

MSD 2023 Datathon

Evaluating substance abuse treatment programs within the US

Presented by Correlation One

Problem Statement

Welcome to the 2023 Datathon! This document explains the topic of the Datathon, important details about the datasets you'll be using, and guidance on how to submit your results.

Background

Substance abuse is a major problem in the United States, with an estimated 21.2 million Americans aged 12 or older struggling with substance use disorder in 2020. Of those individuals, only 4.2 million received any type of treatment for their addiction. There are a variety of substance abuse treatment programs available throughout the country, including outpatient counseling, residential treatment centers, and medication-assisted treatment. In 2021, there were approximately 14,000 specialized substance abuse treatment facilities in the US, serving over 1.3 million people. However, there are still significant barriers to accessing treatment, including lack of insurance coverage, stigma, and limited availability of services in certain areas. Despite these challenges, substance abuse treatment programs remain a critical component of addressing the ongoing addiction crisis in the United States.

Given this, your goal is to evaluate substance abuse treatment across specific states or even the entire United States. You can evaluate them from any lens you want!

Your Task

You are asked to pose your own question and answer it using the available datasets, as well as any supplementary datasets that you find to aid your analysis. Both the creativity of your question, and the quality of your analysis are of paramount importance. You need not be comprehensive; depth of insight is more important than breadth of question posed.

Submissions may be predictive, using machine learning to classify or predict patterns. Submissions may also be illuminating by way of data visualization or sound statistical inference.

Consider exploring one of the sample questions below, or creating your own variation. Creativity in formulating your own question is **encouraged**; however, it should not be at the expense of analytical depth, precision, and rigor, which are far more important.

<u>Sample Question 1:</u> What is the difference in treatment completion rates between different types of substances? Which substances have the greatest differences in the probability of treatment completion and the shortest completion times? Is this consistent across states/statistical areas or are there states/statistical areas that have policies that may suggest why these trends are occurring?

<u>Sample Question 2:</u> Are treatment programs triaging admissions appropriately? For example, do pregnant women have to wait for treatment as long as a non-pregnant person? How do these wait times impact treatment efficacy, are they decreasing the probability of a successful treatment? Are individuals with a co-occurrence of mental illness and substance abuse more likely to commit crimes before admission or discharge from the treatment program? What data driven policy recommendations can you make for triaging admissions and maximizing treatment outcomes? Are these recommendations consistent across the US or are there specific calibrations to the recommendations needed for certain regions, states or statistical areas?

<u>Sample Question 3:</u> What is the impact of self help groups on improving the treatment outcomes of substance abuse? Are there any benefits of being in one before admission or discharge? Which demographics or substance abuse cohorts benefit the most from them? Compare areas that are known to have these more (while adjusting for population)? How do community demographics interact with the efficacy of self help groups?

Sample Question 4: Can you discover any natural experiments between states, counties or statistical areas? If you do, how can they help inform/evaluate substance abuse treatment programs? Can you showcase that their unique demographics are being overlooked or is this demographic less unique than expected and can you suggest how it should inform policy surrounding substance abuse treatment programs? What tactics or strategies for running treatment programs can be uncovered or inspired by this discovery?

<u>Sample Question 5:</u> Identify what states, regions, or statistical areas are neglecting types of substance abuse programs or consistently having failed treatments. Can you demonstrate a data driven understanding of why this is occurring and offer potential and practical solutions to aid in solving their issues.

A note on sample questions: The MSD judges are looking for you to flex your statistical inference skills, not just your ML skills. The ability to showcase novel analyzes through the lens of survival analysis or other causal inference techniques may earn you bonus points (however is not required to produce a winning report, so do not bound yourself by this).

Datasets

The provided datasets are stored in the "Datathon Materials" folder on Google Drive. Your team need only use the data / datasets that are relevant to your chosen question / topic.

<u>treatments_2017-2020.csv</u> — (4,049,136 rows x 81 columns) — size: 1.58 Gb: This dataset is a repeated cross sectional study (pseudo longitudinal study) containing records on admissions of people aged 12 and older, and includes information on admission demographics and substance use as well as information on discharges.

<u>state county cbsa population.csv — (4210 rows x 11 columns) — size: 123.7 Kb:</u> This dataset contains records stat, counties and core based statistical area populations from 2010-2019.

<u>american_community_surveys - 8 datasets - total size: 70.4 Mb:</u> This dataset contains the full american community survey data 1 & 5 year estimates from 2016-2021.

<u>treatment_facilities_2016_2020.csv - (74820 x 97 columns) - size: 19.6Mb : This dataset contains 5 years of data surveys on substance abuse treatment facilities across the US from 2016-2020.</u>

Additional Datasets

Participants are welcome and **encouraged** to scour the internet for their own custom datasets to supplement their analysis. All additional data used should be public and reputable. Additionally, any supplementary datasets should not exceed <2 GB unzipped (consult Correlation One's R&D team if you believe your idea is worthy of an exception).

