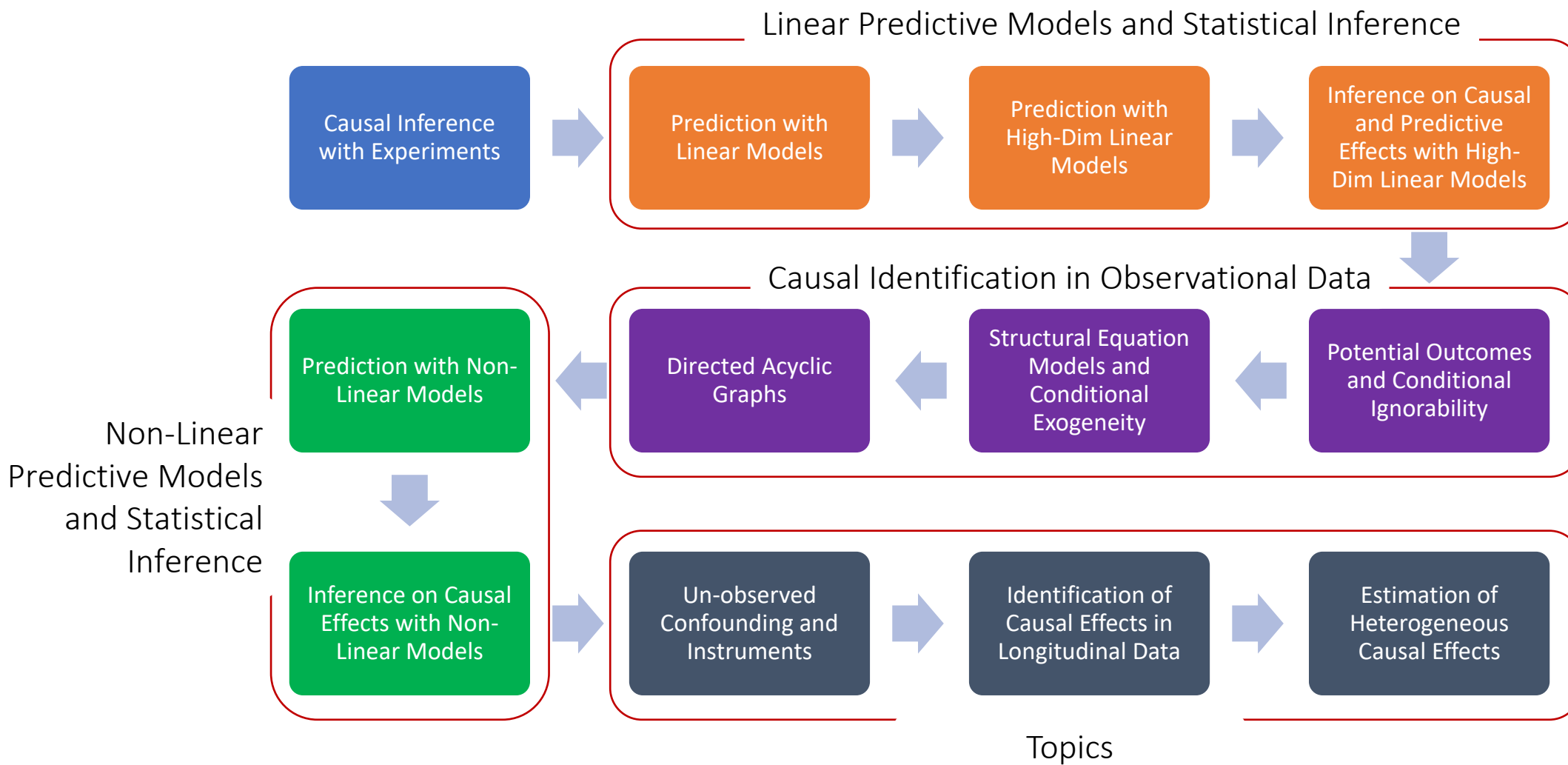
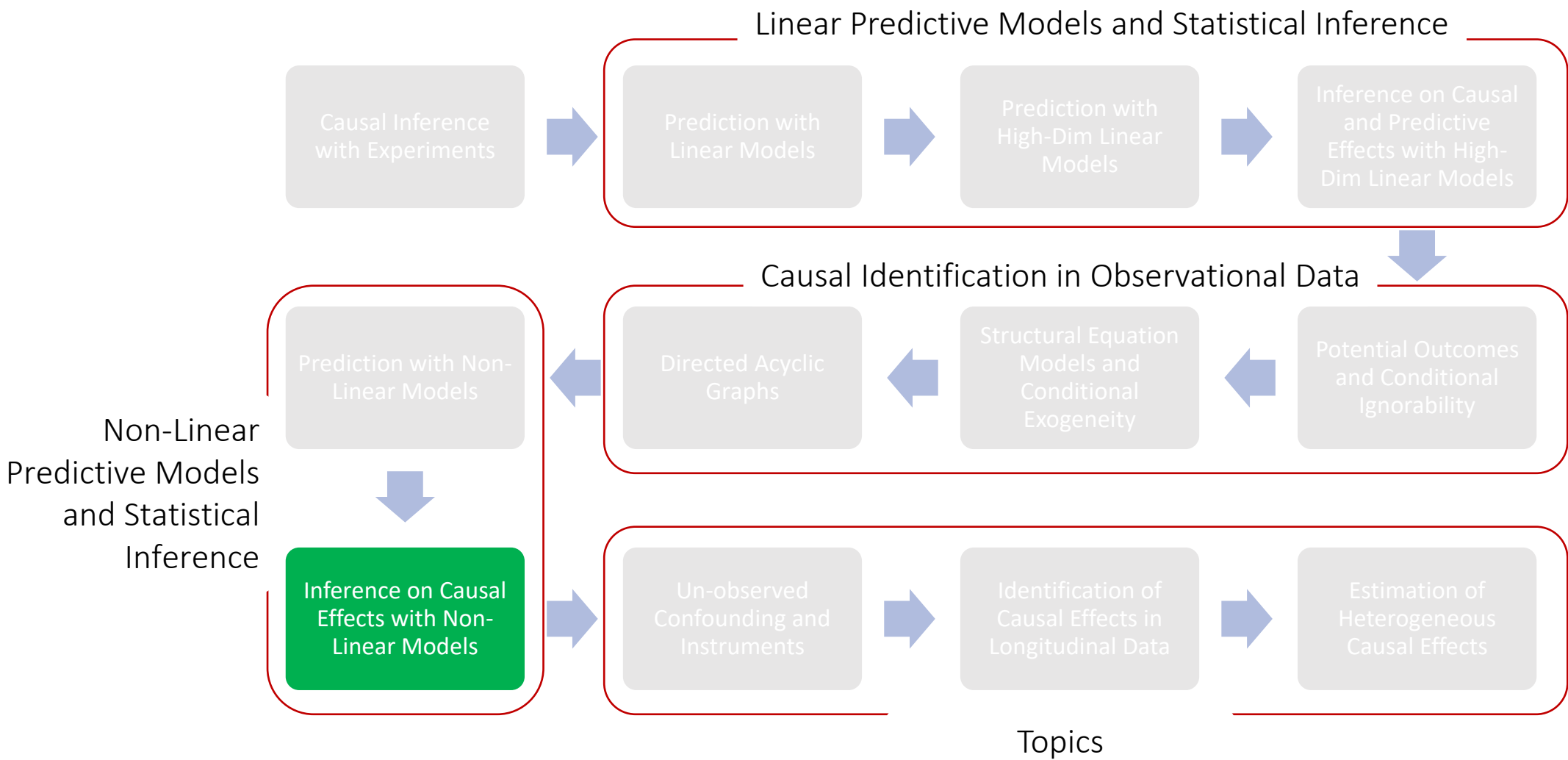


