

Welcome Data Science for the Public Good 2018 Cohort



May 29, 2018

Sallie Keller
Professor of Statistics and Director



BIOCOMPLEXITY INSTITUTE
VIRGINIA TECH

SDAL SOCIAL &
DECISION ANALYTICS
LABORATORY

BIOCOMPLEXITY

INSTITUTE

We are in an ***ALL*** Data Revolution

A new lens for social observing

Infrastructure



- Condition
- Operations
- Resilience
- Sustainability

Environment



- Climate
- Pollution
- Noise
- Flora/ Fauna

People

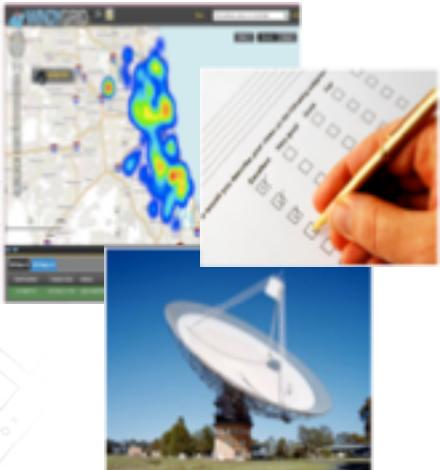


- Relationships
- Location
- Economic Condition
- Communication
- Health

It is time to leverage **ALL** the data sources

Local, State/Provence, and Federal

Designed Data



Administrative Data



Opportunity Data



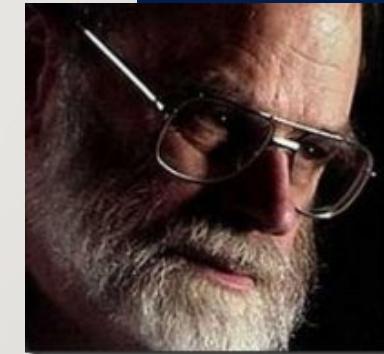
Procedural Data



WHAT IS DATA SCIENCE?

Fourth paradigm

“... change of all sciences moving from observational, to theoretical, to computational and now to the 4th Paradigm - Data-Intensive Scientific Discovery”

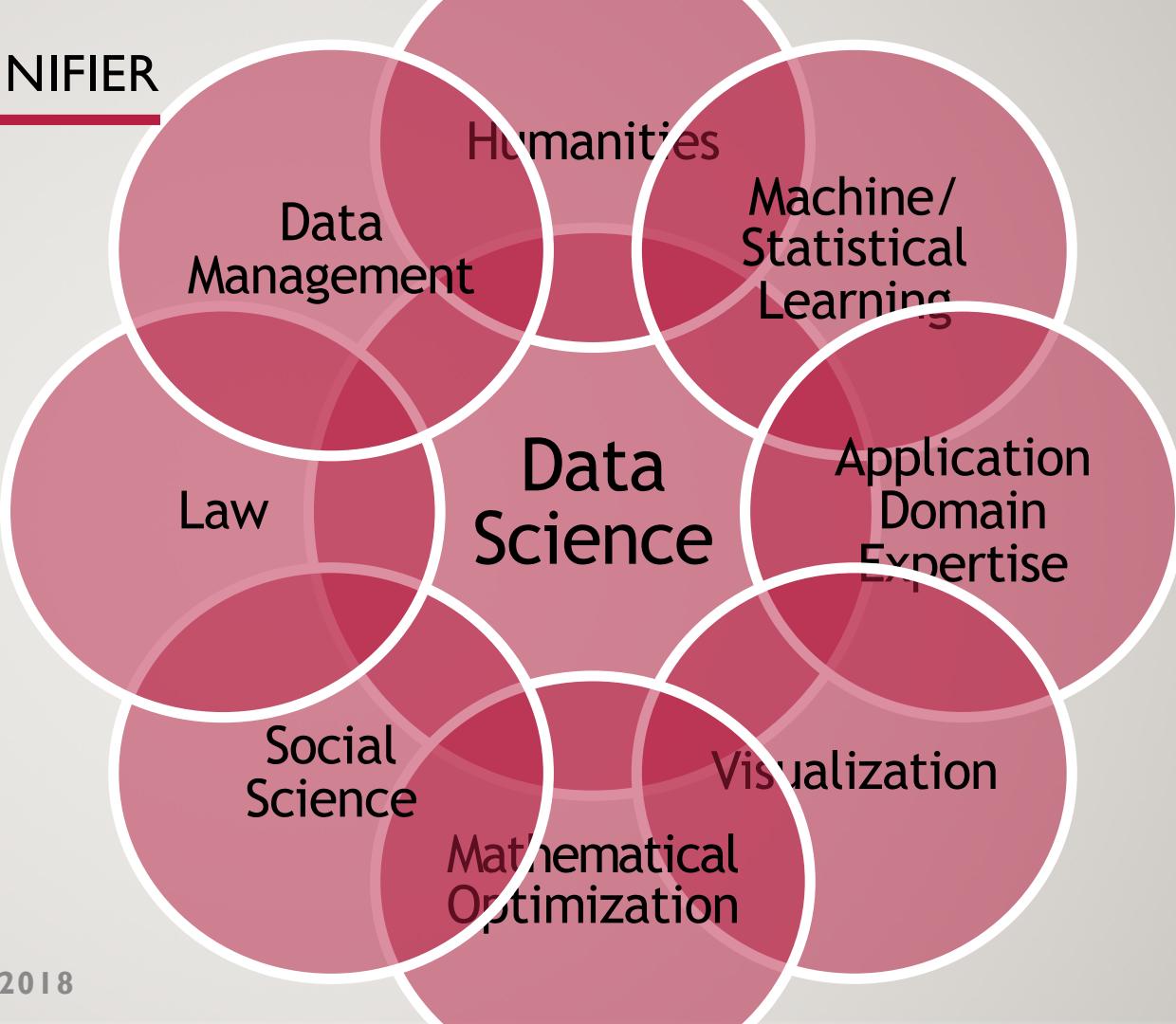


Jim Gray

WHAT IS IMPORTANT?

Need to solve a real problem using data...
No applications, no data science.