Other Materials

We will provide you with the schema for each of the data tables in another packet.

Submissions: Content

Submissions should have two components:

- 1. Report this should have two main sections:
 - a. Non-Technical Executive Summary What is the question that your team set out to answer? What were your key findings, and what are their significance? You must communicate your insights clearly summary statistics and visualizations are encouraged to help explain your thought process
 - b. Technical Exposition What was your methodology / approach towards answering the questions? Describe your data manipulation and exploration process, as well as your analytical and/or modeling steps. Again, the use of visualizations is highly encouraged when appropriate.
- 2. Code please include all relevant code that was used to generate your results. **Although** your code will not be graded, you <u>MUST</u> include it, otherwise your entire submission will be discarded.

Additional information (e.g. roadblocks encountered, caveats, future research areas, and unsuccessful analysis pathways) may be placed in an appendix.

Judges will be evaluating your technical report without your team there to explain it; therefore, **your submission must "speak for itself".** Please ensure that your main findings are clear and that any visualizations are functionally labeled.

Submissions: Evaluations

The competition will have multiple rounds of evaluation. Your Report will judged as follows:

• Technical Executive Summary

o *Insightfulness of Conclusions*. What is the question that your team set out to answer, and how did you choose that question? Are your conclusions precise and nuanced, as opposed to over-generalizations?

• Technical Expositions

- o Wrangling & Cleaning Process. Did you conduct proper quality control and handle common error types? How did you transform the datasets to better use them together? What sorts of feature engineering did you perform? Please describe your process in detail within your Report.
- o *Investigative Depth*. How did you conduct your exploratory data analysis (EDA) process? What other hypothesis tests and ad-hoc studies did you perform, and how did you interpret the results of those tests and analyses? What patterns did you notice, and how did you use these to make subsequent decisions?
- o Analytics & Modeling Rigor. What assumptions and choices did you make, and how did you justify them? How did you perform feature selection? If you build models, how did you analyze their performance, and what shortcomings do they exhibit? If you constructed visualizations and/or conducted statistical tests, what was the motivation behind the particular models you build, and what did you tell you?

Submissions: Formats

Make sure to send your submissions in an email to datathonsubmissions@correlation-one.com.

The email subject line should be in the following format: "Team [n] Submission - MSD Datathon 2023" where [n] is your team number.

Only one of your team members should send us the email (feel free to CC your teammates). You can resend the email before 11:59pm IST / 7:29pm GMT+1 if you would like to make any changes - we will look at your latest submission before the deadline.

Reports can be produced using any tool you prefer (Python Notebook, Shiny Application, Microsoft Office, etc.); however, **your report** <u>MUST</u> be in a universally accessible and readable format (HTML, PDF, PPT, Web link). It must not require dedicated software to open. For example, if your report is a Python Notebook, it should be exported to HTML. If you create a Shiny App, it should be published at an accessible Web link.

Please include the source file used to generate your report. For example, if you submit a PDF with math-type, equations, or symbols please include your LaTeX source file.

Code should be submitted in a single zipped collection of files separate from your report.

Please ensure to give yourselves extra time to send your submissions, especially if the size of your files is large.

Tips and Recommendations

Since this is a single day Datathon, time is of the essence. It will be important to settle on a tractable question early on to ensure that you have time to both thoroughly explore the problem that you pose and write up your results. The outcome of this Datathon, and your overall success, will largely be a product of the quality of the question that you choose to answer and the depth within which you explore that question.

For data engineering, exploration, and modeling, we highly recommend that you install Jupyter Notebook: http://jupyter.org/install.html. Jupyter Notebook is an interactive, real-time development environment that eliminates many pain points of the standard "terminal & text editor" environment, and is compatible with both Python and R.

We also recommend that your team stick to tools and techniques that you have previously used. Learning new skills is certainly valuable, but it can consume a large portion of your available time, leaving less time for completing the task at hand.

We've compiled 3 additional commonalities of successful teams and 3 pitfalls that successful teams will actively avoid. Of course, these may not apply to every team, so we recommend that you and your team apply any tips accordingly.

Tips for Success	Try to Avoid
1. Focus on hypothesis testing when	1. Do not try to exhaust all different models
brainstorming your research question	you know just to yield an ideal cross
	validation accuracy
2. Spend at least 3 hours on your report to	2. Do not violate assumptions of statistical
ensure strong communications through both	models. Sometimes, specific models require
visualizations and writing	specific features so it is best to make sure
	those conditions are sufficiently met
3. Engage in proper causal analysis. Just	3. Do not pick research statements and blindly
because your model passes standard	stick to it trying to get it to work. Often times,
cross-validation checks does not demonstrate	further data exploration will show that it's not
(or even necessarily suggest) causality	true or worthwhile

Ask for Help

Correlation One's Data team is here to help. Let us know about your struggles as early on as you can and we may be able to offer advice on how to best move forward.