

# Economic AI

Matt Taddy



# What do economists do?

320

	JUNE 1980			
	TABLE 2—TOTAL EXPENDITURE AND OWN-PRICE ELASTICITIES			
	Unconstrained $e_{ii}$	Homogeneous $e_i$	First-Differences Model $e_{ii}$	Homogeneous $e_i$
Food	0.21	-0.07	-0.01	0.17
Clothing	2.00	-0.92	0.04	2.92
Housing	0.30	-0.31	1.51	-0.02
Fuel	1.67	-0.28	0.79	0.84
Drink and Tobacco	1.22	-0.60	1.37	1.17
Transport and Communication	1.23	-1.21	-0.48	-0.67
Other goods	1.21	-0.72	-0.16	-0.31
Other services	1.40	-0.93	-0.62	-0.00
			0.04	0.22
			0.83	-0.94
			1.00	-0.31
			1.37	-0.67
			1.14	-1.23
			2.03	-0.52
			1.03	-0.78

the D.W. statistic shows a sharp

heterogeneity is not a new

problem; Ray Byron;

described to a

problem as far

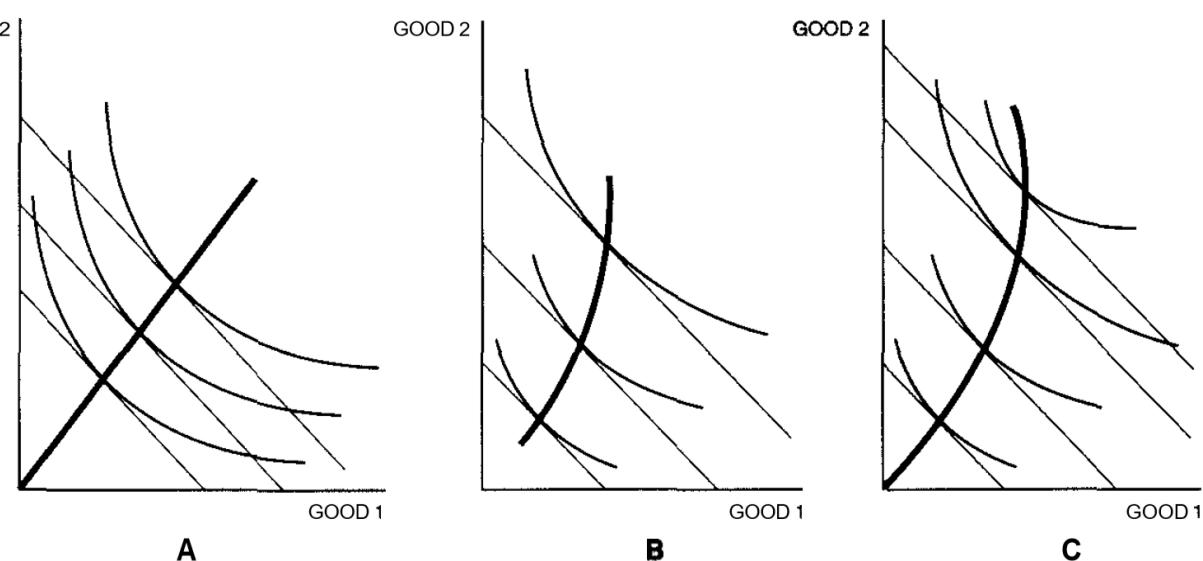
as temporal

A discussion of aggregation above, is that it assume that  $k$ , the distribution of household income structure

Finally, the assumption of separability

is violated.

Temporal



**Income expansion paths.** Panel A depicts unit elastic demands, in panel B good 2 is a luxury good, and in panel C, good 1 is an inferior good.

# What do they need to do today?

Microsoft Azure

New HDInsight Cluster

Cluster Type configuration

Cluster Type configuration

Report a bug

Cluster Type configuration

Learn about HDInsight and cluster versions. [Learn more](#)

Cluster Type: Spark

Cluster Tier ([more info](#)): STANDARD

Operating System: Linux

Version: Spark 2.0.0 (HDI 3.5)

Administration

Manage, monitor, connect

Scalability

On-demand node scaling

99.9% Uptime SLA

toddler shoes

Web Images Videos Maps News Explore

91,800,000 RESULTS Any time

[Shop DSW Kids Shoes | dsw.com](#)  
www.dsw.com/kids · DSW, Inc.

The Latest **Kids** Styles @ Participating DSW Stores Today!  
Find the Best Selection of Infant, Toddler, and Best Brands for Infants, Toddlers, and ...  
taddy@microsoft.com

**Sign Up for DSW® Rewards**  
Earn a \$10 Certificate with Your First Purchase. Free to Enroll!

**Find a Store Near You**  
More Than 480 Locations Available. Shop at a DSW® Near You Today!

Ads

New Balance 150 Slip \$39.99 New Balance

Off Sale - Last Day!

Economic AI breaks complex business  
questions into a structure of ML tasks  
We need to scale-up causal inference

# Causality links actions to their consequences

Good decisions are a direct result of causal understanding (or luck!)

**Pricing:** Ex. How much will sales rise *if I lower prices?*

**Education:** Ex. How does learning change if I *reduce class size?*

**Marketing:** Ex. What is the *causal ROI* from this ad campaign?

This type of reasoning is absent in most AI systems

Without experiments (AB tests)  
causal inference is notoriously difficult

# Traditional methods don't scale

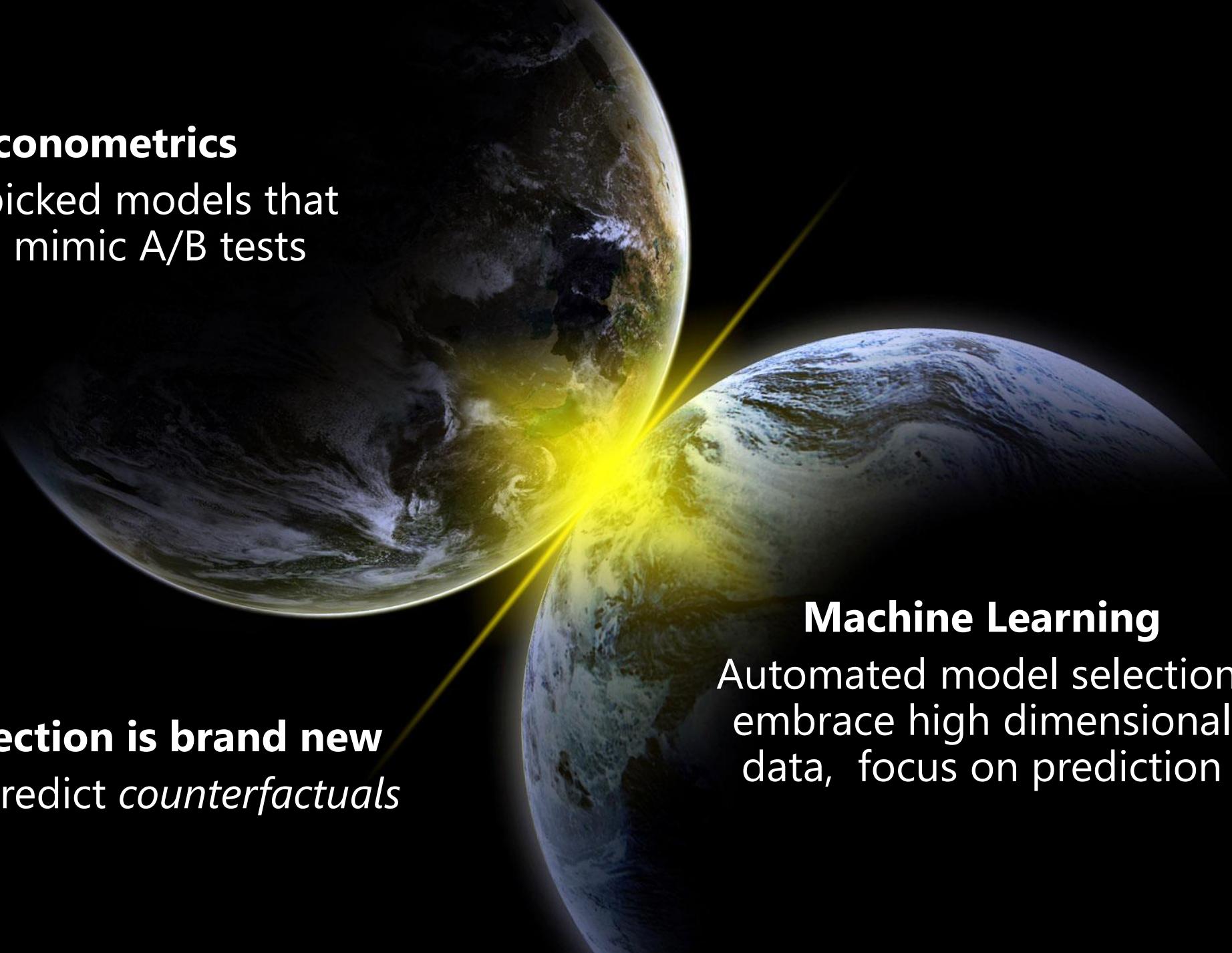
Economists and statisticians look for and analyze **natural experiments**

## Example

- Question: What is the impact of going to charter school on college success?
- Confound: students who seek charter schools are different *to begin with*
- Experiment: compare similar students who live close to catchment boundary

## Limitations

- **Hugely labor intensive**
- Econometric methods limit what data you can work with
- Too cute: these natural experiments are often special scenarios

A photograph of two Earth-like planets in space. One planet is in the foreground, showing its blue oceans and white clouds. A second, larger planet is behind it. A bright yellow beam of light connects the two planets, symbolizing the intersection or connection between them.

