

# Instrumental LLMs

Patrick Power

July 11, 2024

## **Abstract**

In many applied microeconomic contexts, the underlying data is text - think Health Care, Education and Housing. Causal inference in this setting has typically proceeded by hand-selecting numerical representations of the text and estimating the corresponding conditional expectation function assuming that treatment or the instrument is locally randomly assigned. Recent developments in Natural Language Processing/AI though have introduced alternative ways to produce causal estimates from text. In this paper we (1) clarify the general framework for using Large Language Models for causal inference and (2) highlight their relative strengths in the setting of Instrumental Variables.

# 1 Introduction

In many applied microeconomic contexts, the underlying data is text. Think Health Care, Education, and Housing. Causal inference in this setting has typically proceeded by hand-selecting numerical representations of the text and estimating the corresponding conditional expectation function, assuming the treatment or instrument is locally randomly assigned. This vector based approach to casual inference is so fundamental that it’s often one of the first concepts emphasized in an Econometrics course.

Paul Goldsmith-Pinkham writes in his first [lecture](#)– “We say that  $D_i$  is strongly ignorable conditional on a vector  $X_i$ .”

With the development of Natural Language Processing models, though, a natural question to ask is what would it mean to instead condition directly on the underlying text? And as importantly, in what contexts would applied researchers prefer such an approach? We address these two questions, in this paper, in the style of the popular text *Mostly Harmless Econometrics* – meaning we prioritize the conceptual challenges that practitioners must think through rather than the large sample asymptotic properties of the estimators.

We do so first by extending the selection on observable assumption<sup>1</sup> to text. Formally, using the terminology of *Topology* this allows researchers to speak of identification with respect to textual controls - controls such as Emotional well-being, Education, Health, Finances, Skills, Hope, Faith, Social Skills as in [Evans et al. \[2023\]](#). More practically, though, it highlights how LLMs are that much more black-box than even feed-forward neural networks. That is as practitioners, we’re placing no explicit restrictions on the learnt similarity between observations.

We then show that Large Language Models are potentially attractive in the context of instrumental variables with preferential treatment. That is, in the content when the offer for treatment is prioritized among the subset of people who are randomly eligible. As [Figure 1](#) shows, low takeup rates can have an exponential effect on the sampling error of the IV estimator. Small first-stage estimates can lead to noisy IV estimates. In the context of IV with preferential treatment, though, LLM’s are potentially more “efficient” than standard methods because of their representational abilities. They can incorporate information about how the treatment was prioritized making them more sample efficient.

At a high level, of course, perhaps we should have started by saying that there are a lot of good reasons for not wanting to use Large Language Models for causal inference. They are black-box estimators. They are potentially biased. They are expensive to fine-tune. They are harder to fit than linear models. They hallucinate. There is also an expanding way to construct causal estimates using LLMs – from using them for feature selection, to

---

<sup>1</sup>This assumption underlies almost all popular identification strategies in applied econometrics, the paper provides a starting point for researchers interested in running diff-in-diff / instrumental variables / regression-discontinuity-designs on text. In difference-in-difference, we can think of the  $\Delta \hat{Y}_i \perp D_i | X_i$

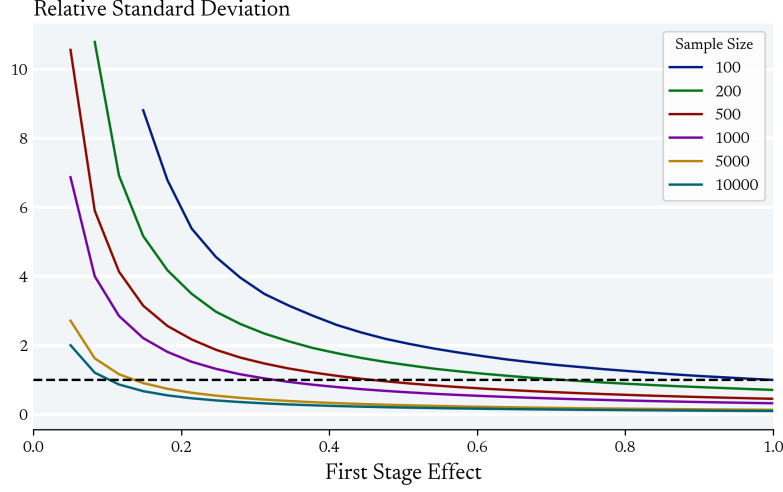


Figure 1: The standard deviation of an Instrumental Variable estimator exhibits exponential like decay with respect to the size of the first stage effect

in-context learning<sup>2</sup>, to fine-tuning (see Figure), which makes an overview almost instantly out of date. The aim of this paper is to provide some insight about the relative tradeoffs bet