MS&E 228: Inference with Modern Non-Linear Prediction

Vasilis Syrgkanis

MS&E, Stanford

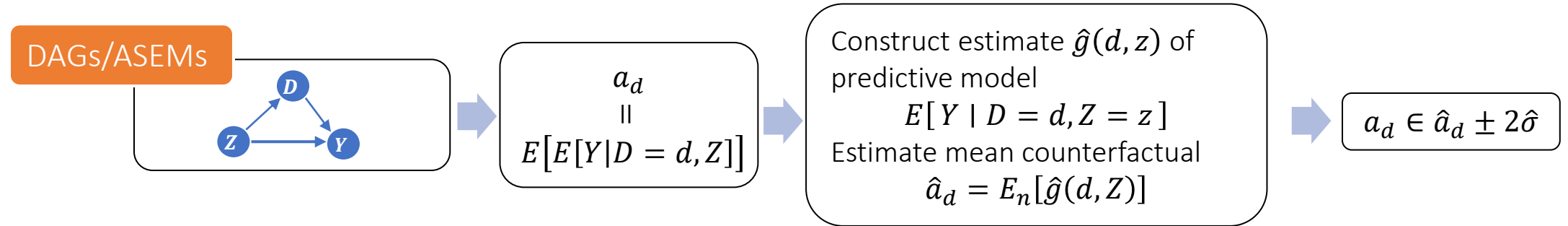




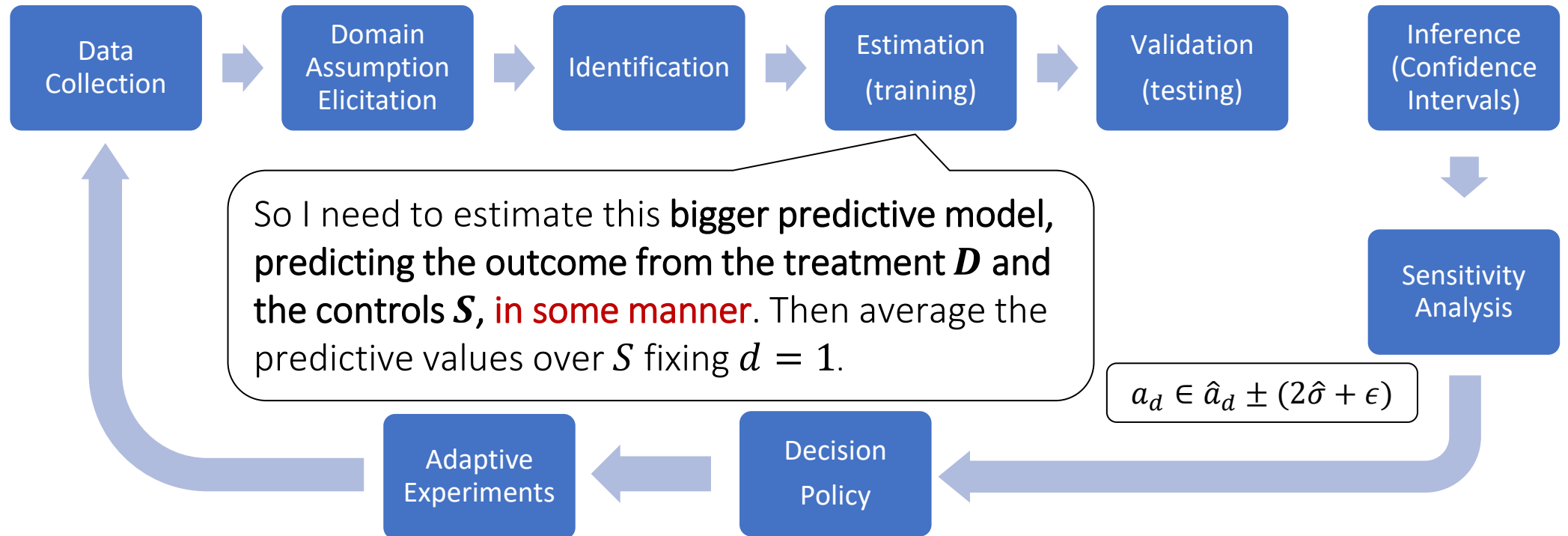
Recap of Last Lecture

Causal Inference Pipeline

Theory

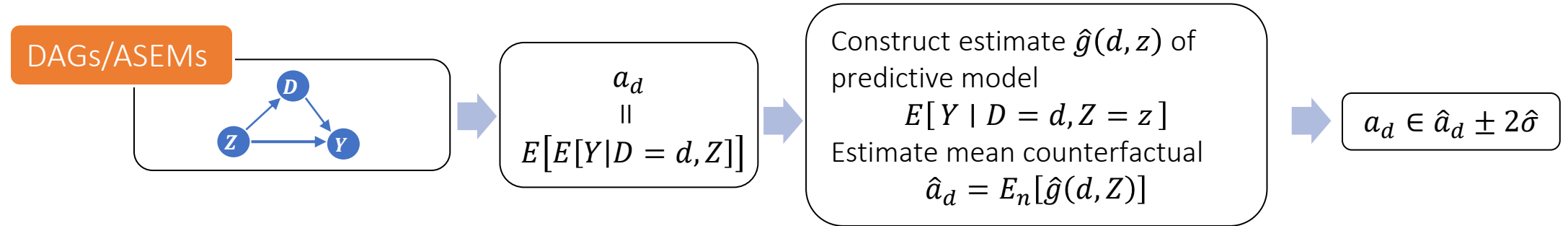


Practice

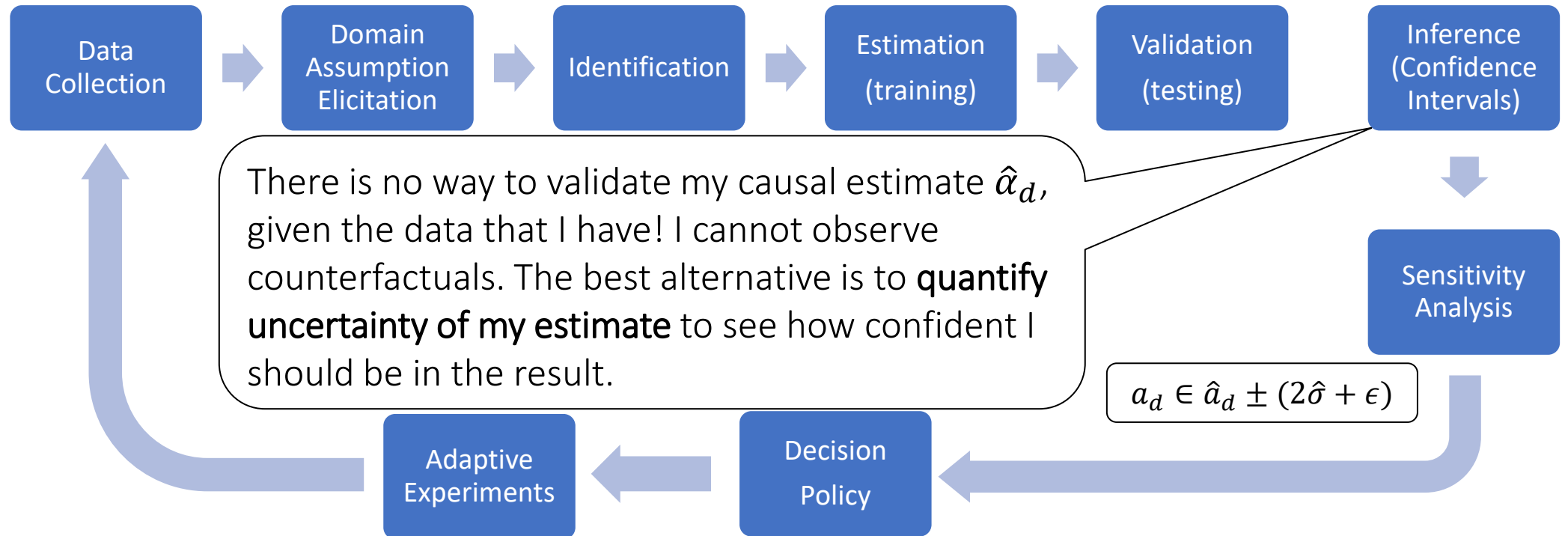


Causal Inference Pipeline

Theory



Practice



Goals for Today

- Methods for Confidence Intervals for ATE with non-linear models
 - General Neyman Orthogonality Framework (Double/Debiased ML)
 - Methods for Confidence Intervals for ATE in a partially-linear model
 - Sample-splitting and cross-fitting
-
- Proof sketch of main theorem*

The Example Problem

Identification under Conditional Ignorability

- Once we condition on enough variables X that affect treatment assignment, remnant variation in D is exogenous (as-if trial)

$$Y^{(d)} \perp D \mid X \quad (\text{conditional ignorability})$$

- Why useful:

$$\begin{aligned} E[Y \mid D = d, X] &= E[Y^{(D)} \mid D = d, X] \\ &= E[Y^{(d)} \mid D = d, X] = E[Y^{(d)} \mid X] \end{aligned}$$

- Average treatment effect is “identified” as (g-formula):

$$\begin{aligned} \theta_0 &= E[Y^{(1)} - Y^{(0)}] = E[E[Y^{(1)} - Y^{(0)} \mid X]] \\ &= E[E[Y \mid D = 1, X] - E[Y \mid D = 0, X]] \end{aligned}$$

Let's take it to data

- We observe n samples Z_1, \dots, Z_n where $Z_i = (X_i, D_i, Y_i)$

- Want to estimate average effect θ_0 , which satisfies:

$$\theta_0 = E[g_0(1, X) - g_0(0, X)]$$

- Where:

$$g_0(D, X) := E[Y \mid D, X]$$

- We want to be able to use ML to learn regression function g_0 !