## Econometrics

Hand picked models that  
try to mimic A/B tests

**The intersection is brand new**  
Use ML to predict *counterfactuals*

## Machine Learning

Automated model selection,  
embrace high dimensional  
data, focus on prediction

# Application: measuring price sensitivity

If I **drop** price 1%, by what % will quantity sold **increase**?

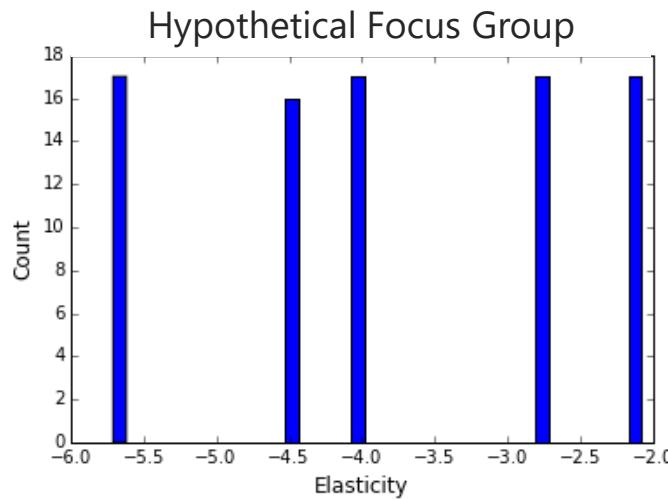
Ex. -3 → drop price 1%, quantity sold goes up 3%

Problem: both prices and sales respond to underlying demand

Need the causal effect of price on sales, not their co-movement

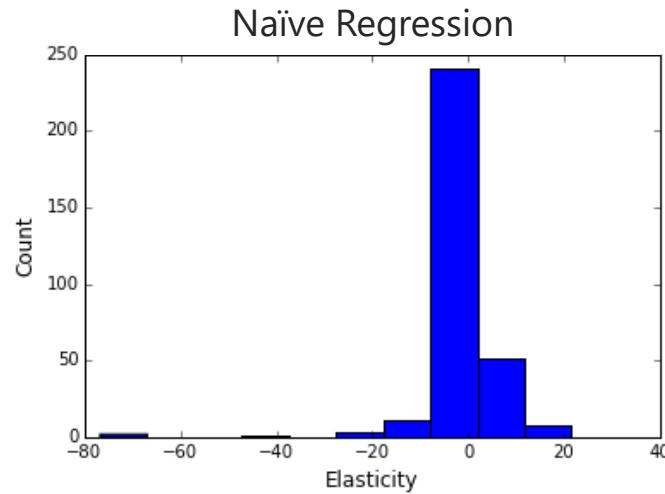
# Buying Beer

## Typical status quo



**Model is too simple, unreliable**

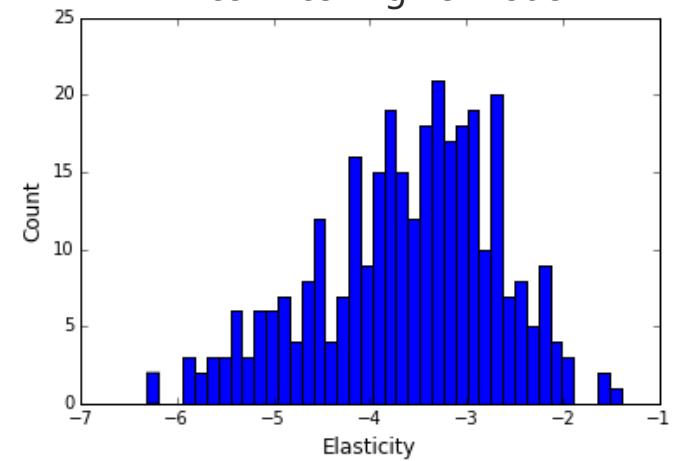
## Old school econometrics



**Model can't handle complexity**



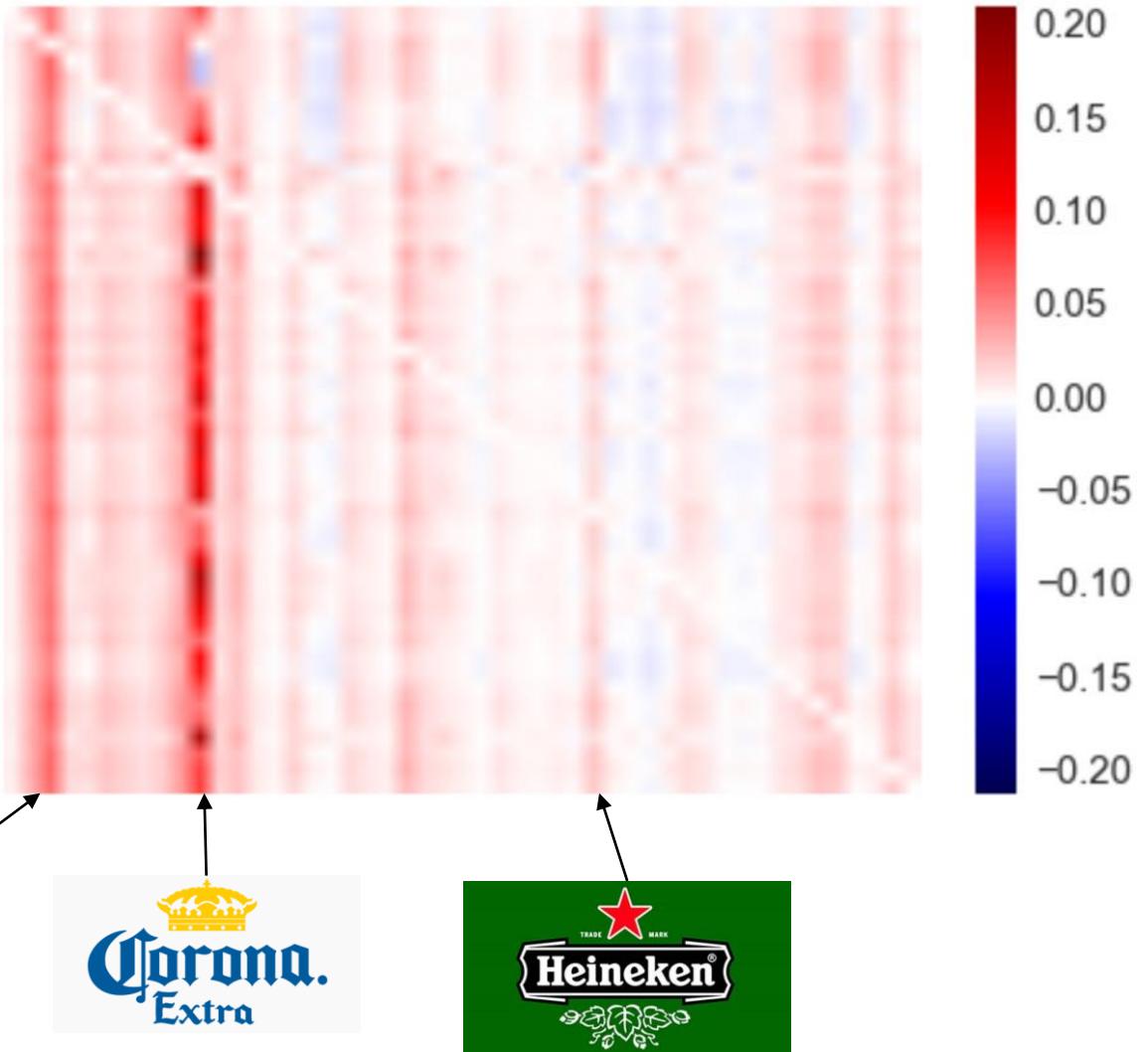
## Alice Price Engine Model



**Selects right level of complexity  
Encourages active experimentation**

# Beer Store cross-product compensated elasticities

*Elasticity matrix (omitting diagonal)*



# How do we do this? How is it automatic?

If we have all of the **price setter's** information about **demand** (product characteristics, demographics, holidays, competitors) then causal inference breaks into two ML tasks:

1. Predict prices from the demand variables:  $p \sim x$
2. Predict sales from the demand variables:  $y \sim x$

The effect of residuals from 1 on those from 2 is causal:

$$(y - \hat{y}(x)) \sim (p - \hat{p}(x))$$

'double ML' by Chernozhukov et al.