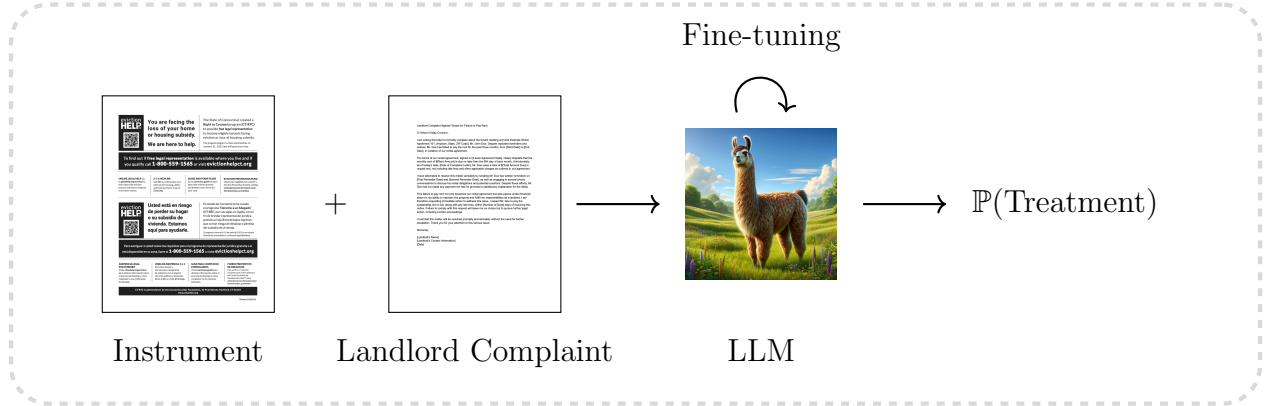


Figure 2: Illustrates how fine-tuning LLMs in the context of instrumental variables. First stage estimates are formed by passing the instrument (a text document) and the controls (another text document) as inputs to a language model which is fine-tuned to output the probability of treatment.

<sup>2</sup>where a frozen language takes the entire dataset in the prompt:  $\mathcal{A} : \{(X_i, Z_i, D_i, Y_i)^n\} \rightarrow X \rightarrow Y$

## 2 Identification

### 2.1 The Fundamental Challenge of Language Models

The fundamental challenge in language models is that the **product topology** breaks down.<sup>3</sup> This is an overly mathematical way to capture the simple idea that context matters which was arguably best expressed by Firth [1957] who wrote, “You shall know a word by the company it keeps.” We take this fantastic line and reinterpret it in the language of topology because topology is the coarsest mathematical construct for talking about casual identification.

A topology is formally a set of subsets of the set of interest that satisfies a series of conditions. For our purposes though, it is a way to encode similarity between elements of a set. For instance, if we think about the set of words, we can define a topology on this set which captures their semantic meaning. One element of the topology may then consist of all words which are related to sports. Another element could contain words related to statistics. And so on. Words which are in the same subsets are in some sense more similar.

We are all implicitly familiar with the standard topology – defined via the open intervals:  $\mathcal{T} = \{(x, y), \forall y > x \in \mathcal{R}\}$ .<sup>4</sup> One useful property of this topology is that it seamlessly extends to n-tuples of real numbers. That is a point  $x := (x_1, x_2)$  is in an element of the topology  $\mathcal{B} := (\mathcal{B}_1 \times \mathcal{B}_2)$  if  $x_1 \in \mathcal{B}_1$  and  $x_2 \in \mathcal{B}_2$ . Which is to say that a point is similar to another point if it is similar across each component of the tuple.

Language doesn’t cooperate with the product topology. Similar sentences have words at each position in the sentence that differ from each other in meaning. Therefore a topology defined on the set of words can not seamlessly scale to a topology on the set of sentences, the way a topology defined on the real numbers scales to a topology defined on vectors.

The recent success of large language models suggests that latently, these models have learned to construct a topology on sentences from an initial word level representations. Or as Radford et al. [2018] notes, they are “captur[ing] higher-level semantics.”

### 2.2 Vector Based Identification

Economists tend to justify identification assumptions by arguing about local variation of the treatment. The typical assumption is that across individuals who share similar features, the treatment (or instrument) can be thought of as good as randomly assigned. That is, locally with respect to the controls there is no selection bias:

$$\mathbb{E}[\tilde{Y}_i(0)|X_i, D_i = 1] \approx \mathbb{E}[\tilde{Y}_i(0)|X_i, D_i = 0] \quad (1)$$

---

<sup>3</sup>Note: The product topology is the coarsest topology under which the projection function,  $p_j : \prod_{i \in J} V_i \rightarrow V_j$  is continuous - [reference](#)

<sup>4</sup>More precisely, the open intervals are a basis for the standard topology

In many areas of applied microeconomics, though, it's common for researchers to have access to the underlying documents. The entire causal inference pipeline, as reflected in Figure 3 then begins with these documents, which are mapped via the Encoder into a vector space, usually a hand-selected feature space. Then a Model, typically a linear model, but potentially a neural network, is fit to this vector space to estimate the corresponding Conditional Expectation Function.