DATA SCIENCE AS A UNIFIER



In our lab data science is policy focused on other people's problems



NCSES

National Center for Science and Engineering Statistics

MITRE

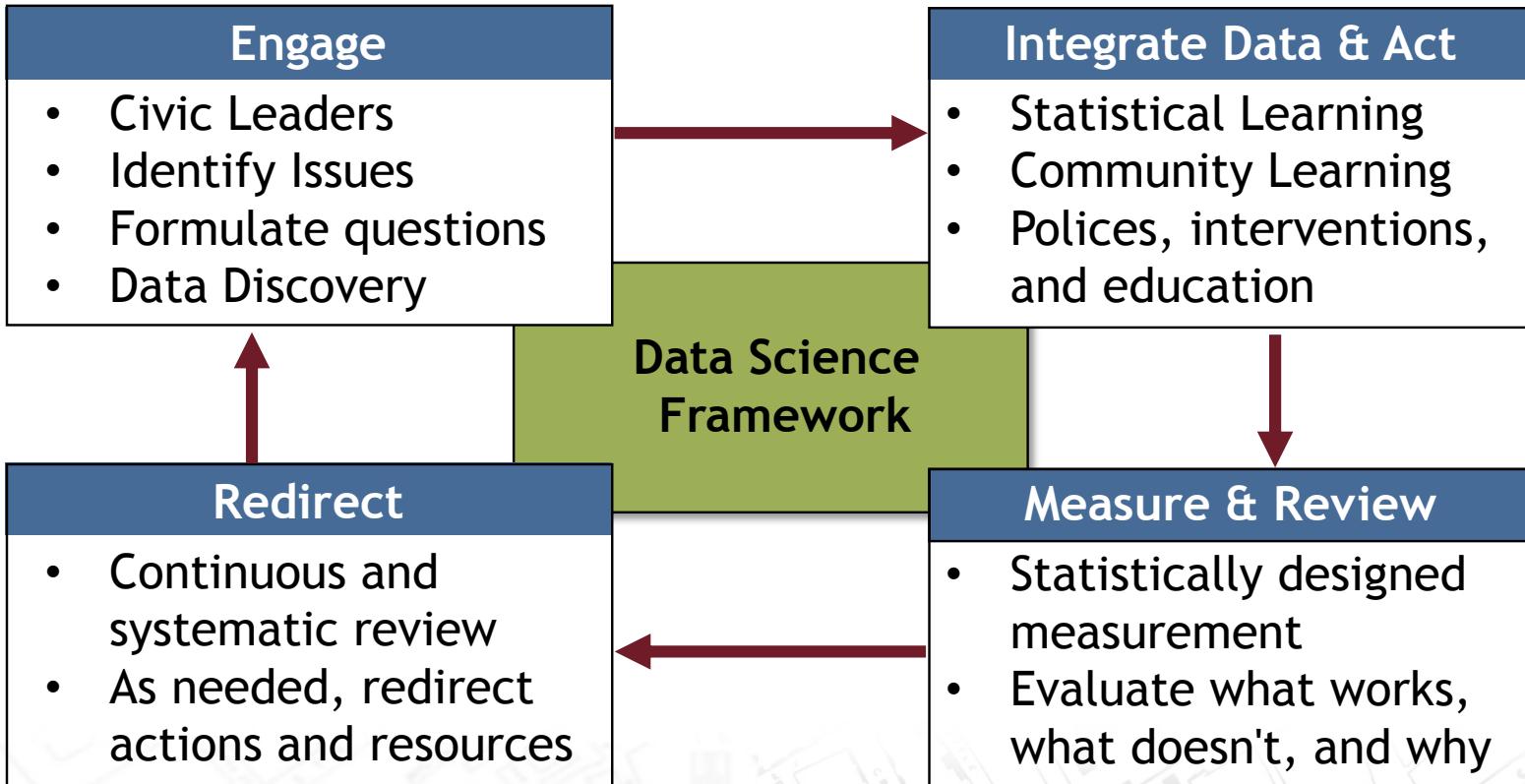


Local / State Government
Federal Statistical Agencies
Department of Defense
Industry

Enhancing Prosperity through Data Science



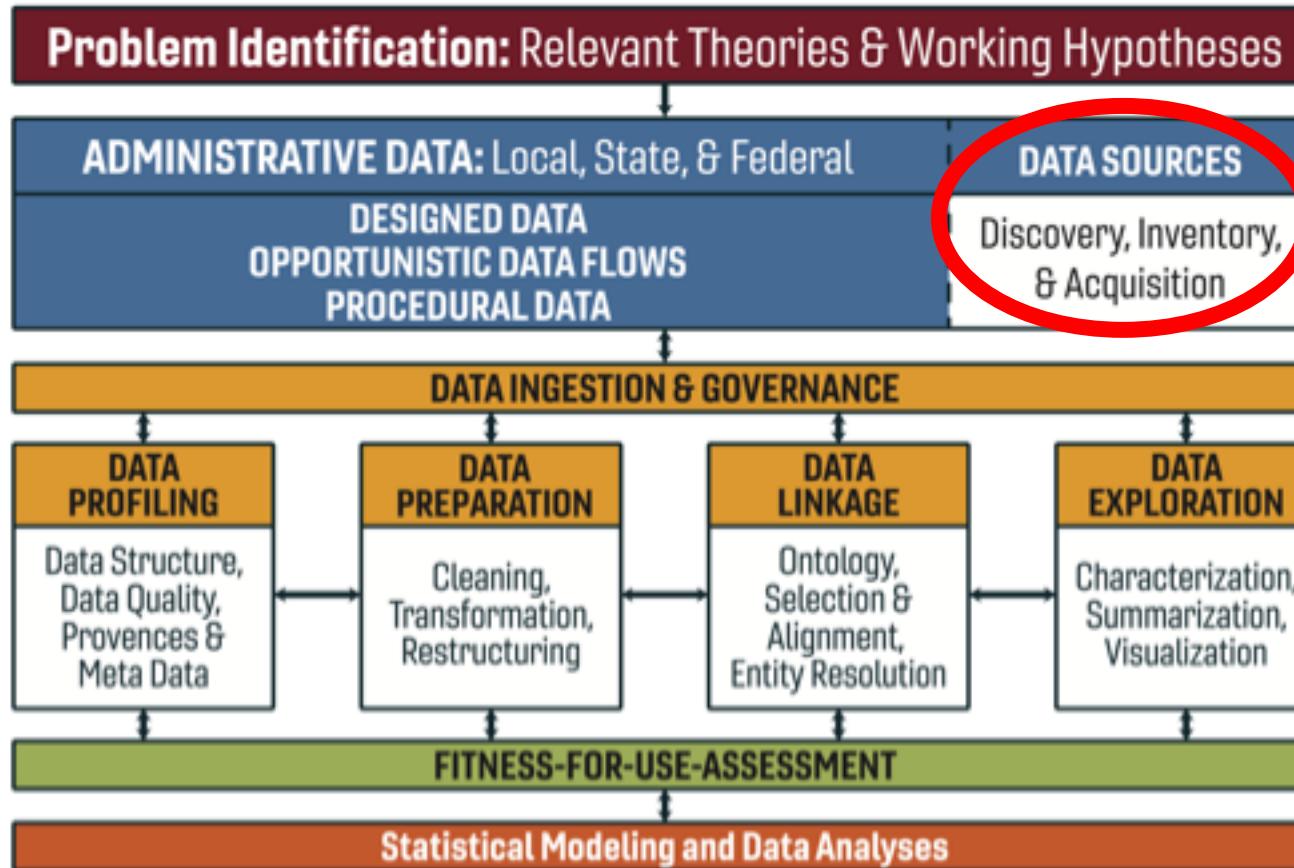
Community Learning through Data-Driven Discovery (CLD3)



