

# Identification of Individuals with Counterfactual Response Types

Ang Li

Florida State University, Department of Computer Science

angli@cs.fsu.com



## Abstract

The unit selection problem aims to identify a set of individuals who are most likely to exhibit a desired mode of behavior or to evaluate the percentage of such individuals in a given population, for example, selecting individuals who would respond one way if encouraged and a different way if not encouraged. Li and Pearl addressed this problem by deriving tight bounds on the “benefit function,” which quantifies the payoff or cost associated with selecting an individual with specific characteristics, using a combination of experimental and observational data. The challenge then becomes designing an algorithm to identify the characteristics that lead this benefit function (when indicating the desired response type) to approach a value of 1. This project aims to explore how machine learning methods can be applied to solve this identification task.

## Introduction

The challenge of identifying individuals with a desired response pattern is pervasive across various industries, including marketing and health science. In customer relationship management, for instance, companies seek to identify customers who might consider leaving but could be persuaded to stay with the right incentive. This incentive must be carefully tailored because customer behavior can be influenced by past incentives offered by the company. In another example, in online advertising, companies are interested in identifying users who would only view an advertisement when prompted. This research aims to address these identification challenges using causal models.

## Terminologies

- **Counterfactuals** are statements or scenarios that express what might have happened if circumstances had been different from what actually occurred.
- **Complier** is a counterfactual response type for individuals who would respond positively if treated and negatively if not treated. For instance, it refers to individuals who would have avoided COVID-19 if they had been vaccinated but would have contracted the virus if they had not been vaccinated.
- **Always-taker** is a counterfactual response type for individuals who always respond positively, regardless of whether they are treated or not. For example, it describes individuals who would have contracted COVID-19 regardless of whether they were vaccinated or not.
- **Never-taker** is a counterfactual response type for individuals who always respond negatively, regardless of whether they are treated or not. For example, it describes individuals who would have avoided COVID-19 regardless of whether they were vaccinated or not.
- **Defier** is a counterfactual response type for individuals who would respond negatively if treated and positively if not treated. For instance, it refers to individuals who would have avoided COVID-19 if they were unvaccinated but would have contracted the virus if they were vaccinated.

## Motivating Example

For the COVID-19 example, according to randomized controlled trials (RCTs) data, Pfizer claims that its vaccine is highly effective, as summarized in Table 1. Nevertheless, despite vaccination, many individuals still contracted COVID-19 infections. The effectiveness of vaccination remains a compelling topic in the field of health science, especially during the initial stages of a pandemic when vaccines are scarce, and RCTs are the only available means for testing vaccine efficacy.

	Comirnaty	Placebo	Vaccine Efficacy
16 years and older	77 positive out of 19993	833 positive out of 20118	91.1%
12 through 15 years	0 positive out of 1057	28 positive out of 1030	100%

Table 1: First COVID-19 Occurrence From 7 Days After Dose 2 by Age Subgroup.

We introduced the unit selection model [2], which signifies a fundamental shift in the treatment evaluation paradigm. This model differs from treatment effects-based and correlation-based decision-making systems in data science and artificial intelligence, enabling us to incorporate data from both experimental and observational studies. The latter is readily available from various sources, is free from selection bias, and, when combined with experimental data, provides insights into the individual (counterfactual) behavior of subjects in the population. More specifically, the unit selection model allows us to estimate the percentage of response types for a population with sufficient data. Based on the same RCT data from Pfizer, our analysis suggests that only 4% of individuals aged 16 and older (or 3% for those aged 12 to 15) would have avoided COVID-19 if they had been vaccinated but would have contracted the virus if unvaccinated (referred to as compliers). We then pose two questions:

1. Who are these compliers? The presence of compliers in each group serves as evidence that vaccination is indeed useful. However, the percentages are small. If we can identify these compliers, we can apply the vaccination more effectively and safely.
2. What is the effectiveness of other age groups that lack RCT data, or can we define age groups more specifically? It is common in reality for RCT data to be insufficient in covering all groups due to time restrictions, ethical issues, and so on.