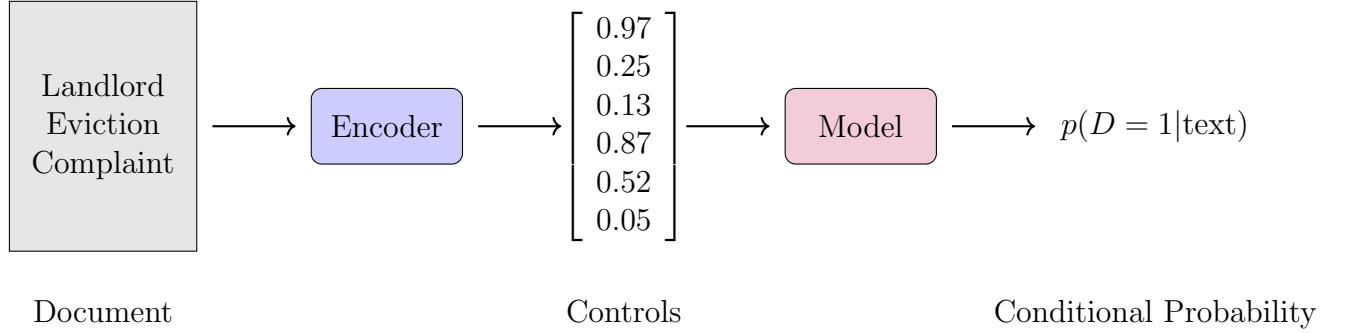


Figure 3: The Standard Pipeline

There are two key aspects in this setup. The first is that the topology on the underlying text can be defined such that the encoder is a continuous function. This allows us to reason about the conditional independence assumption in either the euclidean space defined by the vector representation of controls or more conceptually on the underlying space of text. The second key aspect is that the model is also continuous. Taken together (the composition of continuous functions is continuous), this highlights that the model exploits meaningful variation of the treatment variable. In other words, the causal variation is preserved by the encoder and exploited by the model.

## 2.3 Identification with LLMs

Black-box Large Language models don't expose the encoder or the model separately. This means identification arguments cannot rely on a vector based structure. Everything must be defined or argued about with respect to the underlying textual space. Now, a practitioner may have a fuzzy idea of the underlying topology on the text space under which treatment is as good as randomly assigned but there's no way to pass this information to the language model. And even if there was, a continuous encoder would not preserve the conditional independence.

Therefore, the identification assumption with language models rests on the discrete topology. Or put another way, this means that we're not taking a stand or defining what is local. The argument at the population level is that across the exact same documents treatment is randomly assigned.

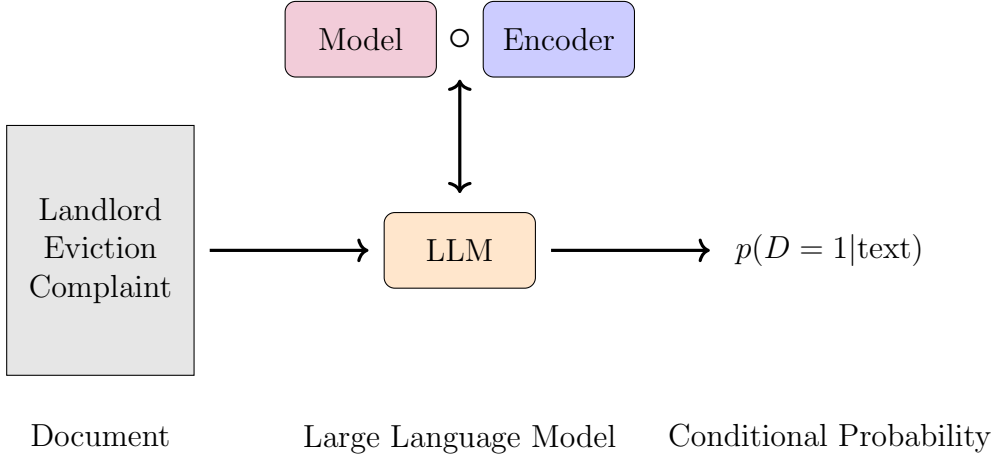


Figure 4: The Black Box Large Language Model Pipeline

## 2.4 Formal Identification

We begin with the underlying probability space which is a tuple of the set containing all possible samples, a  $\sigma$ -algebra, and a probability measure. We subscript these terms by  $n$  to highlight their dependence on the sample size. Note this approach is different from others in the literature which emphasize an observed probability distribution and a counterfactual distribution and express the parameter of interest as a functional of the counterfactual distribution (Kennedy [2022]). With clustered data, which is a staple of empirical economics, the idea that we’re sampling from the observed probability measure doesn’t conceptually work.

$$(\Omega_n, \mathcal{F}_{\Omega_n}, \mathbb{P}_n) \quad (2)$$

On this space, we define the following three random variables of interest: Controls, Treatment, and Potential Outcomes, where for simplicity we assume that both treatment and outcomes are binary.<sup>5</sup>

$$X_i : \Omega_n \rightarrow \mathbb{L}^V \quad (3)$$

$$D_i : \Omega_n \rightarrow \{0, 1\} \quad (4)$$

$$\tilde{Y}_i : \Omega_n \rightarrow \{0, 1\} \rightarrow \{0, 1\} \quad (5)$$

The key aspect is our setup is that the random variable  $X$  transforms the underlying probability space into a probability space defined over text - finite sequence of tokens of length  $L$  from a vocabulary of  $V$  tokens. On this space, we define the discrete topology  $\mathcal{F}_{L^V}$  and assume that conditional on the text, treatment is as good as randomly assigned, which we express as follows.

$$\forall A, B, C \in \sigma(\tilde{Y}_i), \sigma(D_i), \sigma(X_i), \quad \int \mathbb{1}_{A \cap B} d\mathbb{P}_C = \int \mathbb{1}_A d\mathbb{P}_C \int \mathbb{1}_B d\mathbb{P}_C \quad (6)$$

