

# CAUSAL INFERENCE AT HULU

---

Allen Tran

July 17, 2016

Hulu

# INTRODUCTION

Most interesting business questions at Hulu are **causal**

Business: what would happen if we did x instead of y?

- dropped prices for risky subs
- halved our AdWords marketing spend
- bought some piece of Content

# CAUSAL INFERENCE IS NOT SUPERVISED LEARNING

Supervised learning model ( $f : X \rightarrow Y$ )

- Optimized for:  $f(x) \approx y$
- test-train paradigm works
- regularization works

Causal Inference model ( $f : X, T \rightarrow Y$ )

- Optimized to identify treatment effect:  $f(x, T = 1) - f(x, T = 0)$
- there is no test set, need to lean on statistical theory
- naive regularization is a horrible idea

Sacrifice predictive performance to identify treatment effect

1. Potential outcomes framework/treatment effects
2. What we do at Hulu
3. Modern ML approaches to causality (we should do this)

**Goal:** Convince us to think less supervised learning, more causal inference

Note: focus will be on observation studies, not designing randomized trials

# CAUSAL INFERENCE = IDENTIFYING TREATMENT EFFECTS

Let  $x_i$  describe some features (a user), and two potential outcomes,  $Y_1(x_i)$  and  $Y_0(x_i)$ .

$Y_k(x|T = k)$  is observed,  $Y_k(x|T = 1 - k)$  is the unobserved counterfactual.

- individual treatment effect

$$ITE(x_i) = \mathbb{E}_{Y_1 \sim p(Y_1|x_i)} (Y_1(x_i)) - \mathbb{E}_{Y_0 \sim p(Y_0|x_i)} (Y_0(x_i)) \quad (1)$$

- heterogeneous treatment effect

$$HTE(X) = \mathbb{E}_{Y_1 \sim p(Y_1|x \in X)} (Y_1(x)) - \mathbb{E}_{Y_0 \sim p(Y_0|x \in X)} (Y_0(x)) \quad (2)$$

- average treatment effect

$$ATE = \mathbb{E} (Y_1(x) - Y_0(x)) = \mathbb{E}_{x \sim p(x)} (ITE(x)) \quad (3)$$

# ASSUMPTIONS

**Ignorability:** assignment to treatment/control effectively random, conditioning on  $x$

$$(Y_0, Y_1) \perp\!\!\!\perp T \mid x \quad (4)$$

or equivalently,

$$p(Y_0, Y_1, T|x) = p(Y_0, Y_1|x)p(T|x) \quad (5)$$

**Overlap/Common support:** treatment/control across all units

$$p(t|x) > 0 \quad \forall t, x \quad (6)$$

# ESTIMATING AVERAGE TREATMENT EFFECTS

First part of ATE

$$\mathbb{E}_{x, Y_1, T}(Y_1(x)) = \mathbb{E}_x(\mathbb{E}_{Y_1, T|x} Y_1(x, T)) \quad (7)$$

$$= \mathbb{E}_x(\mathbb{E}_{Y_1|x} Y_1(x, T = 1)) \quad (8)$$

Term inside parentheses is observed.

$$ATE = \mathbb{E}_x(\mathbb{E}_{Y_1|x} Y_1(x, T = 1) - \mathbb{E}_{Y_0|x} Y_0(x, T = 0)) \quad (9)$$

But  $x_i$  units we observe are distributed  $p(x_i|T = 1)$  and  $p(x_i|T = 0)$ , not  $p(x)$ .

Randomized trial,  $x \perp\!\!\!\perp T$ , difference in means is sufficient.

In order of simplicity

1. Matching
2. Propensity score/reweighting
3. Covariate adjustment (current strategy at Hulu - incremental retention, portfolio, sub churn)
4. Modern ML methods designed for causal inference



# COVARIATE ADJUSTMENT

Supervised learning with treatment variable as a feature

$$ATE = \frac{1}{N} \sum_i f(x_i, T = 1) - f(x_i, T = 0) \quad (10)$$

Problems

- No goodness of fit metrics. Remember, we are fitting treatment effects, not  $Y$
- Treatment effects wrong if model is misspecified (e.g unconfoundedness does not hold)
  - unobserved demographics behind watch/did not watch (treatment) in incremental retention model
  - specific types of shows (unobserved) get more marketing (treatment) in MMM model
- Regularization? May regularize away our treatment/potentially harmful

# CASUAL TREES/FORESTS

*Athey & Imbens (2015)*

Directly estimate treatment effects via decision trees. Shown to be consistent estimator of treatment effect (i.e is a “correct” model).