But what if you don't have all demand variables? Use upstream randomization

# Endogenous Errors

$$y = g(p, \mathbf{x}) + e \text{ and } \mathbb{E}[p e] \neq 0$$

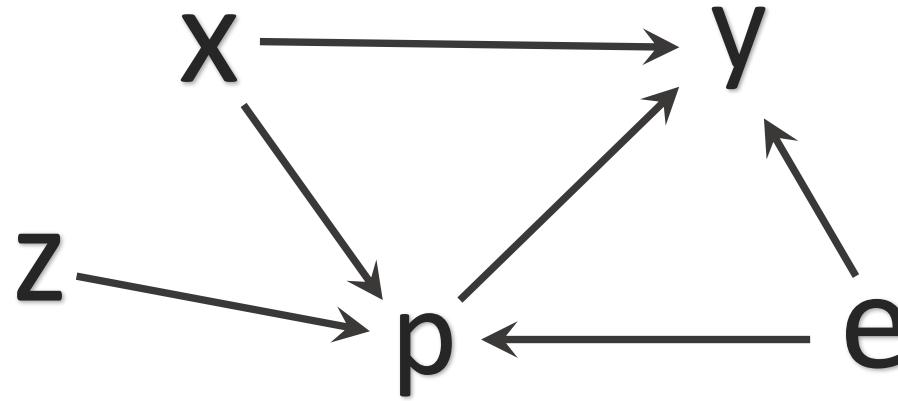
If you estimate this using naïve ML, you'll get

$$E[y|p, \mathbf{x}] = E_{e|p}[g(p, \mathbf{x}) + e] = g(p, \mathbf{x}) + E[e|p, \mathbf{x}]$$

This works for **prediction**. It doesn't work for **counterfactual inference**:

*What happens if I change  $p$  independent of  $e$  ?*

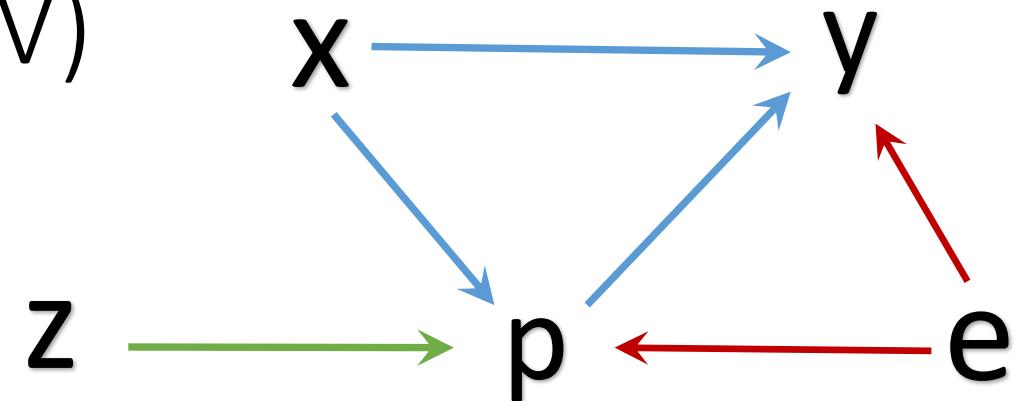
# Instrumental Variables (IV)



In IV we have a special  $z \perp e$  that influences policy  $p$  but not response  $y$ .

- Supplier costs that move price independent of demand (e.g., fish, oil)
- Any source of treatment randomization (intent to treat, AB tests, lottery)

# Instrumental Variables (IV)



The *exclusion structure* implies

$$\mathbb{E}[y|x, z] = \int g(p, x) dF(p|x, z)$$

You can observe and estimate  $\hat{\mathbb{E}}[y|x, z]$  and  $\hat{F}(p|x, z)$

⇒ to solve for *structural*  $g(p, x)$  we have an inverse problem.

cf Newey+Powell 2003

$$\min_{g \in G} \sum \left( y_i - \int g(p, x_i) dF(p|x_i, z_i) \right)^2$$

**2SLS:**  $p = \beta z + \nu$  and  $g(p) = \tau p$  so that  $\int g(p) dF(p|z) = \tau \mathbb{E}[p|z]$

So you first regress  $p$  on  $z$  then regress  $y$  on  $\hat{p}$  to recover  $\hat{\tau}$ .

$$\min_{g \in G} \sum \left( y_i - \int g(p, x_i) dF(p|x_i, z_i) \right)^2$$

Or nonparametric sieves where  $g(p, x_i) \approx \sum_k \gamma_k \varphi_k(p, x_i)$  and

$$\mathbb{E}_F[\varphi_k(p, x_i)] \approx \sum_j \alpha_{kj} \beta_j(x_i, z_i) \text{ (Newey+Powell)}$$

or

$$\mathbb{E}_F[y_i - \sum_k \gamma_k \varphi_k(p, x_i)] \approx \sum_j \alpha_j \beta_j(x_i, z_i) \text{ (BCK, Chen+Pouzo)}$$

*Also Darolles et al (2011) and Hall+Horowitz (2005) for kernel methods.*

But this requires careful crafting and will not scale with  $\dim(x)$

$$\min_{g \in G} \sum \left( y_i - \int g(p, x_i) dF(p|x_i, z_i) \right)^2$$

Instead, we propose to **target the integral loss function directly**  
For discrete (or discretized) treatment

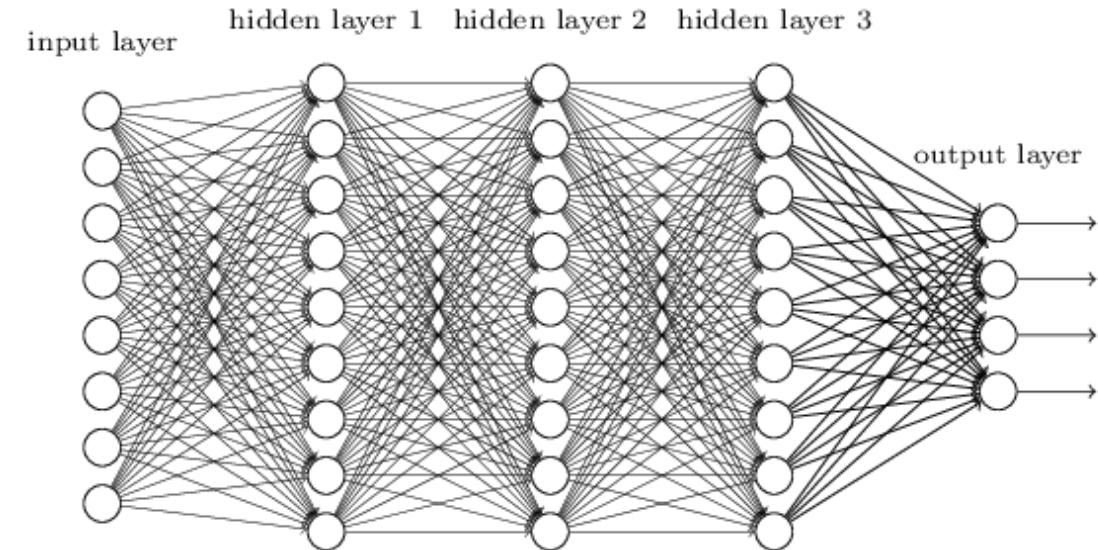
- Fit distributions  $\hat{F}(p|x_i, z_i)$  with probability masses  $\hat{f}(p_b|x_i, z_i)$
- Train  $\hat{g}$  to minimize  $[y_i - \sum_b g(\hat{p}_b, x_i) \hat{f}(p_b|x_i, z_i)]^2$

And you've turned IV into two *generic* machine learning tasks

# Deep Neural Networks

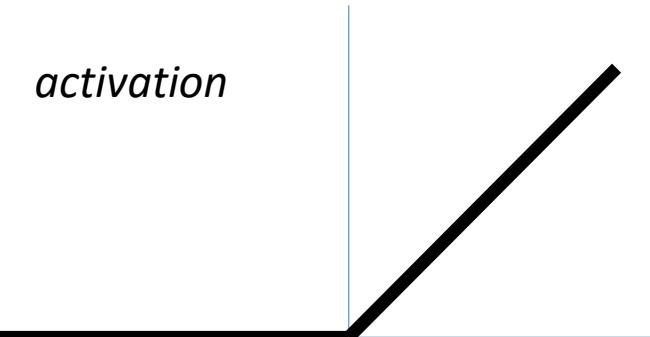
Massive number of parameters,  
mapping output of each layer to  
each node activation in the next

$$\mathbf{z}_i^L \rightarrow h_k(\langle W_k^{L+1}, \mathbf{z}_i^L \rangle)$$



Regularize

- deviance penalties  $\lambda \|W\|$
- dropout training (zeros in grad)