Common CLD3 themes need data science training to address

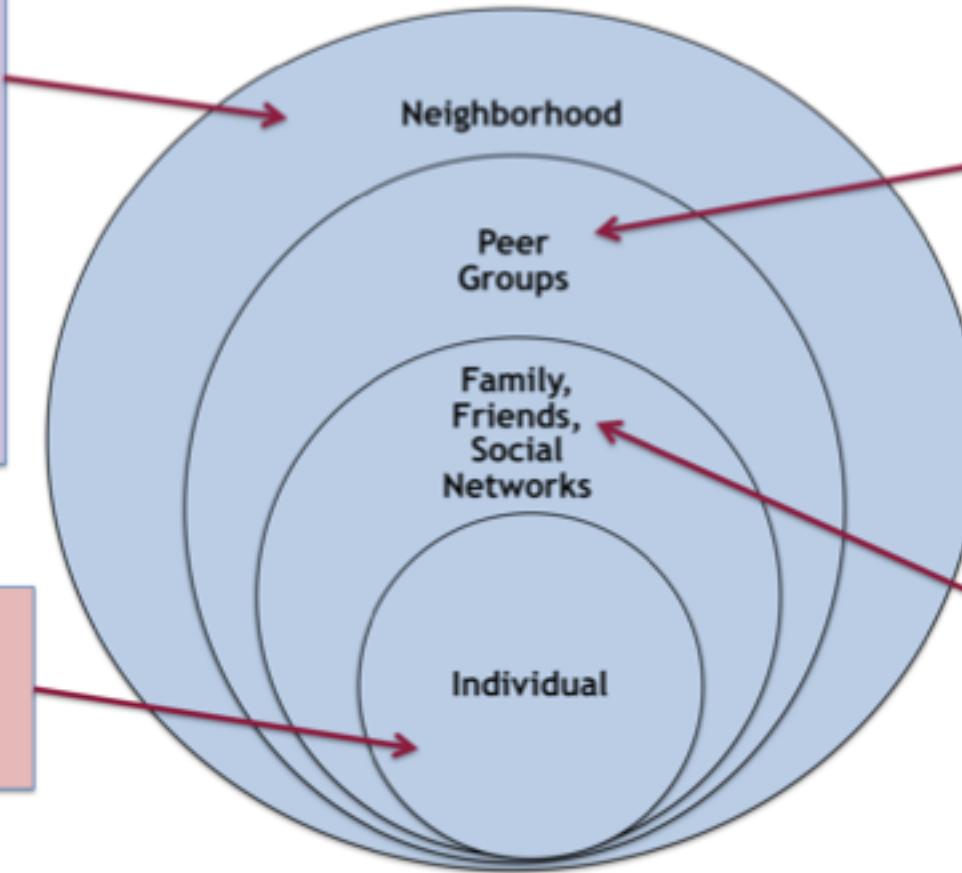
- Locating and describing a population within a community
- Estimating a statistical summary and its margin of error to evaluate its usefulness for the purpose at hand
- Forecasting future needs
- Evaluating a program, policy, or standard operating procedure

Data Science Framework



Local community Data Map

- Access to healthy food - grocery stores, community gardens, farmers markets, restaurants (fast food, other)
- Living Conditions
- Personal Safety
- Engagement
- Support Networks



- Education
- English Literacy
- Health Literacy
- Engagement
- Support Networks

- Behavioral Health
- Physical Health
- Social Wellness
- Support Networks

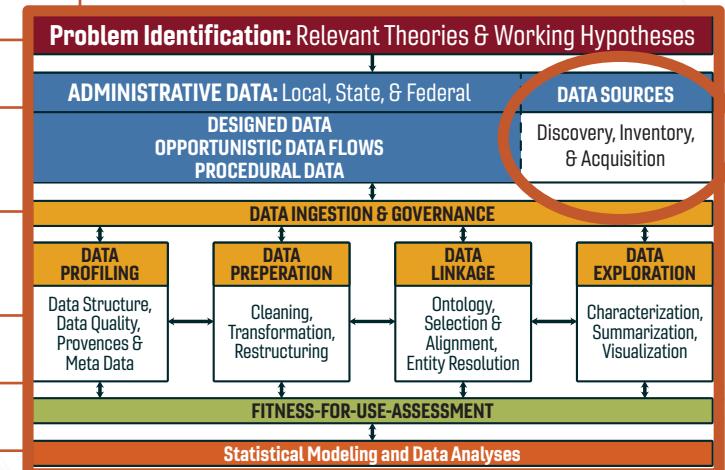
- Family Stability
- Income Stability
- Living Conditions
- Health Literacy
- Support Networks

Data Discovery, Inventory & Acquisition

Data Source	Geography
American Community Survey data (Census), 2011-2015 (updating now to 2012-2016)	Census Tracts and Block Groups
American Time Use Survey (BLS), 2017	National
Youth Risk Behavior Surveillance System, 2015	State
County Health Rankings, 2017	County
Built Environment, e.g., Grocery stores, SNAP retailers, recreation centers, community gardens	Address Level
Fairfax real estate tax assessment data	Address Level
Fairfax Open data: Zoning, Environment, water, Parks, Roads	Shapefiles
Fairfax County Youth Survey, 2016 8 th , 10 th , 12 th graders	High School Attendance Area
Virginia Department of Education, 2017	High School
National Center for Education Statistics, 2014-2015	High School
Center for Disease Control, 2014-2015	High School