Regular decision tree

1. Estimator (of outcome): mean of **outcome** in leaf

$$\hat{y}_l = \frac{1}{N_l} \sum_{i \in l} Y_i \quad (11)$$

2. Splitting criterion: metrics (information gain/Gini/variance reduction) for **outcome** + complexity penalty
3. Score: assess estimates of predicted **outcome** on test set,  $\hat{y}_i$  against  $y_i$

# CASUAL TREES/FORESTS: TRANSFORMING OUTCOMES

Transform outcome  $Y_i$ . Let  $p_i = p(T = 1|x_i)$

$$Y_i^* = \frac{T_i - p_i}{1 - p_i} Y_i \quad (12)$$

Expectation of transformed outcome is average treatment effect!

$$\mathbb{E}_{x,T}(Y_i^*) = \mathbb{E}_x(\mathbb{E}_T(Y_i^*)) \quad (13)$$

$$= \mathbb{E}_x(p_i \mathbb{E}(Y_i^*|T = 1, x) + (1 - p_i) \mathbb{E}(Y_i^*|T = 0, x)) \quad (14)$$

$$= \mathbb{E}_x(p_i Y_1(x_i) - (1 - p_i) Y_0(x_i)) \quad (15)$$

$$= ATE \quad (16)$$

# CASUAL TREES/FORESTS: SPLITTING

Ideal but infeasible criterion, compare treatment effect to ground truth

$$Q(\hat{\tau}) = \mathbb{E} \left( (\tau_i - \hat{\tau}_i)^2 \right) \quad (17)$$

But decompose this:

$$Q(\hat{\tau}) = \mathbb{E} (\tau)^2 + \mathbb{E} (\hat{\tau}_i^2) - 2\mathbb{E} (\tau_i \hat{\tau}_i) \quad (18)$$

First term doesn't involve  $\tau_i$  so can ignore. Second term is sample mean of estimator. Third term can be calculated (after more algebra)

$$\begin{aligned} \mathbb{E} (\tau_i \hat{\tau}_i) &= \mathbb{E}_\ell (\mathbb{E} (\tau_i \hat{\tau}_i | x_i \in \ell)) \\ &= \mathbb{E}_\ell (\hat{\tau}_\ell \mathbb{E} (\tau_i | x_i \in \ell)) \\ &= \mathbb{E} (\hat{\tau}_i^2) \end{aligned}$$

Hence, criterion to continue splitting is

$$Q(\hat{\tau}) = -\mathbb{E}(\hat{\tau}_i^2) \quad (19)$$

Rewards variance in estimates of treatment effects. (Compare to limiting case of no splits, every  $x_i$  has the same treatment effect).

# CAUSAL TREE/FORESTS: IMPLEMENTATION

1. Estimator (of **treatment effect**):  $\hat{\tau}_i$

Mean of transformed outcome (with propensity score weighting as treatment/control distributions in leaf are not representative).

$$\hat{\tau}_i = \frac{1}{N_l} \sum_{i \in \ell} Y_i^* \frac{1}{p(T=1|X_i)} \quad (20)$$

2. Splitting criterion: MSE on **treatment effect** + complexity of number of leaves

$$\frac{1}{N} \sum_i \hat{\tau}_i^2 + \lambda \cdot n_{\text{leaves}} \quad (21)$$

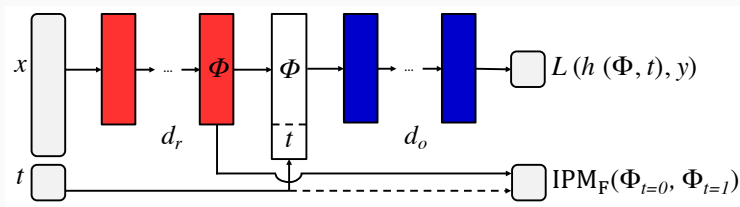
3. Score: out of sample MSE on **treatment effect**

$$\frac{1}{N_{\text{test}}} \sum_{i \in \text{test}} (\hat{\tau}_i - Y_i^*)^2 \quad (22)$$

# DEEP LEARNING/REGULARIZING REPRESENTATIONS

*Shalit, Johansson & Sontag (2016)*

- Learn representations for  $x$ ,  $\Phi(x)$
- Merge representations and treatment to output prediction
- Measure difference in distribution of representations for treatment and control



# DEEP LEARNING/REGULARIZING REPRESENTATIONS

Loss function

$$\min_{\Phi, h} \frac{1}{N} \sum_{i=1} L(h(\Phi(x_i), T_i), Y_i) + \alpha \cdot IPM_F(\{\Phi(x_i)\}_{i|T_i=0}, \{\Phi(x_i)\}_{i|T_i=1}) \quad (23)$$

- $IPM_F$  is some distance metric over distributions
- Penalize differences between treatment and control group
  - mitigates problem that treatment is not randomly assigned (incremental retention: different demographics watch different shows)
  - regularization in the right place
- Theoretical result: if  $IPM_F$  is Wasserstein or Max-Mean Discrepancy, loss function is bound on true counterfactual error



What if we know that we have some unobserved variable that affects outcome?

