Data Commons

# Context

Data powers everything

- policy,
- journalism,
- health,
- science

How do we make it easier?

# Problem is not a shortage of data

Demographics: ACS, Housing Survey, Community Survey,

Economics: BLS, BEA, World Bank, OECD, …

Health: CDC Wonder, CDC Diabetes, County Health, …

Climate: NOAA, Hurricanes, Flooding, …

Genomics: NCBI, ENCODE, …

Too many formats, schemas, …

# Using Data: current model

Forage for data sources, track down assumptions/provenance

Clean it up, compile data sources, figure out storage, …

High upfront costs, sparse ecosystem, few tools, …

# Google maps made satellite imagery part of everyone's life

# We want to do that for data



From **search for dataset, download, clean, normalize, join ...**

to

*Just ask Google*

Journalists  Students  Search, News  Researchers  Cloud

APIS

Aggregated
Knowledge
Graph

Medicare

Sequence
data

CDC

NOAA

FBI

BLS

Census(ACS)

Wikidata

EPA

Landsat

Grid

Base
Wikidata
Schema.org
NOAA
US Census

Country — instanceOf → USA
State — instanceOf → California
Alameda County
3838 Hollis St, Emeryville, CA
The Home Depot — store
Alika Live — startDate 8/11/2017 — location → Starline Social Club
Oakland — event → Live Music: Grails
City
Emeryville — temperature 60F / 65F
$396,900 ← medianHousePrice
6,263   10,080
HousingMeasurement   $109,218
EmploymentMeasurement   94.4%
2/21/2018

# Data Commons vs collection of datasets

Collections of datasets (ala NIH Data Commons). Solves the problem of finding the dataset.

But the remaining problems --- cleaning, joining … remain


Datacommons --- a single KG built by cleaning, normalizing, joining all these datasets

# Data Commons v 0.9

- People, places, …: Integrated view of Census (ACS), CDC, BLS, BEA, FEC, NOAA, DEA, DOL,… on average 5k variables for every state, county, city, zip, school district, congressional district, …
- Education: College Scoreboard, NCES, …
- Disasters: earthquakes, hurricanes, floods, fires


- Scientific Collections: Bronx Botanical Gardens, ENCODE, parts of NCBI, GTEx

# Data Commons v 0.9

Applications for 4 categories

- Researchers
- Students
- Journalists
- Google Users

# APIs

APIs in

- REST
- Python / Python Notebooks
- SPARQL
- SQL against Big Query
- Google Sheets

**Python Notebooks for students & researchers**

CO △ Case Study: Prevalence of Obesity in 500 US Cities ☆

File Edit View Insert Runtime Tools Help

Open in playground

```python
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd

import json
```

## Case Study: Prevalence of Obesity in 500 US Cities

Obesity is well known to correlate with health factors such as high blood pressure, but is also known to correlate with economic factors such as low-income, unemployment, etc [1][2]. The Center for Disease Control (CDC) provides prevalence percentages on health conditions such as obesity, high blood pressure, and high cholesterol for approximately 500 major cities in the US (e.g. San Francisco, New York, and Austin). Meanwhile, the US Bureau of Labor Statistics provides unemployment rates while the US Census provides poverty rates for most cities across the United States.
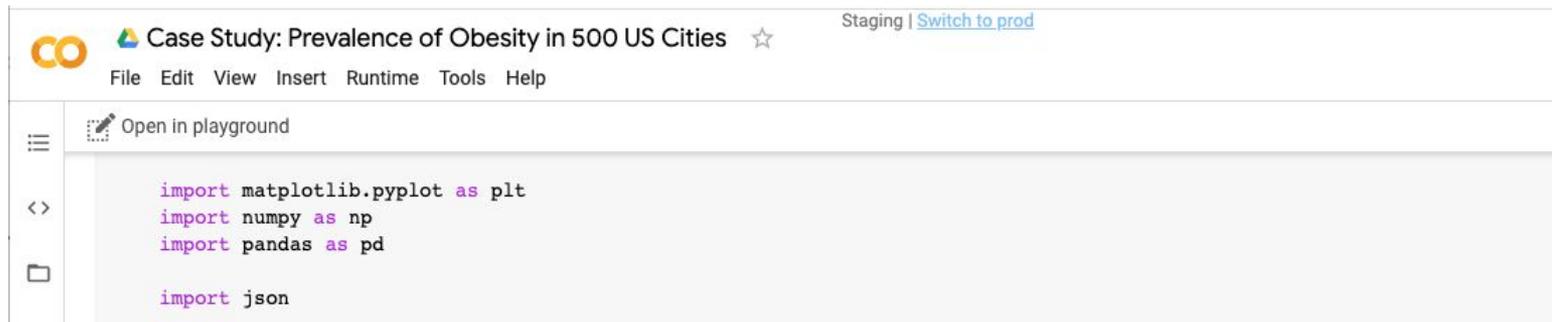
Even though these statistics come from different datasets across different government agencies with different storage formats, Data Commons surfaces each of these in a single, uniform knowledge graph. In fact, you can see this in the browser by looking at the *provenance* column. Let's use the data in Data Commons to create a linear regression model that incorporates variables:

- Prevalence of high blood pressure
- Unemployment rate
- Percent of population living with income below the poverty line

to predict the prevalence of obesity in the 500 cities that the CDC provides data for. One thing you may note is that the US Census also provides employment statistics (you can see this by navigating to the "employment" and "employmentStatus" sections for San Francisco and observing the different provenances). Our choice of using statistics from the Bureau of Labor Statistics is purely demonstrative, but it would be interesting to see if similar results can be reproduced using US Census employment statistics.

