

---

# **Non-Standard Databases and Data Mining**

Intervention

Dr. Özgür Özçep

**Universität zu Lübeck**

**Institut für Informationssysteme**

Presented by: Prof. Dr. Ralf Möller

---

# Structural Causal Models

Slides prepared by Özgür Özçep

## Part II: Intervention

# Literature

---

- J.Pearl, M. Glymour, N. P. Jewell:  
Causal inference in statistics – A primer, Wiley, 2016.  
(Main Reference)
- J. Pearl: Causality, CUP, 2000.

# Intervention

---

- Important aim of SCMs for given data: Where to intervene in order to achieve desired effects.

## Examples

- Data on wildfires: How to intervene in order to decrease wildfires?
  - Data on TV and aggression: How to intervene in order to lower aggression of children?
- 
- How to model “intervention” and associated effects within SCMs and their graphs?

# Randomized Controlled Experiment

---

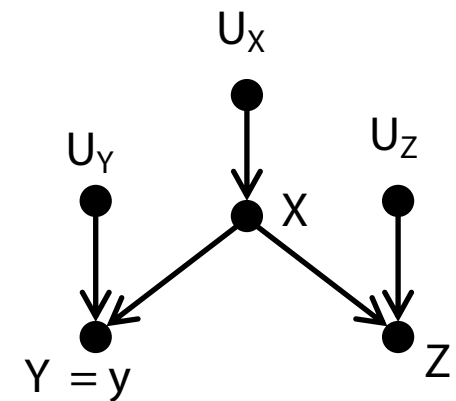
- **Randomized controlled experiment** gold standard
  - Aim: Answer question whether a change in RV  $X$  has indeed an effect on some target RV  $Y$
  - If outcome of experiment is yes,  $X$  is a RV to intervene upon
  - Test condition: all variables different from  $X$  are static (fixed) or vary fully randomly.
- **Problem:** Cannot always set up such an experiment
  - **Example:** cannot control weather in order to test variables influencing wildfire
- **Instead:** use observational data & causal model

# Intervention

## Example (SCM 5; Intervention)

(  $X$  = Temperature,  $Y$  = Ice cream sale,  $Z$  = Crime)

- Would intervention on ice cream sales ( $Y$ ) lead to decrease of crime ( $Z$ )?
- What does it mean to intervene on  $Y$ ?
  - Fix value of  $Y$  in the sense of inhibiting the natural influences on  $Y$  according to SCM (here of  $U_Y$  and  $X$ )
  - Leads to change of the SCM