Under this conditional independence assumption, the corresponding conditional expectation function has a casual interpretation. Note that the left hand side is integrated over the population probability space:  $(\Omega, \mathcal{F}, \mathbb{P})$

$$\int_{\Omega} \mathbb{E}[\tilde{Y}_i(1) - \tilde{Y}_i(0)] d\mathbb{P} = \int_{\Omega_n} \mathbb{E}[Y_i | D_i = 1, X_i] - \mathbb{E}[Y_i | D_i = 0, X_i] d\mathbb{P}_n \quad (7)$$

### 3 Instrumental Variables

The typical instrumental variable setup is two stage least squares where we first predict treatment using a linear model of the controls ( $X$ ) and the instrument ( $Z$ ). Then in the second stage, fit a linear model to the outcome ( $Y$ ) using the predicted treatment ( $\hat{D}$ ) as a control.

$$Y_i = \beta_0 + \beta_1 \hat{D}_i + \beta_2 X_i + \varepsilon_i, \quad \hat{D}_i = \hat{\gamma}_1 X_i + \hat{\gamma}_z Z_i \quad (8)$$

Under the Frish Waugh Lovell Theorem, this two step procedure can be understood as a single variable regression where we regress the outcome variable on a single residualized variable. The residual variable in this context is the difference between the predicted treatment based on the control and the instruments ( $\hat{D}_i$ ) and the predicted treatment based only on the controls ( $\tilde{D}_i$ ). This residualized term (which we can also capture via nonparametric methods) highlights the essence of the instrumental variable strategy: use only the local variation of the treatment variable generated by the instrument.

$$Y_i = \beta_1 (\hat{D}_i - \tilde{D}_i) + \eta_i \quad (9)$$

$$Y_i = \beta_1 (\mathbb{E}[D_i | X_i, Z_i] - \mathbb{E}[D_i | X_i]) + \eta_i \quad (10)$$

From this vantage point it's clear that flexible models are potentially more attractive in this context if they're better able to capture the local variation generated by the instrumental variable. Using real data, Figure 5 highlights how the structure imposed on the linear model can actually generated nonsensical variation: individuals offered the instrument have a negative signed residual.

For rest of the paper, we'll work with the nonparameteric residualized model, where we approximate the conditional expectation functions by fine-funning LLMs. As mentioned previously, in this setting both the instrument, the controls, and the treatment will be text.

$$Y_i = \beta_1 (\mathbb{E}[D_i | X_i, Z_i] - \mathbb{E}[D_i | X_i]) + \eta_i \quad (11)$$

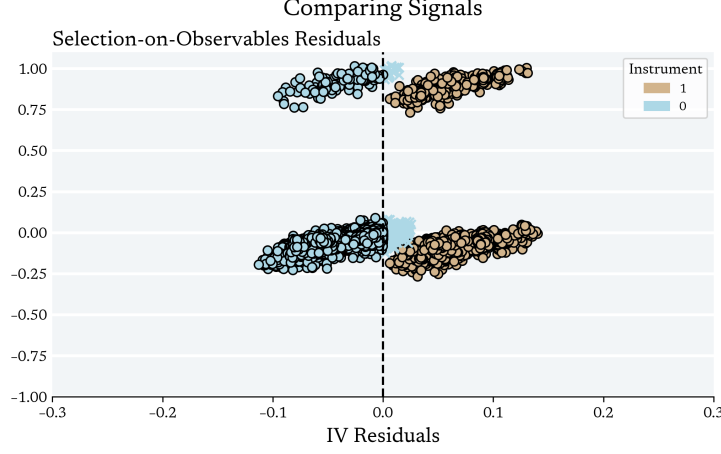


Figure 5: Scatter plot of the residual terms used in OLS and IV colored by the instrumental variable status.

$$\begin{aligned}
\text{Residual}(Z = 1, X) &= \mathbb{E}[D|Z = 1, X] - \mathbb{E}[D|X] \\
&= \mathbb{E}[D|Z = 1, X] - (\mathbb{P}(Z = 1|X)\mathbb{E}[D|Z = 1, X] \\
&\quad + (1 - \mathbb{P}(Z = 1|X)\mathbb{E}[D|Z = 0, X]) \\
&= (1 - \mathbb{P}(Z|X))(\mathbb{E}[D|Z = 1, X] - \mathbb{E}[D|Z = 0, X]) \\
&= \underbrace{(1 - \mathbb{P}(Z|X))(\mathbb{P}(\text{Complier}|X))}_{\geq 0} \\
\text{Residual}(Z = 0, X) &= \underbrace{-\mathbb{P}(Z = 1|X)(\mathbb{P}(\text{Complier}|X))}_{\leq 0}
\end{aligned}$$

## 4 Efficiency

The primary motivation for including controls in the instrumental variable setup is for identification. We would like to have a sufficient set of controls such that the instrumental can be thought of as locally assigned.

$$\tilde{Y}_i \perp Z_i | X_i \quad \tilde{D}_i \perp Z_i | X_i \quad (12)$$

A secondary reason for adding controls is that it can increase the precision of our estimate of the LATE. Intuitively, instrumental variables estimates treatment heterogeneity over a latent variable (the Compliers). One can rightly reason then that perhaps if the compliers are partially observed, we can leverage this information to do better. Using a saturated feature space, we'll show that adding controls which partially differentiate between compliers and non-compliers can (1) increase the variance of the predicted treatment, (2) has no impact



