

# Chapter 18: Variable Selection for Causal Inference

## Contents

<b>1</b>	<b>18.1 The Traditional Approach (pp. 265-267)</b>	<b>1</b>
1.1	Prediction vs Causal Inference . . . . .	2
1.2	Stepwise Selection . . . . .	2
<b>2</b>	<b>18.2 Confounding and Confounders (pp. 267-270)</b>	<b>2</b>
2.1	Backdoor Paths . . . . .	2
2.2	Sufficient Adjustment Sets . . . . .	3
<b>3</b>	<b>18.3 Confounding Adjustment (pp. 270-273)</b>	<b>3</b>
3.1	Variables to Include . . . . .	3
3.2	Variables to Exclude . . . . .	3
3.3	Colliders . . . . .	3
<b>4</b>	<b>18.4 Instrumental Variables and M-bias (pp. 273-276)</b>	<b>4</b>
4.1	M-bias (Butterfly Bias) . . . . .	4
4.2	Instrumental Variables Revisited . . . . .	4
<b>5</b>	<b>18.5 Confounders, Mediators, and Intermediate Confounders (pp. 276-279)</b>	<b>5</b>
5.1	Time-Varying Confounding . . . . .	5
5.2	Intermediate Confounder . . . . .	5
<b>6</b>	<b>18.6 Selecting Variables for Precision (pp. 279-281)</b>	<b>5</b>
6.1	Precision Variables . . . . .	6
6.2	Instruments as Precision Variables? . . . . .	6
6.3	Practical Strategy . . . . .	6
<b>7</b>	<b>18.7 Using Causal Diagrams (pp. 281-282)</b>	<b>6</b>
7.1	Steps for Using DAGs . . . . .	6
7.2	Software Tools . . . . .	7
<b>8</b>	<b>Summary</b>	<b>7</b>

This chapter addresses a critical question in causal inference: **Which variables should we adjust for?** Not all variables that predict the outcome should be included in causal models. Some variables, if adjusted for, can introduce bias rather than remove it. We provide guidance on variable selection using causal diagrams.

This chapter is based on Hernán and Robins (2020, chap. 18, pp. 265-282).

**Central message:** Variable selection for causal inference is fundamentally different from variable selection for prediction. We must use causal reasoning (often encoded in DAGs) rather than purely statistical criteria.

## 1 18.1 The Traditional Approach (pp. 265-267)

---

Traditional variable selection methods are designed for **prediction**, not causal inference.

## 1.1 Prediction vs Causal Inference

**Prediction goal:** Minimize prediction error for  $Y$  given covariates

- Include any variable that improves prediction
- Use criteria like AIC, BIC, cross-validation
- More variables (if not overfitting)  $\rightarrow$  better prediction

**Causal inference goal:** Estimate  $E[Y^a]$  or  $E[Y^{a=1}] - E[Y^{a=0}]$

- Include variables that remove confounding
- Exclude variables that introduce bias
- More variables  $\rightarrow$  better causal estimates

**Why they differ:**

In prediction, we want to capture all associations between covariates and outcome. In causal inference, we want to isolate the causal effect of treatment, which means blocking certain associations and preserving others.

**Example:** A mediator predicts the outcome well, but adjusting for it removes part of the causal effect we want to estimate.

## 1.2 Stepwise Selection

**Traditional approach:** Stepwise regression (forward, backward, or both)

- Add/remove variables based on statistical significance or information criteria
- Maximize  $R^2$  or minimize AIC/BIC

**Problem for causal inference:**

- May exclude important confounders (if weak predictors)
- May include colliders or mediators (if strong predictors)
- Ignores causal structure

**Recommendation:** Do not use stepwise selection for causal inference.

## 2 18.2 Confounding and Confounders (pp. 267-270)

---

What exactly is a confounder, and when should we adjust for it?

**Definition 2.1** (Confounder (Formal Definition)). A variable  $L$  is a **confounder** for the effect of  $A$  on  $Y$  if:

1.  $L$  is associated with treatment  $A$
2.  $L$  is a cause of outcome  $Y$
3.  $L$  is not affected by  $A$  (not a descendant of  $A$  on a causal DAG)

**Causal criterion:**  $L$  is on a backdoor path from  $A$  to  $Y$ .

### 2.1 Backdoor Paths

**Backdoor path:** A path from  $A$  to  $Y$  that starts with an arrow into  $A$

$$A \leftarrow L \rightarrow Y$$

Such paths create non-causal association between  $A$  and  $Y$ .