# Intervention vs. Conditioning

---

- Intervention denoted by  $\text{do}(Y = y)$

$$P(Z = z \mid \text{do}(Y = y)) =$$

probability of event  $Z = z$  on intervening  
upon  $Y$  by setting  $Y = y$

Intervention changes the data generation mechanism

- In contrast

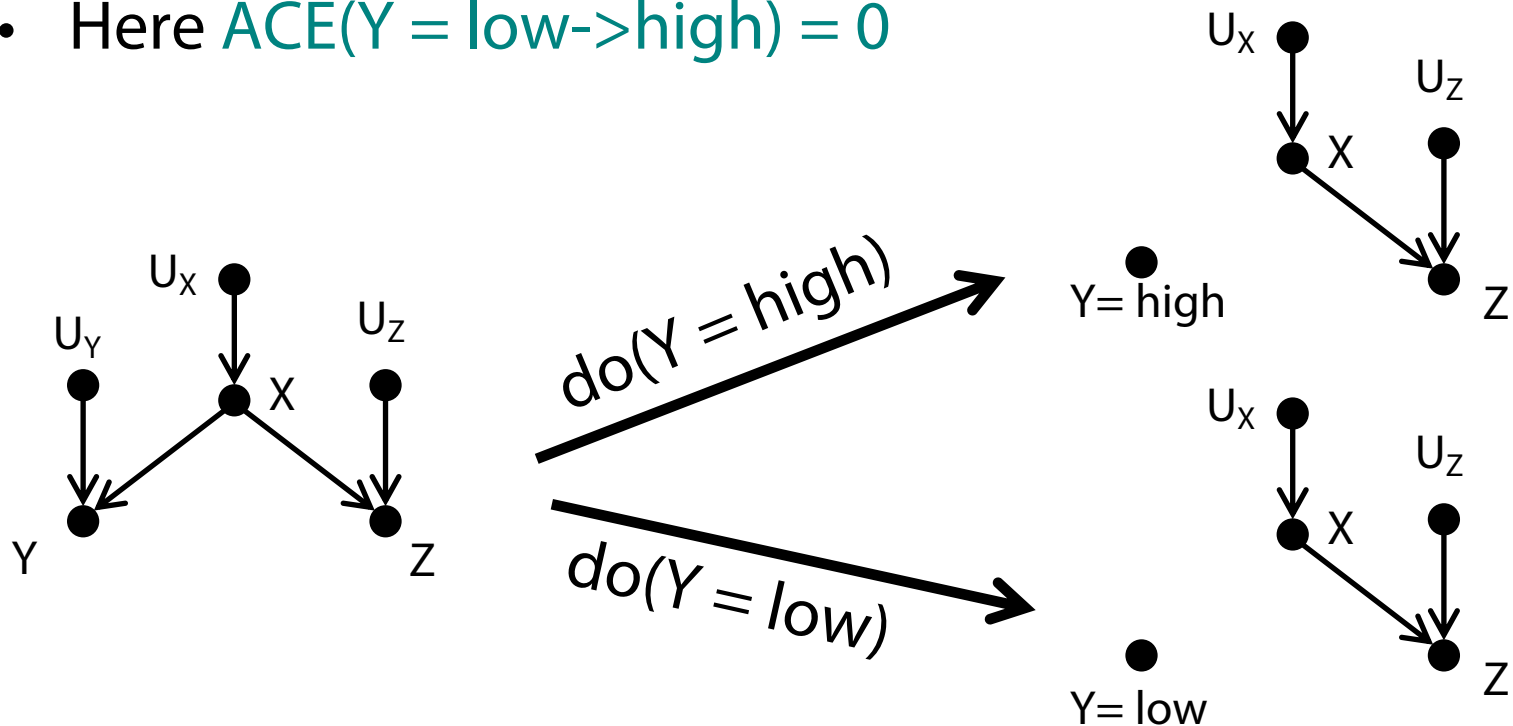
$$P(Z = z \mid Y = y) =$$

probability of event  $Z = z$  when knowing that  $Y = y$

Conditioning only filters on the data

# Average Causal Effect (ACE)

- Would an intervention on ice cream sales (Y) by increasing Y lead to a decrease of crime (Z)?
- Causal Effect Difference/Average Causal Effect (ACE)  
 $P(Z = \text{low} | \text{do}(Y = \text{high})) - P(Z = \text{low} | \text{do}(Y = \text{low}))$
- Here  $\text{ACE}(Y = \text{low} \rightarrow \text{high}) = 0$





# General Causal Effect

- How effective is drug usage for recovery?  
 $ACE = P(Y = 1 \mid \text{do}(X = 1)) - P(Y = 1 \mid \text{do}(X = 0))$
- Need to compute **general causal effect**

## Definition

The **general causal effect** (GCE) of  $X$  on  $Y$  is given by

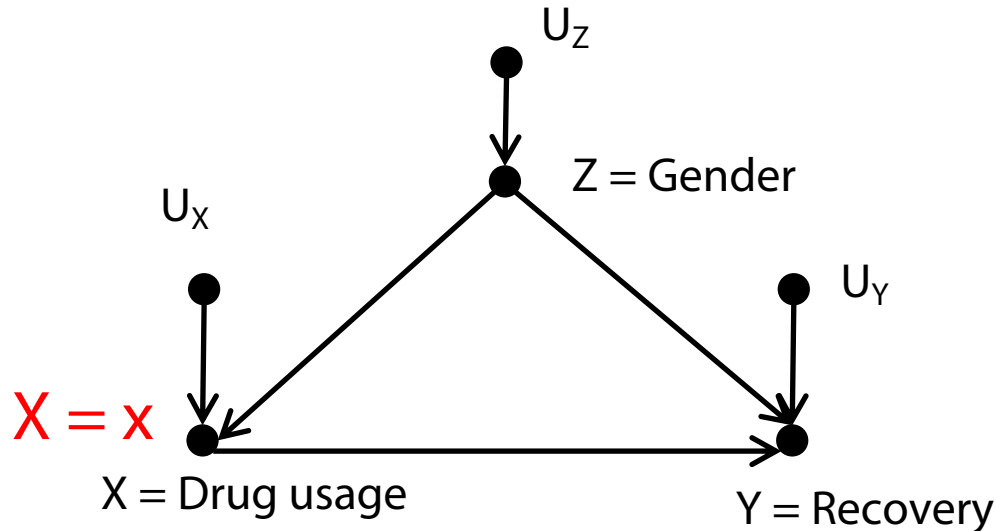
$$P(Y = y \mid \text{do}(X = x)) = P_m(Y = y \mid X = x)$$

= probability in **m**modified graph

# General Causal Effect

## Example (drug-recovery effect)

- How effective is drug usage for recovery?  
 $ACE = P(Y = 1 \mid \text{do}(X = 1)) - P(Y = 1 \mid \text{do}(X = 0))$
- $P(Y = y \mid \text{do}(X = x)) = P_m(Y = y \mid X = x)$



