**National Human Services Data Consortium**

Increasing Capacity & Building Connections:
**Bridging to the Future**

**2019 Spring Conference**
**Nashville, TN**
April 15-17, 2019

**CARES**
OF NY, INC
ENDING HOMELESSNESS

## Dr. Ruth Kassel

Assoc Dir Academic Integration Siena College in Albany NY

Board Member-CARES of NY

(518) 782-6951

rkassel@siena.edu

## Allyson Thiessen

Director of HMIS, serving 13 CoCs (25 Counties) for CARES of NY

Board member-NHSDC

(518) 489-4130 x103

athiessen@caresny.org

2

National Human Services Data Consortium
N-HSDC

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

## In 2016 Siena reached out to CARES of NY for a possible partnership

Siena was in need of big, messy data sets for analysis. Quickly.

CARES of NY had LOTS of big, messy data in dire need of analysis (because HMIS)

*This was either going to go really well or really poorly, but we agreed we were up to giving it a try and hope for The Unicorn.... a mutually beneficial relationship.*

3

**National Human Services Data Consortium**

Increasing Capacity &
Building Connections:
Bridging to the Future

**2019 Spring Conference**
**Nashville, TN**
April 15-17, 2019

# We believe in unicorns!

4

# Academic Community Engagement SPIn Program

**Three Year Plan**

**Y1- Relationship Building**

**Y2- Identity Building**

**Y-3 Sustainability**

10 Years 2008 - 2018
Center for Academic Community Engagement

**CARES OF NY, INC**
ENDING HOMELESSNESS

## TeamBILD
**(Big Issues and Leading-edge Discovery)**
Students and professors across all three schools are collaborating on the same issue and pursuing the same goals by concurrently tackling different parts of the project. Their issue is homelessness. Their mission is to help homeless services better serve their clients.

**cares** ENDING HOMELESSNESS

NECIP DOGANAKSOY AND TRAVIS BRODBECK

**PROFESSORS:**
VERNIZZI GRAZIANO, TING LIU ,
MICHAEL JARCHO, NECIP DOGANAKSOY,
MICHELLE MCCOLGAN, CHINGYEN MAYER,
JENNIFER DORSEY

**STUDENTS:**
MATTHEW JOHNSON, HAMZA MEMON,
LUKE MCKENNA, TIA BROWN, TRAVIS BRODBECK,
NICHOLAS CARPINELLO, LINDSAY CLARKE,
SERNA RIZZO, GORDON MACCAMMON,
AUSTIN SNYDER, THOMAS YAKALIS

TING LIU
HAMZA MEMON
LUKE MCKENNA

National Human Services Data Consortium NHSDC

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# Year 1: Relationship Building (and growing pains)

**Once all the paperwork was out of the way, the first year was spent getting to know each other and the dataset.**

- Throwing ideas around

- Throwing ideas away

- Asking way more questions than could possible have answers for

- Getting to know each other and where our specific needs and skills fit into this partnership

National Human Services Data Consortium NHSDC

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# First and biggest lesson: Stop "justing" people!

*WHY DON'T THEY JUST....*

*IF THEY WOULD ONLY....*

*DON'T THEY REALIZE?????!!!!???*

National
Human Services
Data Consortium

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# Researchers want to know:

## Why don't users "just"...

### fix the data?

### do real-time data entry?

### ask the clients more questions?

National Human Services Data Consortium

N-HSDC

Increasing Capacity & Building Connections: Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# Community Members want to know:

**?????**

## Why can't the HMIS team "just"...

### pull this report "real quick"

### give me a quick answer to a complex question?

### get HUD to change the regs? Programming? EVERYTHING?

9

National Human Services Data Consortium

N-SD

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

*And, of course, the classic:*

**Why does the HMIS team even have an opinion?
They're "JUST" data people…**

**Takeaway: collaboration can not begin until
the "Justs" are out of the way.**

Used Emergency Shelter at least once (Albany county)

Used Emergency Shelter at least once (Dutchess county)

"If I can model the universe I think I can model homelessness in areas of New York"
*Matt Bellis (Physics)*

[15]

**Basic mapping.** Here we show a plot generated by hmis and *folium*, a wrapper to leaflet.js. It shows a marker for all the zip codes that CARES works with. We are interested in performing a more sophisticated mapping analysis to see where help is most needed and how the homeless

National Human Services Data Consortium

Increasing Capacity & Building Connections: Bridging to the Future

# EDA…Data Cleaning



*"This is a lot like problem solving a mechanical issue at GE; a lot of messy data that we need to organize to solve a problem."*
*Necip Doganaksoy & Travis Brodbeck (Accounting)*

