

Linear Regression in a Nonlinear World: The Multivariate Case

Nadav Kunievsy

July 3, 2023

In economics and other social sciences, researchers frequently document the causal effect of a variable of interest on another outcome variable. In non-experimental scenarios, this is usually achieved using a conditional independence assumption¹. To operationalize this assumption when analyzing the data, researchers often employ multivariate linear regression. They regress an outcome variable on the variable of interest and a set of control variables, thereby ensuring the conditional independence needed to identify the causal effect. It is widely known that if the underlying Data Generating Process (DGP) is linear, the structural parameters that govern the outcome can be identified from the population regression coefficients. A natural question arises then how does the population regression coefficient on the variable of interest is related to the underlying marginal causal effect of the variable of interest, when the DGP is not linear. In this note, we demonstrate, using a generalization of Yitzhaki [1996], how the coefficient of the variable of interest is determined by the underlying DGP. Furthermore, we demonstrate how, when the DGP is non-linear, the imposition of linear structure on the model can be thought of as an omitted variable bias and/or a measurement error.

We begin by considering the univariate case. Suppose that the underlying Data Generating Process (DGP) is represented by

$$Y = g(T, \epsilon)$$

where Y is the outcome of interest, T is a continuous variable of interest, ϵ are unobserved variables that influence the outcome, and g is the function that delineates the DGP. A researcher who is interested in the effect of T might choose to estimate the coefficient of T in the following linear model

$$Y = \alpha + \beta_1 T + u \tag{1}$$

In the univariate scenario, Yitzhaki's theorem (Yitzhaki [1996]) provides a method to associate the population regression coefficient with the underlying DGP. Specifically, under certain regularity conditions, Yitzhaki demonstrates

¹It's worth noting that in many instances researchers assume conditional mean independence. However, for the sake of discussion, we focus on conditional independence, which forms the backbone of nonparametric identification arguments Pearl [2009]

that

$$\beta_1 = \int_{-\infty}^{\infty} \frac{\partial E[Y|t]}{\partial t} w(t) dt \quad (2)$$

where $w(t) = \frac{E[T - E[T]|T > t]P(T > t)}{\text{Var}(T)}$ and $\int_{-\infty}^{\infty} w(t) dt = 1$. Moreover, if ϵ is independent of T , then the regression coefficient provides us with a positively weighted average of the average marginal effects of the variable of interest

$$\frac{\partial E[Y|T = t]}{\partial t} = E \left[\frac{\partial g(t, \epsilon)}{\partial t} \right]$$

In this note, we inquire whether Yitzhaki's result extends to the more generalized case where researchers assume that T is independent of the unobserved variables, given a set of control variables, \mathbf{X} . Specifically, suppose the DGP is now represented by

$$Y = g(T, \mathbf{X}, v), \quad (3)$$

$$T = h(\mathbf{X}, \varepsilon), \quad (4)$$

where v and ε are unobserved variables that influence the outcome and the variable of interest's value, respectively, and h and g are the underlying causal functions that govern the DGP. We further suppose that $T \perp\!\!\!\perp v | \mathbf{X}$.²

When researchers aim to utilize the conditional independence assumption and estimate some functional of the causal effect distribution, they often resort to linear regression. Specifically, they estimate the following linear model³:

$$Y = \beta T + X\gamma + \epsilon \quad (5)$$

Often, researchers interpret the coefficient on the variable of interest, β , as an average (or some average) of the marginal causal effect of T . Naturally, one might wonder how Yitzhaki's theorem applies to this multivariate context and how we should interpret β in this scenario.

To answer this question, we first define the "Naive Regression-Weighted Causal Effect" (NRwCE) as

$$NRwCE = E_X \left[\int_{-\infty}^{\infty} \frac{\partial E[Y|t, \mathbf{X}]}{\partial t} w(t, \mathbf{X}) dt \right] \quad (6)$$

²This can also be framed in potential outcomes notation as $T \perp\!\!\!\perp Y_t | X, Y = \int_t \mathbb{1}\{T = t\} Y_t dt$

³Throughout the analysis, we assume that \mathbf{X} contains a constant

where $w(T, \mathbf{X}) = \frac{E[T - E[T]|T > t, \mathbf{X}]P(T > t|\mathbf{X})}{E[\text{Var}(T|\mathbf{X})]}$.