# Intervention (alternatively)

- There are different ways to define intervention (other than by manipulated graph)
- Model intervention  $\text{do}(X=x)$  with force variable  $F$ 
  - $F$  is parent of  $X$ ,
  - $\text{Dom}(F) = \{\text{do}(X=x') \mid x \text{ in } \text{dom}(X)\} \cup \{\text{idle}\}$
  - $\text{pa}'(X) = \text{pa}(X) \cup \{F\}$
  - New ``CPT'' for  $X$

$$P(X=x \mid \text{pa}'(X)) = \begin{cases} P(X=x \mid \text{pa}(X)) & \text{if } F = \text{idle} \\ 0 & \text{if } F = \text{do}(X=x') \text{ and } x \neq x' \\ 1 & \text{if } F = \text{do}(X=x') \text{ and } x = x' \end{cases}$$

$Z$  value not effected by  
intervention on  $x$ :  $f_Z: Z = f(U_Z)$

## Example (drug-recovery effect)

- $P_m(Y = y \mid X = x) = ?$
- Need to reduce to probabilities w.r.t. original graph

1.  $P_m(Z = z) = P(Z = z)$

2.  $P_m(Y = y \mid Z = z, X = x) = P(Y = y \mid Z = z, X = x)$

3. Summing out

$$\begin{aligned} P(Y = y \mid \text{do}(X = x)) &= P_m(Y = y \mid X = x) \\ &= \sum_z P_m(Y = y \mid X = x, Z = z) P_m(Z = z \mid X = x) \\ &= \sum_z P_m(Y = y \mid X = x, Z = z) P_m(Z = z) \\ &= \sum_z P(Y = y \mid X = x, Z = z) P(Z = z) \end{aligned}$$

$Y$  value not effected by intervention  
on  $x$ ,  $f_Y: Y = f(x, z, u_y)$




# Digression

- Conditioning

- $P(Y) = \sum_{z \in Z} P(Y, z) = \sum_{z \in Z} P(Y|z)P(z)$

- $P(Y|X) = P(Y, X) / P(X)$   
 $= \sum_{z \in Z} P(Y, X, z) / P(X)$   
 $= \sum_{z \in Z} P(Y|X, z) P(X, z) / P(X)$   
 $= \sum_{z \in Z} P(Y|X, z) P(z, X) / P(X)$   
 $= \sum_{z \in Z} P(Y|X, z) P(z|X) P(X) / P(X)$   
 $= \sum_{z \in Z} P(Y|X, z) P(z|X)$



Bayes rule is  
your friend

# Adjustment

---

## Definition

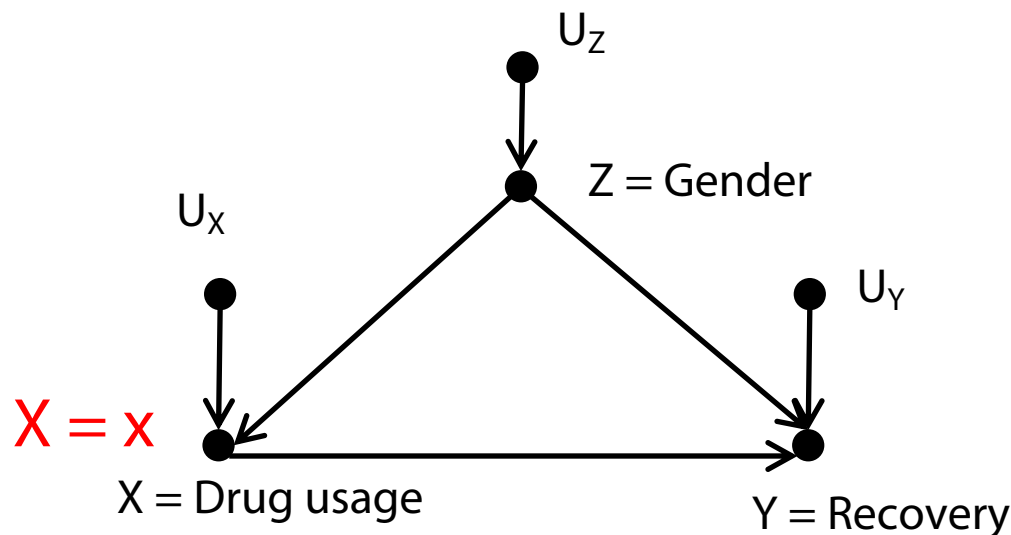
The adjustment formula (for single parent  $Z$  of  $X$ ) for the calculation of the GCE is given by

$$P(Y = y \mid \text{do}(X = x)) = \sum_z P(Y = y \mid X = x, Z=z) P(Z = z)$$

Wording: „Adjusting for  $Z$ “ or „controlling  $Z$ “

# Simpson's Paradox

- How effective is drug usage for recovery?  
 $ACE = P(Y = 1 \mid \text{do}(X = 1)) - P(Y = 1 \mid \text{do}(X = 0))$
- $P(Y = y \mid \text{do}(X = x)) = P_m(Y = y \mid X = x)$



