

# DeepIV: A Flexible Approach for Counterfactual Prediction

Jason Hartford and Kevin Leyton-Brown

University of British Columbia

Greg Lewis\* and Matt Taddy†

Microsoft Research &

\*NBER / † University of Chicago



a place of mind

THE UNIVERSITY OF BRITISH COLUMBIA

Microsoft®

Research

I need a model that predicts the effect of price on ticket sales

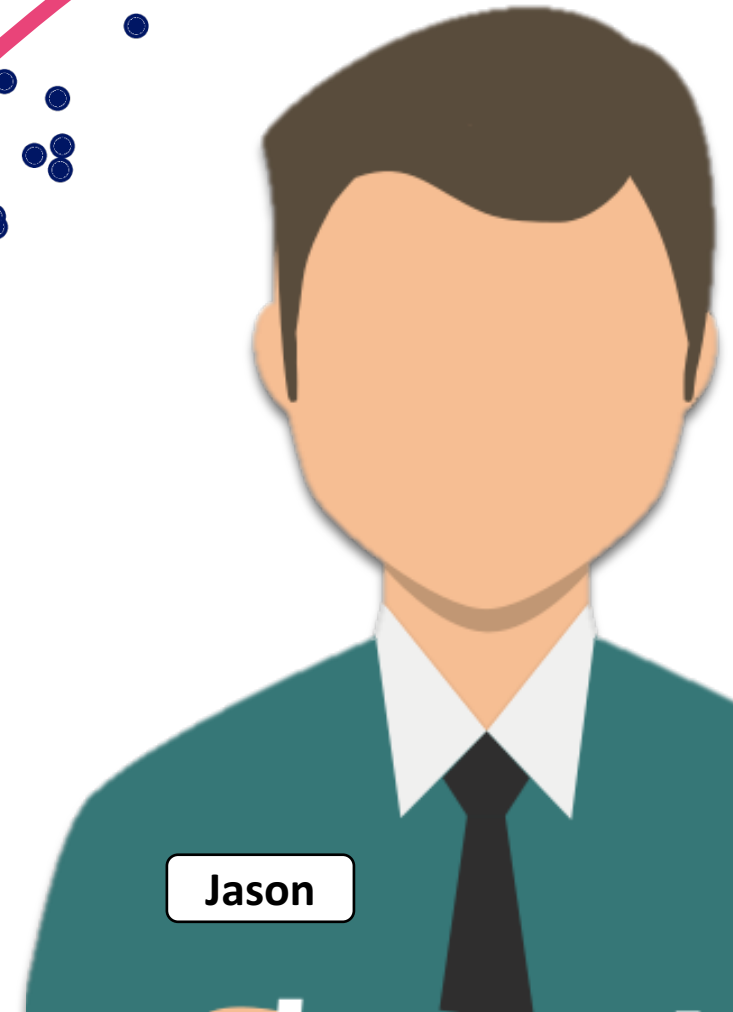
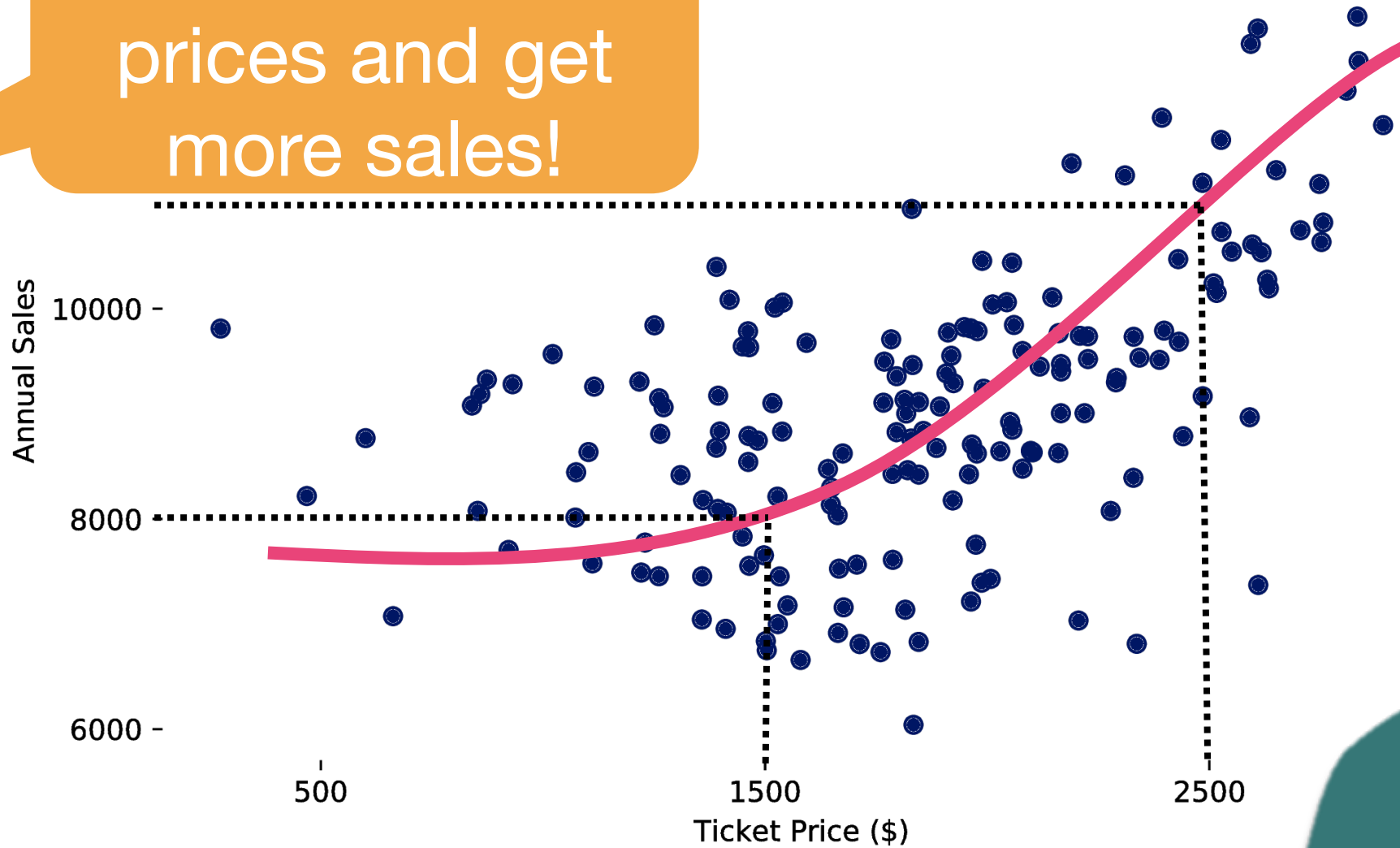


*SkyHighAir*

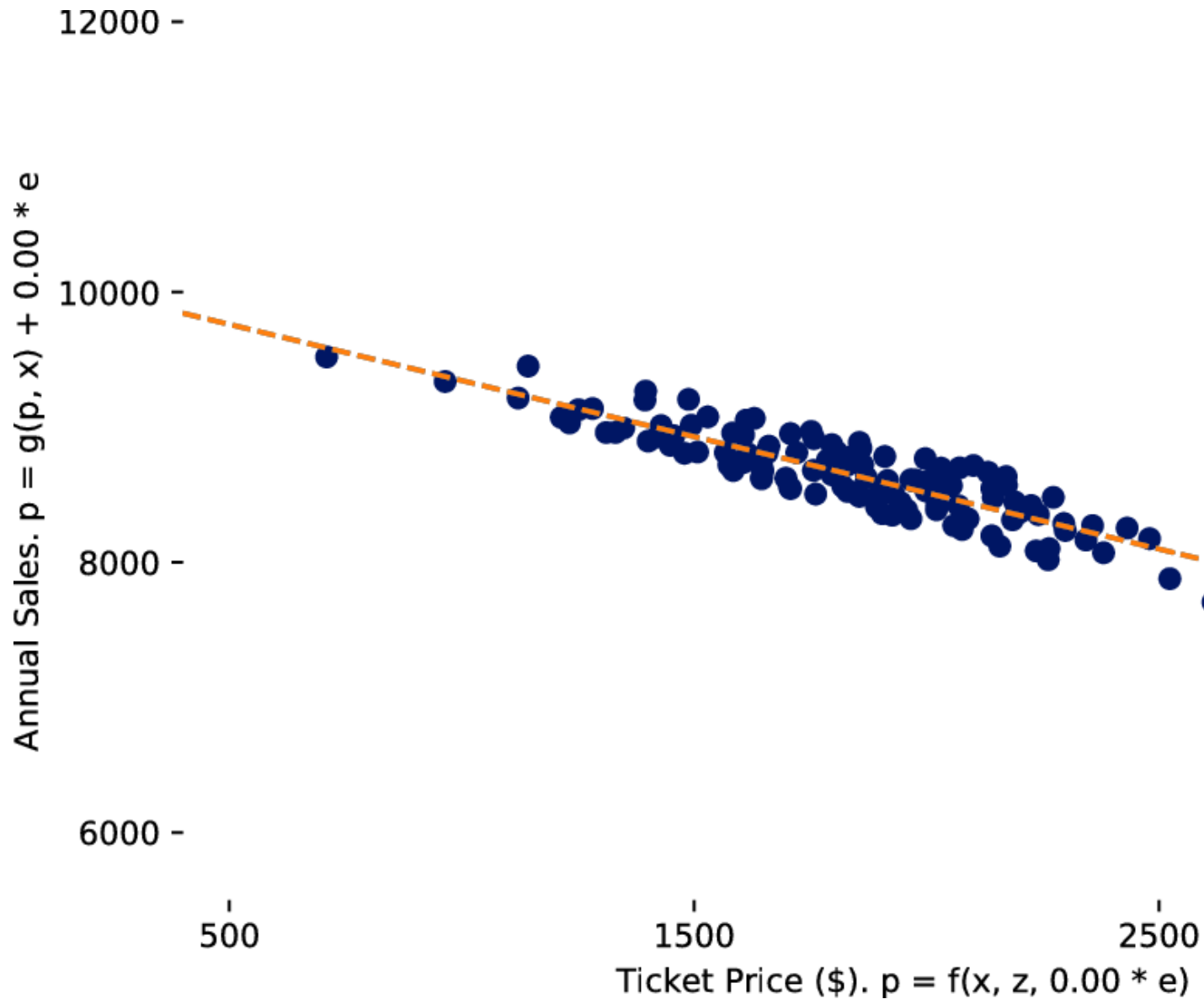


# Prediction with confounding effects

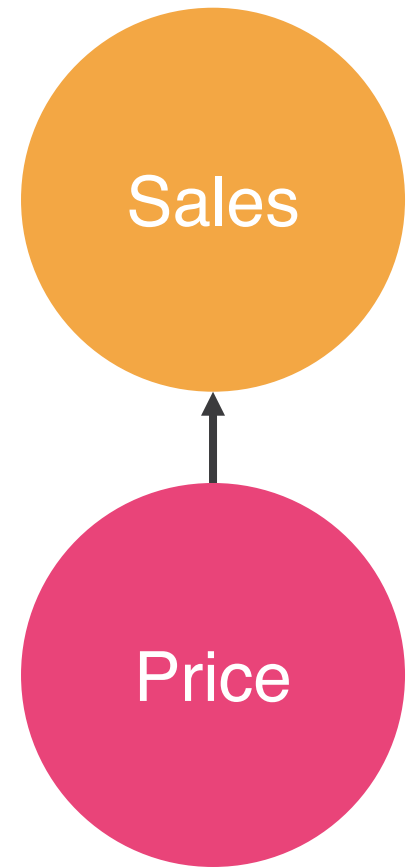
We can raise prices and get more sales!