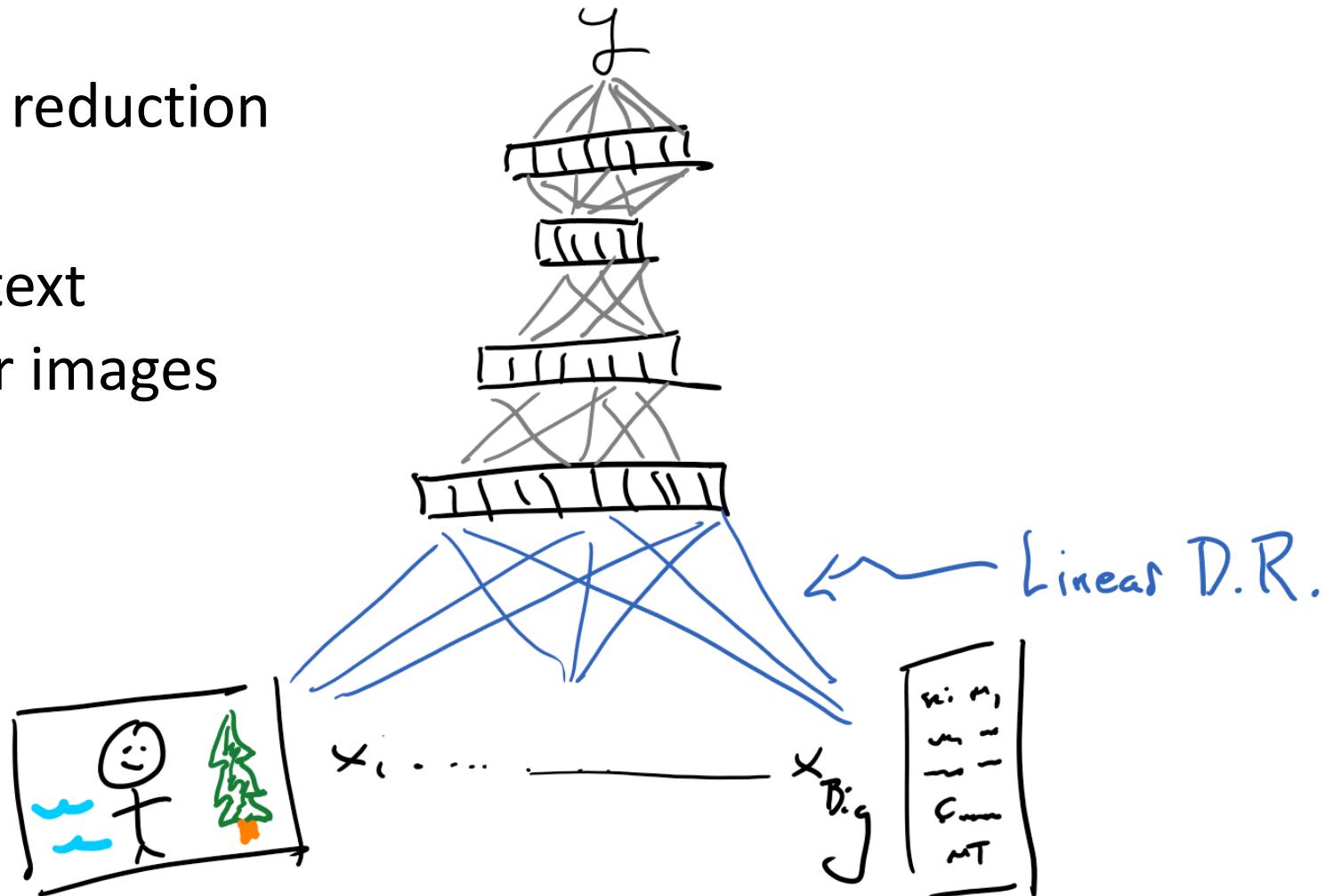
# Deep nets are not really sieves

1<sup>st</sup> layer is a big dimension reduction

e.g.,

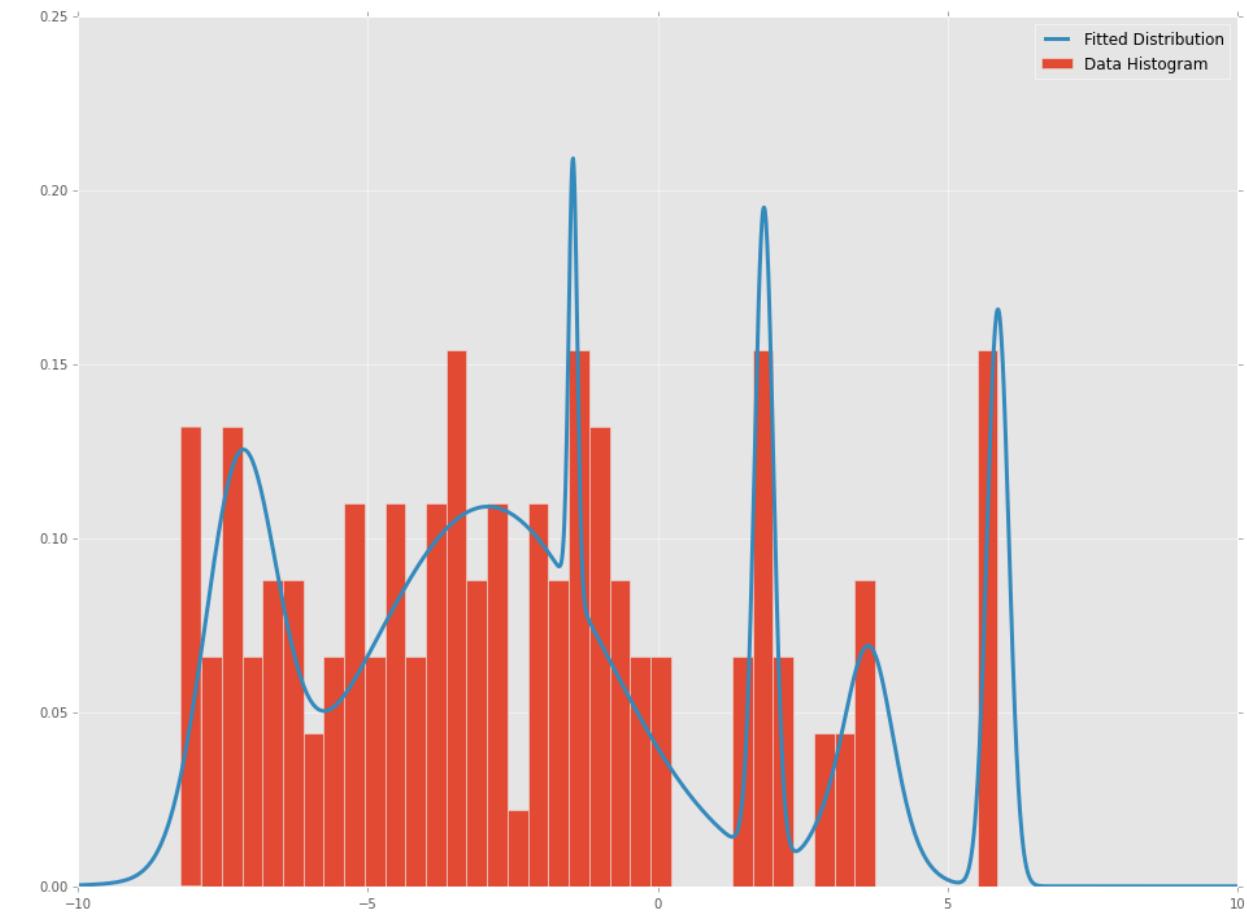
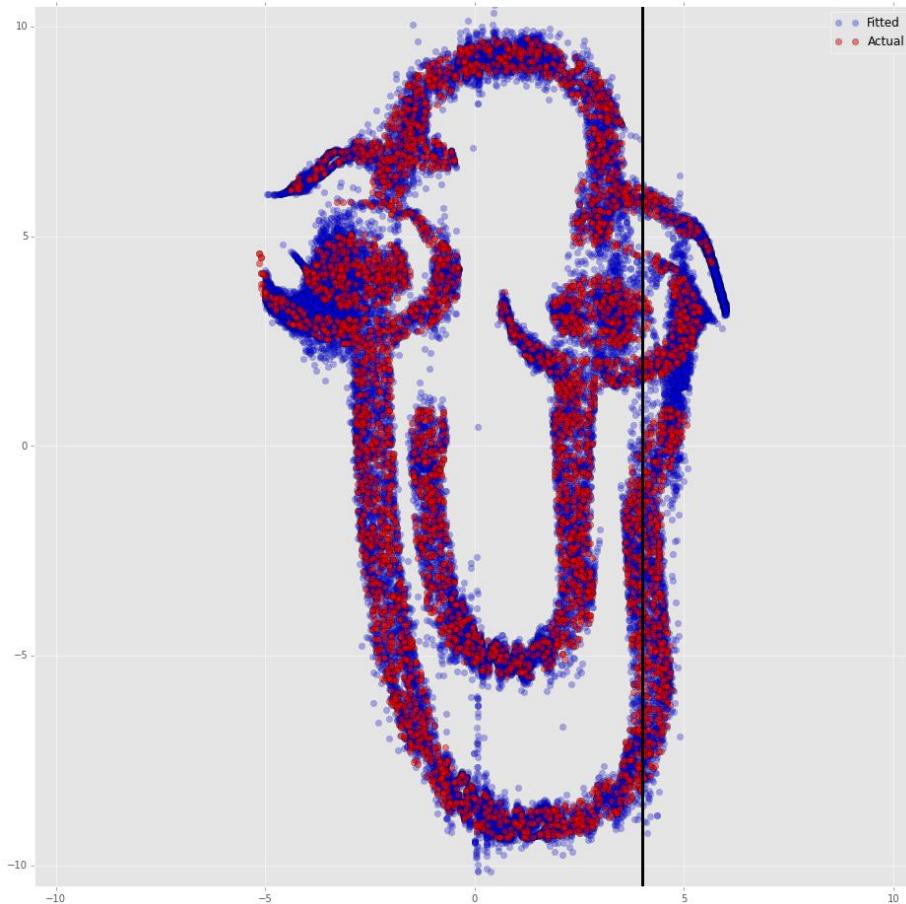
- word embedding for text
- matrix convolution for images

We need to study these...



e.g., first-stage learning for  $F(p|x_i, z_i)$

Bishop 96: Final layer of network parametrizes a mixture of Gaussians



## Stage 2: Integral Loss

The second stage involves an integral loss function

If  $p$  is not discrete or can take many values, not easy!