Initial data sources used with geographic specificity

- All are **updated** as new data are available



Data Discovery, Inventory, & Acquisition

High School

Postsecondary Education

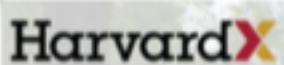
Credentials and Skill-based Training

Work Experience & STEM Occupations

Formal Education



Credentials & Skill-based Training



Community

County Health Rankings & Roadmaps

Building a Culture of Health, County by County

Job Postings & Resumes



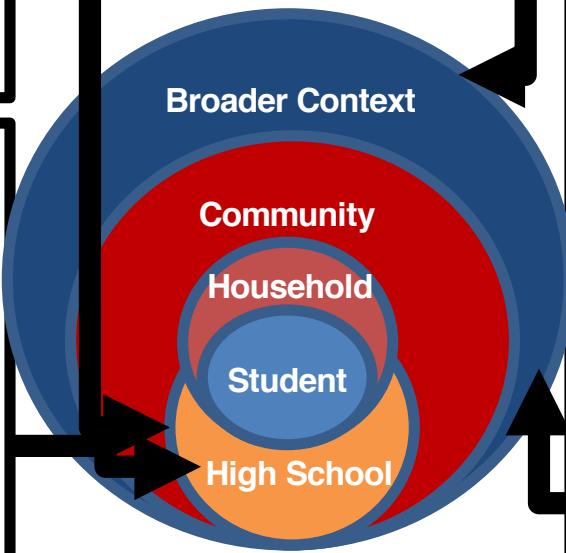
Data Map

Community Characteristics

- % Population w/ Postsecondary Ed (ACS)
- % Households on SNAP (ACS)
- % Households with limited English proficiency (ACS)
- % Employment opportunities by education requirement (Open Data Jobs)
- % Employment opportunities by experience level (Open Data Jobs)

High School “Postsecondary-Going” Culture

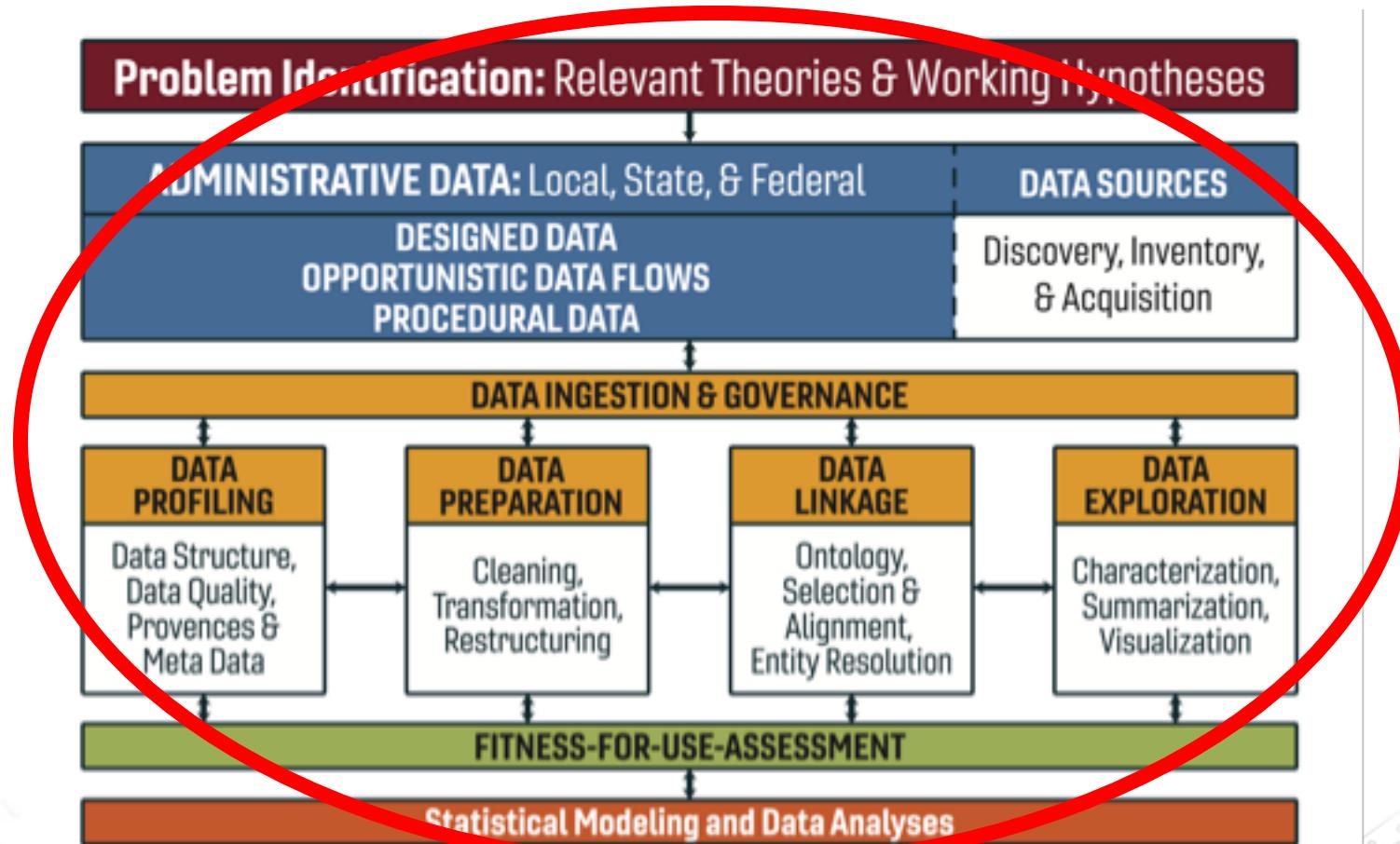
- Graduation rate (VDOE)
- Advanced/regular degree ratio (VDOE)
- % CTE program graduates (VDOE)
- College application rate (SCHEV)
- College acceptance rate (SCHEV)
- % Enrolled in AP classes (VDOE)
- % Passed AP tests (VDOE)
- % in Dual Enrollment courses (VDOE)
- % Teachers w/ graduate degrees (VDOE)
- % Students took the SAT (College Board)
- Mean SAT scores (College Board)
-