# Recap: Simpson's Paradox

- Record recovery rates of 700 patients given access to a drug

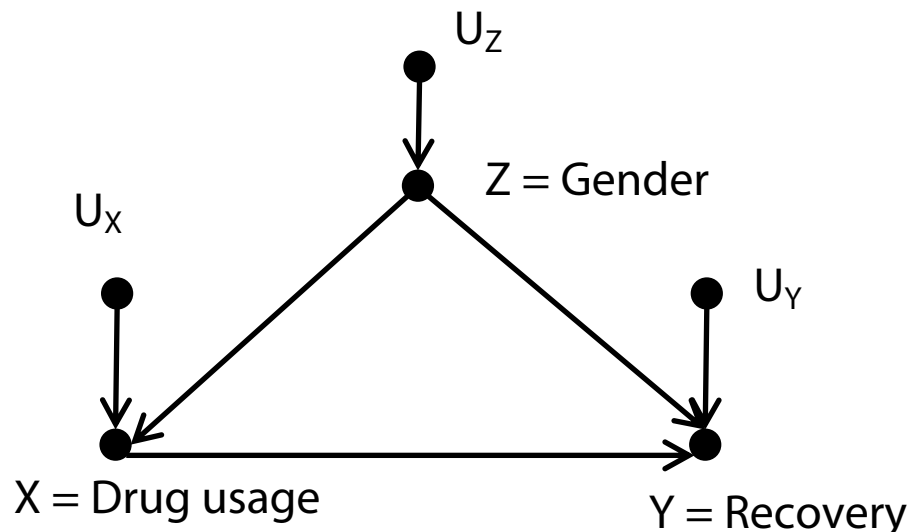
	Recovery rate <b>with</b> drug	Recovery rate <b>without</b> drug
Men	81/87 (93%)	234/270 (87%)
Women	192/263 (73%)	55/80 (69%)
Combined	273/350 (78%)	289/350 (83%)

- Paradox:
  - For men, taking the drug has benefit
  - For women, taking the drug has benefit, too.
  - But: for all persons taking the drug seems to have no benefit



# Resolving the Paradox (Formally)

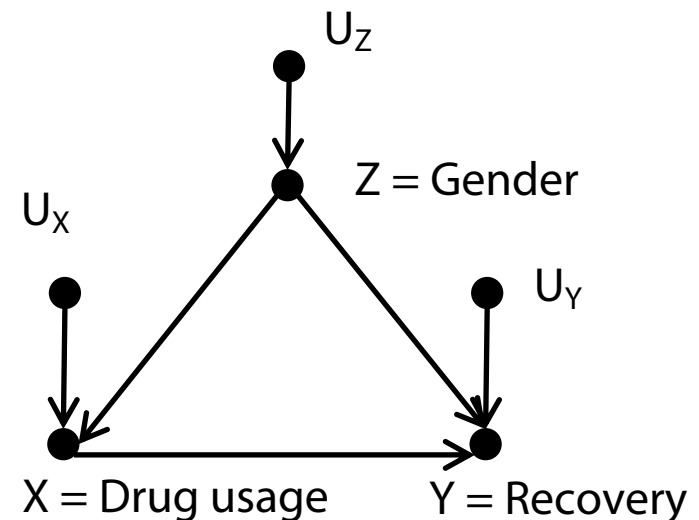
- We have to **understand the causal mechanisms** that lead to the data in order to resolve the paradox
- Formally: What is the general causal effect of drug usage  $X$  on recovery  $Y$ ?
  - $P(Y = y \mid \text{do}(X = x)) = ?$
  - $\text{ACE} = P(Y = 1 \mid \text{do}(X = 1)) - P(Y = 1 \mid \text{do}(X = 0)) = ?$



# Resolving the Paradox (Formally)

- $P(Y = 1 \mid \text{do}(X = 1)) =$  (using adjustment formula)
- $= P(Y = 1 \mid X = 1, Z = 1)P(Z = 1) + P(Y = 1 \mid X = 1, Z = 0)P(Z = 0)$   
 $= 0.93(87 + 270)/700 + 0.73(263 + 80)/700 = 0.832$
- $P(Y = 1 \mid \text{do}(X = 0)) = 0.7818$
- $ACE = 0.832 - 0.7818 = 0.0502 > 0$
- One has to segregate the data w.r.t. Z (adjust for Z)