What do we want from $\hat{\theta}$?

- **Ideally** parametric rates for θ_0 even when we have slower rates for g_0
- **Ideally** construction of confidence intervals for θ_0

- **One approach.** Asymptotic normality

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, \sigma^2)$$

- Implies construction of approximately correct confidence intervals

$$\text{with prob. } \approx 95\%: \theta_0 \in [\hat{\theta} \pm 1.96\hat{\sigma}/\sqrt{n}]$$

Natural Estimation Algorithm

- Estimate \hat{g} of g_0 from data
- Calculate empirical plug-in average:

$$\hat{\theta} := \frac{1}{n} \sum_{i=1}^n \hat{g}(1, X) - \hat{g}(0, X)$$

Natural Algorithm Gone Wrong

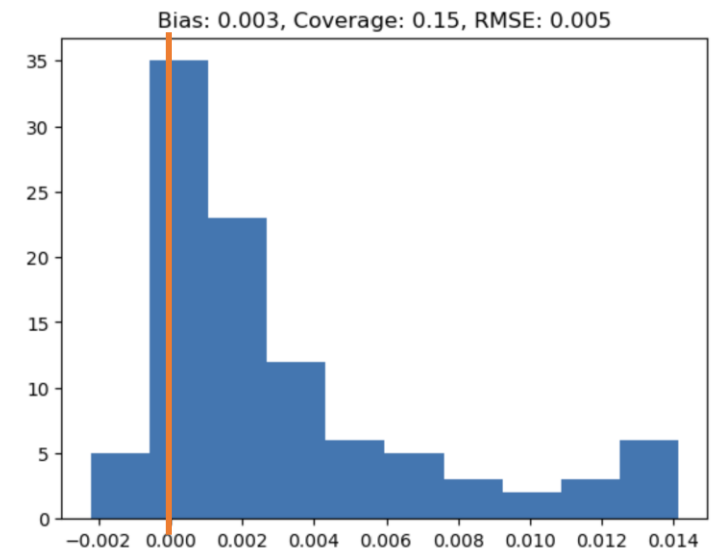
```
def est(X, D, y): # direct non-orthogonal estimator of average effect
    est = RandomForestRegressor(min_samples_leaf=20)
    est.fit(np.hstack([D.reshape(-1, 1), X]), y)
    ones = np.hstack([np.ones((X.shape[0], 1)), X])
    zeros = np.hstack([np.zeros((X.shape[0], 1)), X])
    preds = est.predict(ones) - est.predict(zeros)
    return np.mean(preds), np.std(preds)/np.sqrt(X.shape[0])
```

Simple Example

$$X \sim N(0, I_{20})$$

$$D \sim \text{Binomial}(0.5 + \text{clip}(X_0, -0.4, 0.4))$$

$$y \sim \theta_0 D + X_0 + X_1 + N(0, 1)$$



Natural Estimation Algorithm (Draft 2)

- Split the data in half S_1, S_2
- On first half S_1 , estimate \hat{g} of g_0
- Calculate empirical plug-in average on second half S_2 :

$$\hat{\theta} := \frac{1}{|S_2|} \sum_{i \in S_2} \hat{g}(1, X) - \hat{g}(0, X)$$

Natural Estimation Algorithm (Draft 3)

- Split data in K parts, S_1, \dots, S_K
- For each part k , estimate \hat{g}_k using data from all parts except S_k
- Calculate average over all data:

$$\hat{\theta} = \frac{1}{n} \sum_{k=1}^K \sum_{i \in S_k} \hat{g}_k(1, X) - \hat{g}_k(0, X)$$

Natural Algorithm (Draft 3) Gone Wrong

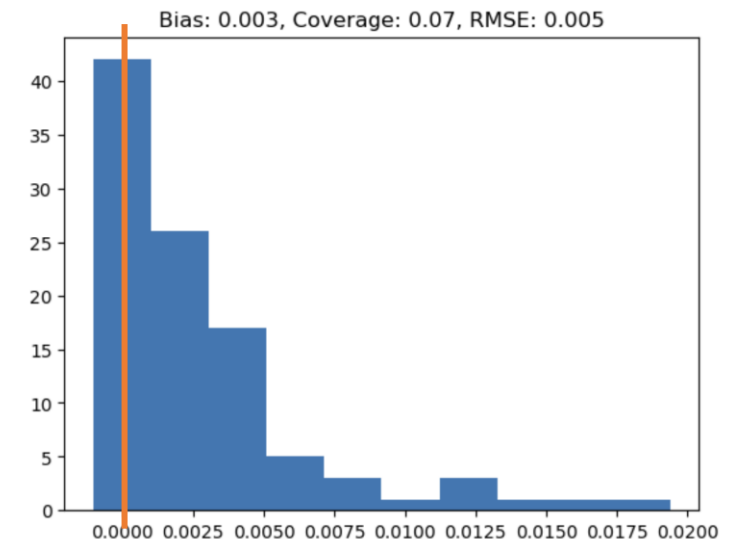
```
def est2(X, D, y): # direct non-orthogonal estimator with sample splitting
    effects = np.zeros(X.shape[0])
    for train, test in KFold(n_splits=3).split(X):
        est = RandomForestRegressor(min_samples_leaf=20)
        est.fit(np.hstack([D[train].reshape(-1, 1), X[train]]), y[train])
        ones = np.hstack([np.ones((X[test].shape[0], 1)), X[test]])
        zeros = np.hstack([np.zeros((X[test].shape[0], 1)), X[test]])
        effects[test] = est.predict(ones) - est.predict(zeros)
    return np.mean(effects), np.std(effects)/np.sqrt(X.shape[0])
```

Simple Example

$$X \sim N(0, I_{20})$$

$$D \sim \text{Binomial}(0.5 + \text{clip}(X_0, -0.4, 0.4))$$

$$y \sim \theta_0 D + X_0 + X_1 + N(0, 1)$$



When is estimate $\hat{\theta}$ \sqrt{n} -asymptotically normal?

