

Project Overview

“How many distinct individuals and families did we serve today?” - Jeff Price, VDSS Research and Planning

Project Title: Unsupervised Learning Data Analysis to Develop Network Visualizations of Family Relationships

Problem Statement: Depending on the program, and the way the program identifies, tracks, and manages people within the program, each program may define ‘household’ and ‘family’ differently

Task: Create proof of concept using data from the National Incidence Study of Child Abuse and Neglect (NIS) that will easily translate for VDSS datasets in the future

1 DATA EXPLORATION

VDSS stakeholder definitions of “family” and “nonfamily” vary significantly. Some are articulate with detail, while others are more vague. Stakeholders have different needs for defining “family.” This creates an issue for the department as a whole as they try to identify who they are serving.

Stakeholders:

- Childcare Subsidy
 - “Any individual(s), adult(s), and/or children related by blood, marriage, adoption or expression of kinship who function as a family unit.”
- TANF
 - No definition of “family.” However, “assistance unit” means “those persons who must participate together as a family unit.” Allows child to receive benefits from non-family caretakers through “Family Cap” provision.
- Child Protective Services, Child Support Enforcement, Foster Care, Medicaid, SNAP

About VDSS

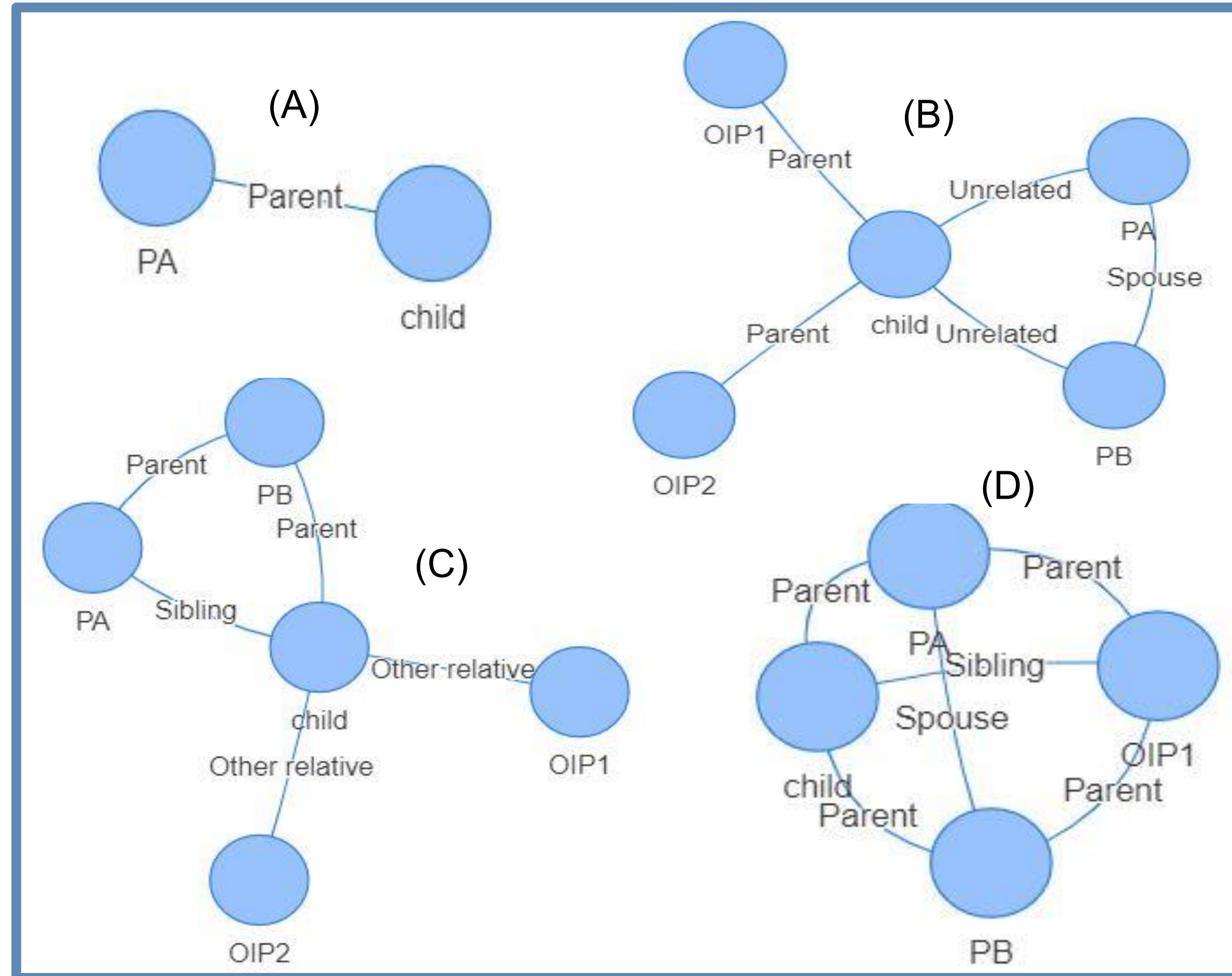
VDSS is dedicated to ensuring the wellbeing of children, families, and households across the state of Virginia. The agency is one of the largest in Virginia, partnered with over 120 social services departments, non-profit and faith-based organizations. Their vast and diverse services benefit children and their families in many aspects, but the variety in their services leads to an agency-wide problem with the definition of family and household. VDSS needs to be able to report on their service accounts of individuals and families. Their current system prohibits them from efficiently completing this task.