	Recovery rate <b>with drug</b>	Recovery rate <b>without drug</b>
Men	81/87 (93%)	234/270 (87%)
Women	192/263 (73%)	55/80 (69%)
Combined	273/350 (78%)	289/350 (83%)



# Simpson Paradox (Again)

- Record recovery rates of 700 patients given access to a drug w.r.t. blood pressure (BP) segregation

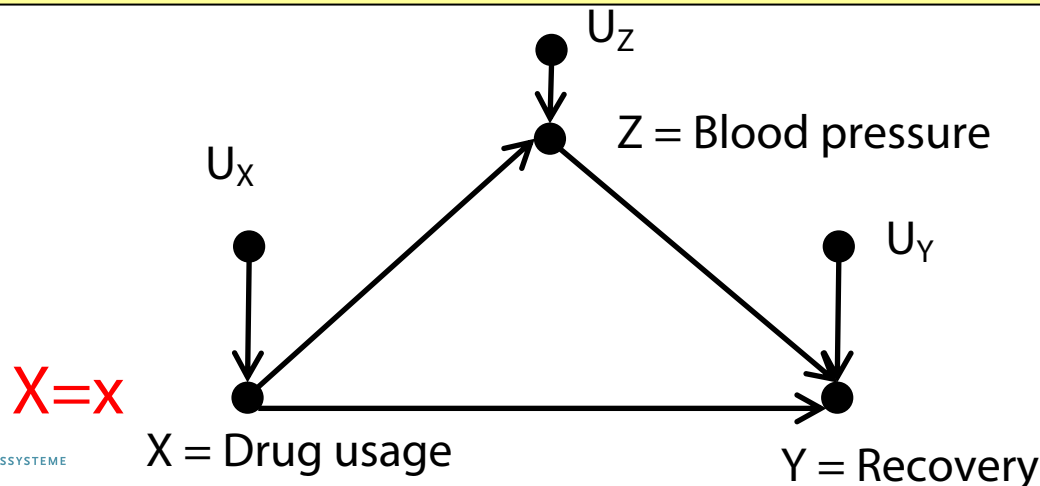
	<b>Recovery rate with drug</b>	<b>Recovery rate without drug</b>
Low BP	234/270 (87%)	81/87 (93%)
High BP	55/80 (69%)	192/263 (73%)
Combined	289/350 (83%)	273/350 (78%)

- BP recorded at end of experiment
- This time segregated data recommends **not** using drug whereas aggregated does

# Resolving the Paradox (Formally)

- We have to **understand the causal mechanisms** that lead to the data in order to resolve the paradox
- Formally: What is the general causal effect of drug usage  $X$  on recovery  $Y$ ?
  - $P(Y = y \mid \text{do}(X = x)) = ?$   
 $= P_m(Y = y \mid X = x) = P(Y = y \mid X = x)$

So: Do not adjust for/segregate w.r.t. any variable



# Causal Effect for Multiple Adjusted Variables

**Rule** (Calculation of causal effect)

$$P(Y = y \mid \text{do}(X = x)) =$$

$$\sum_z P(Y = y \mid X = x, \text{Pa}(X) = z) P(\text{Pa}(X) = z)$$

- $\text{Pa}(X)$  = parents of  $X$
- $z$  = instantiation of all parent variables of  $X$

**Rule** (Calculation of causal effect (alternative))

$$P(Y = y \mid \text{do}(X = x)) =$$

$$\sum_z P(Y = y, X = x, \text{Pa}(X) = z) / P(X = x \mid \text{Pa}(X) = z)$$

# Truncated Product Formula

- Handling of multiple interventions straightforward
- Joint prob. distribution on all other variables  $X_1, \dots, X_n$  after intervention on  $Y_1, \dots, Y_m$

That is, all variables are partitioned in  $X_i$ s and  $Y_j$ s

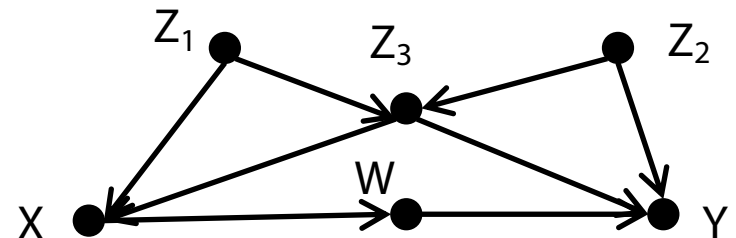
**Definition** (Truncated product formula (g-formula))

$$P(x_1, \dots, x_n \mid \text{do}(Y_1=y_1, \dots, Y_m=y_m)) = \prod_{1 \leq i \leq n} P(x_i \mid \text{pa}(X_i))$$