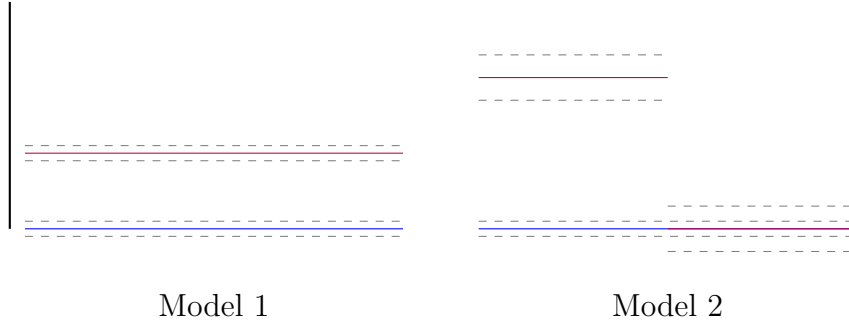


Figure 6: Caption

on the estimator of the First Stage but (3) reduces the variance of the LATE estimator. We'll then show via simulations that these results hold when we relax the saturated feature space assumption, and highlight the language models relative strength.

## 4.1 Saturated Feature Space

Consider the following setup. We' assume that the instrument is randomly assigned, and that we have a categorical variable  $X_i$  which is informative of who is a complier.

We then consider the following two estimators.

$$\hat{\mathbb{E}}[D_i|Z_i] \quad (13)$$

$$\hat{\mathbb{E}}[D_i|X_i, Z_i] \quad (14)$$

As figure 7a highlights, by conditioning on  $X$ , that variance of predictions increase. This isn't surprising. It's more of a sanity check on our simulation. What may be intitially suprising (although it won't be upon further reflection) is that this has no impact on the mean or standard deviation of the first stage estimator as captured in figure 7b.

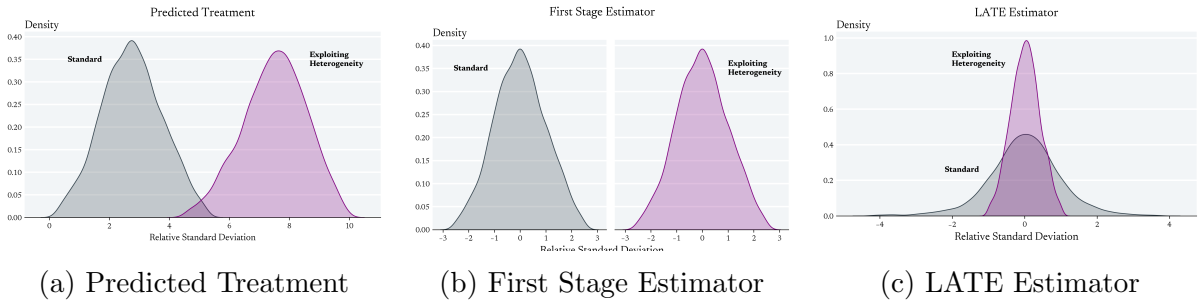


Figure 7

## 4.2 Simulation

The aim of these simulations is to demonstrate that language model – relative to feed-forward neural networks or linear models – *can* exploit characteristics which are associated with compliers. We generate synthetic observations by drawing numerical features:  $x = [x_0, x_1, x_2, x_3, x_4]$ . Using an Anthropic model, and the following prompt, we then map these numerical features into text,  $x \mapsto t^*(x)$  so that we have a numerical representation of the features for the linear and feed-forward model and a textual representation for the LLM.

### Prompt

Task: Write a paragraph description of a tenant in their {age\_group} who is currently {overdue\_phrase} {\$\_}. Mention that they are in relatively {health}, live in a {living\_situation}, have been living there for {months}, and have {pets}. Include some details about their {roommate\_status} who {contribute\_status} to the rent. Also mention somewhere that {additional\_detail}

The instrumental variable is randomly assigned so we don't need to control for any features for identification reasons. A key design choice in this simulation is that (A) the first stage depends heavily on  $x_4$  and (B)  $x_4$  is not passed as input to the linear and feed-forward models. Or put another way, there is a "some information" in the text which is (A) highly indicative of who is a complier and (B) is not a feature that a researcher would have chosen apriori to select as a control variable. We then evaluate a fine-tuned LLM, a feed-forward neural network, a linear model with no controls, and a linear model which does control for contradicting what we said before but we label it the "oracle model" so it's not really a contradiction. We provide greater details of the simulation setup in the accompanying [GitHub repository](#).

In figure 8, we show simulation results for two different first stage functions. Figure 8a shows the sampling distribution for the LATE effect when the treatment is completely determined by the interaction between the instrument and  $x_4$ . Since the LLM has a textual representation of this feature, and the linear and neural network models do not, the LLM is relatively more concentrated around the true parameter. In figure 8b, we simulate all of the control variables are uninformative about the takeup of the treatment variable. With no meaningful signal conveyed by the features, the LLM has greater dispersion relative to the linear model.