NIS Data

In order to build a prototype family structure query and visualization tool and to explore the analysis possibilities of family data network, we obtained data from the Fourth National Incidence Study of Child Abuse and Neglect (NIS-4). The NIS data were collected from 2004–2006. Demographic characteristics of children and their families were collected. Data came from Child Protective Services (CPS) agencies and from various community professionals in health, law enforcement, schools, and other community institutions. Each row of data represented a child. We extracted details about the the child’s primary caregiver and up to five additional individuals involved in each case.

Dashboard

We developed a dashboard that allows users to select the number of parents, the relationship Person A has to the child, the relationship Person B has to the child, and Other Involved Persons. The dashboard returns a table of applicable cases as well as the network relationship image (examples below).



2 EXPLORATORY ANALYSIS

Using NIS data, and at the request of our sponsor, we developed a broad classification tool that anyone can use in the future to identify family types. To do this, we:

- Used R to calculate summary statistics
- Created a network structure of tables (nodes to edges)
- Created network visualizations
- Developed a dashboard for family structure queries and displaying vizualizations.

3 FUTURE RESEARCH

VDSS is interested in continuing the project after the summer. After acquiring access to OASIS and VaCMS data, SDAL can exchange the information in the dashboard from NIS to Oasis and VaCMS. Eventually, we can implement data sharing for confirming PIs (case id, address, last name, and client ID), to validate the proof of concept demonstrated described here.

Clustering & Statistical Summaries

Relationship to Child (%)			Primary Caregiver Relationship to Person B	
Type	Primary Caregiver	Person B		
Spouse			Spouse	32.8
Parent	91	32.8	Unmarried partner	17.7
Grandparent	2.5	4.0	Housemate/roommate	0.32
Step Parent	1.3	8.4	Parent	3.0
Foster Parent	1.2	0.56	Sibling	0.54
Other Relative	1.2	1.4	Son/Daughter	0.60
Unrelated	0.98	10.2	In-law	0.06
Guardian	0.84	0.28	Other relative	1.0
Sibling	0.28	0.46	Other non-relative	0.63

We conducted unsupervised learning on the NIS family structure data. As the data comprised categorical and count data, we first used Random Forest to generate a dissimilarity measure between observations. Unlike most other unsupervised learning methods, Random Forest handles non-continuous features without difficulty. We then used k-medoids clustering as implemented in the *cluster* package in R. The clusters obtained were primarily defined by who filled the primary and secondary caregiver roles, and what their relationship to each other was. Using Random Forest also results in a measure of feature importance. This importance measure confirmed that the primary and secondary caregiver roles and their relationships were the most important features defining family structure.

References/Acknowledgements | Barlow Condensed Bold 30pt

“National Data Archive on Child Abuse and Neglect (NDACAN).” n.d. Accessed August 1, 2018. <https://www.ndacan.cornell.edu/index.cfm>.

Sedlak, A.J., Mettenburg, J., Brown, J., Basena, M., and Madden, K. (2010). Fourth National Incidence Study of Child Abuse and Neglect (NIS-4) Data File. Rockville, MD: Westat, Inc. Available from the National Data Archive on Child Abuse and Neglect web site: <http://www.ndacan.cornell.edu>

“VDSS Manuals.” n.d. Accessed August 2, 2018. <http://www.dss.virginia.gov/about/manuals.cgi>.

“Virginia Administrative Code - Title 22. Social Services.” n.d. Accessed August 2, 2018. <https://law.lis.virginia.gov/admincode/title22/>.

