# Gender Profiling in WikiData

Mar 14, 2024
Queer Data Days
Isabella Lu, Lane Rasberry, Peter Alonzi

DATA SCIENCE

ON
OMIDYAR NETWORK
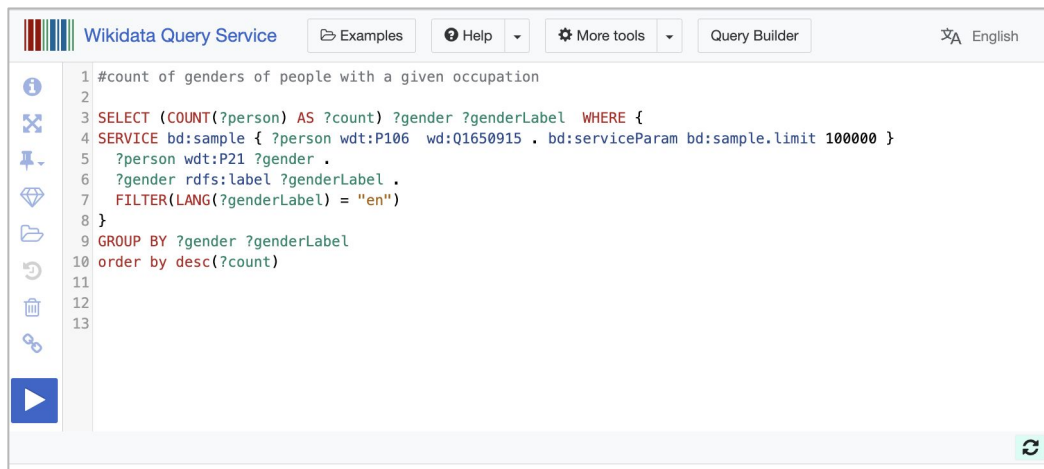
WILLIAM + FLORA
Hewlett Foundation

# Goal

Objective

- Showcase Wiki's capacity to encompass representative data of society
- Profiling occupations by gender

Not:

- Expose demographic trends to advance an agenda

# Wikidata Query

- SPARQL
  - RDF (Resource Description Framework) query programming language
  - Subject-predicate-object triples
- Purpose:
  - Used to retrieve and manipulate RDF datasets
  - Query through Wiki Database
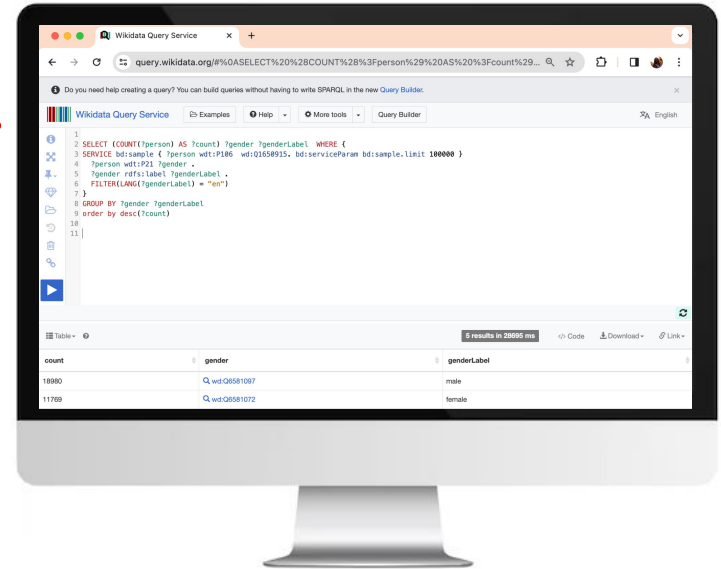- Query Service
  - Limited > Use Samples

# Accessing WikiData

# EXAMPLE SPARQL QUERY:

Count of genders of people with a given occupation

```
SELECT (COUNT(?person) AS ?count) ?gender ?genderLabel  WHERE {

SERVICE bd:sample { ?person wdt:P106  wd:Q1650915. bd:serviceParam
bd:sample.limit 100000 }

   ?person wdt:P21 ?gender .

   ?gender rdfs:label ?genderLabel .

   FILTER(LANG(?genderLabel) = "en")

}

GROUP BY ?gender ?genderLabel

order by desc(?count)
```

TRIPLE PATTERN: outlines what you are querying for

RANDOM SAMPLE: causes limitations

MATCH: values with their respective values

ORDER: determines how you want to order your query

# Gender Distribution of People by Occupation

(Researcher)

female: 38.5%

male: 61.5%

Github

# Application

Prevalence in Education Industry

# Why Is This Important?

Pew Research Center

"**Incomes are rising in the U.S., but the increase is not being felt equally by all Americans. "**

How is income inequality being pictured in Wikidata?

Demographic

# Gender Gap

Forbes : "In 2022, women earned 17% less than men on average."