## 4.3 Leveraging Prompts

$$f(x) = \sum_i a_i 1_{A_i}(x)$$

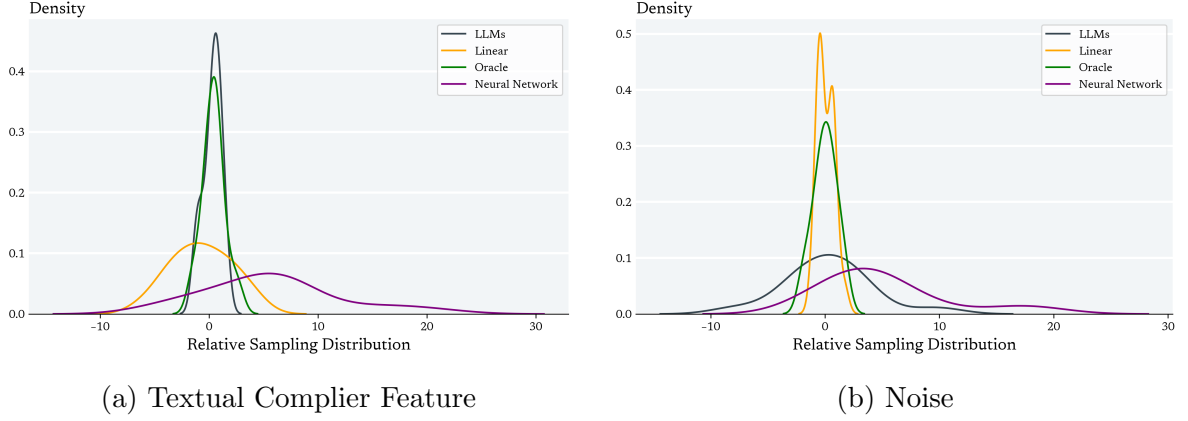


Figure 8: Sampling distributions under two different first stage functions

In certain policy contexts, the offer of a service is (conditionally) randomly assigned. But among those who receive the offer, more individuals would like to take up the offer than can be accommodated (demand exceeds supply). As an example, consider again the roll-out of the Right to Counsel across Connecticut. The offer of legal representation was made available in certain zip codes but not others. But in the zip codes where it was made available, demand for legal representation exceeded the availability of legal aid lawyers.<sup>6</sup> In such settings, it’s typical then for the service to be prioritized. According to one lawyer we talked with, tenants with vouchers and disabilities were often prioritized over other tenants (based on conversations with legal aid lawyers).

The standard approach to leveraging this information would be to fit a linear regression model where we interact the availability of legal aid with a self-constructed variable of whether a tenant has a housing voucher or a disability. There is a wide-number of disabilities, so in practice we would probably like to condition on each type

$$D_i = \beta_0 + \beta_1 \text{Offer}_i + \beta_2 \text{Vulnerable}_i + \beta_3 \text{Offer}_i \times \text{Vulnerable}_i + \varepsilon_i \quad (15)$$

As discussed in section ??, the statistical aim in this setting is to estimate the probability of receiving the offer conditional on covariates and the offer status:  $\mathbb{E} = [D_i | X_i, Z_i]$ . What’s unique, though, is that when service is prioritized, we can perhaps reduce the variance in our estimator by leveraging the additional information of how the offer is prioritized. There

In an instrumental variable setting where the instrument is conditional randomly assigned, and conditional on the offer, demand exceeds supply we may have information about who is likely to take up the treatment. This is most likely in cases where treatment is a result of both an individual accepting the treatment, and the offer of the treatment being prioritized due to limited supply. For example, consider the context of the rollout of the Right to Counsel across Connecticut. Free legal representation was available in only certain zip codes (the instrument). However within these zip codes, given the limited supply of legal aid relative to the demand,

<sup>6</sup>Based on conversations with legal aid lawyers across the state

Language models can leverage this information to improve the precision of the first stage estimate. Mechanically we can do so by prepending each text observation with a description of how treatment is prioritized. As Figure 9 highlights fine-tuned language models with prompting can take advantage of this information – outperforming a standard large language model when the first stage noise is low – and yet is flexible enough to ignore this knowledge if it turns out not to be true – outperforming a fixed large language model when the first stage noise is high.

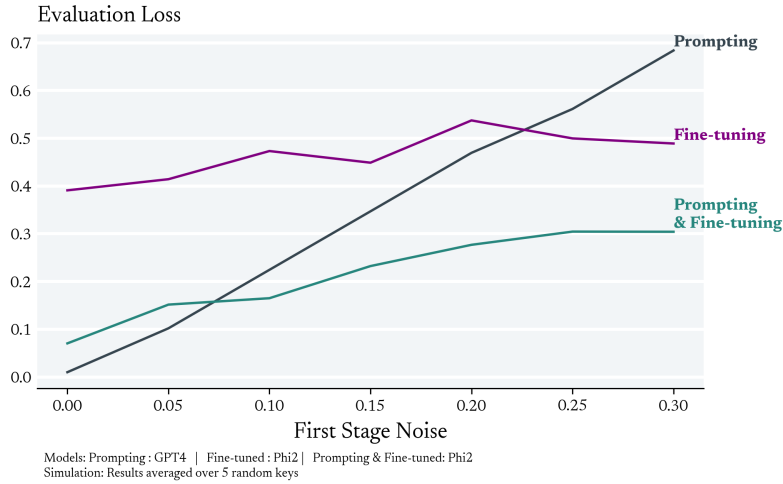


Figure 9: Optimal Evaluation Loss with Early Stopping

#### 4.4 Model Choice

## 5 Interpretability

- “also a challenge for interpretability: how can we hope to understand a function over such a large space, without an exponential amount of time?” [Olah \[2023\]](#)
- “the parameters are a finite description of a neural network, if we can somehow understand them” [Olah \[2023\]](#)