# Simple Charting Tool

**Simple Correlation Tool**

# Google Sheets API

# Biomedical Data Commons

Citizen Scientists

Students

Doctors

Researchers

Cloud

APIS

Aggregated Knowledge Graph

Base
Wikidata
Schema.org
NOAA
US Census

Country — instanceOf — USA

State — instanceOf — California

containedInPlace

3838 Hollis St, Emeryville, CA

Alameda County

8/11/2017

Alika Live — startDate

location — Starline Social Club

The Home Depot — store

containedInPlace

instanceOf — City

Oakland — event — Live Music: Grails

$396,900 — medianHousePrice — Emeryville

hasCensus

temperature — 60F

65F

description — "..." ...

2/21/2018

6,263

10,080

HousingMeasurement

$109,218

population

EmploymentMeasurement

94.4%

instanceOf

Medicare

Sequence data

WHO

CDC

dbSNP

CheMBL

UniProt

Census(ACS)

Entrez Gene

GTEx

MeSH

# Biomedical Data Commons Datasets

- CDC - 500 Cities
- CDC - Diabetes Atlas
- CDC - Wonder
- ChEMBL*
- ClinVar
- dbSNP
- Dartmouth Medicare Atlas

*Source: UCSF SPOKE

- Disease Ontology*
- FDA - Pharmacologic Class*
- GTEx
- ENCODE
- Entrez Gene
- MeSH*
- NY Times - COVID-19 Cases + Deaths

- SIDER*
- SPOKE*
- UCSC Genome Browser
- US Census - SAHIE
- UniProt*
- WHO - COVID-19 Cases + Deaths
- WHO - ICD-10 Codes

# Eg: Extract 3 data points on given variants

Given a list of genetic variants

Find:

1.  Clinical Significance
2.  Functional Category
3.  Significant Gene Associations in Whole Blood

# Current Methods to find information

**Option 1:**
**Site Search**

Use needed database browser and search all variants individually

**Option 2:**
**Download and Analyze Data**

Download data, read it into memory, and parse it for the needed info programmatically

**Option 3:**
**Data Commons**

Run 1 query on Data Commons

# More on Option 2

- I want to find out more about a list of genetic variants
- Download ClinVar and analyze clinical significance of the variants - **a few hours**
- Download dbSNP and identify the type of genomic region these SNPs are located - **a couple of days**
- Download GTEx and identify which genes a genetic variant is significantly associated in a given tissue - **a week**

# Example Data Commons Query

Issue Query to Data Commons

```
SELECT ?gv   ?p
WHERE {
      ?chr dcid "bio/hg38_chr21" .
      ?gv inChromosome ?chr .
      ?gv typeOf GeneticVariant .
      ?gv geneSymbol ?RUNX1 .
      ?gv hg38GenomicPosition ?p
 }
```

Time to run query:  21 seconds

# Many Other Commons ...

- Economics
- Covid
- Energy

….

# Technical Challenges

- Representation
- Infrastructure
- Inference


- General thoughts on KGs …

# Representational Issues

Representing time, geo, measures, provenance, …

Stat. aggregates exacerbate the problem of long predicates: maleLatinoPopulationUnder…

Triples are easy to understand, relatively widely understood

But real systems need more expressiveness

# Infrastructure

Data Commons has about 50B triples

Bio part is another 100B triples

Wide range of queries and latency requirements

- simple queries that require millisecond responses
- complex sparql queries
- everything in between

# Data Commons approach

Multiple 'backing' stores, that support different classes of queries with different latencies

- Big Query … in memory hash tables
- KG is a view on top of more native encodings
- DLG + function terms + contexts + simple kinds of inference is the real 'language' --- serves as the 'epistemological/knowledge level'

# Serving System

Datalog (Sparql, Graphql, …)

Graph Nav APIs, Time Series, etc.

Translator

Mixer

SQL

KG Schema <->
Relational Schema maps

Cache Builder

Big Query

Triples

Places

Observations

Populations

Weather

Cache

# Ending thoughts

- Vocabulary creep: Google KG, Cyc, Wikidata all have many tens of thousands of 'schema' terms. Human language manages with few thousand terms. How do we bring the compositionality of NL to KR?


- The problems of old AI haven't been solved. They just can't be expressed clearly in today's ML formalisms. e.g., World's tallest mountain yesterday?