# Prediction with confounding effects

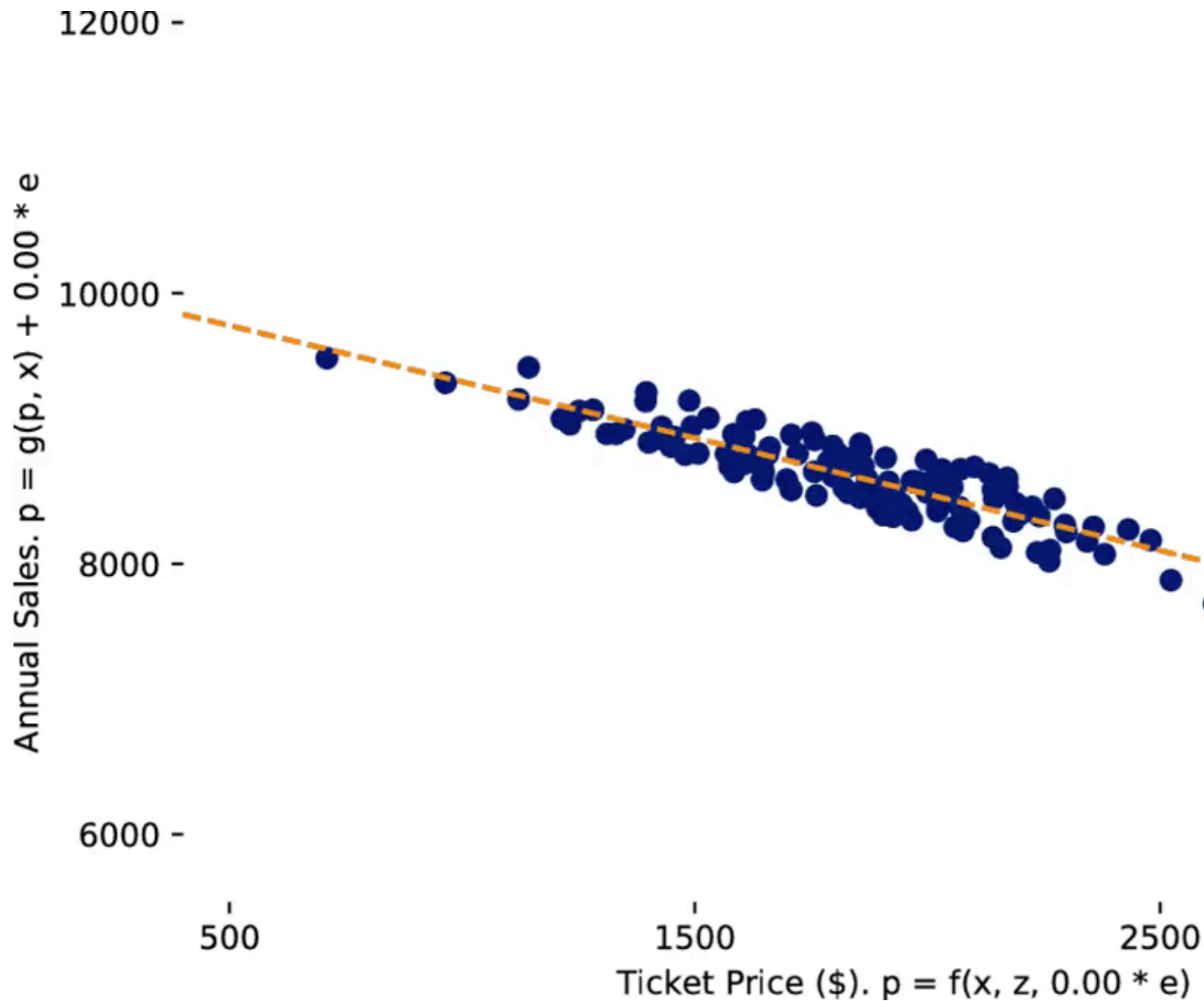


$$y = g(p)$$

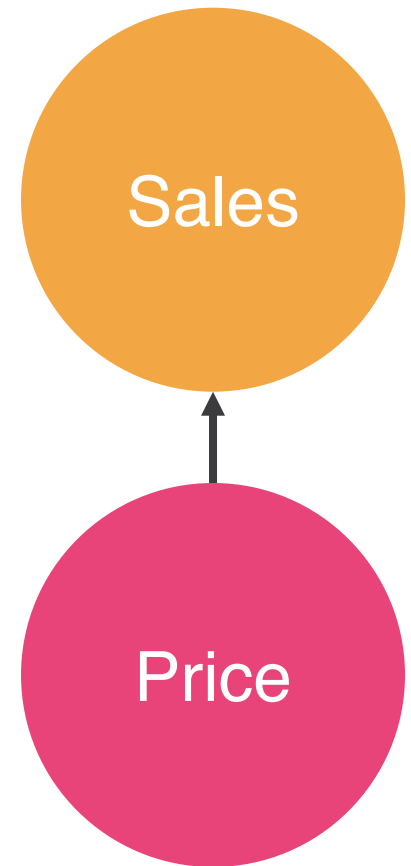


$p$

# Prediction with confounding effects

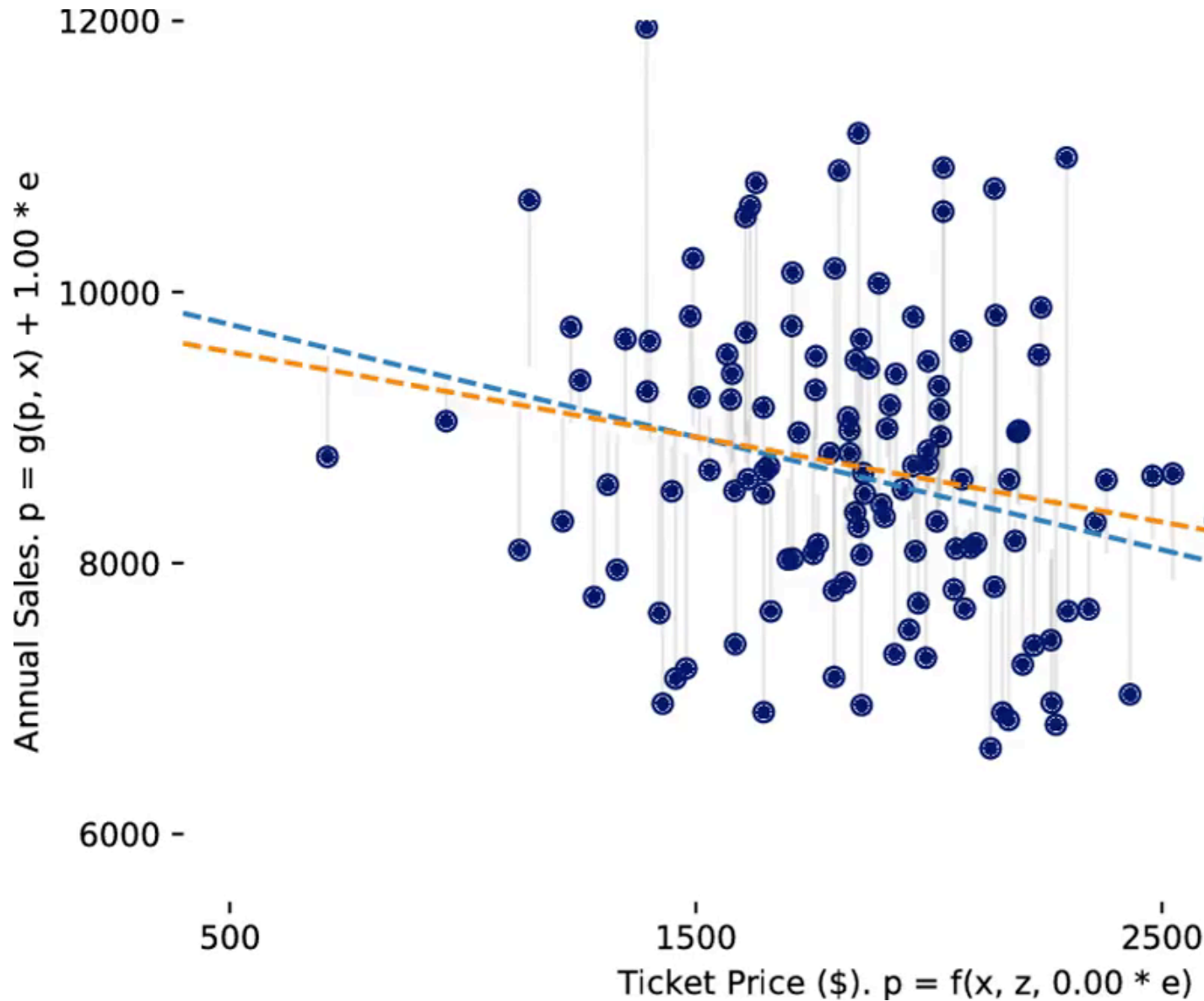


$$y = g(p, e)$$