Perception of Postsecondary Availability

- Number of vocational schools, colleges, and universities in geographic area (IPEDS)
- Cost (tuition, fees, room and board, financial aid) of colleges in geographic area (IPEDS)
- Acceptance rate/college selectivity of colleges (IPEDS/SCHEV)
- College “choice set” of peers (SCHEV)
- College enrollment rates of students within school district (SCHEV)

Data Science Framework



Parks and Recreation Participation

Locating populations and Evaluating a Policy



Issue: In **Arlington Virginia**, the Department of Parks & Recreation (DPR) policy provides fee reductions to households that meet certain economic criteria

Goal: DPR wants to ensure that **fees are not a barrier** to program participation by targeting outreach

The DATA

Household Level

- Department of Parks & Recreation 2016 enrollment data for all programs
- Geocoded household locations

Census Tract Block Group Level

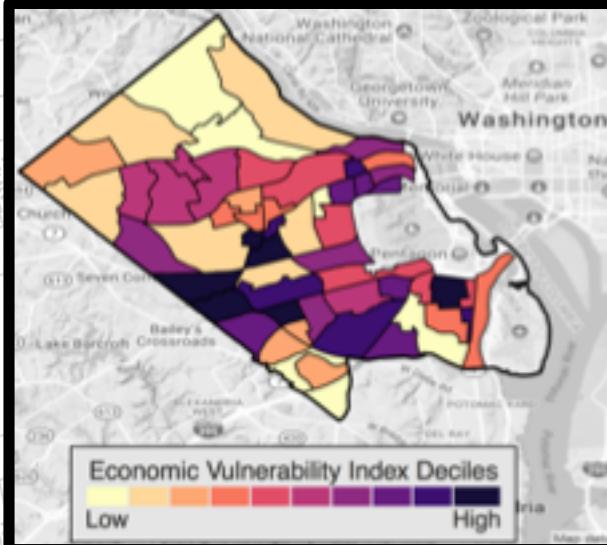
- 5-year 2015 American Community Survey - household level measures of economic vulnerability

School Level

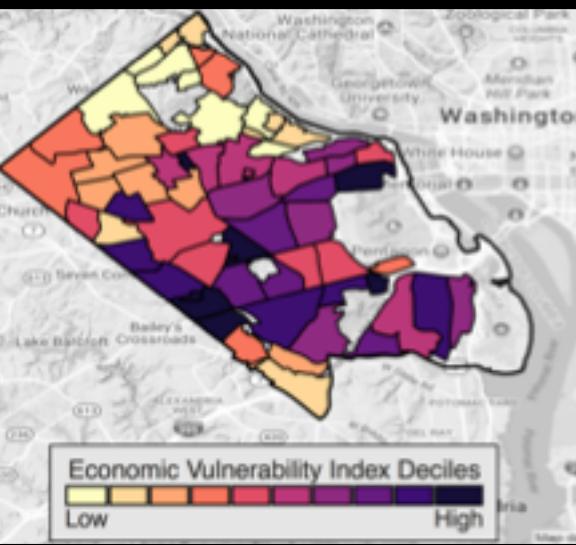
- National Center for Education Statistics
- Center for Disease Control 2014-15
- Virginia Department of Education - student level measures of economic vulnerability aggregated by school

Arlington County Vulnerability Indicators

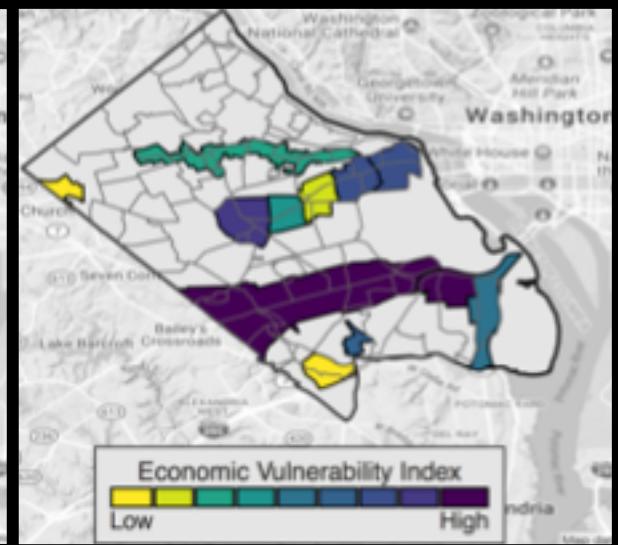
Census Tracts



Civic Association Neighborhoods



High-Density Planning Regions



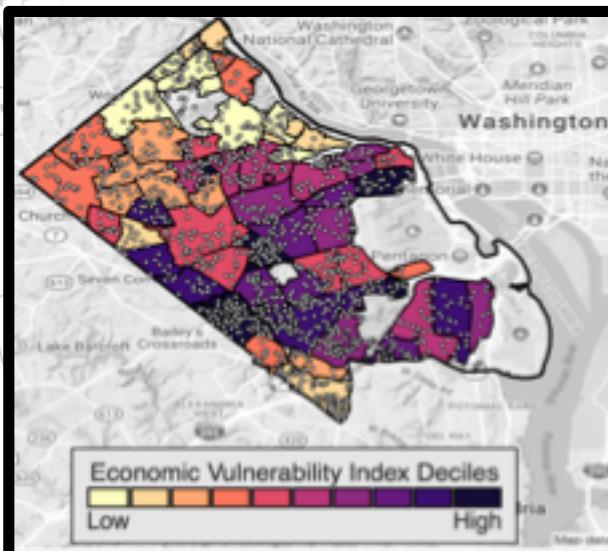
Based on a **statistical combination** of the percentage of Households with:

- housing burdens > 50% of Household income
- no vehicle
- receiving Supplemental Nutrition Assistance Program (SNAP)
- in poverty

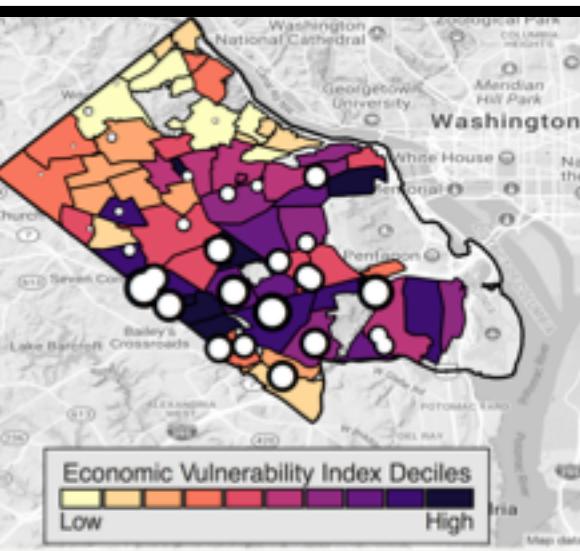
Source: American Community Survey 2012-20156 aligned to geographic areas using **SDAL Synthetic Technology**.