13

National Human Services Data Consortium
N-SDC

Increasing Capacity &
Building Connections:
Bridging to the Future

**2019 Spring Conference**
**Nashville, TN**
April 15-17, 2019

# Build an Online Homeless Database

**SIENA**college
*The education of a lifetime*

**Dr. Ting Liu[1], Luis Concepcion-Bido[1], Caleb Ryor[1], and Travis Brodbeck[2]**
[1]Computer Science, Siena College, Albany, NY  [2]Accounting, Siena College, Albany, NY

### Abstract

Organizations that help for homeless people usually are not willing to share the data because of sensitive personal information. Their data isn't stored the most optimally either. We plan to build a nationwide database that can be shared with organizations/researchers to help find new approaches to the issue of homelessness through proper storage of data and analysis.

### Introduction to the Siena Homelessness Project

The main hurdle with the data we were faced with was the sheer amount of it needed to be handled. Everything was being stored in CSV files which provided plenty of function based issues as well as space issues causing Excel to perform slowly. All the personal information given to us was double hashed so we could still track information regarding a given individual. The data itself consists of various site and client data from homeless shelters all across the Albany area all of which is collected and managed by CARES, Inc., our data provider and community partner.

### Solution: Database

Utilizing a SQL database for the sake of storing all the information provided to us gave us much more freedom as it enabled there to no longer be a limit on space as well as having more flexible options in how we could record this data. The next steps that needed to be made before the database could be complete however involved optimizing where certain data was being stored for even more ease of access.
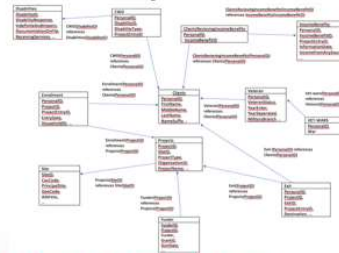
*Example of Data Entries in Table*

### Building an SQL Homeless Database

As pictured below, the first step to creating said database involved us forming an ER diagram to visualize what was being created. Each of the 9 CSV files provided became its own table with the exception of the Clients table being split to hold Veteran and VeteranStatus information elsewhere for the sake of storage optimization. Important strides were made to create code to deduplicate data within specific tables within the database to make queries provide more accurate analysis as well as just further improving upon the database itself. Fortunately the data itself was sorted well enough for there not to be much more resorting to be done on our part, but the small changes that we made were paramount to the data involved to be analyzed properly at all.

*ER Diagram of database*

### SQL Queries for Statistical Analysis

With the creation of a SQL database we are able to run various queries to parse through the data to find various figures for all kinds of issues. Here are a few examples of what we are able to currently do:
- Compare Veteran and Clients tables to find percentage of veterans
- Sort IDs within database to find percentage of youth, adults, and elderly within the system
- Find percentage of those with a given recorded disability

More statistical analysis is planned but much of what we can do is reliant on optimizing our tables for less complicated queries as well as continuing further with cleaning and deduplicating the data we store.

### Develop a Web Interface

**Goals for web interface development**
- Help visitors understand issues in homelessness data through short articles
- Serve as hub connecting organizations and researchers to share data
- Generate and showcase statistics from current datasets using SQL queries
- Be a center for Siena College data analysis tools with examples to bring in more clients

*Screenshot of current frontpage*

Siena College Homelessness Project

### Future Work

- Generate more code for data deduplication
- Work with more companies to get them involved in the project for more datasets
- Talk with other parts of project on campus to showcase/host more tools on webpage

### Acknowledgements

16

**National Human Services Data Consortium**

Increasing Capacity &
Building Connections:
Bridging to the Future

**2019 Spring Conference**
Nashville, TN
April 15-17, 2019

**But…whenever you want to use software, there's a cost……**



18

National Human Services Data Consortium — N-HSDC

Increasing Capacity & Building Connections:
**Bridging to the Future**

**2019 Spring Conference**
**Nashville, TN**
April 15-17, 2019

## Siena:

- Creates community among faculty
- Re-energizes faculty and students
- Elevates the reputation for community engagement
- Helps get community engagement classifications

## CARES of NY :

- Identifies and addresses gaps in data quality
- Research opportunities using HMIS data
- Excitement about USING data
- Educating on true state of homelessness
- Connects community organizations with students

National Human Services Data Consortium

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# Year 2: Identity Building

20

National Human Services Data Consortium

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

## Accomplishments

- Clear delineation of roles
- Better understanding of commitment
- First publications
- Data for Good Exchange poster
- Redefining wants and expectations

21

# Our first publication!

Journal of Open Source Software (JOSS) article on the HMIS Visualization suite developed in Year 1.