Brute force just samples from  $\hat{F}(p|x_i, z_i)$  and you take gradients on

$$\frac{1}{N} \sum_i \left( y_i - \frac{1}{B} \sum_b g(p_b, x_i; \theta) \right)^2, \quad p_b \sim \hat{F}(p|x_i, z_i)$$

This is what economists usually do, but this is super inefficient

# Stochastic Gradient Descent

You have loss  $L(\mathbf{D}, \theta)$  where  $\mathbf{D} = [\mathbf{d}_1 \dots \mathbf{d}_N]$

In the usual GD, you iteratively descend

$$\theta_t = \theta_{t-1} - \mathbf{C}_t \nabla L(\mathbf{D}, \theta_{t-1})$$

In SGD, you instead follow *noisy* but *unbiased* sample gradients

$$\theta_t = \theta_{t-1} - \mathbf{C}_t \nabla L(\{\mathbf{d}_{t_b}\}_{b=1}^B, \theta_{t-1})$$

# SGD for integral loss functions

Our one-observation stochastic gradient is

$$\nabla L(d_i, \theta) = -2 \left( y_i - \int g_\theta(p, x_i) d\hat{F}(p|x_i, z_i) \right) \int g_\theta'(p, x_i) d\hat{F}(p|x_i, z_i)$$

Do SGD by pairing each observation with *two independent* treatment draws

$$\nabla \hat{L}(d_i, \theta) = -2(y_i - g_\theta(\dot{p}, x_i)) g'_\theta(\ddot{p}, x_i), \quad \dot{p}, \ddot{p} \sim \hat{F}(p|x_i, z_i)$$

So long as the draws are independent,  $\mathbb{E} \nabla \hat{L}(d_i, \theta) = \mathbb{E} \nabla L(d_i, \theta) = L(\mathbf{D}, \theta)$

# Validation and model tuning

We can do OOS *causal validation*

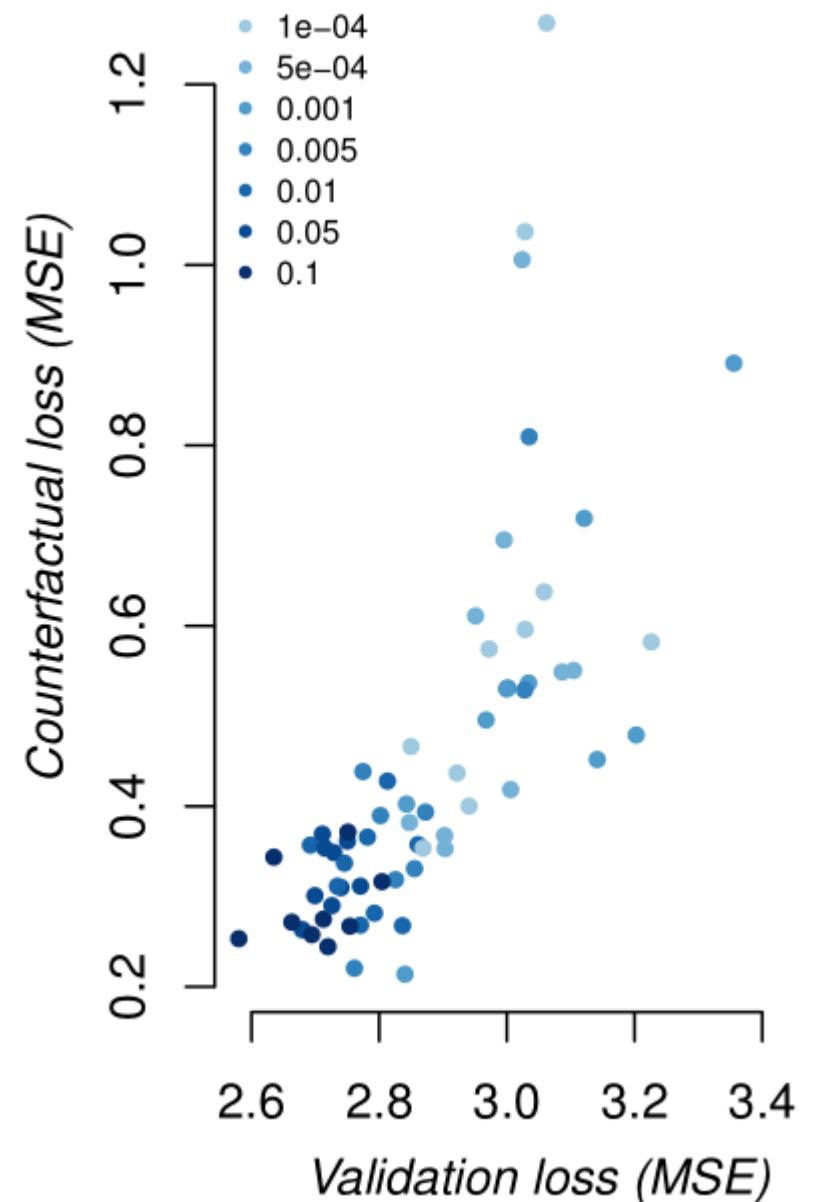
Leave-out deviance on first stage

$$\sum_{i \in LO} -\log \hat{f}(p|x_i, z_i)$$

Leave-out loss on second

$$\sum_{i \in LO} (y_i - \int g_\theta(p, x_i) d\hat{F}(p|x_i, z_i))^2$$

You want to minimize both of these (in order).