This parameter intuitively extends Yitzhaki's interpretation of β from the univariate case to the multivariate case. Firstly, for each \mathbf{X} -cell, the numerator of the weights assigns the same weight that Yitzhaki's weights would assign in a regression of the outcome variable on T at this particular value of \mathbf{X} . The denominator of the weights is simply the average over the conditional variance of T , which ensures that the weights sum to 1, in a manner similar to Yitzhaki's original weights. Secondly, just like in the univariate case, if ϵ is independent of T given \mathbf{X} , the *NRwCE* provides us with a weighted average of causal effects.

$$E_X \left[\int_{-\infty}^{\infty} \frac{\partial E[Y|t, \mathbf{X}]}{\partial t} w(t, \mathbf{X}) dt \right] = E_X \left[E \left[\int_{-\infty}^{\infty} \frac{\partial g(t, \mathbf{X}, v)}{\partial t} w(t, \mathbf{X}) dt \middle| \mathbf{X} \right] \right]$$

The generalized form of Yitzhaki's theorem, detailed below, reveals that the population regression coefficient, β , is equivalent to the naive regression weighted causal effects only when the relationship between the controls and the variable of interest is linear. In other cases, it often yields a biased estimate in relation to the Naive Regression-Weighted Causal Effect.

Theorem 1 (Generalized Yitzhaki's Theorem). Denote by π the coefficients of \mathbf{X} in the population regression of T on \mathbf{X} . Denote by $\mu(\mathbf{X}) = E[T|\mathbf{X}]$ and denote the prediction error by $\Delta(\mathbf{X}) = \mu(\mathbf{X}) - \pi\mathbf{X}$. Assume the first and second moments and conditional moments exist, then the treatment coefficient, β , from the population regression, $Y = T\beta + \mathbf{X}\gamma + \epsilon$, is given by

$$\begin{aligned} \beta &= \underbrace{\frac{\text{Cov}(Y, (T - \mu(\mathbf{X})))}{\text{Var}(T - \mu(\mathbf{X})) + \text{Var}(\Delta(\mathbf{X}))}}_{\text{Causal Component}} + \underbrace{\frac{\text{Cov}(Y, \Delta(\mathbf{X}))}{\text{Var}(T - \mu(\mathbf{X})) + \text{Var}(\Delta(\mathbf{X}))}}_{\text{Bias}} \\ &= \underbrace{E_{\mathbf{X}} \left[\int_{-\infty}^{\infty} \frac{\partial E[Y|T = t, \mathbf{X}]}{\partial t} w(t, \mathbf{X}) dt \right]}_{\text{Causal Component}} + \underbrace{\frac{\text{Cov}(Y, \Delta(\mathbf{X}))}{\text{Var}(T - \mu(\mathbf{X})) + \text{Var}(\Delta(\mathbf{X}))}}_{\text{Bias}} \end{aligned}$$

where

$$w(t, x) = \frac{E[T - \mu(\mathbf{X})|T > t|\mathbf{X}]P(T > t|\mathbf{X})}{E_{\mathbf{X}}[\text{Var}(T|\mathbf{X})] + \text{Var}(\Delta(\mathbf{X}))} \geq 0$$

The proof, detailed in the appendix, applies the Frisch-Waugh-Lovell Theorem and integration by parts⁴. The

⁴A similar theorem for the case where T is discrete and the researcher conducts a fully saturated regression can be found at [Angrist and Krueger \[1999\]](#)

initial observation from Theorem 1 is that, in general, β does not equal the *NRwCE* parameter. The weights, while they resemble the weights of the *NRwCE*, don't sum to 1, and there's an additional bias term driven by changes in the controls and variation in T . To better understand this expression, we evaluate different DGPs.