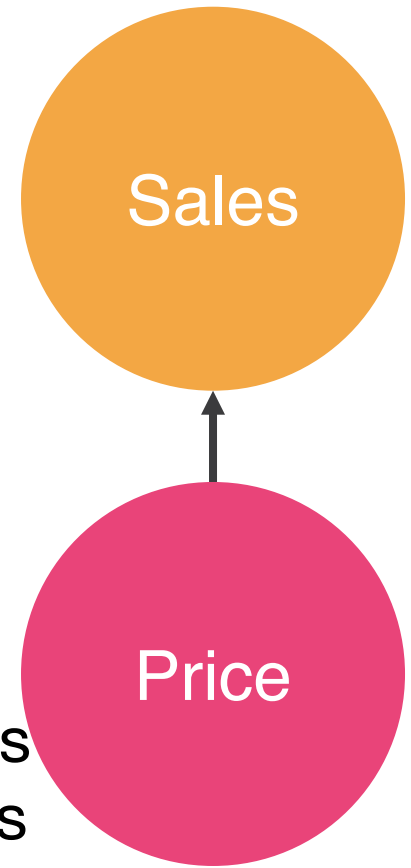


$p$

# Prediction with confounding effects



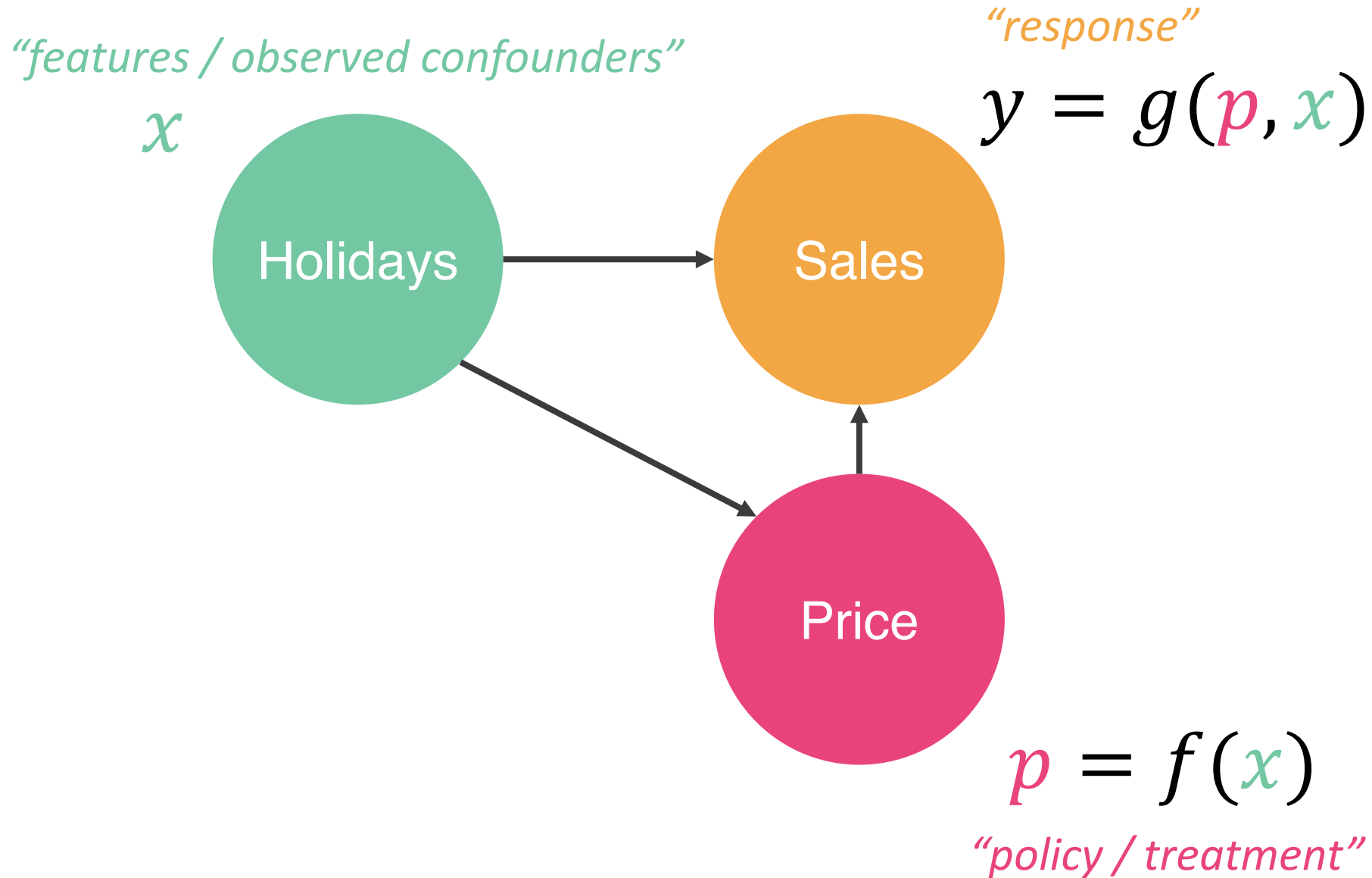
$$y = g(p, e)$$



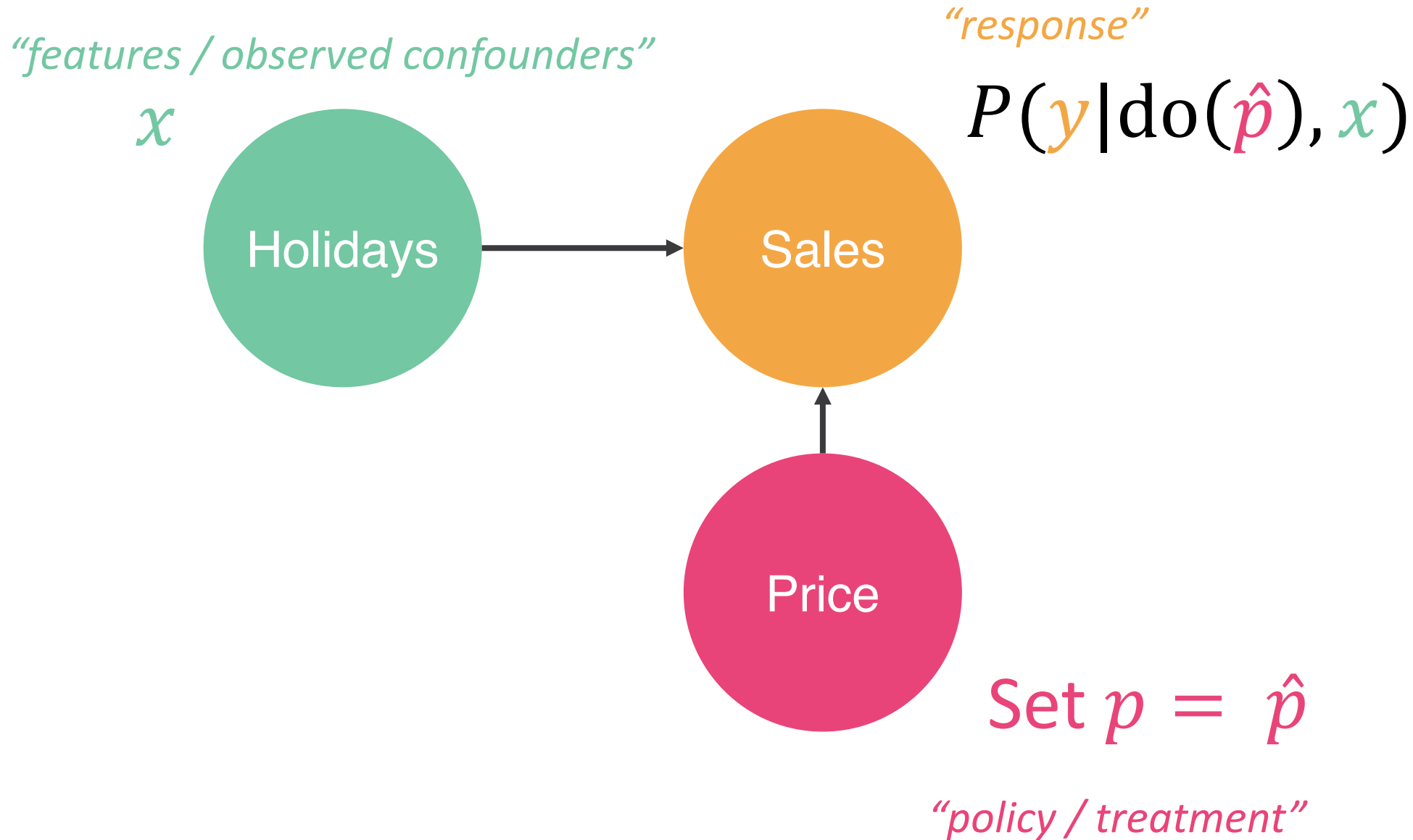
Automated pricing engine increases prices as the plane fills