When is estimate $\hat{\theta}$ \sqrt{n} -asymptotically normal?
We need to change the moment we use

Debiased Machine Learning

Average Causal Effect Example

- We observe n samples Z_1, \dots, Z_n where $Z_i = (X_i, D_i, Y_i)$
- Want to estimate average effect θ_0 , which satisfies:

$$\theta_0 := E[g_0(1, X) - g_0(0, X)]$$

- Where:

$$g_0(D, X) := E[Y \mid D, X]$$

- The identification formula for θ_0 is sensitive to variations in g
- Any bias or error in g propagates to bias or error in moment and $\hat{\theta}$
- Can we add a correction that corrects the biases of \hat{g}

Better Formula for ATE

Regression residual is a proxy that g is biased

- **Key Idea.** Add a debiasing correction

$$M(g, a) = E[g(1, X) - g(0, X)] + E[a(D, X) (Y - g(D, X))]$$

- What is a_0 ?
- Insensitivity: Take derivative with respect to g at θ_0, g_0, a_0 in any direction $v \in G$

$$\left. \frac{\partial}{\partial t} M(g_0 + t v, a_0) \right|_{t=0} = E[v(1, X) - v(0, X)] - E[a(D, X) v(D, X)] = 0$$

- If this holds then if g is very wrong but a is correct:

$$\begin{aligned} \theta &= E[a_0(D, X)Y] = E[a_0(D, X)E[Y | D, X]] \\ &= E[a_0(D, X)g(D, X)] = E[g(1, X) - g(0, X)] \end{aligned}$$

Inverse Propensity Weighting (IPW)

- The following works: inverse propensity scoring

$$a_0(D, X) = \frac{D}{\Pr[D = 1|X]} - \frac{1 - D}{\Pr[D = 0|X]}$$

- Sketch:

$$\begin{aligned} E \left[\frac{D}{\Pr[D = 1|X]} g(D, X) \right] &= E \left[\frac{D}{\Pr[D = 1|X]} g(1, X) \right] \\ &= E \left[\frac{E[D|X]}{\Pr[D = 1|X]} g(1, X) \right] \\ &= E[g(1, X)] \end{aligned}$$

New Formula is Insensitive

$$M(g, a) = E[g(1, X) - g(0, X)] + E[a(D, X) (Y - g(D, X))]$$

- Take derivative with respect to g at g_0, a_0 in any direction $v \in G$

$$\left. \frac{\partial}{\partial t} M(g_0 + t v, a_0) \right|_{t=0} = E[v(1, X) - v(0, X)] - E[a(D, X) v(D, X)] = 0$$

Take derivative with respect to a at g_0, a_0 in any direction $v \in A$

$$\left. \frac{\partial}{\partial t} M(g_0, a_0 + t v) \right|_{t=0} = E[v(D, X) (Y - g_0(D, X))] = 0$$

Asymptotic Normality of De-biased Estimate

$$\hat{\theta} := E_n[\hat{g}(1, X) - \hat{g}(0, X) + \hat{a}(D, X) \cdot (Y - \hat{g}(D, X))]$$

- Assume that propensities are bounded away from 0 and 1 (strict overlap)
- Assume \hat{g}, \hat{a} estimated on separate sample (or cross-fitting), are consistent and:

$$\sqrt{n} E[(a_0(D, X) - \hat{a}(D, X))(\hat{g}(D, X) - g_0(D, X))] \rightarrow_p 0$$

- Assume random variables $Y, a(D, X), g(D, X)$ have bounded fourth moments
- Then:

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N(0, \sigma^2), \quad \hat{\sigma} = Var_n(\hat{g}(1, X) - \hat{g}(0, X) + \hat{a}(X) \cdot (Y - \hat{g}(X)))$$

Python Pseudocode

```
cv = KFold(n_splits=nfolds, shuffle=True, random_state=123)
yhat0, yhat1 = np.zeros(y.shape), np.zeros(y.shape)
# we will fit a model  $E[Y | D, X]$  by fitting a separate model for  $D=0$ 
# and a separate model for  $D=1$ .
for train, test in cv.split(X, y):
    # train a model on training data that received zero and predict on all test data
    yhat0[test] = modely.fit(X[train][D[train]==0], y[train][D[train]==0]).predict(X[test])
    # train a model on training data that received one and predict on all test data
    yhat1[test] = modely.fit(X[train][D[train]==1], y[train][D[train]==1]).predict(X[test])
# prediction for observed treatment
yhat = yhat0 * (1 - D) + yhat1 * D
# propensity scores
Dhat = cross_val_predict(modeld, X, D, cv=cv, method='predict_proba', n_jobs=-1)[: , 1]
Dhat = np.clip(Dhat, trimming, 1 - trimming)
# doubly robust quantity for every sample
drhat = yhat1 - yhat0 + (y - yhat) * (D/Dhat - (1 - D)/(1 - Dhat))
point = np.mean(drhat)
var = np.var(drhat)
stderr = np.sqrt(var / X.shape[0])
return point, stderr, yhat, Dhat, y - yhat, D - Dhat, drhat
```