First, let's assume $g(T, \mathbf{X}, v) = \beta T + \gamma \mathbf{X} + v$, and allow h to be unrestricted. In this case, the outcome equation is correctly specified. By using the standard argument from the consistency of the OLS, we understand that the population regression coefficient equals the structural β , and, therefore, is trivially equal to the *NRwCE* parameter⁵. Next, let's consider the scenario in which h , the function governing the intensity of the variable of interest, is linear in \mathbf{X} . In this case, $\Delta(\mathbf{X}) = 0$ for all \mathbf{X} , and Theorem 1 shows that both the bias term and $\text{Var}(\Delta(\mathbf{X}))$ equal zero, which implies that β equals the naive regression weighted causal effect.

Finally, let's consider the situation where both $g(T, \mathbf{X}, v)$ and $h(\mathbf{X}, \epsilon)$ are non-linear in \mathbf{X} . In this scenario, the population regression coefficient doesn't yield a weighted average of treatment effects. Instead, it gives a weighted average of the marginal causal effect of T and a bias term, which is introduced by the correlation between the outcome variable Y and $\Delta(\mathbf{X})$. These $\Delta(\mathbf{X})$ terms represent deviations of the conditional expectations from their best linear approximation⁶. If these deviations are systematically correlated with the outcome variable, our causal estimate will be biased. If, however, we find that $\text{Cov}(Y, \Delta(\mathbf{X})) = 0$ but $\text{Var}(\Delta(\mathbf{X})) \neq 0$, then the bias term in theorem 1 becomes zero, but we encounter a measurement error issue that attenuates β downwards. Even when control variables enter the outcome equation linearly, there can still be a bias if the treatment variable, T , has a nonlinear effect on the outcome. Specifically, consider $g(T, \mathbf{X}, v) = f(T) + \mathbf{X}\gamma + v$. In this scenario, β would not, in general, equal *NRwCE*. To illustrate this, let's again assume that $T \perp v | \mathbf{X}$ and $E[v | \mathbf{X}] = 0$. Then, we have:

$$\begin{aligned} \text{Cov}(Y, \Delta(\mathbf{X})) &= \text{Cov}(f(T) + \mathbf{X}\gamma + v, T - \pi\mathbf{X} - (T - \mu(\mathbf{X}))) \\ &= \text{Cov}(f(T), \Delta(\mathbf{X})) \neq 0 \end{aligned}$$

where the second equality stems from the mean zero assumption, and the fact that the residuals are not correlated with \mathbf{X} . This demonstrates that when considering including control variables in the regression, we should concentrate

⁵To see this through the lens of theorem 1, notice that:

$$\text{Cov}(Y, T - \mu(\mathbf{X})) + \text{Cov}(Y, \mu(\mathbf{X}) - \pi\mathbf{X}) = \text{Cov}(Y, T - \pi\mathbf{X}) = \beta \text{Cov}(T, T - \pi\mathbf{X})$$

Divide by the denominator to get β .

⁶Recall that the coefficients provide the best linear approximation to the conditional expectations. See, for example, (Angrist and Pischke [2009])

on their interaction with the variable of interest, and not on their potential impact on the outcome variable. In the accompanying Jupyter Notebook, we demonstrate that the bias and attenuation effects can be significant, such that the regression coefficient may provide limited-to-no information on the effect of the variable of interest.

Therefore, we can see that imposing a linear structure on control variables, while the function $h(\mathbf{X}, \epsilon)$ that regulates the intensity of the variable of interest is non-linear, can lead to issues analogous to omitted variable bias and/or measurement error. A different perspective on the resulting bias can be achieved by expressing $\Delta(\mathbf{X})$ as the difference between residuals:

$$\Delta(X) = \mu(X) - \pi X = \underbrace{T - \pi X}_{\substack{\text{Unexplained Due to} \\ \text{Linearity Restrictions}}} - \underbrace{T - \mu(X)}_{\substack{\text{Fundamentally} \\ \text{unexplained}}}$$

This illustrates that bias is induced by the portion of the variation in T that could potentially be explained using the control variables, but is not explained because we restrict the model to be linear. If these remaining, unexplained components are correlated with the outcome variable, we receive additional bias because we are unable to control for these components. If they are not correlated, we get attenuation bias due to the fact that we mismeasure the control variable function.