**Goal:** Block all backdoor paths to eliminate confounding.

## 2.2 Sufficient Adjustment Sets

**Definition 2.2** (Sufficient Adjustment Set). A set of variables  $L$  is **sufficient for confounding adjustment** if conditioning on  $L$  blocks all backdoor paths from  $A$  to  $Y$ .

Equivalently:  $(Y^a \perp\!\!\!\perp A \mid L)$  for all  $a$  (conditional exchangeability).

**Multiple sufficient sets:** There may be many sufficient adjustment sets. We want to choose one that:

1. Blocks all backdoor paths (necessary)
2. Doesn't introduce new bias (important)
3. Is measurable and measured

**DAG-based approach:**

1. Draw a causal DAG representing your subject-matter knowledge
2. Identify all backdoor paths from  $A$  to  $Y$
3. Find a set  $L$  that blocks all backdoor paths
4. Adjust for  $L$  (and only  $L$ )

This is superior to traditional approaches because it's based on causal structure, not statistical associations.

## 3 18.3 Confounding Adjustment (pp. 270-273)

---

When we adjust for a sufficient set, we remove confounding. But be careful about adjusting for too much.

### 3.1 Variables to Include

**Confounders:** Variables on backdoor paths

- Include to block backdoor paths
- These are causes of both treatment and outcome (or proxies thereof)

**Example DAG:**

$$A \leftarrow L \rightarrow Y$$

Adjust for  $L$  to block the backdoor path.

### 3.2 Variables to Exclude

**Mediators:** Variables on the causal path from  $A$  to  $Y$

- Do NOT adjust (would remove part of the causal effect)

**Example DAG:**

$$A \rightarrow M \rightarrow Y$$

If we adjust for  $M$ , we block the causal path through  $M$ .

**Descendants of treatment:** Variables affected by  $A$

- Usually do NOT adjust (may induce bias)

### 3.3 Colliders

**Definition 3.1** (Collider). A **collider** on a path is a variable with two arrows pointing into it.

**Example:**

$$A \rightarrow C \leftarrow U \rightarrow Y$$

$C$  is a collider on the path  $A \rightarrow C \leftarrow U \rightarrow Y$ .

**Property:** This path is **blocked** by default (without conditioning on  $C$ ).

**Danger:** If we condition on  $C$  (or its descendants), we **open** the path, creating **collider bias**.

**Rule:** Do NOT adjust for colliders (unless necessary to block other paths).

**Collider bias example:**

- $A$ : Athletic ability
- $C$ : Being on a sports team (affected by both  $A$  and parental encouragement  $U$ )
- $U$ : Parental encouragement (also affects academic performance  $Y$ )
- $Y$ : Academic performance

If we condition on being on a sports team ( $C = 1$ ), we induce a negative association between athletic ability and parental encouragement. This can bias estimates of the effect of  $A$  on  $Y$ .

**Practical implication:** Including “selection variables” in regression can introduce bias.

## 4 18.4 Instrumental Variables and M-bias (pp. 273-276)

---

Some variables should not be adjusted for even if they’re associated with both treatment and outcome.

### 4.1 M-bias (Butterfly Bias)

**DAG structure:**

$$\begin{array}{ccc} U1 & \rightarrow & L & \leftarrow & U2 \\ \downarrow & & & & \downarrow \\ A & & & & Y \end{array}$$

**Properties:**

- $L$  is associated with both  $A$  and  $Y$  (through  $U1$  and  $U2$ )
- $L$  is a **collider** on the path  $A \leftarrow U1 \rightarrow L \leftarrow U2 \rightarrow Y$
- This path is blocked by default
- But if we adjust for  $L$ , we **open** this path!

**Result:** Adjusting for  $L$  introduces bias even though  $L$  is associated with both  $A$  and  $Y$ .

**Practical example:**

- $A$ : Smoking
- $Y$ : Lung cancer
- $U1$ : Genetic variant affecting smoking propensity
- $U2$ : Different genetic variant affecting cancer risk
- $L$ : Being in a genetic study (selected based on  $U1$  and  $U2$ )

In the genetic study sample, adjusting for study participation  $L$  induces collider bias.

**Lesson:** Don’t adjust for a variable just because it’s associated with treatment and outcome. Check the causal structure!