## Challenges

When it comes to the first question, the underlying reason we should be concerned about identifying compliers is that the randomness in an RCT is insufficient for decision making. Compliers and always-takers represent latent heterogeneity within the treatment group, while always-takers and defiers represent latent heterogeneity within the control group. Consider two scenarios:

1. An RCT on a population with 50% always-takers and 50% never-takers, as depicted in Figure 1, results in a positive rate of 0.5 for both the treatment and control groups.
2. An RCT on a population with 50% compliers and 50% defiers, as shown in Figure 2, also yields a positive rate of 0.5 for both the treatment and control groups.

In both cases 1 and 2, there is no difference in the RCT results. However, it's evident that these two cases should lead to different decisions. The treatment is entirely ineffective in case 1, and the primary goal in case 2 should be to identify who the compliers are. In reality, there are countless cases that correspond to the same RCT results.

The main difficulty in identifying compliers (or other response types) stems from the fact that the desired response type is not observed directly but, rather, it is defined counterfactually in terms of what the individual would do under hypothetical unrealized conditions.

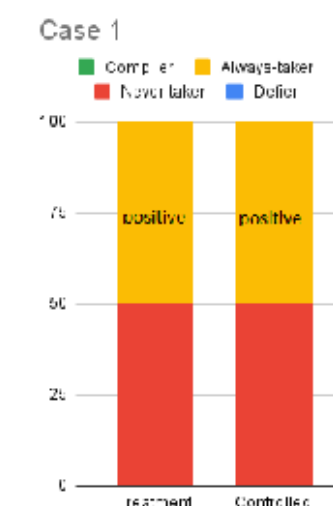


Figure 1: Population consists always-takers and never-takers only.

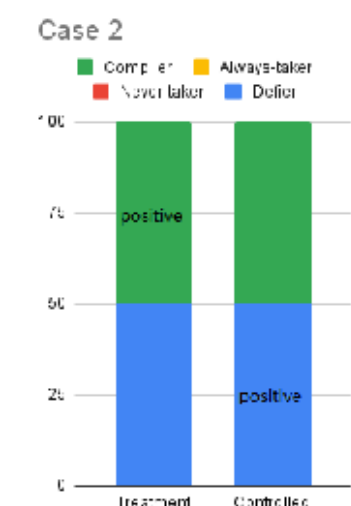


Figure 2: Population consists compliers and defiers only.

As for the second question, the straightforward answer is that if you have sufficient data and apply the unit selection model to each age group, you can derive more precise estimations. However, this task is particularly challenging due to the wealth of characteristics available for each individual (resulting in a large number of groups) and the sparsity of data for each group.

We found that precise estimates require about 1,500 experimental and observational samples per age group, making it impractical to have sufficient data for all age groups. The challenge escalates with multiple characteristics like gender and family history, requiring over 49 million samples each for experimental and observational studies across groups with 15 binary characteristics. This analysis also assumes uniform distribution, ignoring the rarity or non-existence of some groups, which makes evaluation impossible.

## Methods

We hypothesize that an individual's counterfactual behavior is determined by their personal characteristics, establishing a mapping between the probabilities of causation that express this behavior (e.g., Probability of necessity and sufficiency (PNS) represents the complier) and these characteristics. Machine learning is recognized as the most effective tool for learning such mappings and predicting accordingly. Figure 3 best represents this component. Once we have the mapping, we can identify which set of characteristics leads to a population with nearly 100% probability of causation. Figure 4 best represents this component.

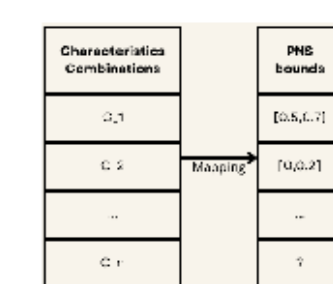


Figure 3: Learning the mapping using characteristics with known bounds, such as  $C_1$  and  $C_2$ , and predicting the PNS bounds for unknown ones, such as  $C_n$ .