Arlington County Neighborhood Insights

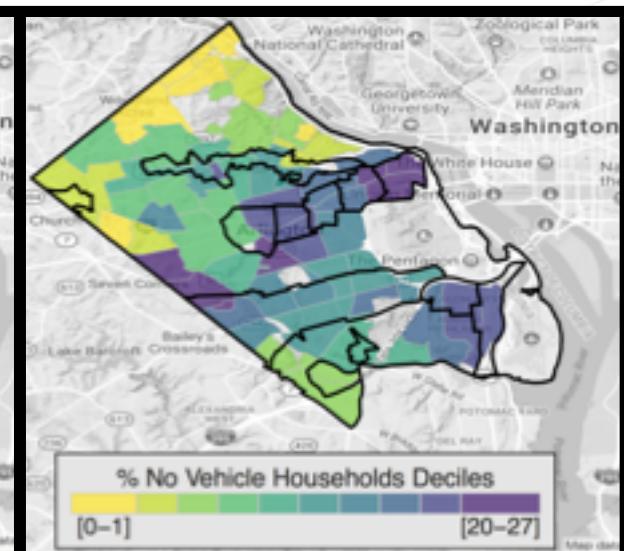
Households receiving
subsidies from Department
of Parks and Recreation



School and neighborhood
vulnerability indices

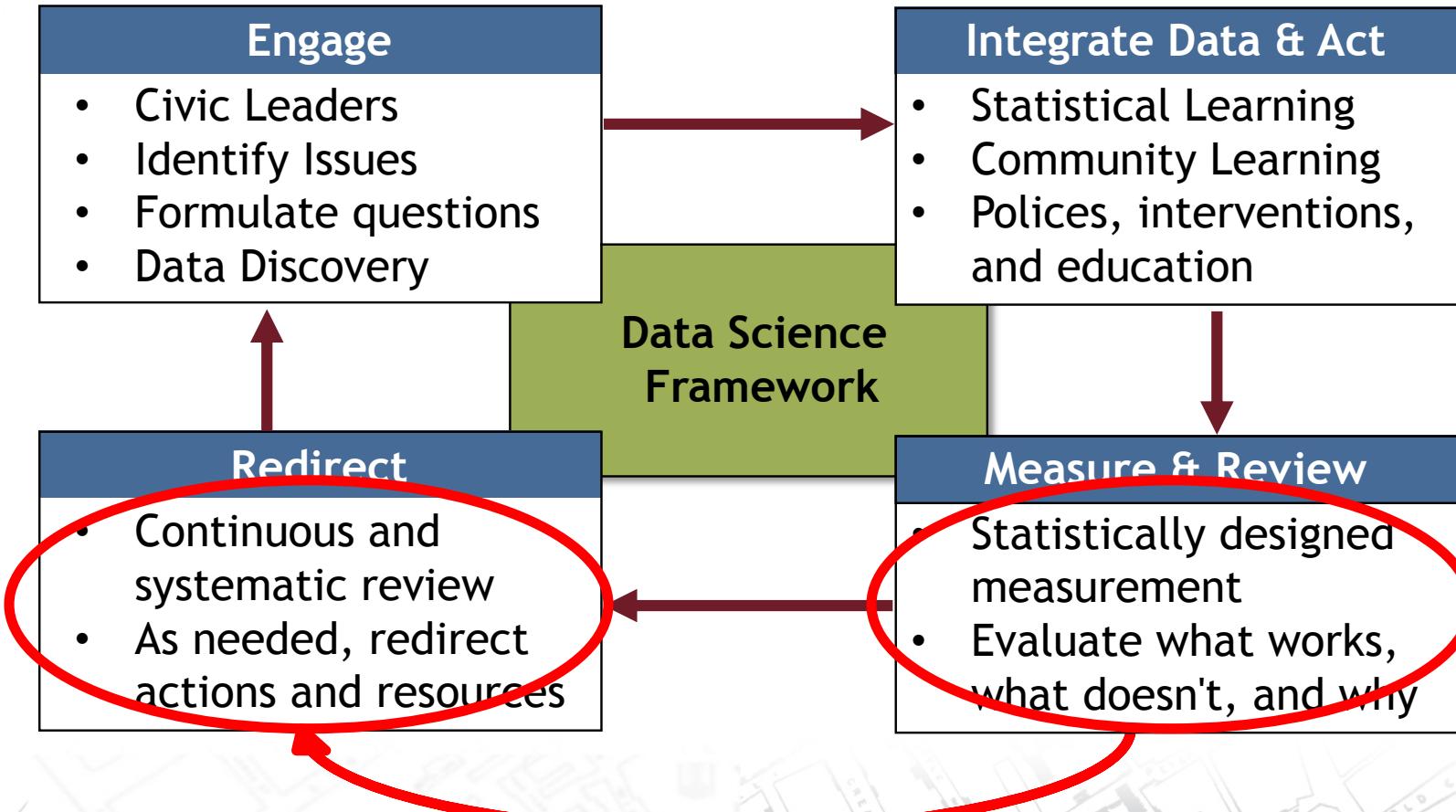


High-Density Planning
Regions with % households
with no vehicles



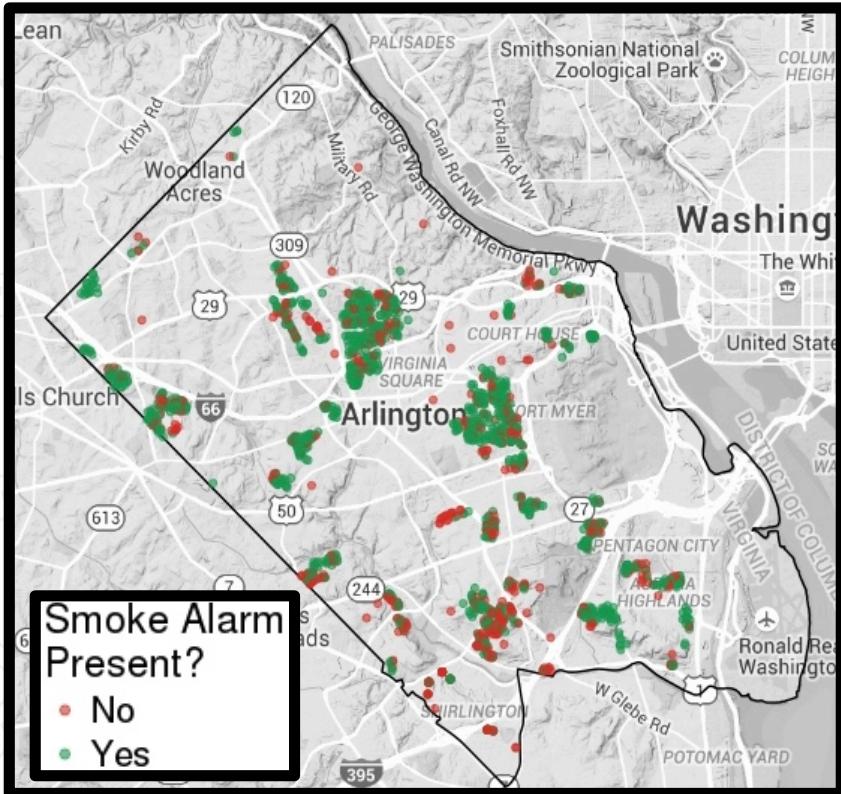
Sources: ACS 2012-2016; NCES, CDC, and VDOE 2014-2015; Arlington County Department of Parks & Recreation 2016.

Next Steps



Arlington Operation FireSafe

Evaluating and Informing a Program



Issue: The Fire Department wants to improve the efficiency of their Operation FireSafe program;

Out of 5,623 visits to single family homes only 1,799 had working smoke detector