$$p = f(e)$$

# The **observational** distribution

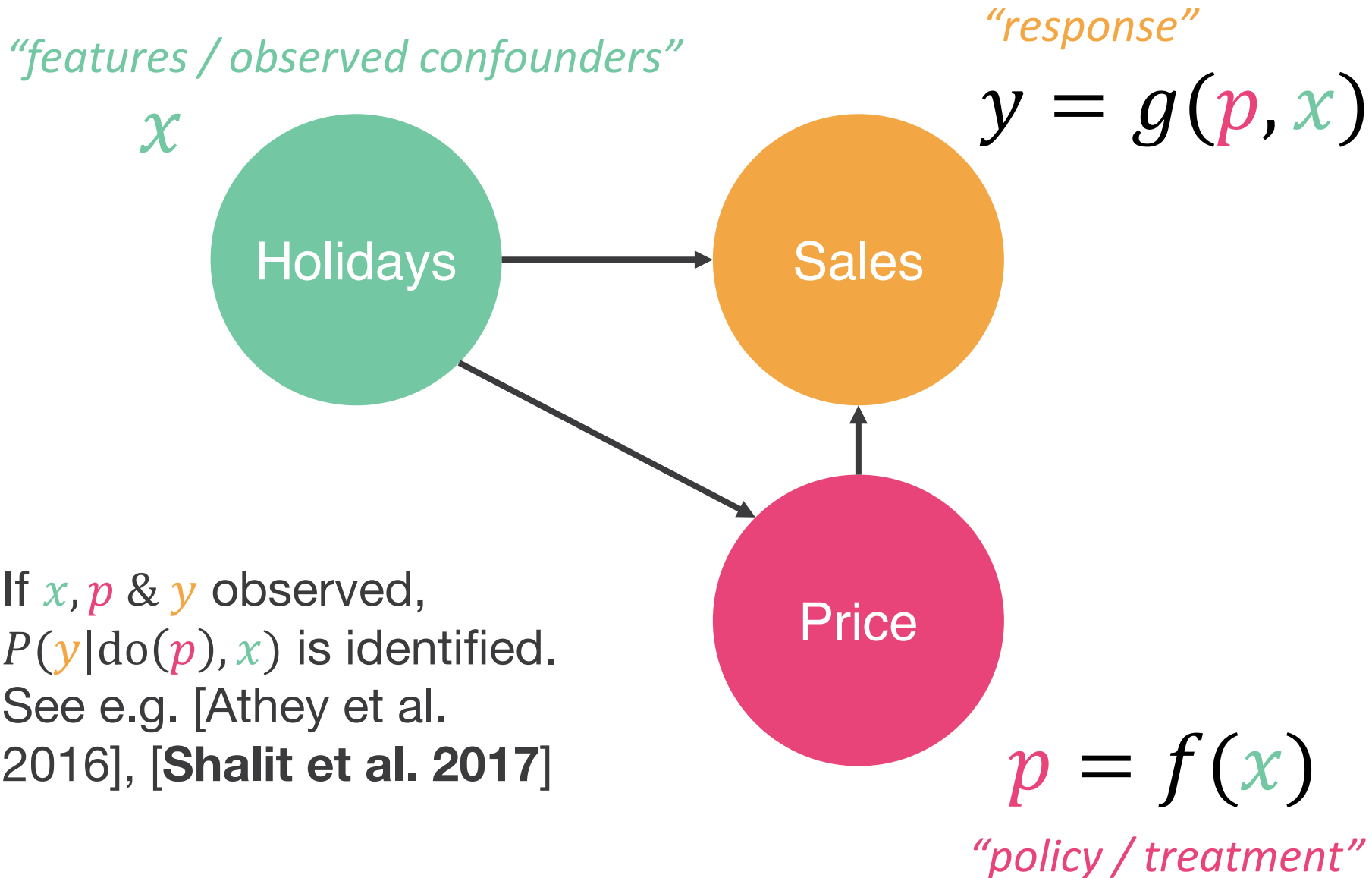


# The *interventional* distribution

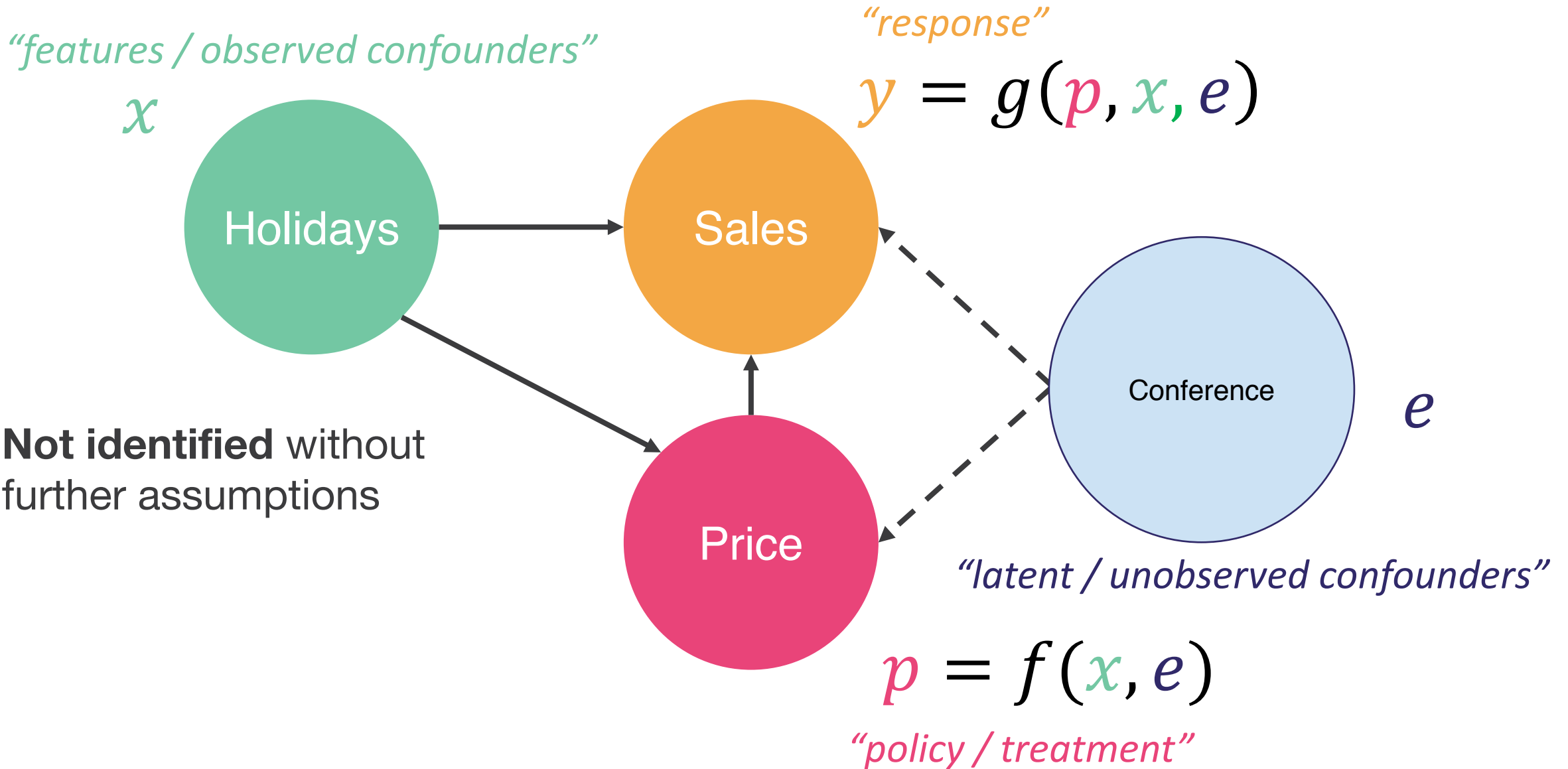




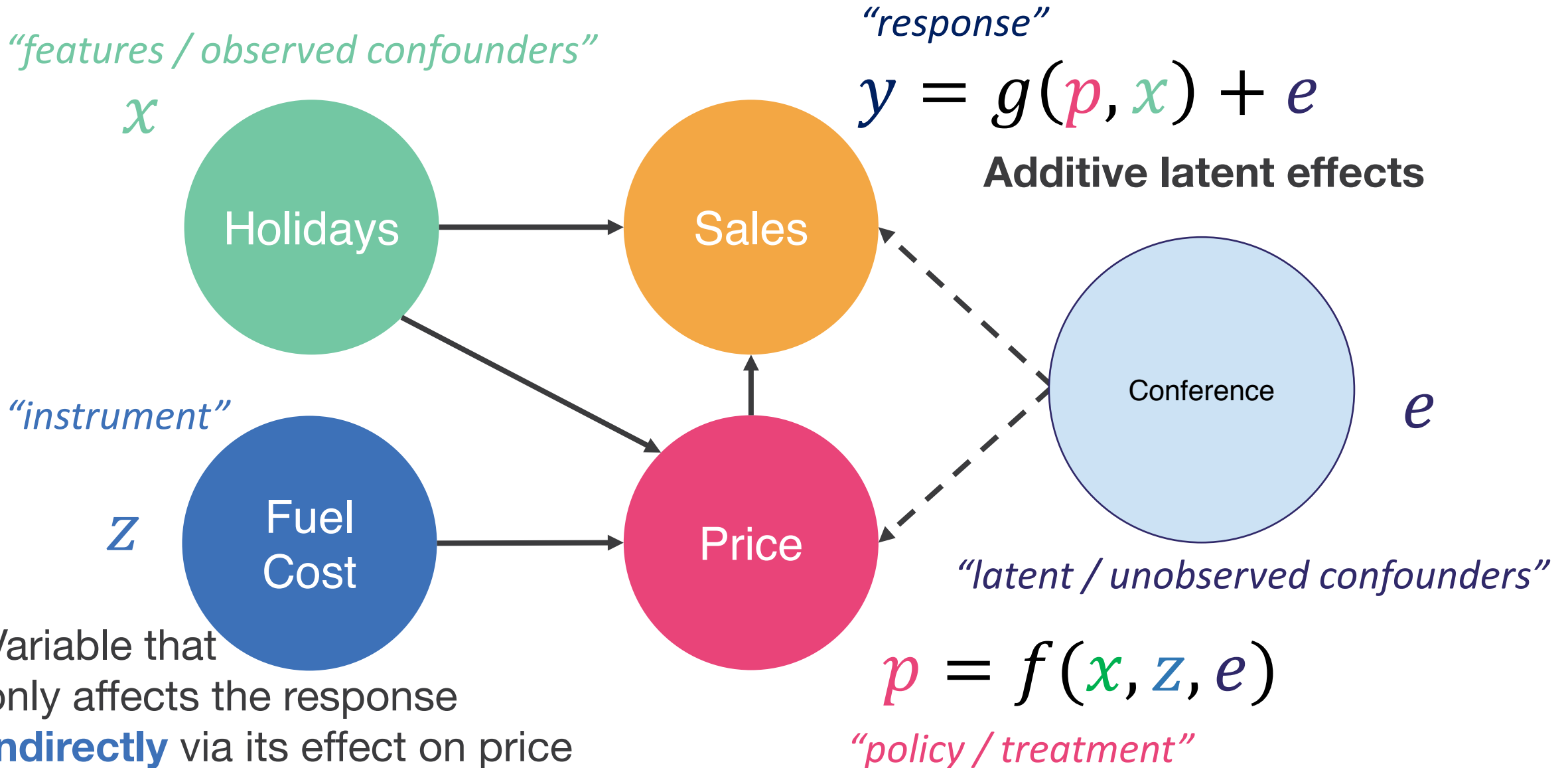
# Identification of causal effects



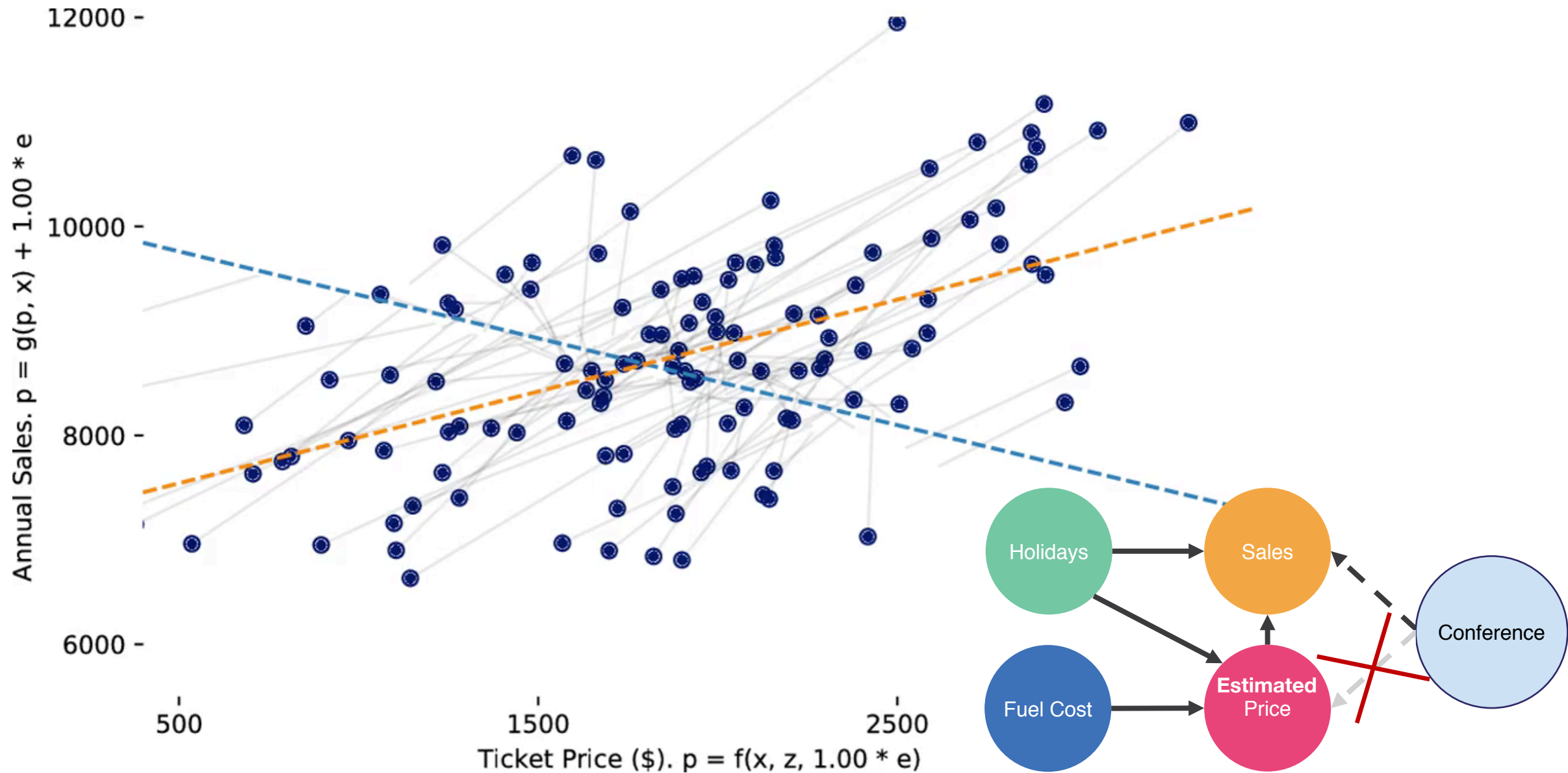
# Identification of causal effects



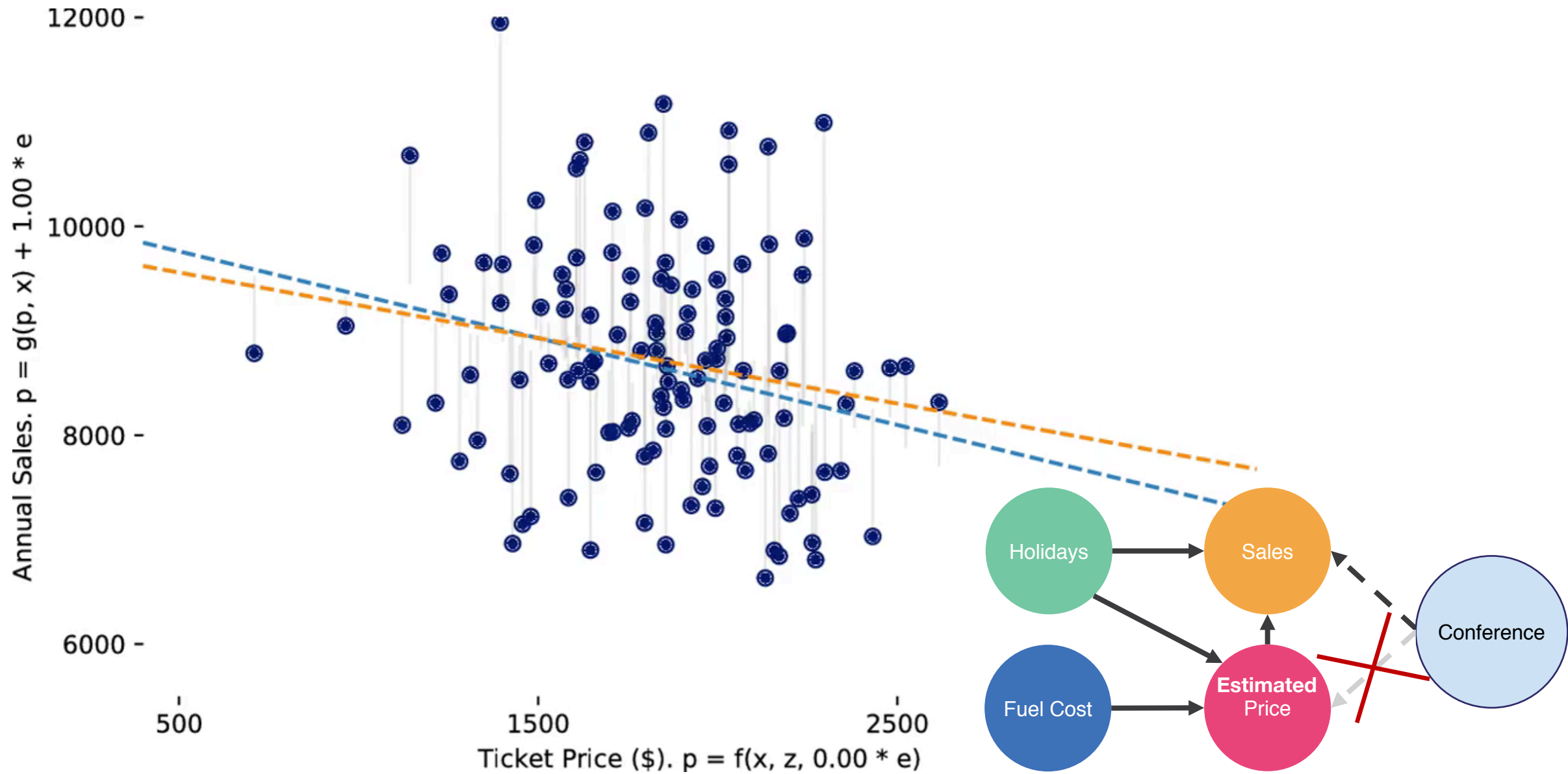
# Identification of causal effects



# Simulate a world without latent effects on price



# Simulate a world without latent effects on price



# The learning problem

These assumptions imply the following **identity**<sup>1</sup>,

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

So we can recover  $g(p, x)$  solve the implied **inversion problem**...

$$\min_{g \in G} \sum_{t=1}^n \left( y_t - \int g(p, x_t) dF(p|x, z) \right)^2$$

1. This holds if  $E[e|x] = 0$ . In general we recover  $g(p, x)$  up to a constant wrt  $p$  – see paper for details.