### 4.2 Instrumental Variables Revisited

An **instrumental variable**  $Z$  satisfies:

$$Z \rightarrow A \rightarrow Y$$

with no backdoor paths from  $Z$  to  $Y$ .

**Should we adjust for  $Z$ ?**

- If using IV methods: NO (use  $Z$  as instrument)
- If using standard methods and  $Z$  is not a confounder: NO (unnecessary, may hurt efficiency)
- If  $Z$  confounds some other relationship of interest: MAYBE

## 5 18.5 Confounders, Mediators, and Intermediate Confounders (pp. 276-279)

---

Time-varying treatments create new challenges for variable selection.

### 5.1 Time-Varying Confounding

**Setting:** Treatment varies over time ( $A_0, A_1, \dots$ ), as do confounders ( $L_0, L_1, \dots$ )

**Time-varying confounder:**  $L_1$  is a confounder for the effect of  $A_1$  on  $Y$

**Problem:** If  $A_0$  affects  $L_1$ , then:

- $L_1$  is a confounder (need to adjust)
- $L_1$  is a mediator (should not adjust in standard regression)

### 5.2 Intermediate Confounder

**Definition 5.1** (Intermediate Confounder (Time-Dependent Confounder Affected by Prior Treatment)). A variable  $L_1$  is an **intermediate confounder** if:

1.  $L_1$  is a confounder for the effect of  $A_1$  on  $Y$
2.  $L_1$  is affected by prior treatment  $A_0$

**DAG:**

$$\begin{array}{ccc} A_0 & \rightarrow & L_1 & \rightarrow & Y \\ \downarrow & & \downarrow & & \\ A_1 & \rightarrow & & & Y \end{array}$$

**Standard regression fails:** Cannot correctly adjust for  $L_1$  using standard methods.

**Solutions:**

- **G-methods:** Parametric g-formula, IP weighting, g-estimation (Part III)
- These methods properly handle time-varying confounders affected by prior treatment

**Why standard regression fails:**

If we adjust for  $L_1$ , we block the indirect effect  $A_0 \rightarrow L_1 \rightarrow Y$ . If we don't adjust, we have confounding of  $A_1 \rightarrow Y$ .

**Example:** HIV treatment and CD4 count

- $A_0, A_1$ : Antiretroviral therapy at times 0 and 1
- $L_1$ : CD4 count at time 1 (affected by  $A_0$ , affects treatment choice  $A_1$ , affects outcome  $Y$ )
- Standard regression cannot handle this correctly

**Part III solution:** Marginal structural models with IP weighting or g-formula can properly estimate effects in this setting.

## 6 18.6 Selecting Variables for Precision (pp. 279-281)

---

After ensuring confounding is addressed, can we include additional variables to improve precision?

## 6.1 Precision Variables

**Definition:** Variables associated with the outcome but not with treatment (after accounting for confounders).

**Example DAG:**

$$A \rightarrow Y \leftarrow V$$

$V$  is associated with  $Y$  but not with  $A$  (no arrow from  $V$  to  $A$  or shared causes).

**Effect of adjustment:**

- Does NOT affect bias (no confounding)
- DOES improve precision (reduces residual variance)

**Recommendation:** Include precision variables to improve efficiency.

## 6.2 Instruments as Precision Variables?

**Question:** Should we include instrumental variables in outcome models?

**Answer:** Generally NO.

- Instruments are associated with  $A$  but (by exclusion) not directly with  $Y$
- Including them in outcome models doesn't improve precision
- May slightly worsen precision due to additional parameters

## 6.3 Practical Strategy

1. **First priority:** Include all variables needed to block backdoor paths (confounders)
2. **Second priority:** Exclude colliders, mediators, and descendants of treatment
3. **Third priority:** Consider including precision variables if they:
  - Strongly predict the outcome
  - Are not affected by treatment
  - Don't introduce collinearity issues

**Balance precision and bias:**

- Including precision variables:  $\uparrow$  efficiency, but more complex models
- Excluding precision variables:  $\downarrow$  efficiency, but simpler models

With large samples, precision gains may be modest. With small samples, precision improvements can be valuable.

**Practical consideration:** In many applications, confounders ARE strong predictors of the outcome, so adjusting for them serves both purposes (removes bias and improves precision).

