# Advanced ML in Google Cloud (2)



Abhay Agarwal (MS Design '19)

# Agenda

- 'Productizing' analytics

- Data wrangling

- Data fundamentals

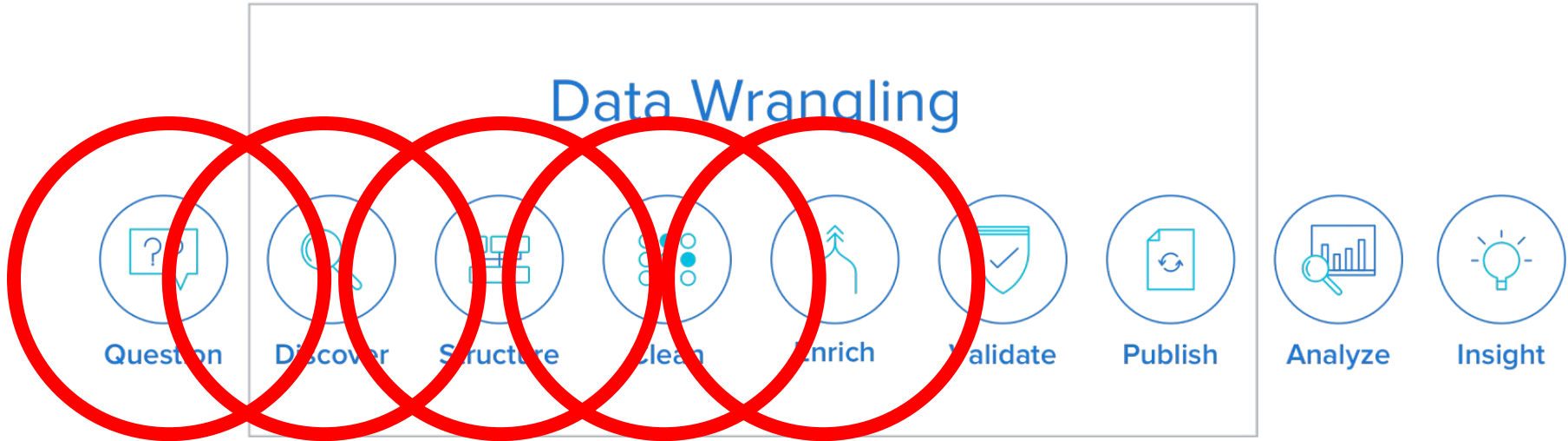- Data studio vs datalab vs colab

# 'Productizing'

- What does it mean to 'productize' your ML?

# Pitfalls in Productizing

- My algorithm has a 95% accuracy -- is it ready for production?

- My algorithm has a 95% accuracy and 95% precision -- is it ready for production?

- My  algorithm has a 95% accuracy, 95% precision, and my training data is roughly sampled from real examples -- is it ready for production?

- My algorithm has a 95% accuracy, 95% precision, training data sampled from real examples, and my algorithm tests hypotheses that match the use cases -- is it ready for production?

# Data wrangling



Data Wrangling

Question · Discover · Structure · Clean · Enrich · Validate · Publish · Analyze · Insight

# DATA COLLECTION FUNDAMENTALS

# Kev Concepts
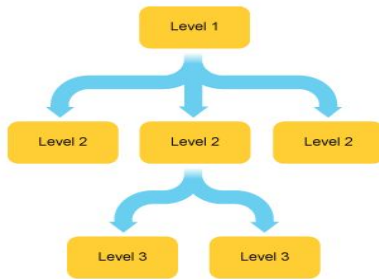
**Quantity**

**Quality**

**Cost**

**Structure**

**Freshness**

# Quantity

- Breadth
  - Number of entities or observations
  - E.g., People, companies, stars, shopping trips,…
  - Ideally: comprehensive

- Depth
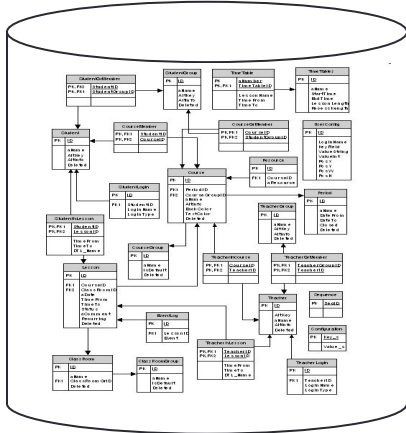  - Data gathered on each entity or observation