Continuous Treatments under Partial Linearity

Partially Linear Model

- Relevant in many applications: dose-response curve in healthcare, effect of price on demand, return-on-investment
- Assume conditional exogeneity

$$Y^{(d)} \perp D \mid X$$

- Assume partially linear response

$$g_0(D, X) = E[Y \mid D, X] = \theta_0 D + f_0(X)$$

- Parameter of interest θ_0 is constant marginal effect of treatment

Generalization of FWL Theorem

- Let's define a slight variant of residualization

$$\tilde{V} = V - E[V|X]$$

- Generalization of FWL theorem to partially linear models

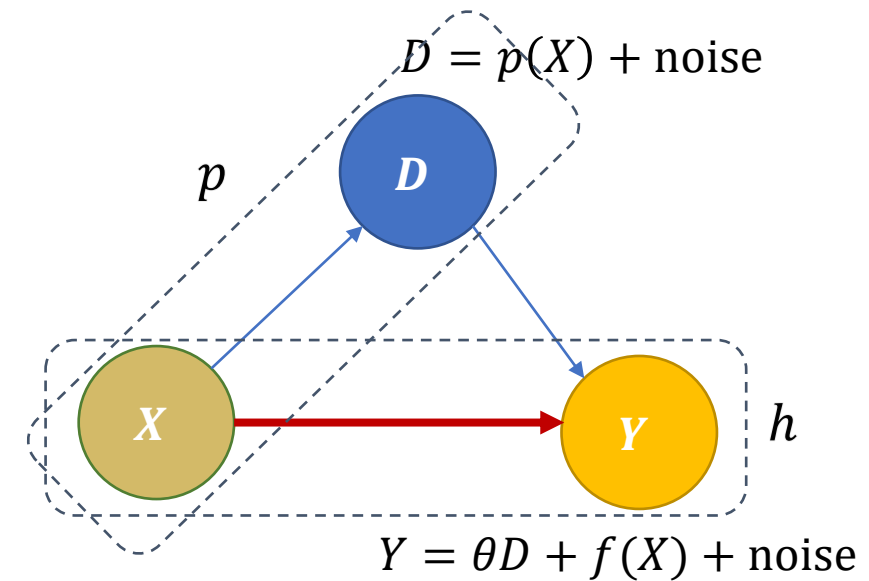
$$\tilde{Y} = \theta_0 \tilde{D} + \epsilon, \quad E[\epsilon|\tilde{D}] = 0$$

- Let's consider the residual outcome

$$\begin{aligned}\tilde{Y} &= Y - E[Y|X] \\ &= \theta_0 D + f_0(X) + \epsilon - E[\theta_0 D + f_0(X) + \epsilon|X] \\ &= \theta_0 D + f_0(X) + \epsilon - \theta_0 E[D|X] - f_0(X) \\ &= \theta_0 (D - E[D|X]) + \epsilon\end{aligned}$$

Orthogonal Method: Double ML

- **Double ML.** Split samples in half
 - Regress $Y \sim X$ with ML on first half, to get estimate $\hat{h}(S)$ of $E[Y|X]$
 - Regress $D \sim X$ with ML on first half, to get estimate $\hat{p}(S)$ of $E[D|X]$



Orthogonal Method: Double ML

- **Double ML.** Split samples in half

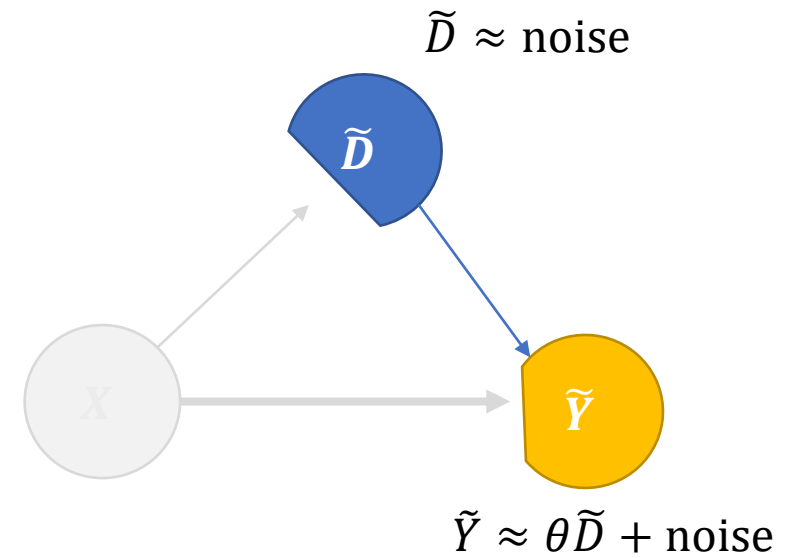
- Regress $Y \sim X$ with ML on first half, to get estimate $\hat{h}(S)$ of $E[Y|X]$
- Regress $D \sim X$ with ML on first half, to get estimate $\hat{p}(S)$ of $E[D|X]$
- Construct residuals on other half, $\tilde{D} := D - \hat{p}(X)$ and $\tilde{Y} := Y - \hat{h}(X)$
- Run OLS on residuals: $\tilde{Y} \sim \tilde{D}$ to get $\hat{\theta}$

- OLS equivalent to solving moment condition:

$$E[(\tilde{Y} - \theta \tilde{D})\tilde{D}] = 0$$

- Orthogonal Moment condition:

$$M(\theta, h, p) = E \left[\left(Y - h(X) - \theta (D - p(X)) \right) (D - p(X)) \right]$$



General Theory

Semi-Parametric Moment Restrictions

- Observe samples Z_1, \dots, Z_n i.i.d. from data distribution D

- Distribution D satisfies vector of moment restrictions

$$M(\theta_0, g_0) := E_{Z \sim D}[m(Z; \theta_0, g_0)] = 0$$

- $\theta_0 \in \mathbb{R}^d$ finite dimensional target parameter of interest
- $g_0 \in G$ potentially infinite dimensional parameter we don't care (nuisance)
- g_0 is un-known and needs to be estimated from data
- Examples:

$$\begin{aligned} m(Z; \theta, g, a) &= g(1, X) - g(0, X) - a(D, X)(Y - g(D, X)) - \theta \\ m(Z; \theta, h, p) &= (Y - h(X) - \theta(D - p(X)))(D - p(X)) \end{aligned}$$

Natural Estimation Algorithm (sample-splitting)