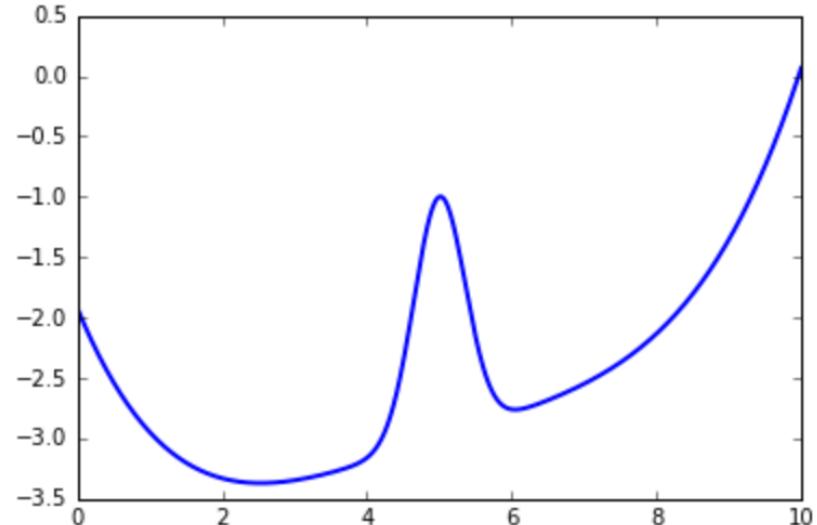


# heterogeneous price effects

$$y = 100 + s\psi_t + (\psi_t - 2)p + e,$$

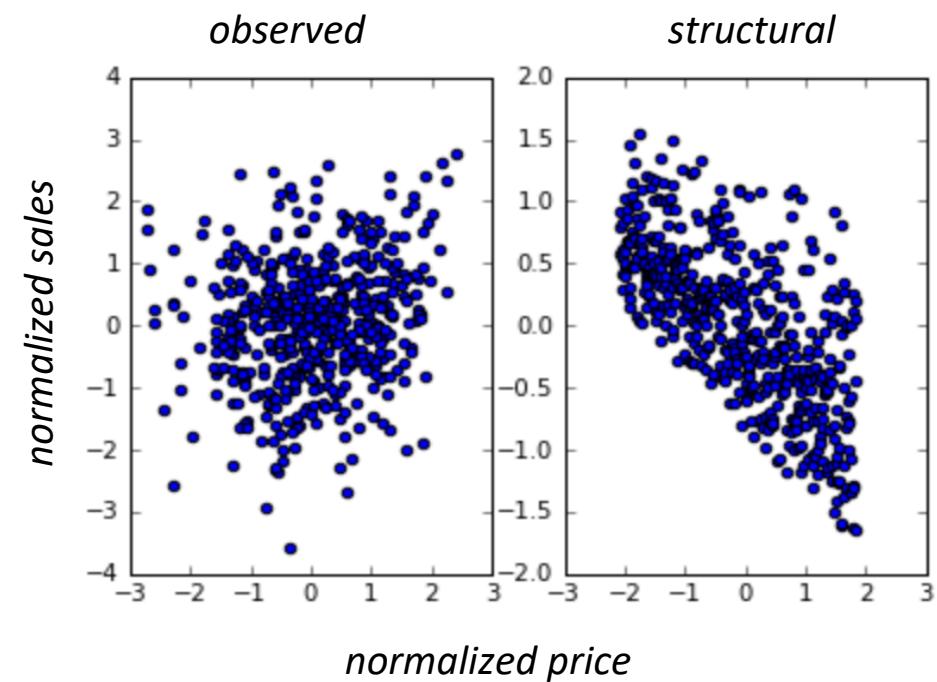
$$p = 25 + (z + 3)\psi_t + v$$

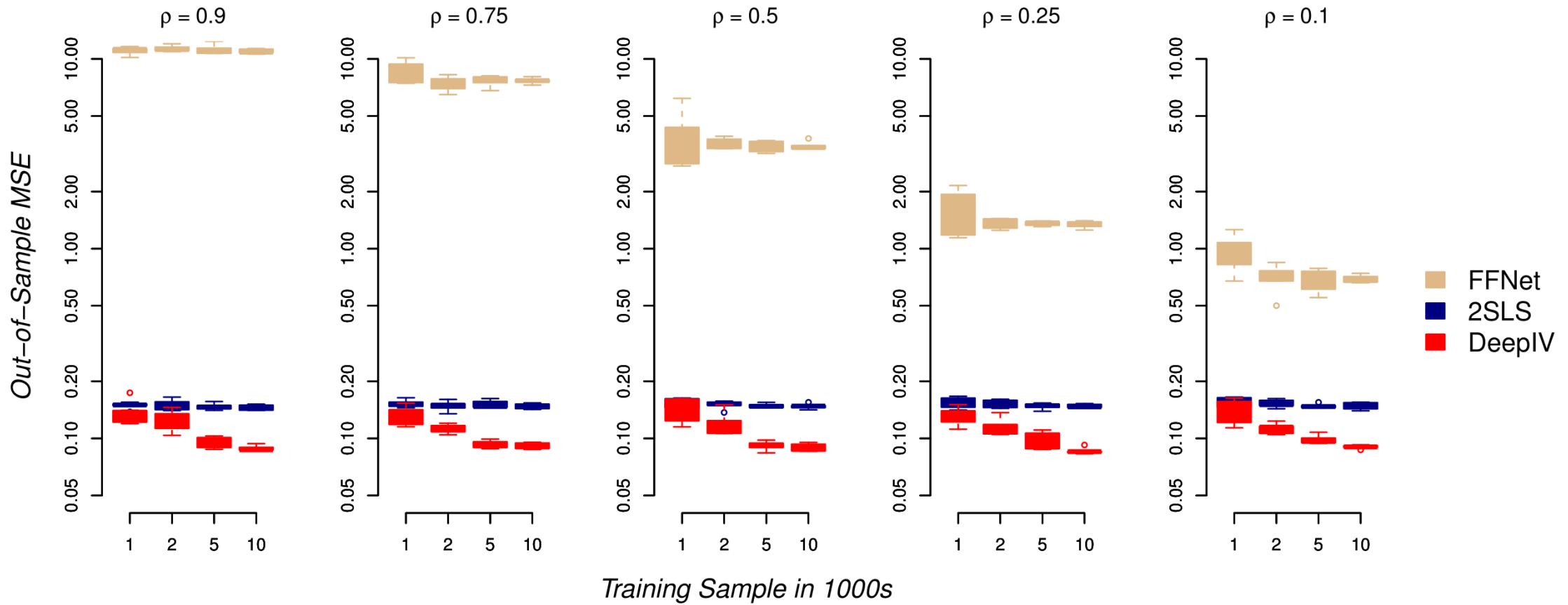
$$z, v \sim N(0, 1) \text{ and } e \sim N(\rho v, 1 - \rho^2),$$



'time' dependent prices, sensitivity, utility

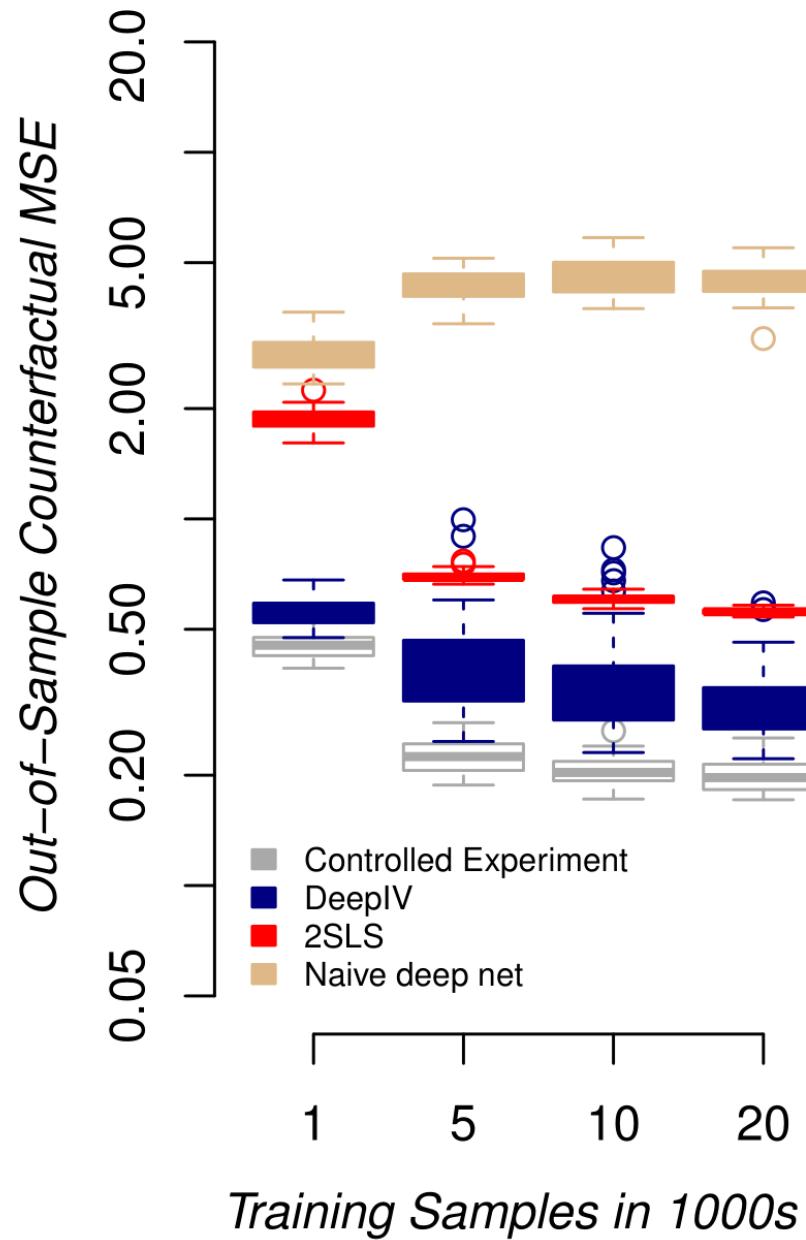
Customer 'type' 1-7 impacts demand





Make it  
harder...

1 1 5 4 3  
7 5 3 5 3  
5 5 9 0 6  
3 5 2 0 0



Inference? Good question

# Frequentist inference

**Data split!** Get top node values and averages on left-out data:

$$\eta_{ik} = \eta_k(x_i, p_i) \text{ and } \bar{\eta}_{ik} = \mathbb{E}_{\hat{F}(p|x_i, z_i)} \eta_k(x_i, p)$$

Stack as instruments  $\bar{H} = [\bar{\eta}_1 \cdots \bar{\eta}_L]'$  and treatments  $H = [\eta_1 \cdots \eta_L]'$

Post-net 2SLS coefficients are  $\hat{\beta} = (\bar{H}' H)^{-1} \bar{H}' y$  with variance  $V_\beta$  and

$$\text{var}[\hat{g}(x, p)] = \boldsymbol{\eta}'(x, p) V_\beta \boldsymbol{\eta}(x, p)$$

# Bayesian inference

Variational Bayes: fit  $q$  to minimize  $\mathbb{E}_q[\log q(W) - \log p(W|D)]$

Diversion... in training we use **dropout**:

At each SGD update, calculate gradients against  $W^l = \Xi^l \Omega^l$  at layer  $l$  where

$$\Xi^l = \text{diag}(\xi_{l1} \dots \xi_{lK_l}), \quad \xi_{kj} \sim \text{Bern}(c)$$

i.e., dropout randomly drops *rows* of each layer's weight matrix

## Dropout is Variational Bayes!

VB minimizes  $KL(q) \propto \mathbb{E}_q[\log q(W) - \log p(\mathbf{D}|W) - \log p(W)]$

If  $q(W) = \prod_l \prod_k (c\mathbb{1}_{[W_k^l = \Omega_k^l]} + (1 - c)\mathbb{1}_{[W_k^l = 0]})$  and  $w \sim N(0, \lambda^{-1})$ ,

$$KL(q) \propto \mathbb{E}_q L(\mathbf{D}|W) + c\lambda|W|_2 + K [c \log c + (1 - c) \log(1 - c)]$$

(see also Gal and Ghahramani)

