Handling Missing Values in Joint Sequence Analysis Alexandra Ballow¹, Alina Lazar¹, Alex Sim², Kesheng Wu² ¹Youngstown State University, ²Lawrence Berkeley National Laboratory

~Option 1~

-10 -5

-10

ABSTRACT

This study focuses on developing methodologies to minimize the effects of incomplete data. Specifically, it hopes to reduced the noise and bias caused in categorical sequence data by data gaps. Some strategies investigated include choosing a substitution "cost" to replace missing values and deleting the missing values at the end of a sequence. Cluster validity metrics are used to determine the accuracy of the unsupervised clustering algorithms and t-SNE is employed to visualize clusters and age biases. It became clear that deleting missing values provided the best results, but all data sets are different. Thus this study recommends employing the studied procedures before conducting analysis on longitudinal sequence data to ensure the results are unbiased. After these tests optimize the data, clustering is conducted to understand the correlation between a person's state in life and their travel.

DATA

- Data acquired from the WholeTraveler survey
- Consists of demographic and travel information from nearly a thousand Bay Area citizens
- This study focuses on the binary variables including: Whether they walked or used a bicycle
- Having a child
- Having a partner
- Use of public transport
- Employment Status
- Is full of missing values, as demonstrated by the family size sequences (right)



RESEARCH QUESTION

What effect does the treatment of missing values in categorical sequence data have on cluster analysis results?



due to its lack of noise and evaluation of outliers







general survey pool

CONCLUSIONS

Recently released data detailing travel tendencies was wrought with large groups of missing values, typically found in longitudinal data. This project was intended to investigate and mitigate the effect of missing data on clustering results on this new data. After testing two different procedures to minimize the biases created by missing values, it is determined that removing missing values from the end of the sequences and normalizing the distance measure produced the best results. This procedure generated more interpretable results, with less parameter tuning and less error. Proceeding with the results acquired during this step, it was possible to create clusters and sequences plots which will be evaluated in future studies.

ACKNOWLEDGMENTS

This work was prepared in partial fulfillment of the requirements of the Berkeley Lab Undergraduate Faculty Fellowship (BLUFF) Program, managed by Workforce Development & Education at Berkeley Lab. I want to thank Ling Jin, Annika Todd, and Anna Spurlock from the Energy Analysis Environmental Impacts Group for providing datasets and guidance for the project.







WORKFORCE **& EDUCATION**

Youngstown STATE UNIVERSITY.

ENERGY

Office of Science