# A two-stage solution

$$\min_{g \in G} \sum_{t=1}^n \left( y_t - \int g(p, x_t) dF(p|x, z) \right)^2$$

**Stage 1:** fit  $\widehat{F}_\phi(p|x, z)$  using the model of your choice.

**Stage 2:** train network  $\widehat{g}_\theta$  using **stochastic gradient descent** with **monte-carlo integration**.

We use **mixture density networks** [Bishop 94]


$$\widehat{F}_\phi(p|x, z)$$

At each SGD iteration

Sample

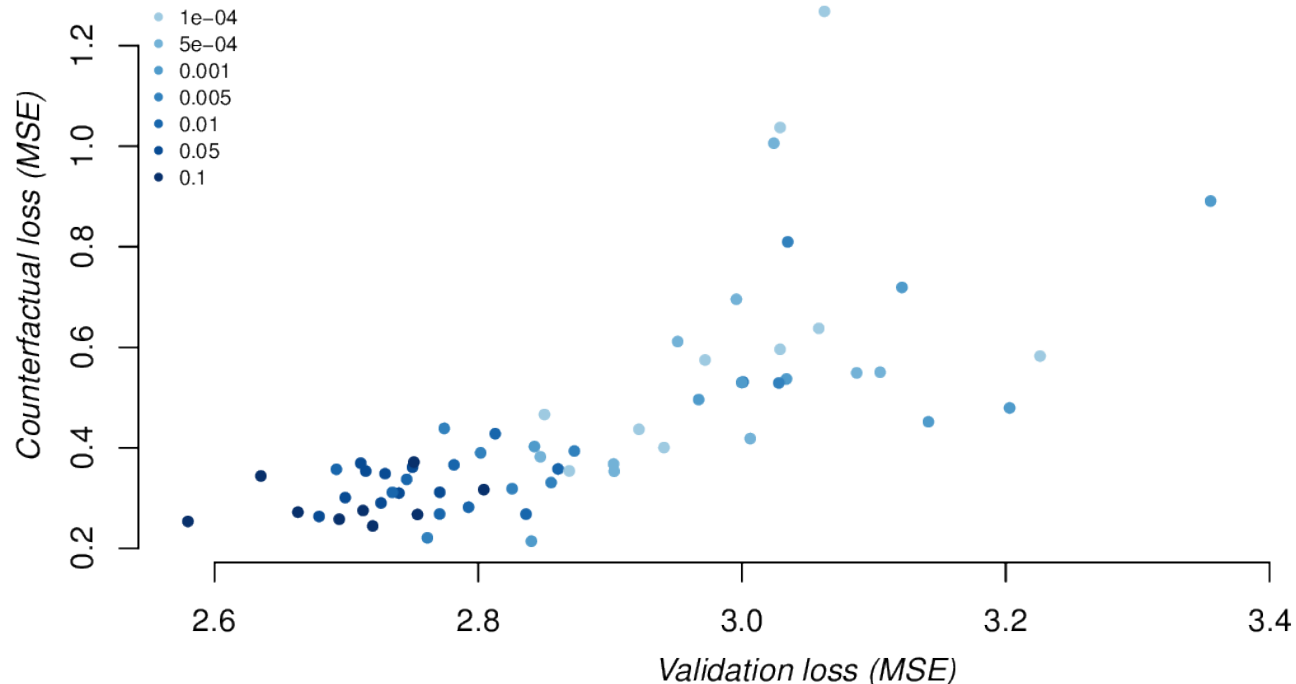
$$\nabla L(\theta) = -2 \left( y_t - \frac{1}{|\dot{p}_1|} \sum_{p_1 \sim \widehat{F}(p|x, z)} \widehat{g}(\dot{p}_1, x_t) \right) \times$$

$$\left( \frac{1}{|\dot{p}_2|} \sum_{p_2 \sim \widehat{F}(p|x, z)} \nabla_\theta \widehat{g}(\dot{p}_2, x_t) \right)$$

$$\widehat{g}_\theta(p, x)$$

# Causal Validation

- In general, out-of-sample validation causal models is **challenging / impossible**...
- But... both our losses depend only on **observable** quantities *and* reflect causal loss, so we can simply use **standard validation sets**.