## 6 Implementation

- “use a quantized base model, which is not changed at all by the training, and add trainable LoRA adaptors that are not quantized.” [Reference](#)
- The right model can depend on the structure of your text: [Reference](#)

We offer suggestions based on our own experience of fine-tuning these models. For information regarding the details of actually configuring these models, we refer to reader to more technical material ([Chunked Cross Entropy Forward](#)).

When the probability of receiving an offer is small ( $p(Z = 1) \leq 0.3$ ), and the conditional take-up rate is small ( $p(D = 1|Z = 1) \leq 0.15$ ), fine-tuned LLMs can suffer from mode collapse: the model appears to learn –low validation loss – but the histogram of the models outputs looks like an empty figure. To counter this problem, we (1) augment the cross entropy loss function proposed by scaling the negative log probability the the relative inverse of the probability of receiving the instrument (2) use a relatively long warm-up ratio 0.5 and a lower learning rate  $1e^{-5}$  v.s.  $1e^{-4}$ , and (3) track the average recall rate over the class  $D = 1$ .

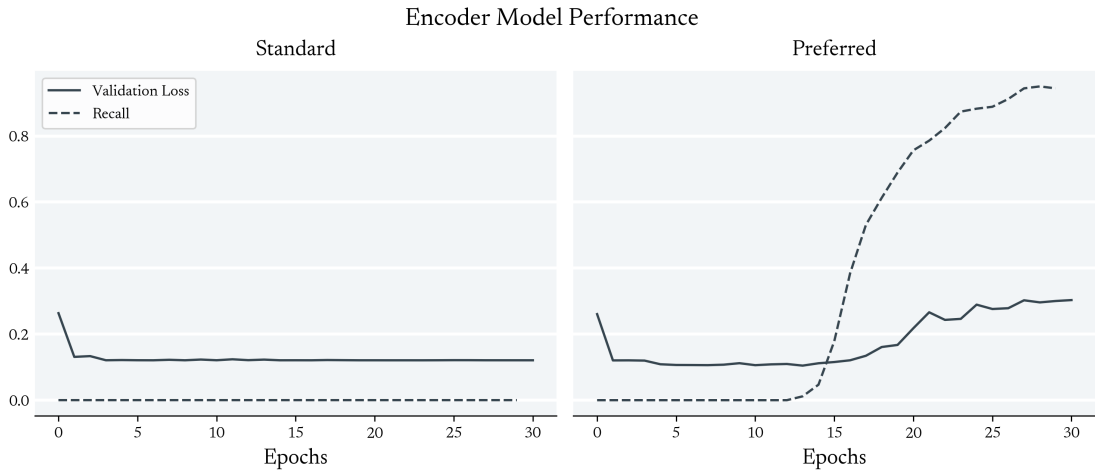


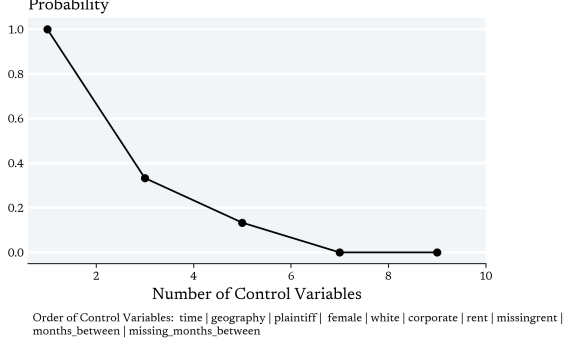
Figure 10: Importance of considering metrics in addition to the validation loss

## 7 Conclusion

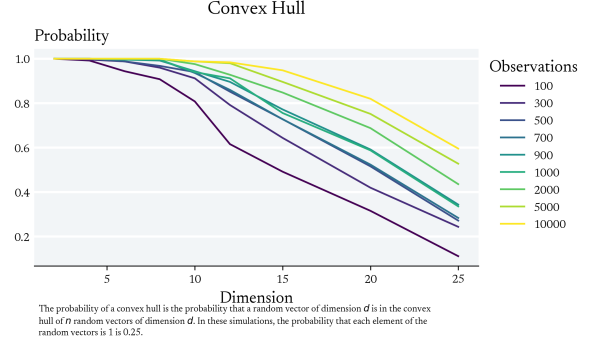
In the finite sample, the success of these language models (as with all models) depends on their ability generalize ([Balestrierio et al. \[2021\]](#)). Despite having within variation of the treatment at the population level, the need to generalize in the finite sample is clear. The first is that we don’t have literal within covariate variation of the treatment. In most cases, the function space that we’re searching over is flexible enough such that we can perfectly fit the treatment variable as a function of the controls. In figure ??, we perfectly interpolate the Right to Counsel indicator based on controls.

