

Causal Inference in Statistics: An Overview

Lorenzo Lagos
February 7, 2017

I found Pearl's distinction between standard statistical analysis and causation quite enlightening. The crucial difference is that causal analysis does not only infer beliefs or probabilities under static conditions, but it must also make claims on how changing conditions affect these parameters. Therefore, **dynamics** is the crucial point of differentiation since such a concept cannot be defined from observed distributions alone, but rather requires a set of assumptions. One thing that was puzzling for me is the difference between assumptions and models. Pearl defines **identifiability** as the condition when any two models (M_1 and M_2) with the same set of identifying assumptions (A) such that the joint distribution of both models are the same ($P(M_1) = P(M_2)$) imply that the parameter/quantity of interest is the same ($Q(M_1) = Q(M_2)$). Hence, what matters is the assumptions rather than the models. However, isn't a model simply a set of assumptions? Especially, if we rule out models where $P(M_1) \neq P(M_2)$ then all assumptions that distinguish M_1 from M_2 seem superfluous.

Pearl's presentation also clarified the distinction between structural equations and regression equations. The former is meant to represent a physical reality (depicted in causal graphs), while the latter are artifacts of analysis. Hence, a missing link in a causal graph may posit the assumption that $Cov(U_Y, U_X) = 0$ whereas residual terms in regression are already assumed to be uncorrelated with the regressors. Thus, a structural β can be **proven** to be a causal effects, whereas a regression β can only be **interpreted** as a causal estimate.

Regarding the use of **do-calculus**, I was intrigued by the fact that questions that cannot be answered with experimental studies cannot be expressed in do-notation. For example, my research question (what would be the average change in separations had there been no probationary period in place?) falls under this category since no separation case can be tested twice. Nonetheless, such questions can be answered with a probabilistic analysis of counterfactuals, e.g., "Y would be y has X been x in situation $U = u$ ". The **unit-level counterfactual** would be $Y_{M_x}(u)$ where Y is separations, M_x is the model under probationary period length x , and u are the characterizing attributes that vary across units of analysis. The bulk of my future work will be on developing the model and selecting the characterizing attributes. Doing so will allow me to **1)** Define the excess separations around the probationary period as a function that can be computed from any model $Q(M)$; **2)** Formulate causal assumptions and represent them in a causal graph; **3)** Determine if the target quantity is identifiable; and **4)** Estimate or approximate the target quantity.