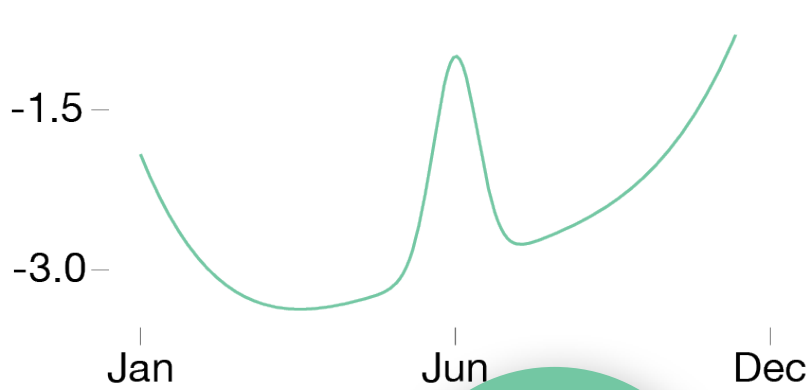


# Evaluation

Simulation & Bing Ads Experiments

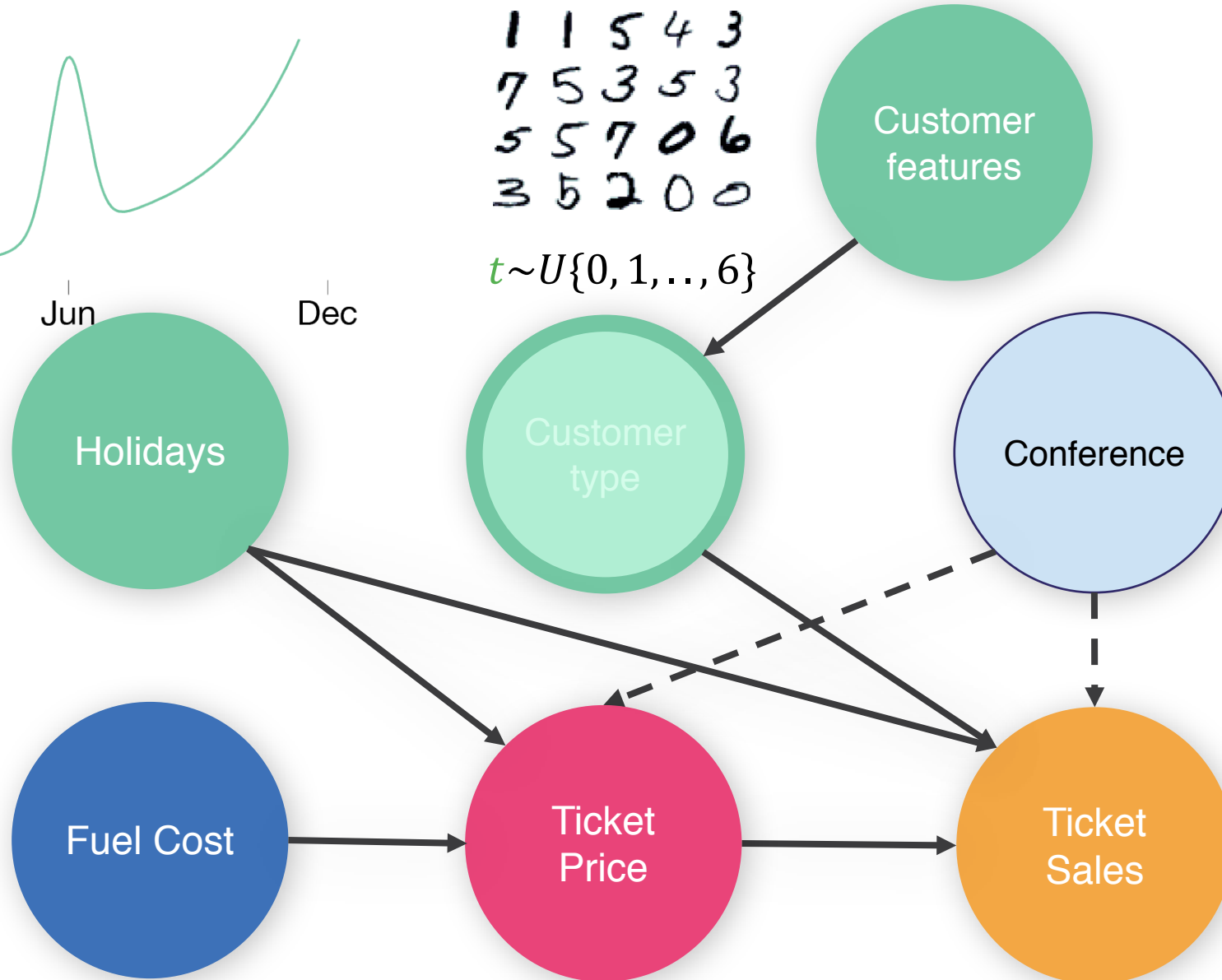
# Simulation Experiments

Price Sensitivity



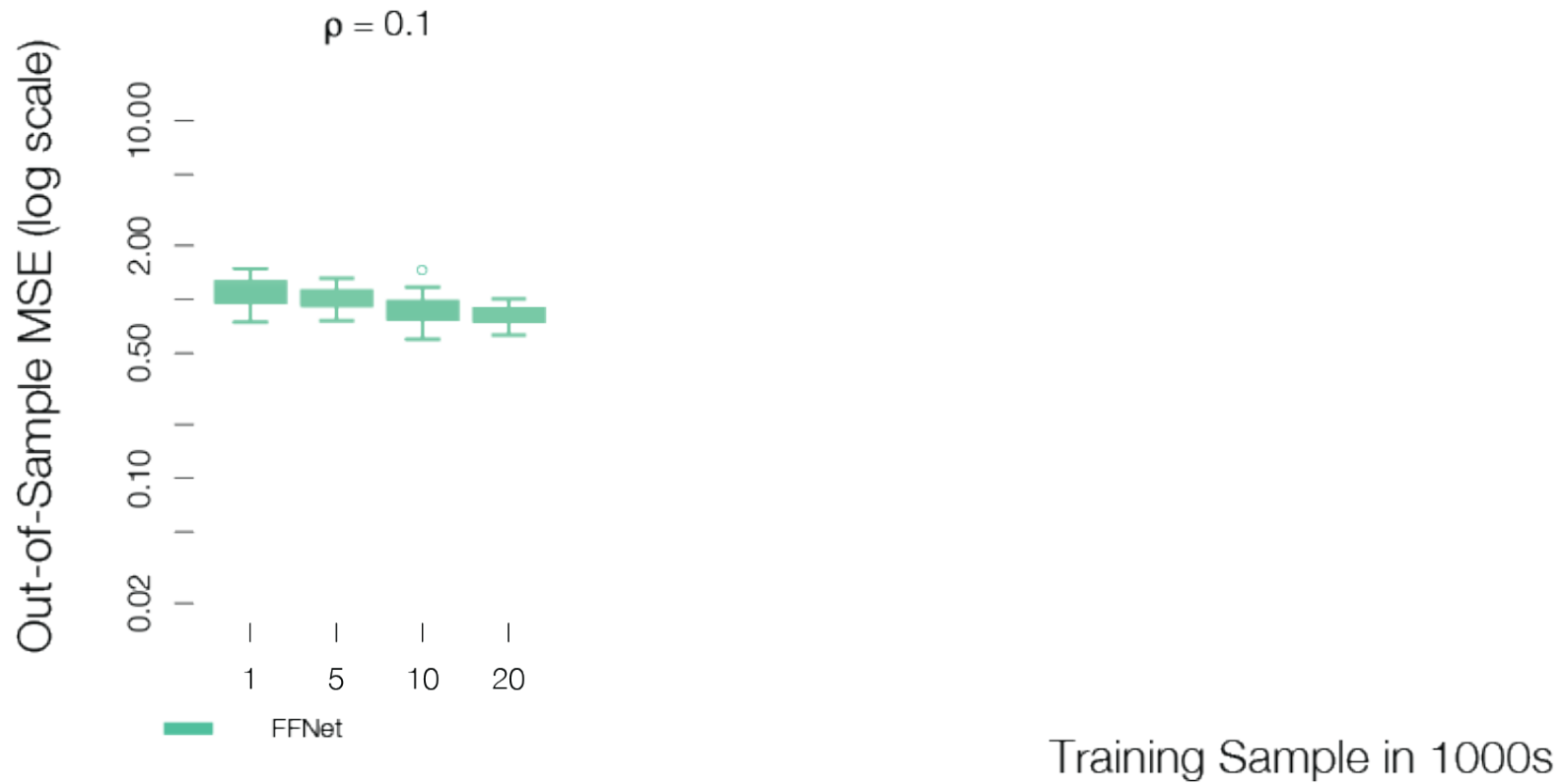
1	1	5	4	3
7	5	3	5	3
5	5	7	0	6
3	5	2	0	0

$t \sim U\{0, 1, \dots, 6\}$