# Breadth and Depth

**Depth**

**Breadth**

| | Population | Surface area | Population density | Gross national income, Atlas method | Gross national income per capita, Atlas method | Purchasing power parity gross national income | | Gross domestic product | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | per capita | | per capita |
| | millions | sq. km thousands | people per sq. km | $ billions | $ | $ billions | $ | % growth | % growth |
| | 2017 | 2017 | 2017 | 2017 | 2017 | 2017 | 2017 | 2017 | 2017 |
| Afghanistan | 35.5 | 652.9 | 54 | 20.2 | 570 | 71.1 | 2,000 | 2.6 | 0.1 |
| Albania | 2.9 | 28.8 | 105 | 12.4 | 4,320 | 34.8 | 12,120 | 3.8 | 3.9 |
| Algeria | 41.3 | 2,381.7 | 17 | 163.5 | 3,960 | 621.9 | 15,050 | 1.7 | -0.1 |
| American Samoa | 0.1 | 0.2 | 278 | .. | .. | .. | .. | -2.6 | -2.7 |
| Andorra | 0.1 | 0.5 | 164 | .. | .. | .. | .. | 1.9 | 2.3 |
| Angola | 29.8 | 1,246.7 | 24 | 99.1 | 3,330 | 180.5 | 6,060 | 0.7 | -2.6 |
| Antigua and Barbuda | 0.1 | 0.4 | 232 | 1.4 | 14,170 | 2.3 | 22,980 | 3.3 | 2.3 |
| Argentina | 44.3 | 2,780.4 | 16 | 577.1 | 13,040 | 897.2 | 20,270 | 2.9 | 1.9 |
| Armenia | 2.9 | 29.7 | 103 | 11.7 | 4,000 | 29.5 | 10,060 | 7.5 | 7.3 |
| Aruba | 0.1 | 0.2 | 585 | .. | .. | .. | .. | .. | .. |
| Australia | 24.6 | 7,741.2 | 3 | 1,263.5 | 51,360 | 1,160.1 | 47,160 | 2.0 | 0.3 |
| Austria | 8.8 | 83.9 | 107 | 400.3 | 45,440 | 462.5 | 52,500 | 3.0 | 2.2 |
| Azerbaijan | 9.9 | 86.6 | 119 | 40.2 | 4,080 | 164.2 | 16,650 | 0.1 | -1.0 |
| Bahamas, The | 0.4 | 13.9 | 39 | 11.5 | 29,170 | 11.8 | 29,790 | 1.4 | 0.4 |
| Bahrain | 1.5 | 0.8 | 1,936 | 30.2 | 20,240 | 64.1 | 42,930 | 3.9 | -0.8 |
| Bangladesh | 164.7 | 147.6 | 1,265 | 242.8 | 1,470 | 664.5 | 4,040 | 7.3 | 6.2 |
| Barbados | 0.3 | 0.4 | 664 | 4.4 | 15,540 | 5.1 | 17,830 | 1.7 | 1.4 |

# Structure

**Structured**
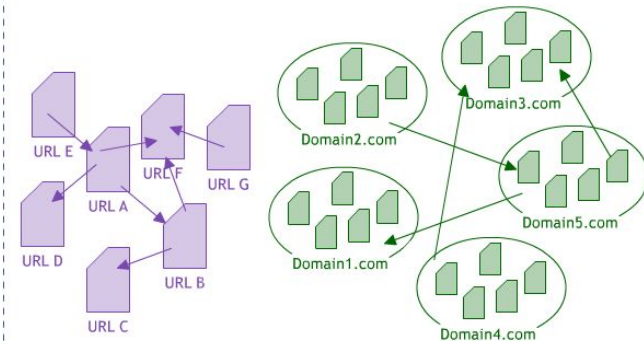
**Semi-structured**

**Unstructured**

# Graph Data

Graphs arise naturally in many settings

Many interesting techniques

e.g., Page Rank, community detection



Page vs. Domain-Level Link Graph

The page-level link graph counts links between individual URLs and values pages based on these links.

The domain-level link graph aggregates the links that exist on every page of a domain and considers only links that pass between unique sites to create domain-wide metrics.

# Data Quality

- Errors
  - E.g., human labeling mistakes
- Missing data
  - E.g., missing addresses in customer records
- Bias
  - Sample bias, measurement bias, prejudice/stereotype

# Data Quality: Sample Bias



Day Driving vs Night Driving

Tank recognition

# Data Quality: Prejudice/Stereotype Bias Algorithmic Law Enforcement


Satoshi Kambayashi



**His honour the machine**
Prisoners released on bail*
%

| | |
|---|---|
| Chosen by judges | 18.6 |
| Suggested by algorithm | 14.9 |

— *of which:* re-offend†

*From a representative sample of the US Department of Justice database 1990-2009
†Failure to appear in court and re-arrest before trial

Source: Jens Ludwig, University of Chicago

Economist.com

But what about perpetuating bias against minorities?

# Data Quality: Measurement Bias

# Data Freshness

Rate of data collection must match rate of change of underlying phenomenon

# Data manipulation in Google Cloud

- Data Studio

- Datalab

- Colab

- (offline!)

# Data Studio

- Data Studio - glorified spreadsheets with a few integrations to Google Cloud to pull data

- Use cases: excel-like functions, simple visualizations (e.g. geographic)

# Datalab

- Datalab - hosted Jupyter instance with preset libraries

- Use cases: python scripting, visualization, ML pipelining, some long-running scripting, versioned scripts and models

# Colab

- Colab - Shared, no-setup version of Datalab that is designed around sharing

- Use cases: creating publicly accessible work, collaboration, but no long-running scripting