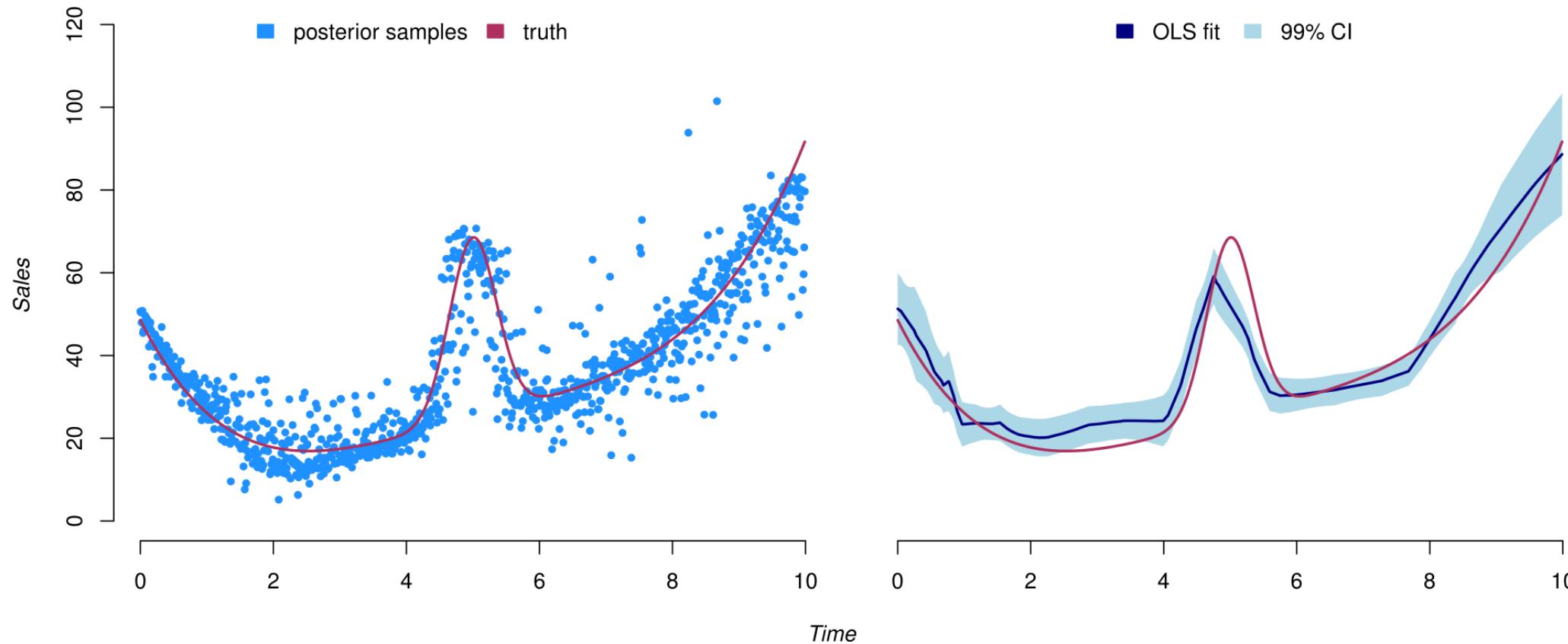
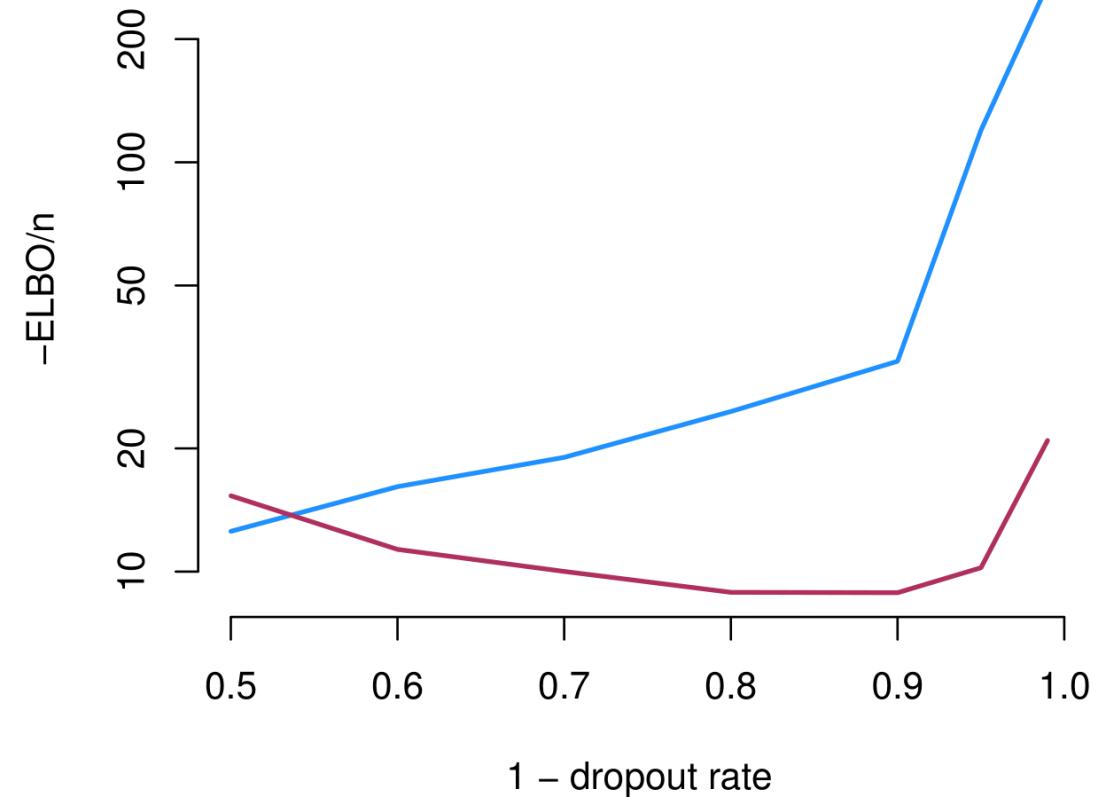
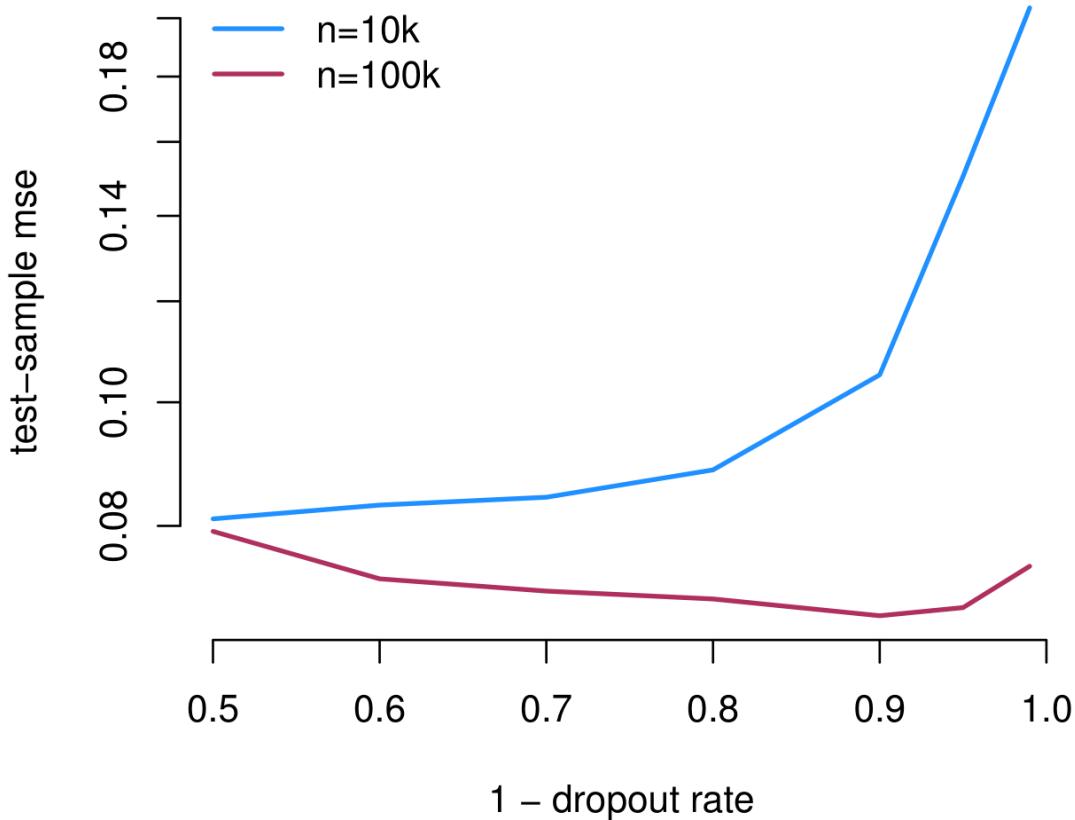


Figure 3: Bayesian (left) and Frequentist (right) inference for a central slice of the counterfactual function, taken at the average price and in our 4<sup>th</sup> customer category. Since the price effect for a given customer at a specific time is constant in (27), the curves here are a rescaling of the customer *price sensitivity* function.

Tuning the dropout rate is like treating it as a variational parameter



# Ads Application

Taken from Goldman and Rao (2014)