Split the data in half (S_1, S_2)

- On first half S_1 , estimate \hat{g} of g_0
- On second half S_2 , $\hat{\theta}$ is solution w.r.t. θ of empirical plug-in moment equation:

$$M_n(\theta, \hat{g}) := \frac{1}{n_2} \sum_{i \in S_2} m(Z_i; \theta, \hat{g}) = 0$$

Neyman Orthogonality

Moment $M(\theta, g)$ is **Neyman orthogonal** if for any $v \in G - g_0$:

$$D_g M(\theta_0, g_0)[v] := \left. \frac{\partial}{\partial t} M(\theta_0, g_0 + t v) \right|_{t=0} = 0$$

Main Theorem

If moment is Neyman orthogonal and RMSE of \hat{g} is $o_p(n^{-1/4})$, plus regularity conditions

$$\sqrt{n} (\hat{\theta} - \theta_0) \rightarrow N \left(0, J_0^{-1} \Sigma (J_0^{-1})^\top \right)$$

where $J_0 := \nabla_\theta M(\theta_0, g_0)$ and $\Sigma := E[m(Z; \theta_0, g_0) m(Z; \theta_0, g_0)^\top]$

Practical Variants of Sample-Splitting

Cross-fitting and semi-cross-fitting

Cross-fitting

- Sample splitting is statistically lossy
- Only half of the data are used for the final parameter estimation
- Can we utilize all the data?
- *Cross-fitting*: analogous to cross-validation
- Use the second half to train g and predict on first half
- Then calculate parameter using all the data

Cross-fitting Estimation Algorithm

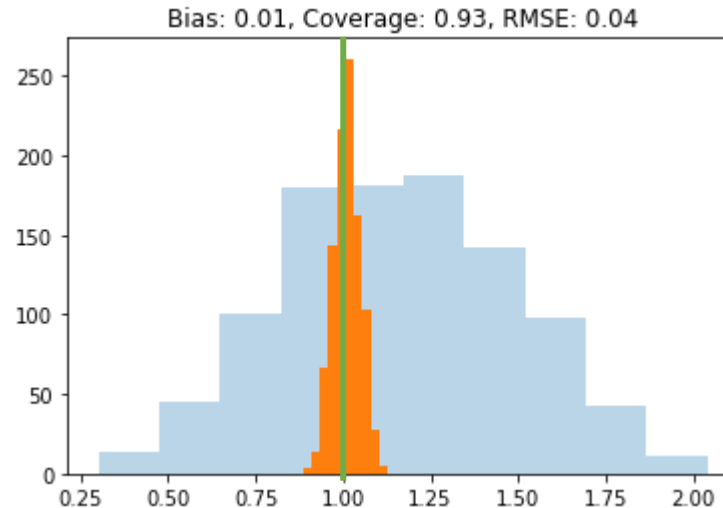
- Split the data in half
- On first half, estimate \hat{g}_1 of g_0 and predict on second half
- On second half, estimate \hat{g}_2 of g_0 and predict on first half
- On all data, solution $\hat{\theta}$ to empirical plug-in moment equation:

$$M_n(\hat{\theta}, \hat{g}) := \frac{1}{n} \sum_{i \in S_1} m(Z_i; \hat{\theta}, \hat{g}_2) + \frac{1}{n} \sum_{i \in S_2} m(Z_i; \hat{\theta}, \hat{g}_1) = 0$$

- In practice do this with $K \approx 3$ to 5 folds: for each fold k train on all other folds and predict on fold k

Natural Algorithm (Draft 3) Gone Right

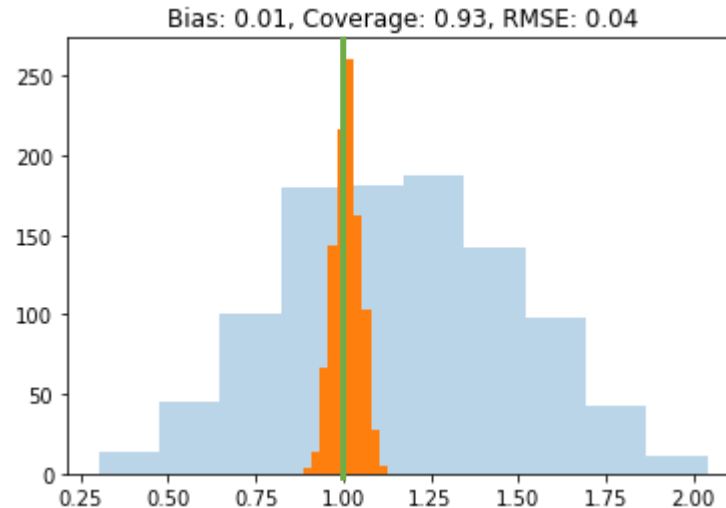
```
def dml2(X, D, y): # orthogonal dml with sample-splitting
    est_y = RandomForestRegressor(min_samples_leaf=20)
    yres = y - cross_val_predict(est_y, X, y, cv=3)
    est_t = RandomForestRegressor(min_samples_leaf=20)
    Dres = D - cross_val_predict(est_t, X, D, cv=3)
    theta = np.mean(yres * Dres) / np.mean(Dres**2)
    var = np.mean((Dres**2) * (yres - theta*Dres)**2) / np.mean(Dres**2)
    stderr = np.sqrt(var / X.shape[0])
    return theta, stderr
```



Natural Algorithm (Draft 3) Gone Right

```
from econml.dml import LinearDML

dml = LinearDML(model_y=RandomForestRegressor(min_samples_leaf=20),
                 model_t=RandomForestRegressor(min_samples_leaf=20))
est.fit(y, D, W=X).effect_inference()
```



Stacking and Model Selection

- If we want to choose among many models or perform stacking, we can just use a stacked or automl model in place of each ML model