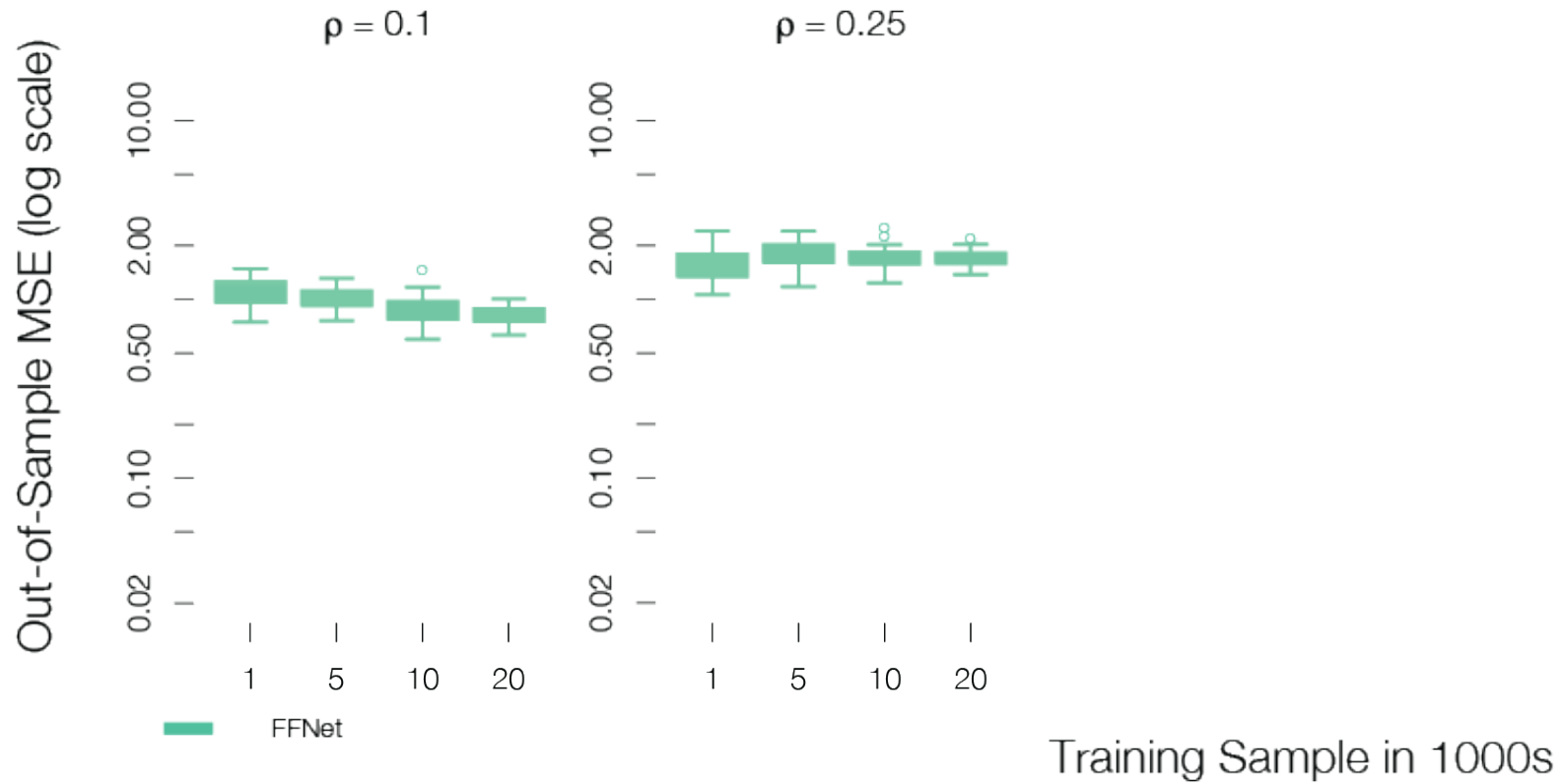


$\rho$  lets us smoothly vary the correlation between sales and price

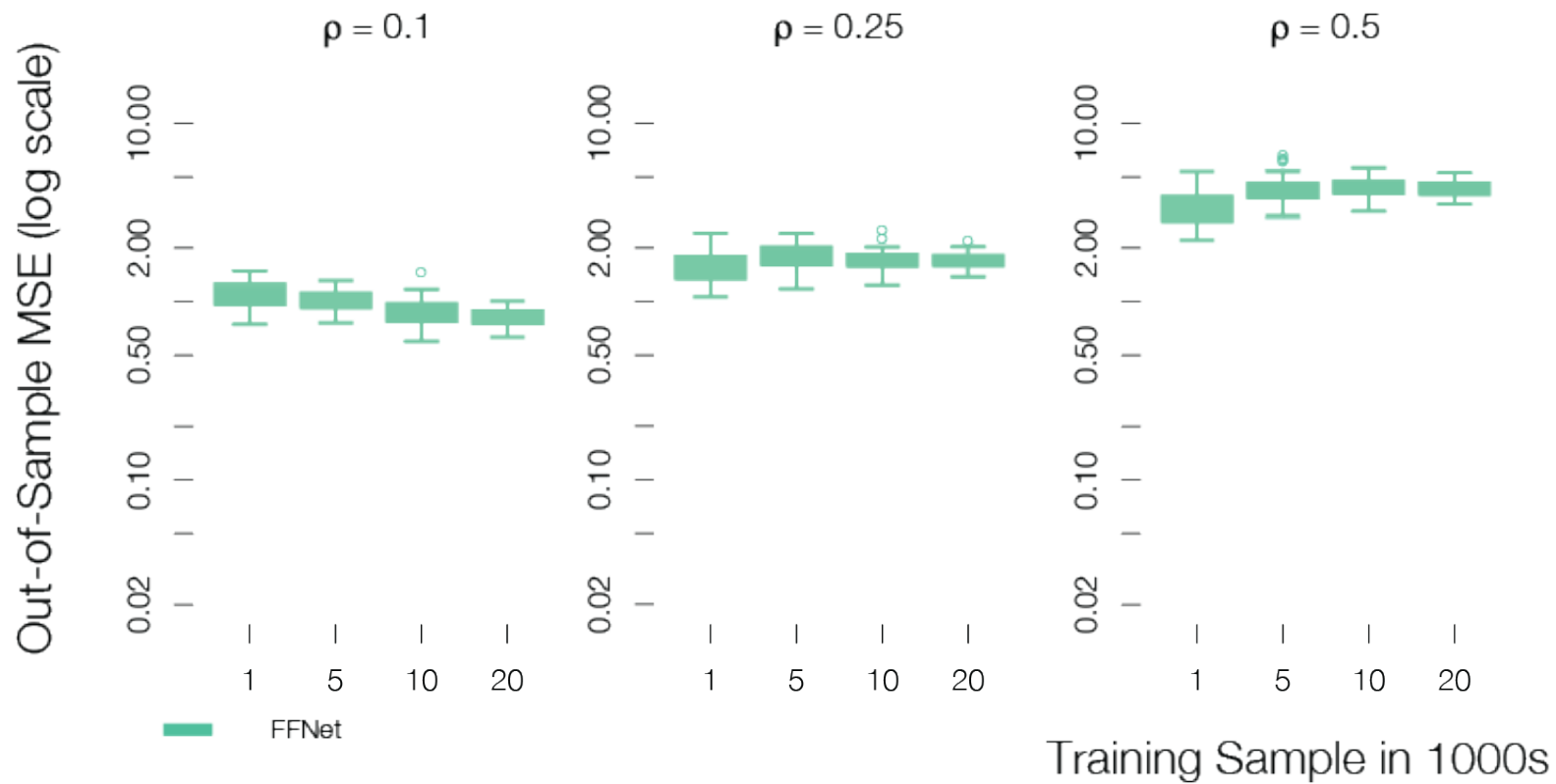
# Simulation – low dimensional feature space



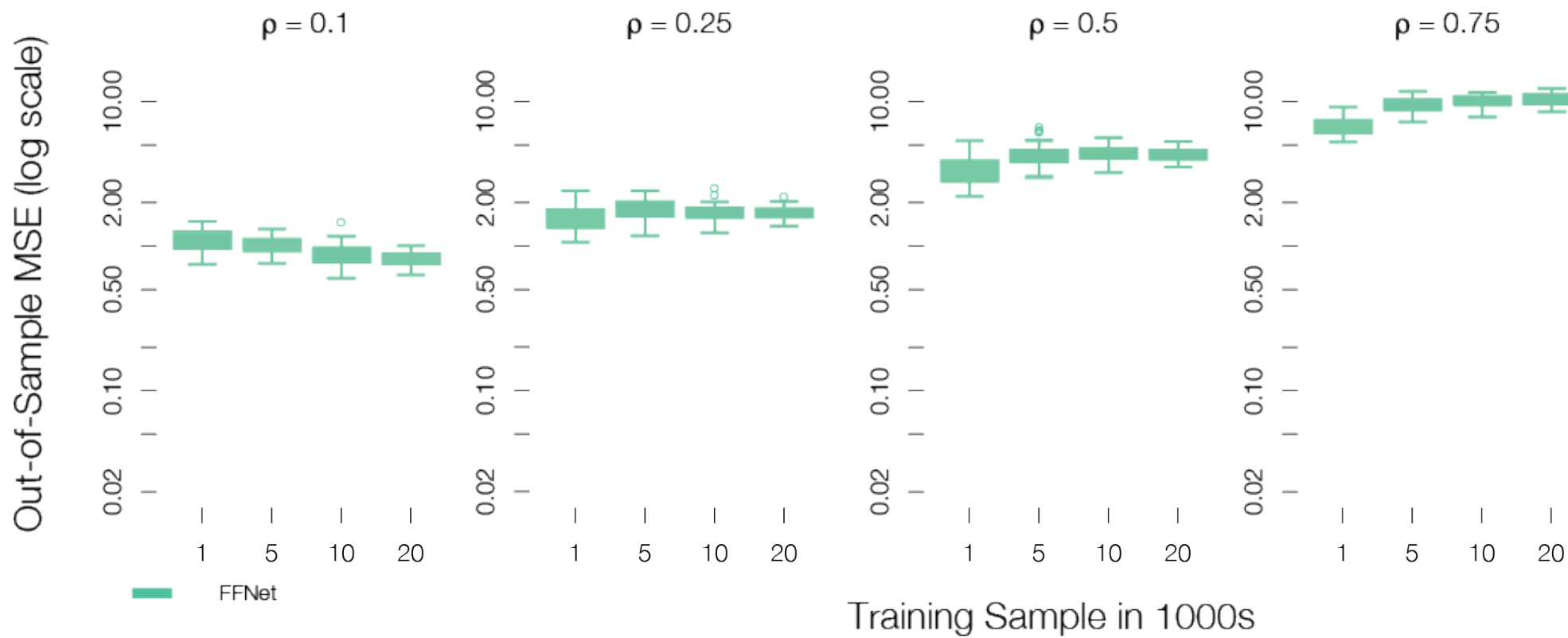
# Simulation – low dimensional feature space