$\text{pa}(X_i)$  = sub-vector of  $(x_1, \dots, x_n, y_1, \dots, y_m)$  constrained to parents of  $X_i$

## Example 1

$$P(z_1, z_2, w, y \mid \text{do}(X=x, Z_3=z_3)) \\ = P(z_1)P(z_2)P(w|x)P(y|w, z_3, z_2)$$



# Truncated Product Formula

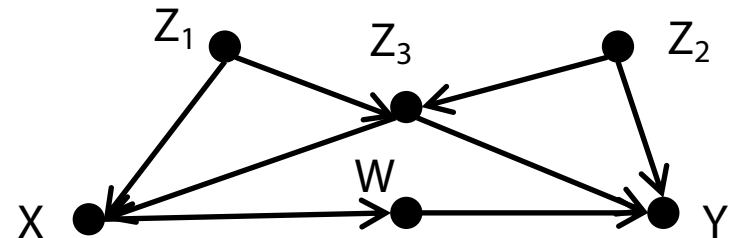
**Definition** (Truncated product formula (g-formula))

$$P(x_1, \dots, x_n \mid \text{do}(Y_1=y_1, \dots, Y_m=y_m)) = \prod_{1 \leq i \leq n} P(x_i \mid \text{pa}(X_i))$$

**Example 2** (summing out)

$$\begin{aligned} &P(w, y \mid \text{do}(X=x, Z_3=z_3)) \\ &= \sum_{z_1, z_2} P(z_1)P(z_2)P(w|x)P(y|w, z_3, z_2) \end{aligned}$$

Can check that this formula is compatible with the adjustment formula



# Backdoor Criterion (Motivation)

---

- Intervention on  $X$  requires adjusting parents of  $X$
- But sometimes those variables are not measurable (though perhaps represented in graph)
- Need more general criterion to identify adjustment variables
  1. Block all spurious paths between  $X$  and  $Y$
  2. Leave all directed paths from  $X$  to  $Y$  unperturbed
  3. Do not create new spurious paths



# Backdoor Criterion (Formulation)

## Definition

Set of variables  $Z$  satisfies **backdoor criterion** relative to a pair  $(X,Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- Can adjust for  $Z$  satisfying backdoor criterion

$$P(Y = y \mid \text{do}(X = x)) = \sum_z P(Y = y \mid X = x, Z = z)P(Z=z)$$

# Backdoor Criterion (Intuition)

## Definition

Set of variables  $Z$  satisfies **backdoor criterion** relative to pair  $(X,Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- Ad 1.: Descendants are effects of  $X$ , should not be conditioned on

(compare drug usage  $X$  and blood pressure  $Z$ )

- Ad 2.: One is interested in effects of  $X$  on  $Y$ , not vice versa. Effects of  $Y$  on  $X$  should be blocked.

# Backdoor Criterion Generalizes Adjustment

## Definition

Set of variables  $Z$  satisfies **backdoor criterion** relative to pair  $(X,Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- $Z = \text{Pa}(X)$
- For any  $W$  in  $Z$  both conditions fulfilled
  - $W$  is not a descendant (as **DAG**)
  - $Z$  blocks every path as every path into  $X$  must go through a parent of  $X$

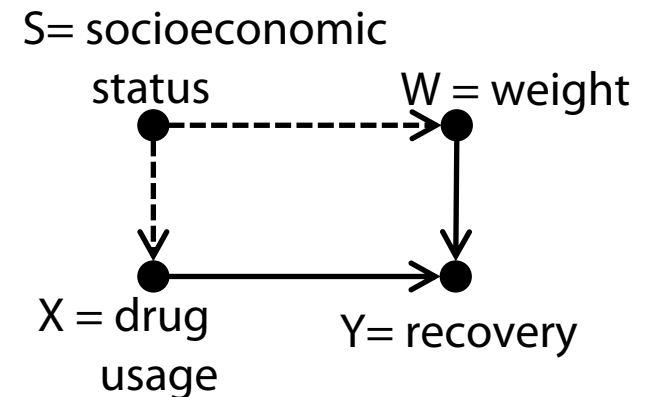
# Backdoor Criterion (Example 1)

## Definition

Set of variables  $Z$  satisfies **backdoor criterion** relative to pair  $(X, Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- Causal effect of  $X$  on  $Y$ ?
- $S$  is not recorded in the data
- $\{W\}$  for  $Z$  fulfills backdoor criterion
  - $W$  not descendant of  $X$
  - Blocks **backdoor path**