Furthermore, in a “typical” empirical settings, observations from the treated group lie outside the convex hull formed from observations in the control group (figure 11).

One motivation for linear models then is that as kernel methods, we understand how they generalize because we’ve specified the similarity function between observations (the inner product). Conversely neural networks, and large language models are black box because



(a) Real Data



(b) Simulated Data

Figure 11: (Left) With a sample size greater than 10,000 observations, with just a small number of controls, the probability that an observation is within the convex hull of the other observations diminishes significantly. (Right) Using simulated data, but focused on dimensions rather than control variables, the same pattern occurs.

the similarity function is learnt during the training phase of the model (Domingos [2020]):  $K(x, x') = \int_{c(t)} \nabla_{\theta} f_{\theta}(x)^T \nabla_{\theta} f_{\theta}(x') dt$ .

One reason for not relying on linear models is that they can generalize poorly. As in the case of residualized IV (explained above), those with an instrument = 1 should have positive residuals while those with the instrument set to 0 should have negative residuals. As figure ?? illustrates, the linear model, by underfitting the data, doesn't necessarily satisfy this condition. Several individuals with the instrument set to 0 have positive residuals and vice-versus.<sup>7</sup> The non-linear language model, meanwhile has no problem satisfying the monotonicity assumption. It cleanly exploits the meaningful variation generated by the instrument.

$$\hat{D}_i = \beta_0 + \sum_{j=1}^k \beta_j x_{ij} + \varepsilon_i \quad (16)$$

$$\hat{D}'_i = \beta'_0 + \sum_{j=1}^k \beta'_j x_{ij} + \beta'_{k+1} z_i + \eta_i \quad (17)$$

$$\text{residual}(x, z = 1) = \hat{\beta}_{k+1} + \sum_{j=0}^k (\hat{\beta}_j - \hat{\beta}'_j) x_j \quad (18)$$

The initial application of neural networks to empirical microeconomics asked us to re-think generalization (Zhang et al. [2021]). The adoption of Large Language Models motivates us to reconsider the meaning of controls. In this paper, we clarify the conceptual framework of causal inference with language models, and highlight that they may be especially attractive in the case of instrumental variables with preferential treatment.

<sup>7</sup>Using data from [The Right to Counsel at Scale](#)

# References

- Randall Balestriero, Jerome Pesenti, and Yann LeCun. Learning in high dimension always amounts to extrapolation. *arXiv preprint arXiv:2110.09485*, 2021.
- Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. A neural probabilistic language model. *Advances in neural information processing systems*, 13, 2000.
- Amanda Bertsch, Maor Ivgi, Uri Alon, Jonathan Berant, Matthew R Gormley, and Graham Neubig. In-context learning with long-context models: An in-depth exploration. *arXiv preprint arXiv:2405.00200*, 2024.
- Vicki Boykis. What are embeddings, 2024. URL [https://vickiboykis.com/what\\_are\\_embeddings/](https://vickiboykis.com/what_are_embeddings/). Accessed: 2024-01-07.
- Ryan Cotterell, Anej Svete, Clara Meister, Tianyu Liu, and Li Du. Formal aspects of language modeling. *arXiv preprint arXiv:2311.04329*, 2023.
- Pedro Domingos. Every model learned by gradient descent is approximately a kernel machine. *arXiv preprint arXiv:2012.00152*, 2020.
- William N Evans, Shawna Kolka, James X Sullivan, and Patrick S Turner. Fighting poverty one family at a time: Experimental evidence from an intervention with holistic, individualized, wrap-around services. Technical report, National Bureau of Economic Research, 2023.
- John Rupert Firth. Ethnographic analysis and language with reference to malinowski’s views. *Man and Culture: an evaluation of the work of Bronislaw Malinowski*, pages 93–118, 1957.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Edward H Kennedy. Semiparametric doubly robust targeted double machine learning: a review. *arXiv preprint arXiv:2203.06469*, 2022.
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- Graham Neubig. Introduction to NLP. YouTube, 2024.
- Chris Olah. Mechanistic interpretability, variables, and the importance of interpretable bases. <https://transformer-circuits.pub/2022/mech-interp-essay/index.html>, 2023. An informal note on some intuitions related to Mechanistic Interpretability.

- Omar Osanseviero. Understanding sentence embeddings. [https://osanseviero.github.io/hackerllama/blog/posts/sentence\\_embeddings/](https://osanseviero.github.io/hackerllama/blog/posts/sentence_embeddings/), 2024. Accessed: 02-08-2024.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Jason Wei, Najoung Kim, Yi Tay, and Quoc V Le. Inverse scaling can become u-shaped. *arXiv preprint arXiv:2211.02011*, 2022a.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022b.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Jian-Qiao Zhu, Haijiang Yan, and Thomas L Griffiths. Recovering mental representations from large language models with markov chain monte carlo. *arXiv preprint arXiv:2401.16657*, 2024.
- Bruno Zimmermann. Ictp diploma - topology. YouTube video, 2016. URL <https://www.youtube.com/watch?v=28BluiBRdUk&t=1599s>. Featured playlist of 20 videos. ICTP Mathematics.