**Women's Estimated Average Pay for Every Dollar Earned by Men, by Industry, 2021**

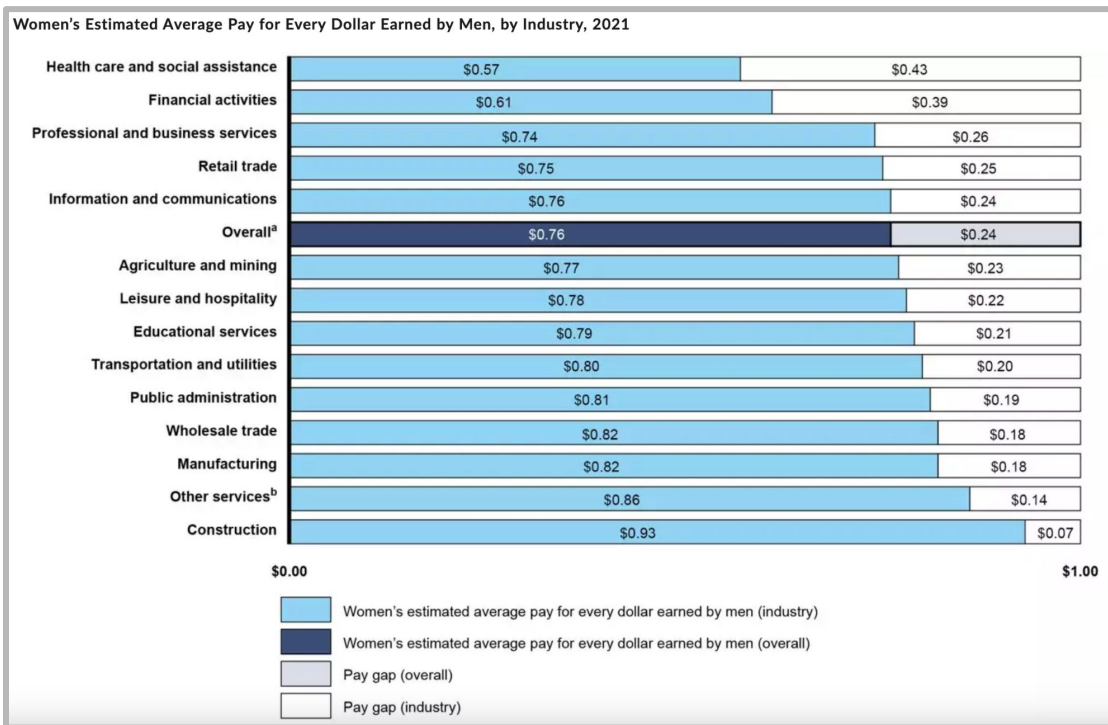| Industry | Women's pay | Pay gap |
|---|---|---|
| Health care and social assistance | $0.57 | $0.43 |
| Financial activities | $0.61 | $0.39 |
| Professional and business services | $0.74 | $0.26 |
| Retail trade | $0.75 | $0.25 |
| Information and communications | $0.76 | $0.24 |
| Overall[a] | $0.76 | $0.24 |
| Agriculture and mining | $0.77 | $0.23 |
| Leisure and hospitality | $0.78 | $0.22 |
| Educational services | $0.79 | $0.21 |
| Transportation and utilities | $0.80 | $0.20 |
| Public administration | $0.81 | $0.19 |
| Wholesale trade | $0.82 | $0.18 |
| Manufacturing | $0.82 | $0.18 |
| Other services[b] | $0.86 | $0.14 |
| Construction | $0.93 | $0.07 |

$0.00                                                    $1.00

- ▢ Women's estimated average pay for every dollar earned by men (industry)
- ▢ Women's estimated average pay for every dollar earned by men (overall)
- ▢ Pay gap (overall)
- ▢ Pay gap (industry)

GAO

Highlight:

In Educational Services, women make $0.79 for every dollar earned by men on average