Strategies

1. Causal graphs (Judea Pearl et al.) - assess whether it is a problem and potential fixes
2. Instrumental variables - find “instruments” that induce random-like properties into our  $x$  variables

*Bishop, Chapter 8: Pattern Recognition and Machine Learning*

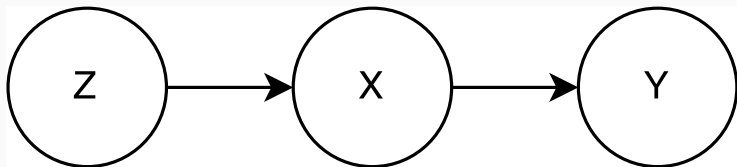
Represent casual effects via directed graphs.

Theoretical results (“d-separation”): infers conditional independence from directed graphs. We want

$$Y \perp\!\!\!\perp Z \mid X \tag{24}$$

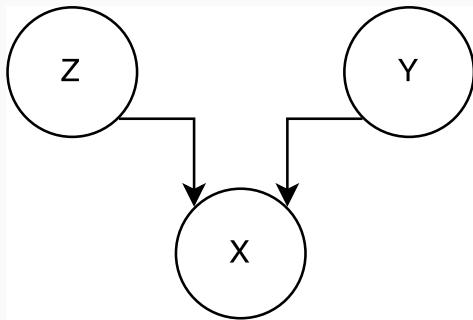
where  $Y$  is our outcome,  $Z$  is unobserved and  $X$  is observed

## CAUSAL GRAPHS: HEAD-TO-TAIL



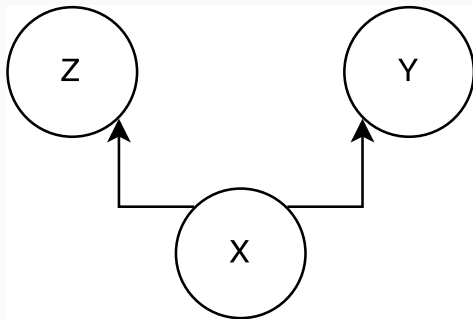
Observing  $X$  “blocks” the path from  $Z$  to  $Y$ , making  $Z$  independent of  $Y$ , conditional on  $X$ .

## CAUSAL GRAPHS: HEAD-TO-HEAD



Observing  $X$  “blocks” the path from  $Z$  to  $Y$ , making  $Z$  independent of  $Y$ , conditional on  $X$ .

## CAUSAL GRAPHS: TAIL-TO-TAIL

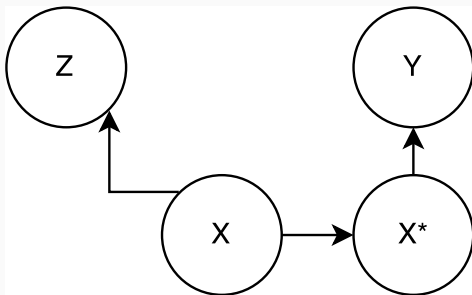


Observing  $X$  “**unblocks**” the path from  $Z$  to  $Y$ , making  $Z$  dependent of  $Y$ , conditional on  $X$ .

We can **not** use  $X$  as a dependent variable

# CAUSAL GRAPHS: FIXING TAIL-TO-TAIL

We can fix the prior situation, by finding something else to condition on.



Conditioning on  $X^*$  “**blocks**” the path from Z to Y, making Z dependent of Y, conditional on X.

# INSTRUMENTAL VARIABLES

Often (nearly all observational data) our  $X$  variable is correlated with unobserved variables making it impossible to identify the effect of  $X$  alone.

Idea: extract the random component on  $X$  uncorrelated with those other variables.

# INSTRUMENTAL VARIABLES: IMPLEMENTATION

Suppose you have an instrument  $Z$  that is correlated with  $X$  but not  $Y$ .

Two Stage Least Squares

1. OLS of  $X$  on  $Z$  gets you  $X_{IV}$
2. Second stage OLS of  $X_{IV}$  on  $Y$

But linear regression is not very powerful, how to improve?

- What if I want to use a non-linear estimator (deep learning)?



# INSTRUMENTAL VARIABLES: EXAMPLES

Note: there is a whole branch of economics dedicated to finding super fun/crazy instruments

- Effect of education on wages: education correlated with unobserved ability. Use exogenous districting/local government policies that force more/less education on people
- Watching show A on retention: but watch behavior correlated with unobserved demographics/preferences. Induce (some) randomness in reco strategy (change probability of exposure) and use exposure probabilities as instruments.

# CONCLUSION

A lot of our questions are causal

- We try to answer these causal questions in a supervised learning heavy manner
- Supervised learning/covariate adjustment is not great at causal questions
- There are alternatives that make use of modern ML algorithms/big data and directly identify treatment effects