# 7 18.7 Using Causal Diagrams (pp. 281-282)

---

**Causal DAGs** (directed acyclic graphs) are invaluable tools for variable selection.

## 7.1 Steps for Using DAGs

1. **Draw the DAG:**
  - Represent your causal assumptions about relationships between variables
  - Include treatment, outcome, all measured covariates, and key unmeasured variables
  - Draw arrows representing direct causal effects
2. **Identify backdoor paths:**
  - Find all paths from  $A$  to  $Y$  that start with an arrow into  $A$
  - These are sources of confounding

### 3. Find sufficient adjustment sets:

- Identify sets of variables that block all backdoor paths
- Avoid inducing collider bias
- Use algorithms (e.g., `dagitty` R package) if DAG is complex

### 4. Choose an adjustment set:

- Select a sufficient set that is measured
- Prefer simpler sets (fewer variables) when multiple options exist
- Check for practical considerations (measurement error, missing data, etc.)

## 7.2 Software Tools

### R package `dagitty`:

- Define DAGs
- Find adjustment sets automatically
- Check conditional independencies implied by the DAG
- Visualize DAGs

### Example:

```
library(dagitty)
dag <- dagitty('dag {
  A -> Y
  L1 -> A
  L1 -> Y
  L2 -> Y
  U -> A
  U -> Y
}')
adjustmentSets(dag, exposure = "A", outcome = "Y")
```

### DAGs encode assumptions:

- What you include in the DAG (and what you exclude) represents your causal knowledge
- DAGs make assumptions explicit and falsifiable
- Subject-matter expertise is crucial for drawing valid DAGs

### Limitations:

- DAG is only as good as your causal knowledge
- Cannot test all assumptions (especially about unmeasured variables)
- Must think carefully about what to include

**Recommendation:** Draw multiple plausible DAGs, find adjustment sets for each, conduct sensitivity analyses.

## 8 Summary

---

### Key principles for variable selection in causal inference:

1. Use **causal reasoning**, not statistical criteria
2. **Adjust for confounders** (variables on backdoor paths)
3. **Don't adjust for mediators** (variables on causal paths from treatment)
4. **Don't adjust for colliders** (unless necessary to block other paths)
5. **Don't adjust for descendants of treatment** (generally)
6. **Consider precision variables** (if they don't introduce bias)
7. Use **causal DAGs** to guide selection

### Variables to include:

- Confounders (on backdoor paths from  $A$  to  $Y$ )
- Precision variables (predict  $Y$ , not affected by  $A$ , not colliders)
- Proxies for unmeasured confounders

**Variables to exclude:**

- Mediators (on causal path from  $A$  to  $Y$ )
- Colliders (except if needed to block backdoor paths)
- Descendants of treatment (except special cases)
- Instruments (when using standard methods)
- Variables that induce M-bias

**Special cases:**

- **Intermediate confounders:** Time-varying confounders affected by prior treatment
  - Cannot be handled by standard regression
  - Require g-methods (Part III)

**Tools:**

- **Causal DAGs:** Represent causal structure graphically
- **Backdoor criterion:** Identify sufficient adjustment sets
- **Software:** dagitty, ggdag (R), DAGitty (web interface)

**Common mistakes:**

1. Using stepwise selection or AIC/BIC for variable selection
2. Adjusting for all predictors of the outcome
3. Adjusting for mediators
4. Conditioning on colliders
5. Ignoring time-varying confounding

**Practical workflow:**

1. Draw a causal DAG based on subject-matter knowledge
2. Identify backdoor paths from treatment to outcome
3. Find sufficient adjustment sets using the backdoor criterion
4. Choose a set that is measured and practical
5. Fit causal model adjusting for that set (using appropriate method)
6. Conduct sensitivity analyses with alternative DAGs/adjustment sets

**Bottom line:**

Variable selection for causal inference requires causal thinking. Statistical significance, prediction accuracy, and  $R^2$  are not appropriate criteria. Instead:

- Think about causal structure
- Use DAGs to formalize assumptions
- Apply the backdoor criterion
- Adjust for the right variables, not just any variables

**Looking ahead:** Part III extends these ideas to longitudinal settings with time-varying treatments and confounders, where variable selection becomes even more critical and complex.

Hernán, Miguel A, and James M Robins. 2020. *Causal Inference: What If*. Boca Raton: Chapman & Hall/CRC. <https://miguelhernan.org/whatifbook>.