JOSS | 10.21105/joss.00384



JOSS — The Journal of Open Source Software

## hmis: A python tool to visualize and analyze HMIS data

Sara Mahar[1] and Matthew Bellis[1]

1 Siena College

### Summary

Many organizations that work to combat homelessness receive funds from the US Department of Housing and Urban Development (HUD). These organizations might be overnight shelters or transitional housing or somewhere in between the Continuum of Care (CoC) provided by these groups. Since 2004, HUD has mandated that groups that receive these funds collect data on the homeless individuals that make use of these services. As such, there is a wealth of data that has been collected all over the country from a variety of organizations. Organizations have some freedom in how they collect and store these data, often making use of 3rd-party software solutions, but the data format is the same everywhere.

This variety of data storage tools means that is is difficult for a data scientist at any of these organizations to dig into this data using standard, open-source computing tools like python or R. These groups can download the data in a standardized "HMIS data dump", which results in 12 separate .csv files, but this still does not make any initial analysis any easier, a priori. These files have information about individuals's name (hashed as a personal ID number), date of birth, prior living, disabilities, jail time, etc.

This module contains a suite of python functions to allow for analysis and visualization of the data collected by the various partners across the CoC. Visualization includes time-series plots, and mapping of the locations of the programs individuals have entered. Analysis can be done with these visualizations and with the ability to withdraw individuals who share a common character. For example, the analyst can withdraw all of the individuals who have visited more than 25 programs and then visualize them.

We have developed these tools to work with the standard HMIS data dump in the RHY (Runaway and Homeless Youth) data format, that produces 12 .csv files in which personal identifying information is de-identified through a hashing algorithm. Because of this standardization, any other tools that leverage this software package can be used by similar networks across the country.

The definitions of the information in the HMIS data can be found on HUD's website (HUD 2017 ; HUD 2016).

This software project started thanks to the help and assistance from members of CARES NY (http://caresny.org/), a group committed to ending homelessness and who applies for grants from HUD and administers them with partners across the CoC. We would like to particularly acknowledge CARES members, Maureen Burns, Terry O'Brien, and Allyson Thiessen, who explained to us the need for tools like this and the data formats themselves. We also acknowledge members of the Siena College community, Ruth Kassel and Paul Thurston, for the initial connects with CARES, and their strong and continued support of this project.

Mahar et al., (2017). hmis: A python tool to visualize and analyze HMIS data. Journal of Open Source Software, 2(18), 384, doi:10.21105/joss.00384

22

National Human Services Data Consortium · NHSDC

2019 Spring Conference
Nashville, TN
April 15-17, 2019

Increasing Capacity &
Building Connections:
Bridging to the Future

# 2ⁿᵈ Publication: Bloomberg Poster

National Human Services Data Consortium
NHSDC

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# 3rd Publication – Council of Undergrad Research Journal

Issue Theme:
"Big Data as a Tool to Promote Undergraduate Research"
*Editor-in-Chief*: James LaPlant
*Issue Editors*: Laurie Gould, Janice DeCosmo
*Proposal Deadline: June 1, 2018*

The theme of the spring and summer 2019 issues of *SPUR: Scholarship and Practice of Undergraduate Research* (formerly *CUR Quarterly*) will focus on big data as a tool to promote undergraduate research. Five to six articles from a wide range of disciplines are sought that explore how the applications and use of big data serve to facilitate undergraduate research in a variety of educational and professional contexts. In addition, vignettes (maximum 300 words) are welcomed that offer concrete, creative suggestions with regard to the connections between big data and undergraduate research. Examples of topics of interest include the following:

BIGData
& UNDERGRADUATE RESEARCH

24

# Year 3: Sustainability

- Pilot Project
- Economic Research
- Case Notes Analysis
- Deep Dive for DSS

National Human Services Data Consortium

N-SDC

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

## How to SPIn your own partnership

**Understand who the connectors are:**

- Undergrad research center
- Community engagement center
- Outreach and volunteer center
- Don't limit yourself to the social sciences



26

National Human Services Data Consortium

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# Be prepared to:

- Give a sample of the data (to evaluate: structure, size, messiness)
- Answer ALL THE QUESTIONS
- Get excited!  It's contagious.
- Accommodate academic timelines
- Be flexible in what you want/need/expect
- Give as much as you get: time, energy, and enthusiasm
- Let this be student led; you'll be AMAZED at what you get (good and bad!)
- Work on projects that are low urgency but high importance

National Human Services Data Consortium

Increasing Capacity &
Building Connections:
Bridging to the Future

2019 Spring Conference
Nashville, TN
April 15-17, 2019

# Do you believe in unicorns?