# Backdoor Criterion (Example 1 (cont'd))

## Definition

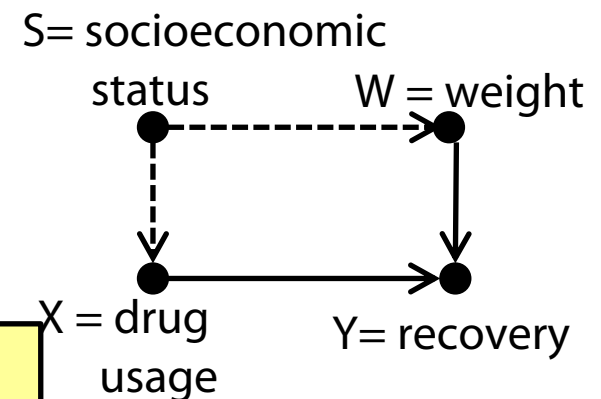
Set of variables  $Z$  satisfies **backdoor criterion** relative to pair  $(X,Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- Causal effect of  $X$  on  $Y$ ?

$$\begin{aligned} P(y \mid \text{do}(x)) &= \sum_w P(Y=y \mid X=x, W=w) P(W=w) \\ &= \sum_s P(Y=y \mid X=x, S=s) P(S=s) \end{aligned}$$

Conditioning on different variables  $S$  vs.  $W$   
with same effect calculation



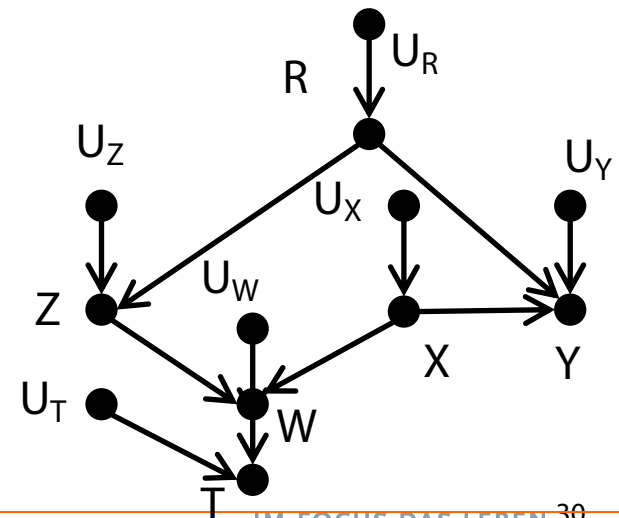
# Backdoor Criterion (Example 2a)

## Definition

Set of variables  $Z$  satisfies **backdoor criterion** relative to pair  $(X, Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- Causal effect of  $X$  on  $Y$ ?
- No backdoor paths
  - Can use  $Z = \{$
  - $P(y \mid \text{do}(x)) = P(y \mid x)$



# Backdoor Criterion (Example 2b)

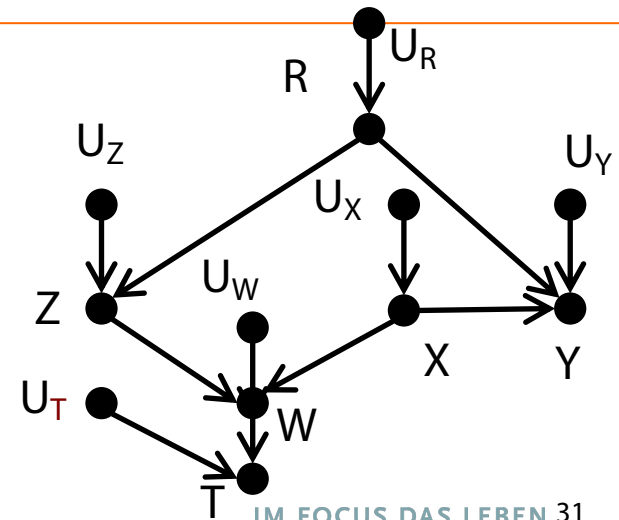
## Definition

Set of variables  $Z$  satisfies **backdoor criterion** relative to pair  $(X, Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- Causal effect of  $X$  on  $Y$ ?
- No backdoor paths
- Can one adjust for  $W$ ?
  - No, then collider  $W$  not blocking