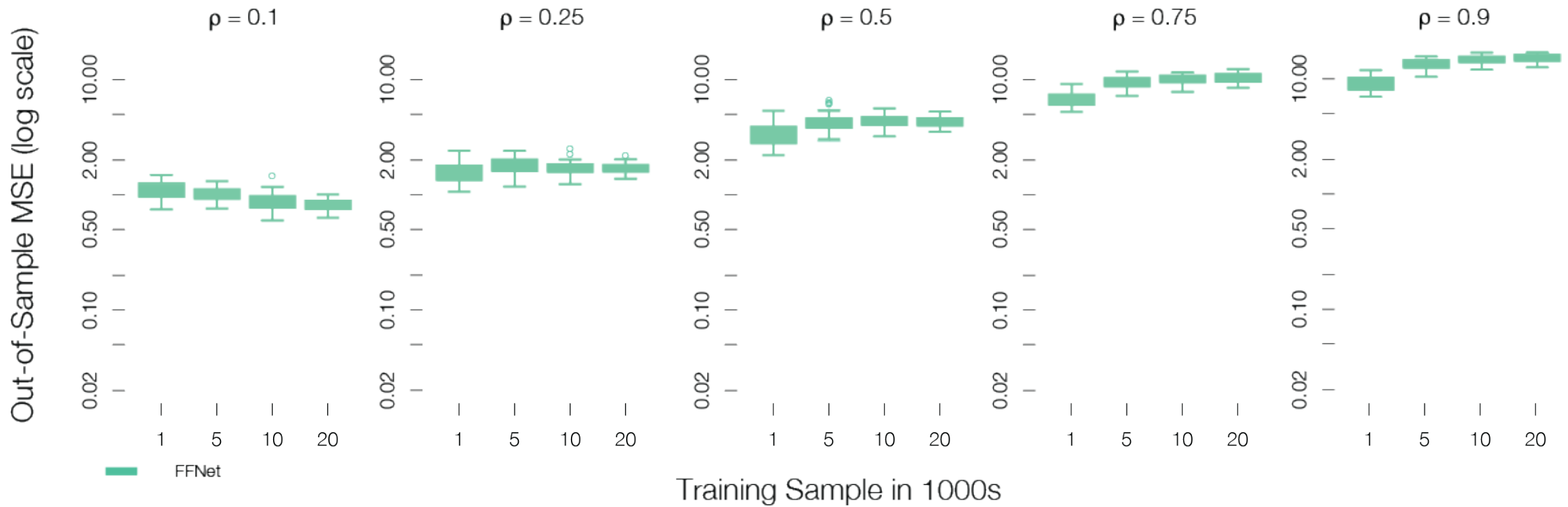
# Simulation – low dimensional feature space



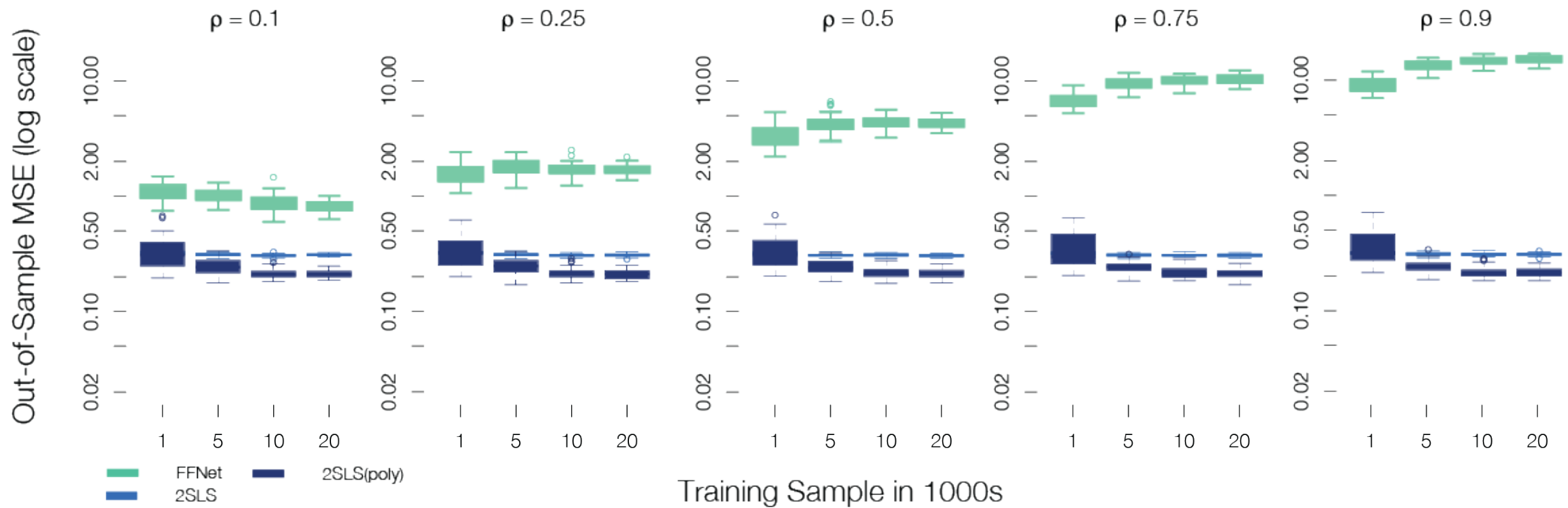
# Simulation – low dimensional feature space



# Simulation – low dimensional feature space

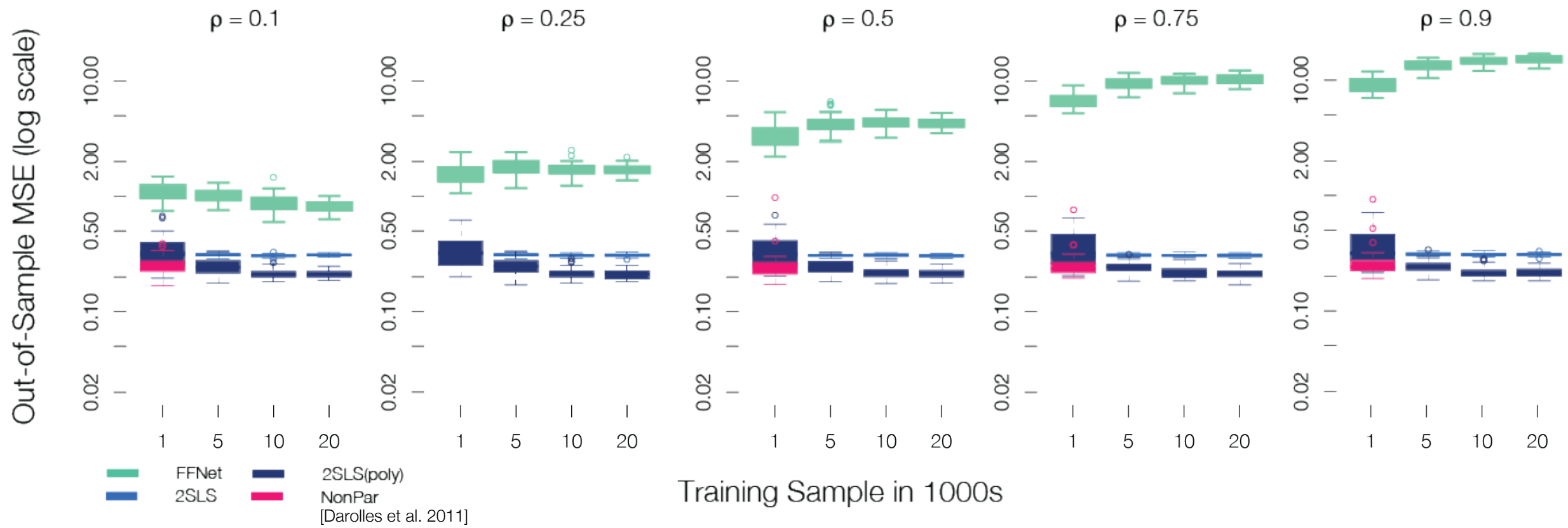


# Simulation – low dimensional feature space

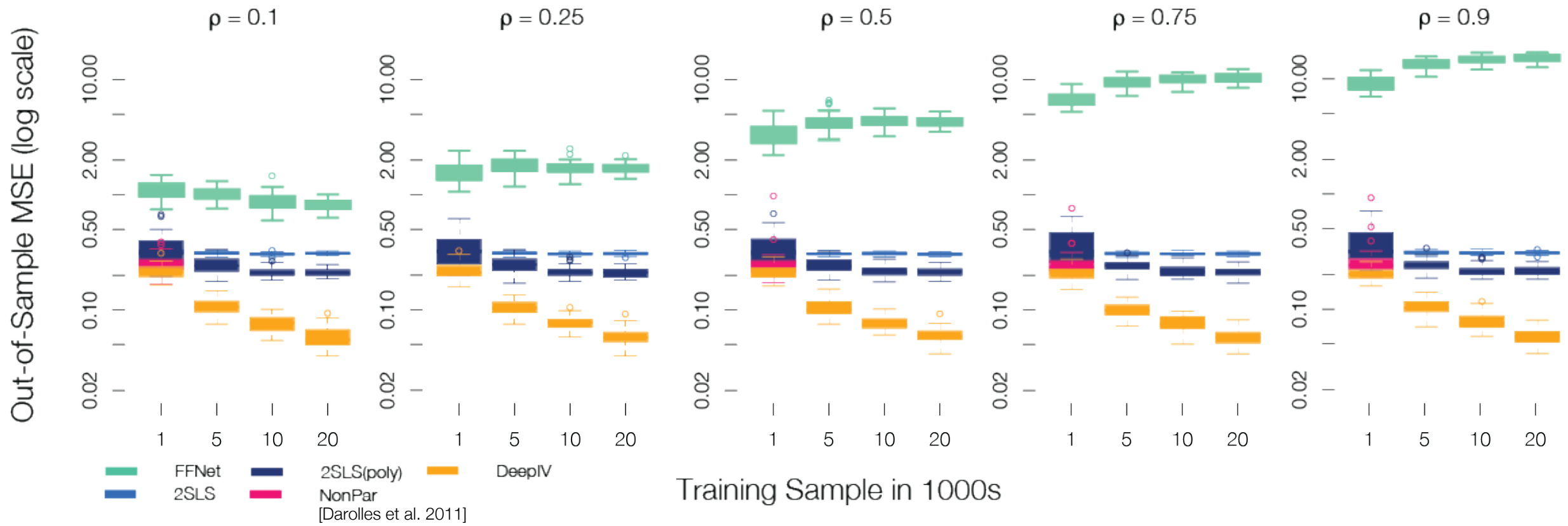




# Simulation – low dimensional feature space



# Simulation – low dimensional feature space



# Implications and future directions

- We recover **heterogeneous treatment effects** in settings with **unobserved confounding effects** for both **discrete** and **continuous** variables... and SGD **scales** naturally to **very large datasets**.
- Can leverage the flexibility of deep nets for rich data types. E.g. **raw text** in our Bing ads application experiments / **images** in simulation.

Future work:

- Methods for **uncertainty** estimates over predictions.