To conclude, the Generalized Yitzhaki theorem emphasizes that merely assuming conditional independence and estimating the causal effect using regression is insufficient to ensure the attainment of some functional form of the treatment effect distribution. However, it also provides guidance on how researchers can address these biases when interested in the $NRwCE$ parameter. The simplest approach to obtain an unbiased estimate of $NRwCE$ is to include an estimate of $E[T|\mathbf{X}]$ as a control variable in the regression⁷. To compute $E[T|\mathbf{X}]$, one can either estimate it nonparametrically [Ullah and Pagan \[1999\]](#), or use prior knowledge on the assignment mechanism ([Borusyak and Hull \[2021\]](#)).

However, this raises a question - if one can estimate $E[T|\mathbf{X}]$ nonparametrically, why choose to estimate the causal effect using regression instead of estimating the entire model nonparametrically? In many cases, researchers opt for linear regression due to its efficiency and stability, two properties that do not always characterize nonparametric estimators, especially when the dimensions of \mathbf{X} are large. Consequently, if researchers wish to use regression, it

⁷This is as the conditional expectation, given the conditional expectation, is a linear function of the conditional expectation $E[T|E[T|\mathbf{X}]] = E[T|\mathbf{X}]$

would be insightful to include in their analysis a discussion on the relationship between the control variables and the variable of interest. This can be done by plotting $E[T|x_j]$ for different components of \mathbf{X}_j . Estimating the conditional expectations for one variable is relatively straightforward and could provide an indication that the linearity assumption is not unreasonable.

Researchers should also bear in mind that the relationship between the treatment and control variables is not invariant to monotonic changes of the treatment. For instance, researchers should be cautious when estimating a linear model where T enters the regression linearly, and a similar model where they use $\log(T)$. Both models are unlikely to obtain a weighted average of treatment effects, and at least one of them is likely to suffer from misspecification bias. Hence, researchers should be more conscious of how they model their control variables.

Appendix

The following is a generalization of [Yitzhaki \[1996\]](#), for the multivariate case. Denote by π the coefficients from linear projection of T on \mathbf{X} . Denote by $f(T)$ and $f(T|\mathbf{X})$ the density and conditional density of T . Using Frisch-Waugh-Lovell (theorem [Frisch and Waugh \[1933\]](#)), we have that

$$\beta = \frac{\text{Cov}(Y(T - \mathbf{X}\pi))}{\text{Var}(T - \pi\mathbf{X})}$$

We start by focusing the numerator. Denote the conditional expectation of T , condition on \mathbf{X} by $\mu(\mathbf{X})$. Then we can express the the numerator as

$$\begin{aligned} \text{Cov}(Y, T - \mathbf{X}\pi) &= \text{Cov}(Y, (T - \mu(\mathbf{X})) + \mu(\mathbf{X}) - \mathbf{X}\pi) \\ &= \text{Cov}(Y, (T - \mu(\mathbf{X}))) + \text{Cov}(Y, \mu(\mathbf{X}) - \mathbf{X}\pi) \end{aligned}$$

Using the law of iterated expectations and integration by parts we can re-express the first term as