9

# Data Outside of Wiki

NCES: National Center for Education Statistics

- Women are the majority



Figure 1. Percentage distribution of teachers in public elementary and secondary schools, by instructional level and sex: School year 2020–21
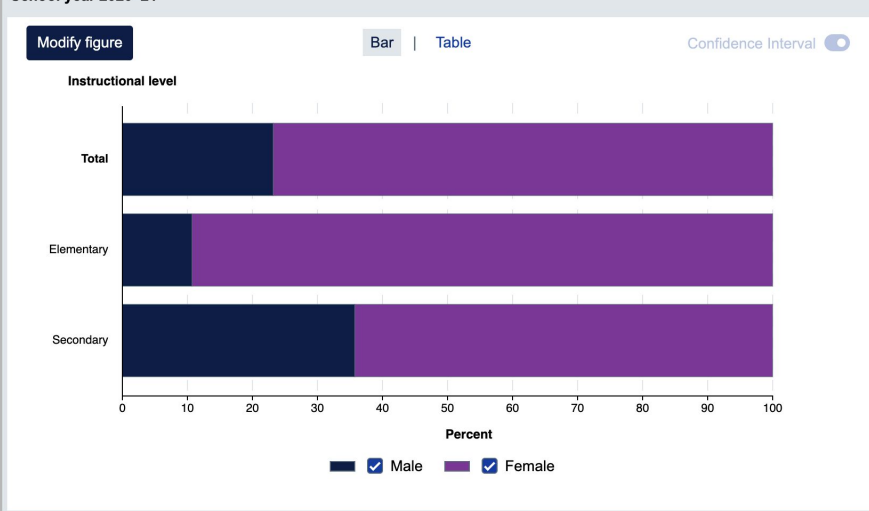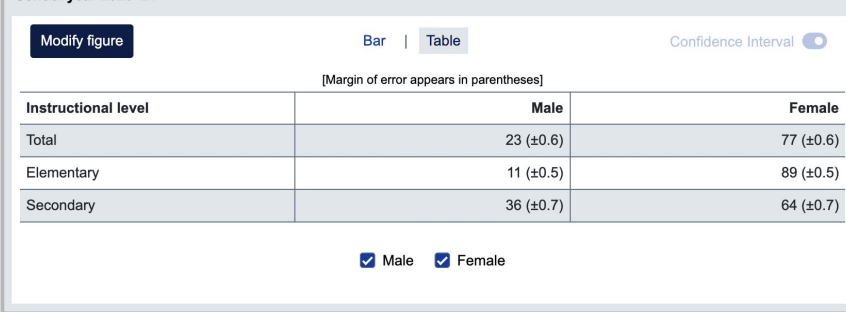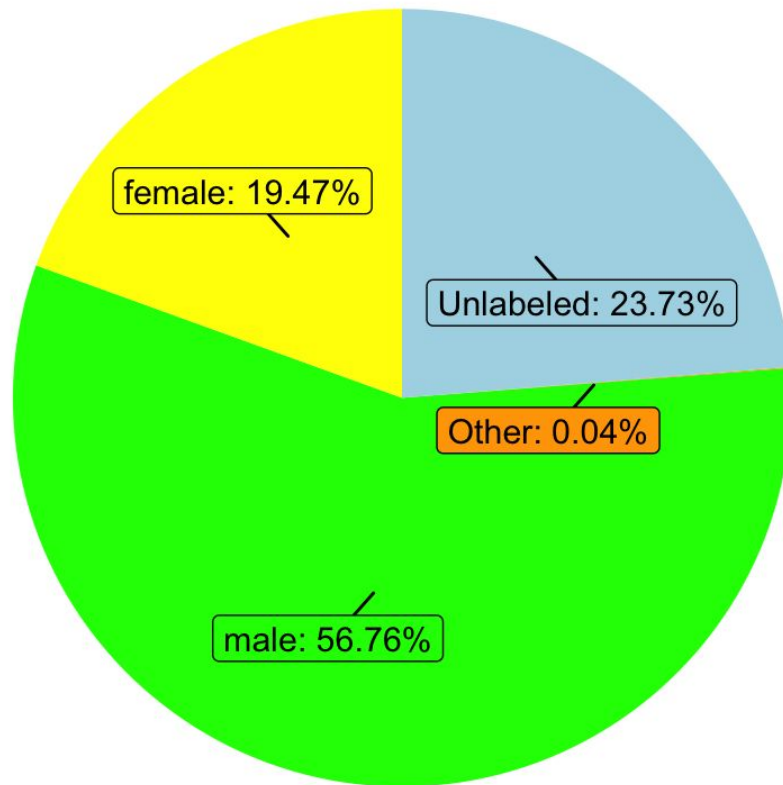
| Instructional level | Male | Female |
|---|---|---|
| Total | 23 (±0.6) | 77 (±0.6) |
| Elementary | 11 (±0.5) | 89 (±0.5) |
| Secondary | 36 (±0.7) | 64 (±0.7) |

[Margin of error appears in parentheses]

# Gender Distribution of Teachers

Demographic Statistics:

- Women are NOT the majority
  - Men: 56.76% ± .1%
  - Women: 19.47% ± .1%

Data Query Details:

- Sample size 100,000
- Random sample of WikiData
- Extra: Those with no gender label

Github



female: 19.47%

Unlabeled: 23.73%

Other: 0.04%

male: 56.76%

# Wikidata Possible Instability

Takeaways:

- Challenges arise with data gaps within WikiData
- Caution is needed when interpreting statistical analysis from data
  - Misinterpretation
- Data may not be fully representative of the entire population

Challenge Image

# Moving Forward

Next Steps:

- Run with larger sample sizes for more accurate analysis
- Create more analyses on different demographic areas of interest
- Filter more succinctly: obtain accurate timely data



Target Image

# Thanks