We have 74 mil click-rates over 4 hour increments for 10k search terms

Treatment: ad position 1-3

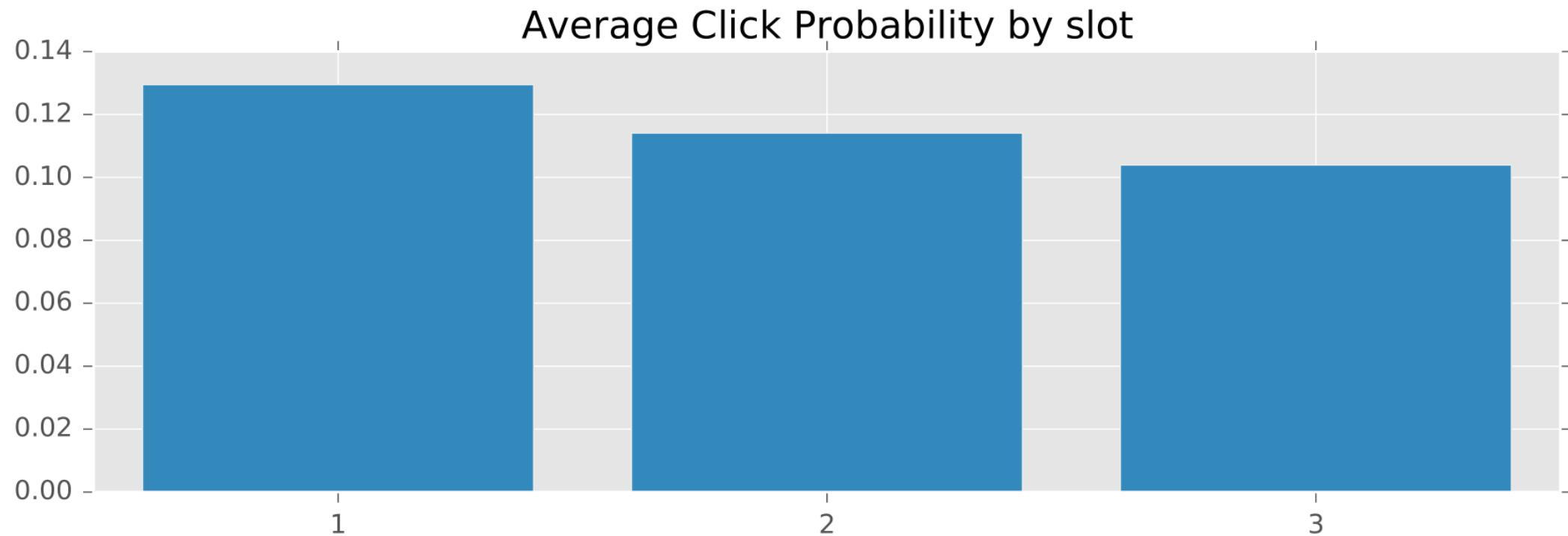
Instrument: background AB testing (bench of ~ 100 tests)

Covariates: advertiser id and ad properties, search text, time period

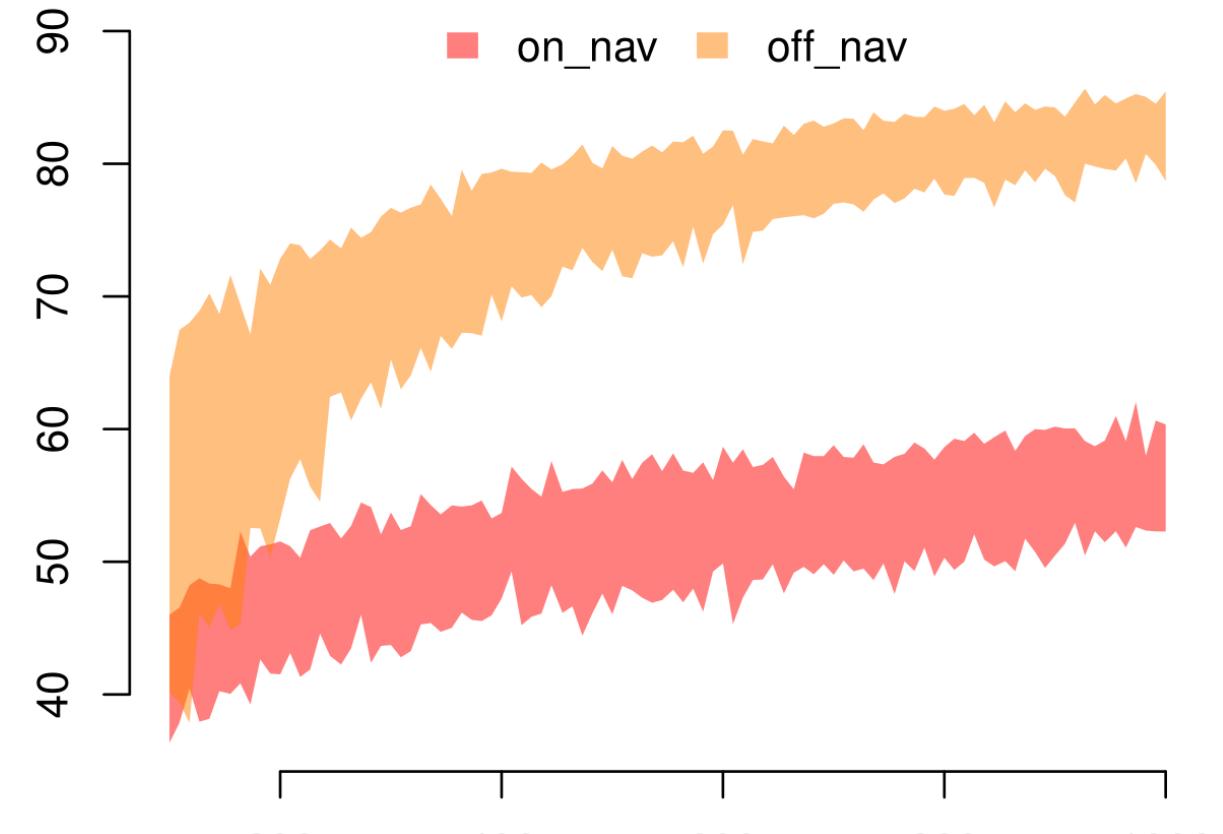
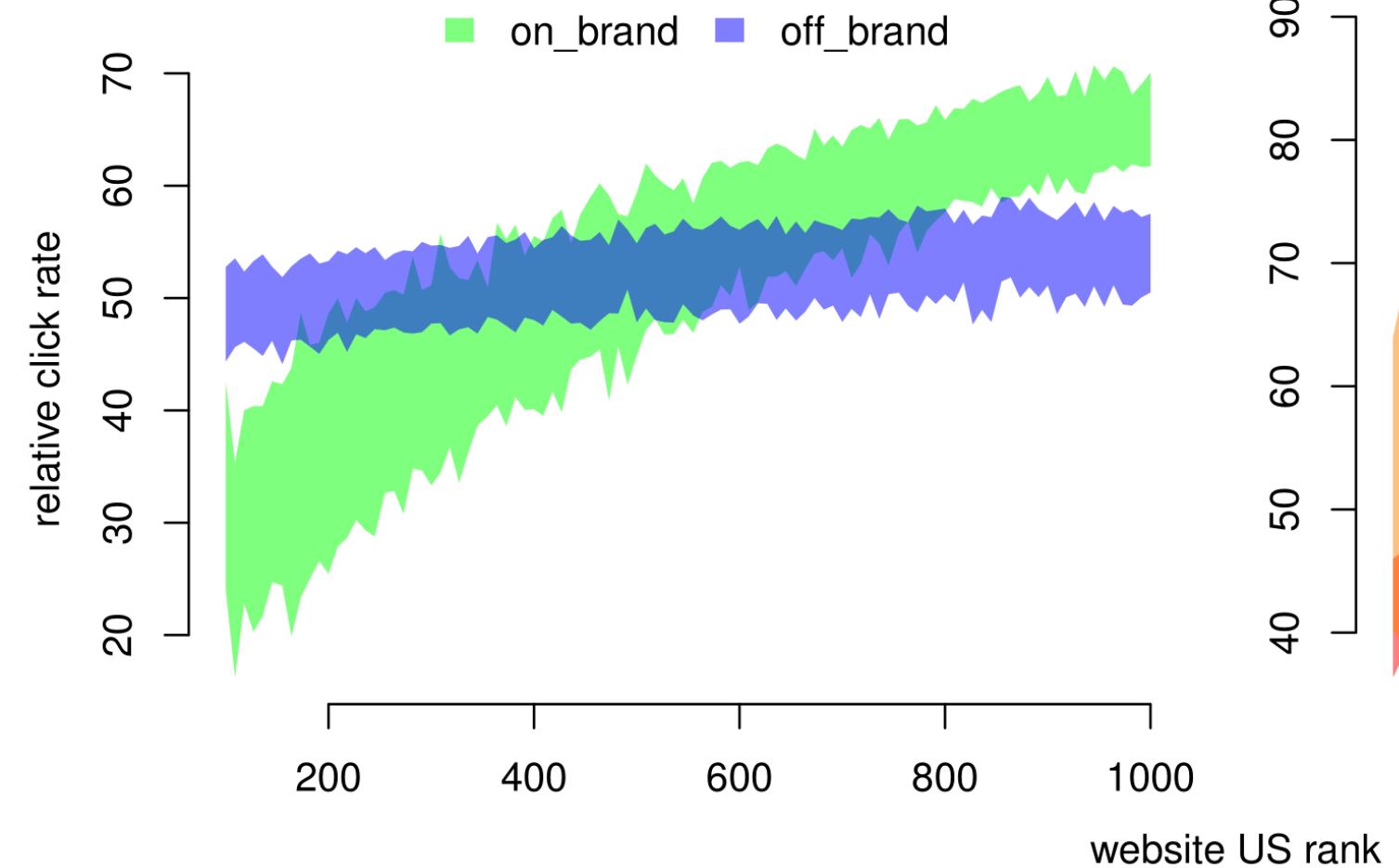
Ads in slot 1 are inherently different (better?) than those in slot 2

We need causal inference for the effect of position on clicks

# Average Treatment Effects



These compare to observed click probabilities of 0.33, 0.1, and 0.05.



Heterogeneity across advertiser and search

# How do we do this?

The ML behind ad click-probability is regularly AB tested

This gives **upstream randomization** – many natural experiments

We can then break causal inference into

1. Model ad position changes with huge number of AB tests
2. Model how ad click rates change with the same tests

And again: you can combine 1+2 to get causal position effects

# Alice

Established: December 5, 2016



## Automated Learning and Intelligence for Causation and Economics

We use economic theory to build systems of tasks that can be addressed with deep nets and other state-of-the-art ML.

This is the construction of systems for *Economic AI*