$$\begin{aligned}
\text{Cov}(Y, (T - \mu(\mathbf{X}))) &= E[Y(T - \mu(\mathbf{X}))] \\
&= E_{\mathbf{X}}[E[Y(T - \mu(\mathbf{X}))|\mathbf{X}]] \\
&= E_{\mathbf{X}}[E[E[Y|T, \mathbf{X}](T - \mu(\mathbf{X}))|\mathbf{X}]] \\
&= E_{\mathbf{X}}\left[\int_{-\infty}^{\infty} E[Y|T, \mathbf{X}](u - \mu(\mathbf{X}))f(u|\mathbf{X})du\right] \\
&= E_{\mathbf{X}}\left[\left[E[Y|u, \mathbf{X}] \int_{-\infty}^t (u - \mu(\mathbf{X}))f(u|\mathbf{X})dt\right]_{t=-\infty}^{\infty} - \int_{-\infty}^{\infty} \frac{\partial E[Y|t, \mathbf{X}]}{\partial t} \int_{-\infty}^t (u - \mu(\mathbf{X}))f(u|\mathbf{X})dudt\right] \\
&= E_{\mathbf{X}}\left[\int_{-\infty}^{\infty} \frac{\partial E[Y|t, \mathbf{X}]}{\partial t} - \int_{-\infty}^t (u - \mu(\mathbf{X}))f(u|\mathbf{X})dudt\right] \\
&= E_{\mathbf{X}}\left[\int_{-\infty}^{\infty} \frac{\partial E[Y|t, \mathbf{X}]}{\partial t} E[T - \mu(\mathbf{X})|T > t, \mathbf{X}]p(T > t|\mathbf{X})dt\right]
\end{aligned}$$

Where the last equality follows from the fact that

$$E[T - \mu(\mathbf{X})|T > t, \mathbf{X}]P(T > t|\mathbf{X}) + E[T - \mu(\mathbf{X})|T \leq t, \mathbf{X}]P(T \leq t|\mathbf{X}) = 0$$

Therefore the numerator is given by

$$\text{Cov}(Y, T - \mathbf{X}\pi) = E_{\mathbf{X}}\left[\int_{-\infty}^{\infty} \frac{\partial E[Y|t, \mathbf{X}]}{\partial t} E[T - \mu(\mathbf{X})|T > t, \mathbf{X}]p(T > t|\mathbf{X})dt\right] + \text{Cov}(Y, \mu(\mathbf{X}) - \pi\mathbf{X})$$

Next, we turn to the denominator. We can re-express it as

$$\begin{aligned}
\text{Var}(T - \pi\mathbf{X}) &= \text{Var}(T - \mu(\mathbf{X}) + \mu(\mathbf{X}) - \pi\mathbf{X}) \\
&= \text{Var}(T - \mu(\mathbf{X})) + \text{Var}(\mu(\mathbf{X}) - \pi\mathbf{X}) + 2\text{Cov}(T - \mu(\mathbf{X}), \mu(\mathbf{X}) - \pi\mathbf{X}) \\
&= \text{Var}(E[T - \mu(\mathbf{X})|\mathbf{X}]) + E[\text{Var}(T|\mathbf{X})] + \text{Var}(\mu(\mathbf{X}) - \pi\mathbf{X}) \\
&= E[\text{Var}(T|\mathbf{X})] + \text{Var}(\mu(\mathbf{X}) - \pi\mathbf{X})
\end{aligned}$$

where we used the law of total variance and the fact that $\text{Cov}(T - \mu(\mathbf{X}), \mu(\mathbf{X}) - \pi\mathbf{X}) = E[(T - \mu(\mathbf{X}))(\mu(\mathbf{X}) - \pi\mathbf{X})] = 0$, due to the law of iterated expectations, which concludes the proof.

References

- Joshua D. Angrist and Alan B. Krueger. Empirical strategies in labor economics. In *Handbook of labor economics*, volume 3, pages 1277–1366. Elsevier, 1999.
- Joshua D. Angrist and Jörn-Steffen Pischke. *Mostly Harmless Econometrics: An Empiricist’s Companion*. Princeton University Press, 2009.
- Kirill Borusyak and Peter Hull. Non-random exposure to exogenous shocks. *Working Paper*, 2021.
- Ragnar Frisch and Frederick V. Waugh. Partial time regressions as compared with individual trends. *Econometrica: Journal of the Econometric Society*, pages 387–401, 1933.
- Judea Pearl. *Causality*. Cambridge University Press, 2009.
- Aman Ullah and Adrian Pagan. *Nonparametric Econometrics*. Cambridge University Press, Cambridge, 1999.
- Shlomo Yitzhaki. On using linear regressions in welfare economics. *Journal of Business Economic Statistics*, 14(4): 478–486, 1996.