeasured, dct#theme | 20.90% |
| authors | so#author, so#creator | 14.12% |
| same_as | so#sameAs, rdf#same_as | 12.72% |
| date created | so#dateCreated | 9.62% |
| alternate name | so#alternateName, rdf-schema#label | 3.40% |
| is accessible for free | so#isAccessibleForFree | 3.04% |

# Outline

- Schema.org
- Job search
- Dataset search
- Data commons
- Questions

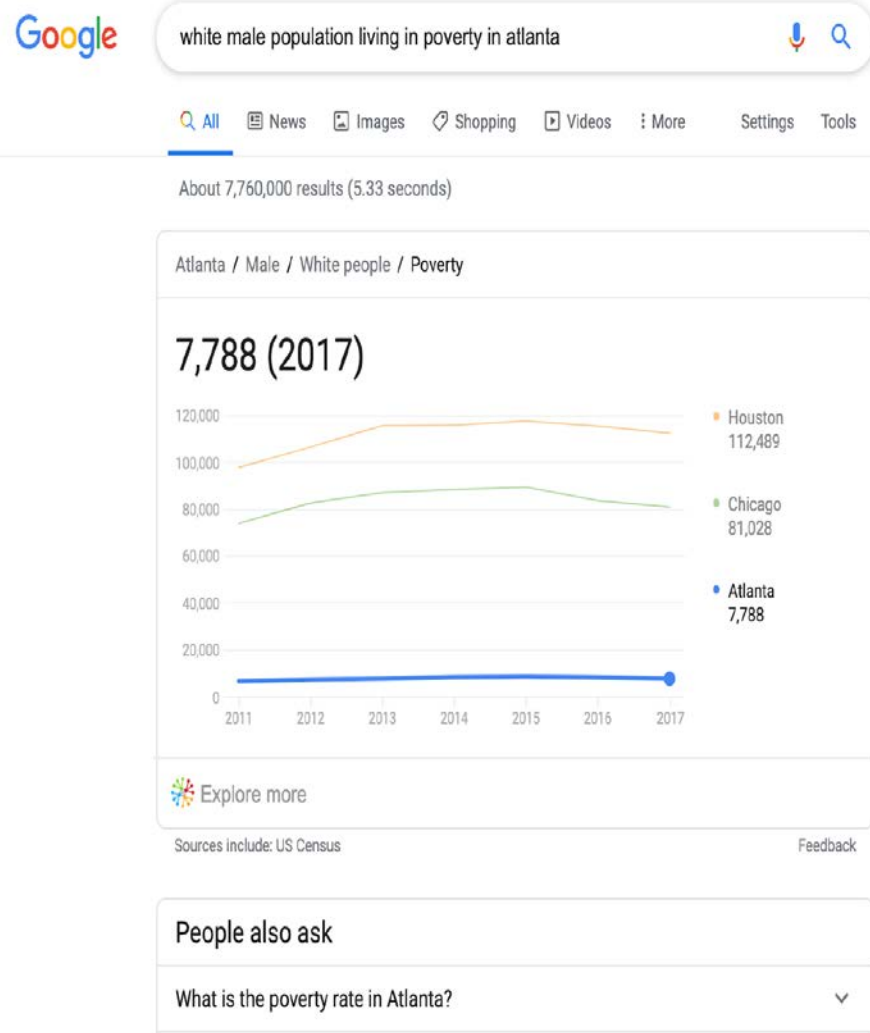# We want to make it easier to find and analyze data



From search for dataset, download, clean, normalize, join …

to

*Just ask Google*

# Coming soon: In Google Search

- **Stats seeking queries** for O(1000) variables from Demographics, Education, Employment, Health, Commerce, Disasters, Crime, and Housing

- **Line charts**

- **US data** from federal government; international data from World Bank, OECD
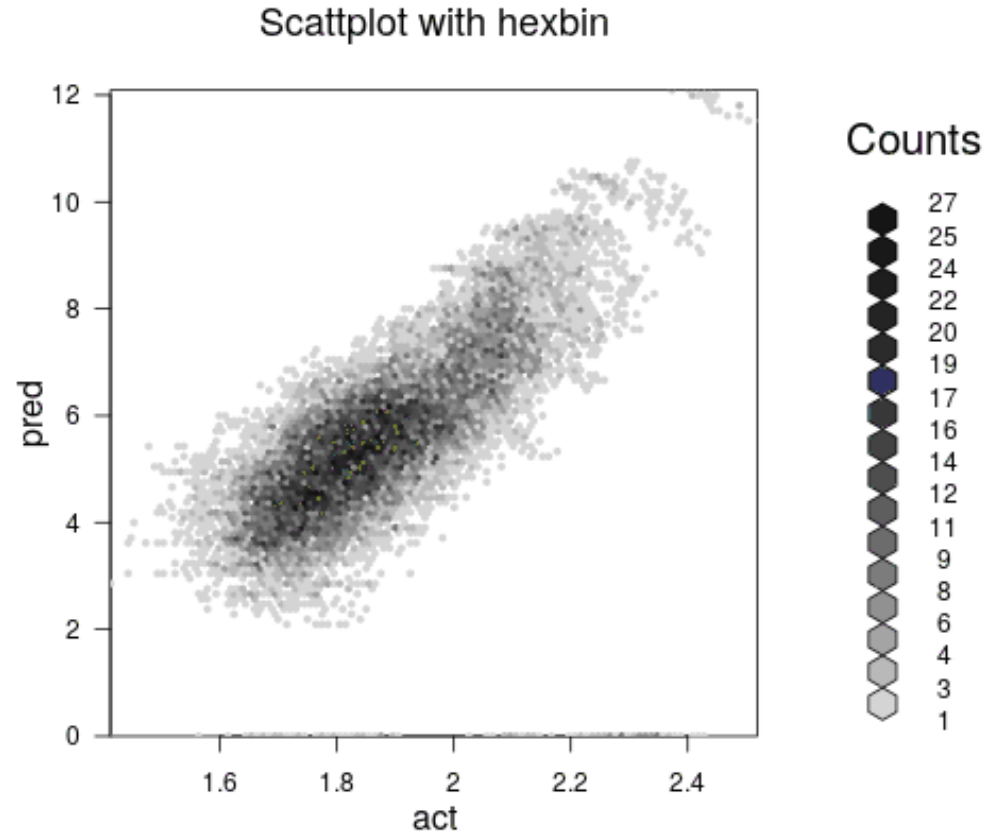
# Important goal: Joining data across providers

Sites like FRED, Census, etc. allow joining data across their data sets by hosting the data and using common standards. But it is the common standards that really matter, since this allows joins across different distinct providers.

Start with dates, move on to geos.

E.g. I recently joined "county health rankings" (125 datasets) with "JHU covid confirmed cases by county" to build a predictive model at the county level.



Scattplot with hexbin

# Outline

- Schema.org
- Job search
- Dataset search
- Data commons
- Questions

# Questions

1. What are best practices wrt to populating dataset metadata?
2. We are updating schema.org. What additional schemas do you want?
3. Who are important international partners?
4. Would "Job Trends" or "Recipe Trends" (modeled on Google Trends) be useful?