Stacking ML Models

```
def dml2(X, D, y): # orthogonal dml with sample-splitting
    est_y = StackingRegressor([rf, nnet, gbf, lasso])
    yres = y - cross_val_predict(est_y, X, y, cv=3)
    est_t = StackingRegressor([rf, nnet, gbf, lasso])
    Dres = D - cross_val_predict(est_t, X, D, cv=3)
    theta = np.mean(yres * Dres) / np.mean(Dres**2)
    var = np.mean((Dres**2) * (yres - theta*Dres)**2) / np.mean(Dres**2)
    stderr = np.sqrt(var / X.shape[0])
    return theta, stderr
```

AutoML Models

```
from flaml import AutoML

def dml2(X, D, y): # orthogonal dml with sample-splitting
    est_y = AutoML()
    yres = y - cross_val_predict(est_y, X, y, cv=3)
    est_t = AutoML()
    Dres = D - cross_val_predict(est_t, X, D, cv=3)
    theta = np.mean(yres * Dres) / np.mean(Dres**2)
    var = np.mean((Dres**2) * (yres - theta*Dres)**2) / np.mean(Dres**2)
    stderr = np.sqrt(var / X.shape[0])
    return theta, stderr
```

Stacking and Model Selection

- If we want to choose among many models or perform stacking, we can just use a stacked or automl model in place of each ML model
- Model selection or stacking done many times within each training fold
- Computationally expensive and statistically lossy
- Can we use all the data to at least select among models?

Semi-Cross-fitting Estimation Algorithm

- Split the data in half (*in practice K folds*)
- On first half, estimate $\hat{g}_1^{(1)}, \dots, \hat{g}_1^{(L)}$ of g_0 and predict on second half
- On second half, estimate $\hat{g}_2^{(1)}, \dots, \hat{g}_2^{(L)}$ of g_0 and predict on first half
- Choose the model $\ell \in \{1, \dots, L\}$ that optimizes out-of-sample RMSE
- On all data, solution $\hat{\theta}$ to empirical plug-in moment equation:

$$M_n(\hat{\theta}, \hat{g}) := \frac{1}{n} \sum_{i \in S_1} m(Z_i; \hat{\theta}, \hat{g}_2^{(\ell)}) + \frac{1}{n} \sum_{i \in S_2} m(Z_i; \hat{\theta}, \hat{g}_1^{(\ell)}) = 0$$

Semi-Crossfitting

```
def dml2(X, D, y): # orthogonal dml with semi-crossfitting
    # cross val predict with many models
    est_y = [rf, gbf, lasso]
    yres = np.array([y - cross_val_predict(est, X, y, cv=3) for est in est_y])
    est_d = [rf, gbf, lasso]
    Dres = np.array([D - cross_val_predict(est, X, D, cv=3) for est in est_d])
    # select models with best out of fold performance
    best_y = np.argmin(np.mean(yres**2, axis=1))
    best_d = np.argmin(np.mean(Dres**2, axis=1))
    yres = yres[best_y]
    Dres = Dres[best_d]
    # go with their corresponding residuals
    theta = np.mean(yres * Dres) / np.mean(Dres**2)
    var = np.mean((Dres**2) * (yres - theta*Dres)**2) / np.mean(Dres**2)
    stderr = np.sqrt(var / X.shape[0])
    return theta, stderr
```

Semi-Crossfitting

- If the number of models L is small, then “spillover” is ok and approach still works. For practical purposes L should be thought as constant.
- Under further regularity, provably asymptotic normality holds if

$$\sqrt{\log(L)} = o(n^{1/4})$$

Semi-Cross-fitting with Stacking

- Split the data in half (*in practice K folds*)
- On first half, estimate $\hat{g}_1^{(1)}, \dots, \hat{g}_1^{(L)}$ of g_0 and predict on second half
- On second half, estimate $\hat{g}_2^{(1)}, \dots, \hat{g}_2^{(L)}$ of g_0 and predict on first half
- Construct weights $\alpha_1, \dots, \alpha_\ell$ on the models using all the data (stacking)

- On all data, solution $\hat{\theta}$ to empirical plug-in moment equation:

$$M_n(\hat{\theta}, \hat{g}) := \frac{1}{n} \sum_{i \in S_1} m(Z_i; \hat{\theta}, \hat{g}_2^{(\ell)}) + \frac{1}{n} \sum_{i \in S_2} m(Z_i; \hat{\theta}, \hat{g}_1^{(\ell)}) = 0$$

Semi-Crossfitting with Stacking

```
def dml2(X, D, y): # orthogonal dml with semi-crossfitting and stacking
    # cross val predict with many models
    est_y = [rf, gbf, lasso]
    ypreds = np.array([cross_val_predict(est, X, y, cv=3) for est in est_y]).T
    est_d = [rf, gbf, lasso]
    Dpreds = np.array([cross_val_predict(est, X, D, cv=3) for est in est_d]).T
    # calculate stacked residuals by finding optimal coefficients
    # and weighing out-of-sample predictions by these coefficients
    yres = y - LinearRegression().fit(ypreds, y).predict(ypreds)
    Dres = D - LinearRegression().fit(Dpreds, D).predict(Dpreds)
    # go with the stacked residuals
    theta = np.mean(yres * Dres) / np.mean(Dres**2)
    var = np.mean((Dres**2) * (yres - theta*Dres)**2) / np.mean(Dres**2)
    stderr = np.sqrt(var / X.shape[0])
    return theta, stderr
```

Semi-Crossfitting

- If the number of models L is small, then “spillover” is ok and approach still works. For practical purposes L should be thought as constant.
- Under further regularity, provably asymptotic normality holds if
$$\sqrt{L} = o(n^{1/4})$$

Equivalent view of cross-fitting with stacking (lens of FWL theorem)

- Construct out of fold predictions based on many ML models
- Use these predictions as engineered features X in a simple OLS regression on D, X
- Use the coefficient and standard error of D from this final OLS