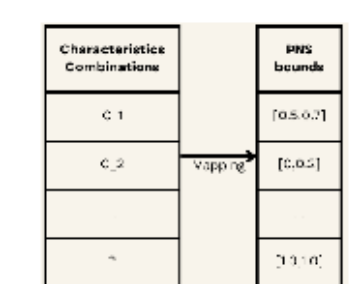


Figure 4: Search for the characteristics to achieve PNS bounds equal to  $[1.0, 1.0]$ .

## Simulated Study

we used neural networks as our estimator to estimate the lower and upper bounds of the PNS. We illustrated this by the results from [1]. Because the simulated data is tabular in nature, we just use a multilayer perceptron (MLP) in the following experiments. Our experimental environment is an AWS p3.2xlarge instance. The key parameters of our model are: embeddings dimension as 128, training epochs as 600, and learning rate as 0.01.

To show the performance of estimation, we randomly selected 200 subpopulations (among 32768) and compared their learned bounds of PNS with the true bounds of PNS computed from the informer data. The results are shown in Figures 5 and 6. The learned bounds of PNS for subpopulations are a good fit for the true PNS bounds. The average error of the learned lower bound among 32768 subpopulations is 0.0775, and the average error of the learned upper bound among 32768 subpopulations is 0.1371. They are both acceptable errors given we have only 423 training size to learn 32768 subpopulations.

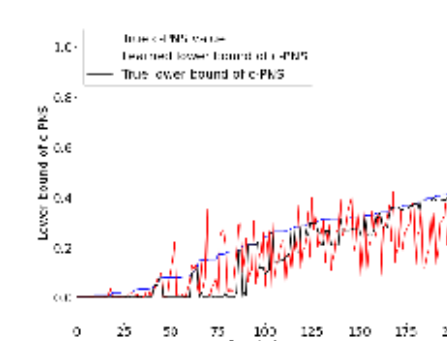


Figure 5: Lower bound of the PNS for 200 sets of characteristics. These 200 sets of characteristics are randomly selected from 32768 of them.

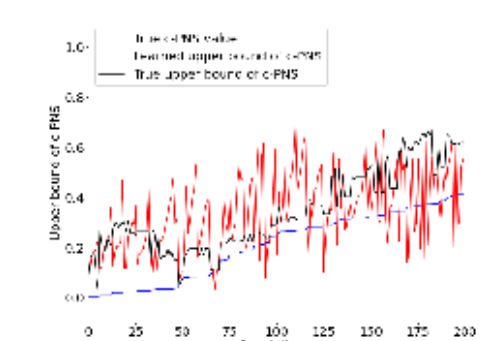


Figure 6: Upper bound of the PNS for 200 sets of characteristics. These 200 sets of characteristics are randomly selected from 32768 of them.

## Application

We utilized a portion of our findings to examine the impact of artificial sweetener consumption on health, particularly its relationship to obesity [3]. Our study revealed that between 20% and 50% of individuals, especially those with poor dietary habits, are more likely to gain weight from consuming Diet Coke. Conversely, in groups such as young females with healthier diets, only a small proportion experience weight gain due to Diet Coke. Our next step is to apply our proposed machine learning methods to identify these individuals.

## References

- [1] Ang Li, Song Jiang, Yizhou Sun, and Judea Pearl. Learning probabilities of causation from finite population data. Technical report, Department of Computer Science, University of California, Los Angeles, CA, 2022.
- [2] Ang Li and Judea Pearl. Unit selection based on counterfactual logic. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 1793–1799. International Joint Conferences on Artificial Intelligence Organization, 7 2019.
- [3] Yicheng Qi and Ang Li. Causality in the can: Diet coke's impact on fatness. *arXiv preprint arXiv:2405.10746*, 2024.