Goal: Construct a model to estimate for each single family home the probability it has adequate smoke detectors

The DATA

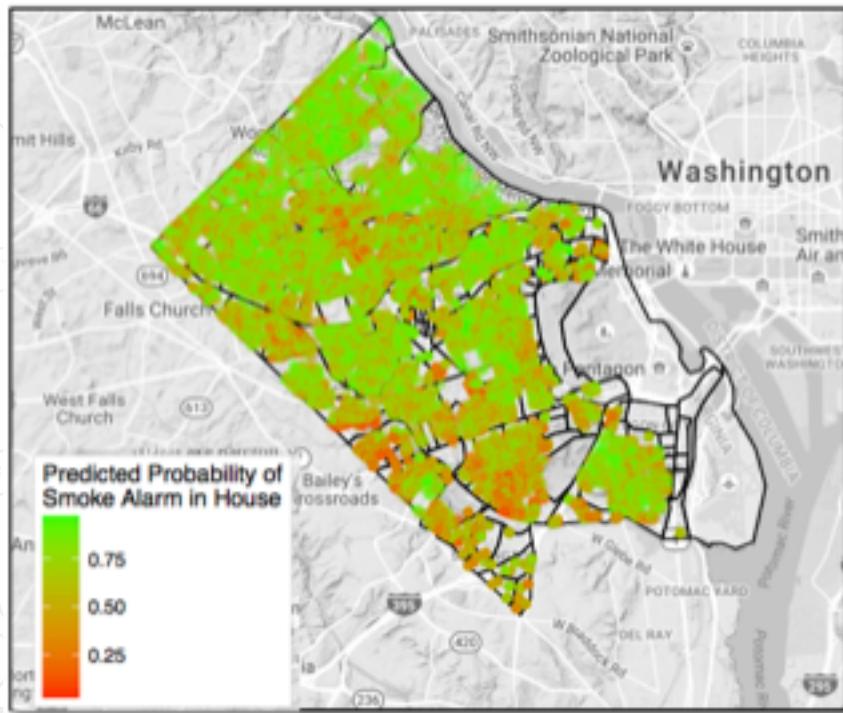
Household Level

- Operation FireSafe data on the location of the 5,623 single family homes visited, time and date of the visit, and the outcome installed smoke detector; working smoke detector in place
- Real estate tax assessments for 60,343 single family homes which includes tenure, home age, value, size, and number of bedrooms
- Geocoded the single family home locations

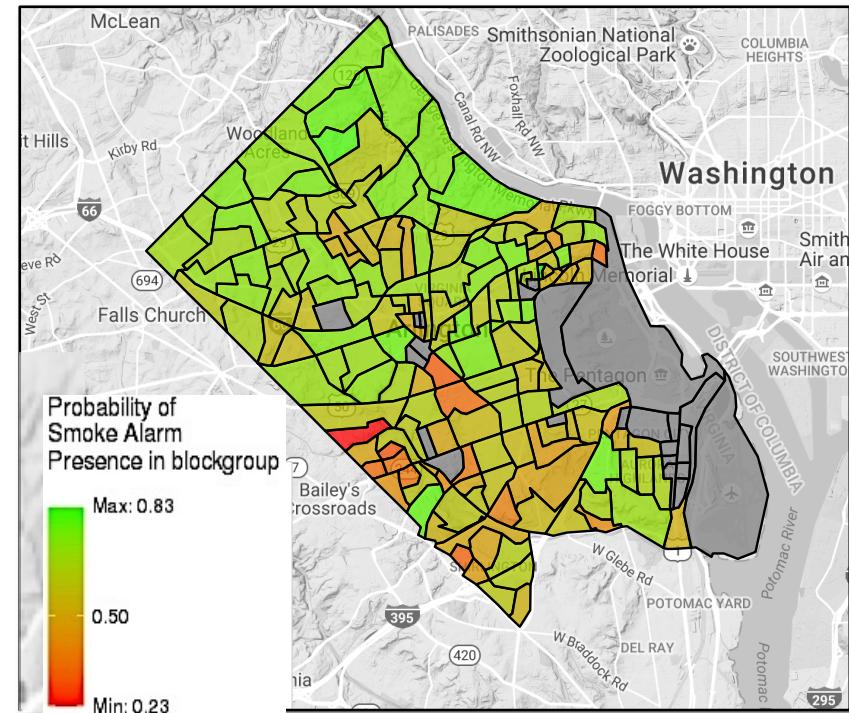
Census Tract Block Group

- 5-year 2015 American Community Survey - household level demographic and socioeconomic data

Probability of Having a Home Smoke Detector

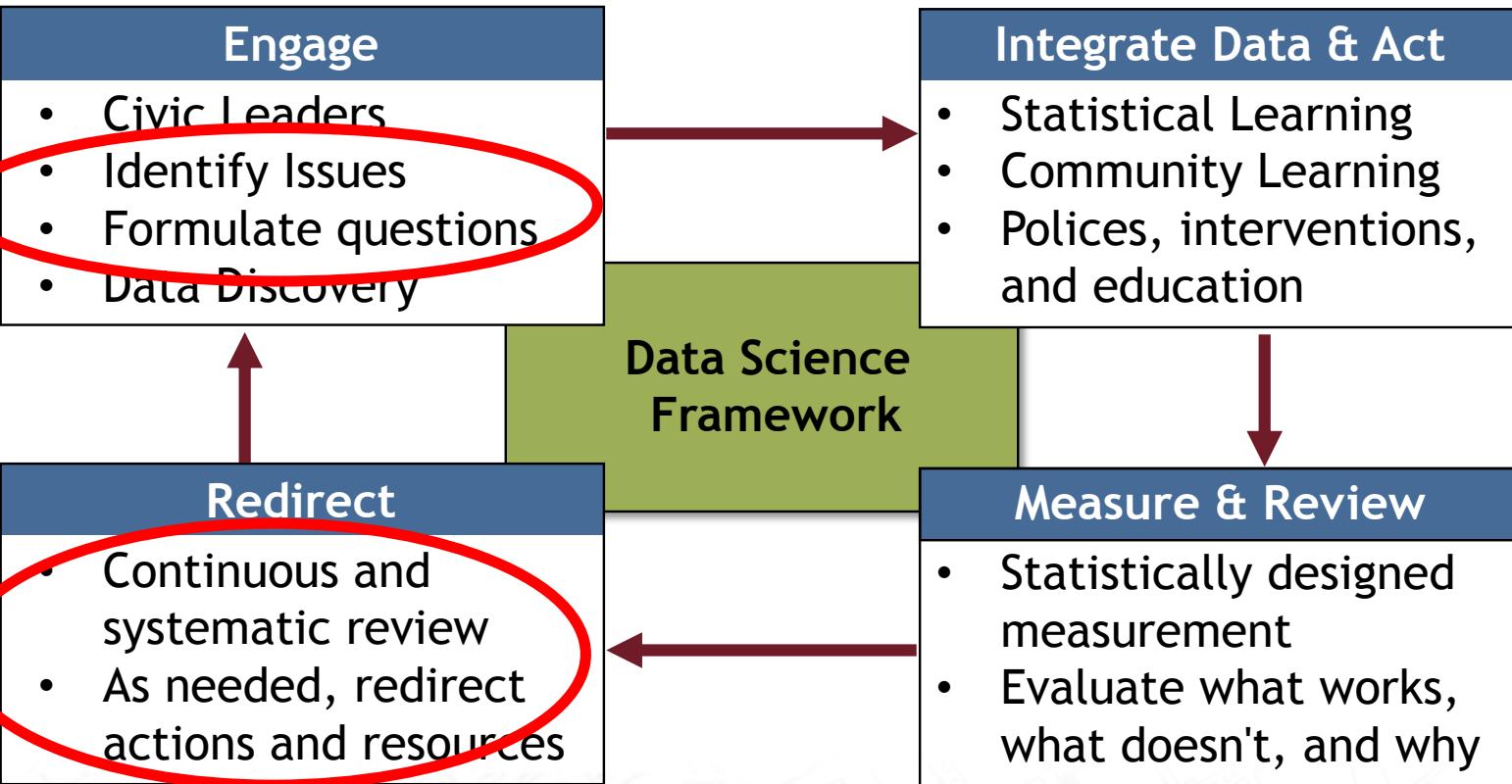


Housing Unit Level Predictions



Census Block Level Predictions

Next Steps



Data Science for the Public Good 2017



IDENTIFYING STEM EDUCATION PATHWAYS
Sponsor: The Bridges Project, The National Center for Interdisciplinary Research and Education

EXPLORING MENTAL HEALTH SERVICES FOR FAIRFAX COUNTY YOUTH
Sponsors: Fairfax County Sheriff's Office, Juvenile and Family Court, Northern Health and Human Services

RESIDENTIAL SMOKE ALARM NEED IN ARLINGTON COUNTY
Sponsor: Arlington Fire Department, Arlington County Fire Department

HOW DO EVENTS AFFECT CRIME?
Sponsor: Captain Bruce Benson and McLucy, Arlington County Police Department

MODELING THE IMPACT OF OPEN SOURCE SOFTWARE: NETWORK OF R PACKAGES
Sponsor: The National Center for Interdisciplinary Research and Education, The National Center for Engineering Education

DISCOVERING NON-TRADITIONAL DATA SOURCES FOR BUSINESS INNOVATION
Sponsor: Microsoft's Cloud for the Small Business, Amazon Web Services, Amazon.com, Inc., Cisco Systems, Google, Oracle, IBM, Dell, and SAP

A STUDY ON WMATA BUS FARE EVASION
Sponsor: Jennifer L. Crampton, Washington Metropolitan Area Transit Authority

ANALYZING THE ECONOMIC IMPACT AND SOCIAL INTEGRATION OF REFUGEES IN ROROHO, VIRGINIA
Sponsor: Kristen Hwang, USA, Kyle Morgan, MTS, Craig Horner, MTS, Harrison Brinkley, MTS, Adrienne Rogers, MTS, with Mark Orr, Stephanie Siragusa, and Monica Price, MSBIS

MODELING RESPONSE TIME FOR STRUCTURE FIRES
Sponsor: National Institute of Standards and Technology, Virginia Tech

PROFILE OF NEW KENT, VA

Sponsor: Parks, Arsenault, Morris, Beaufort, Lewis, Lewis, Moore, Long, Virginia Tech, City of New Kent, Commonwealth of Virginia, Northern Virginia Community College, and 20+ representatives

CREATING SYNTHETIC DATA FOR VIRGINIA LONGITUDINAL DATA SYSTEM
Sponsor: Virginia Department of Education, State Corporation Control Board, Virginia Department of Health, and Virginia Department of Higher Education

DEFINING AND MEASURING EQUITY IN ALEXANDRIA, VA
Sponsor: Living Justice, City of Alexandria

PROFILING ARMY BASES

Sponsor: Identifying publicly available data about military bases, their locations, their sizes, demographic information, and other quantitative profiles of army bases and their surrounding areas; identify relevant variables for use in statistical models

Thank You