**spurious path**



# Backdoor Criterion (Example 2c)

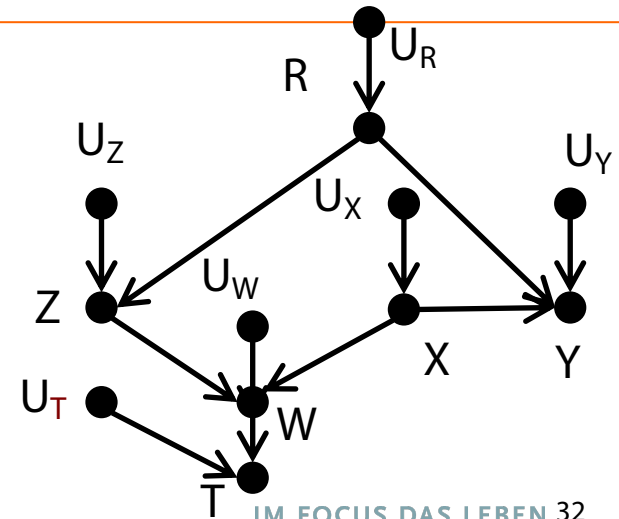
## Definition

Set of variables  $Z$  satisfies **backdoor criterion** relative to pair  $(X,Y)$  of variables iff

1. No node in  $Z$  is a descendant of  $X$  and
2.  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$

- From 2b we know: effect of  $X$  on  $Y$  not via conditioning on  $W$ .
- But how to calculate **w-specific causal effect**:

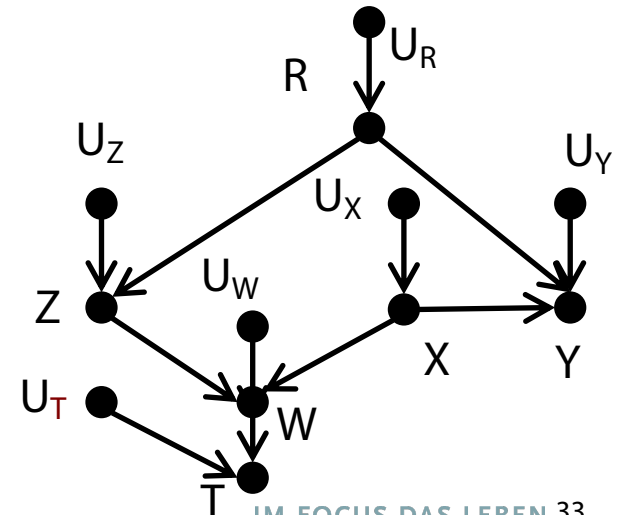
$$P(Y = y \mid \text{do}(X = x), W = w) = ?$$





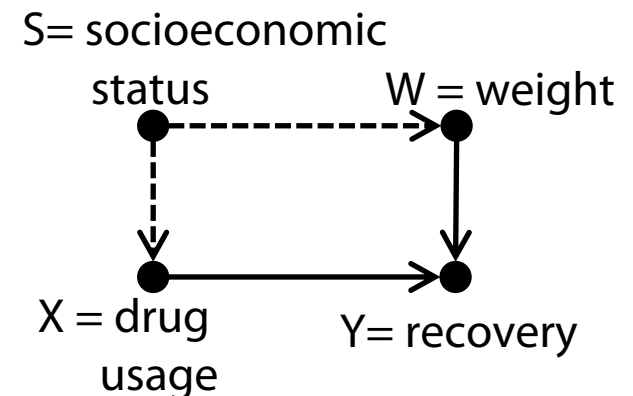
# Backdoor Criterion (Example 2c (cont'd))

- **W-specific causal effect**  $P(Y = y \mid \text{do}(X = x), W = w) = ?$
- Use fork **R** to condition on  
 $P(Y = y \mid \text{do}(X = x), W = w) =$   
$$\sum_r P(Y=y|X=x,W=w,R=r)P(R=r|X=x,W=w)$$
- Degree to which causal effect of **X** on **Y** is modified by values of **W** is called **effect modification** or **moderation**



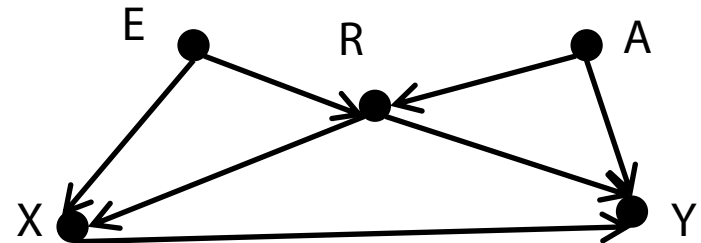
# Backdoor Criterion (Example 3)

- What is effect modification for  $X$  on  $Y$  by  $W$  in drug example?
- Compare  $P(Y = y \mid \text{do}(X = x), W = w)$  and  $P(Y = y \mid \text{do}(X = x), W = w')$
- Here: As  $W$  blocks backdoor
  - $P(Y = y \mid \text{do}(X = x), W = w) = P(Y = y \mid X = x, W = w)$
  - $P(Y = y \mid \text{do}(X = x), W = w') = P(Y = y \mid X = x, W = w')$



# Backdoor Criterion (Example 4)

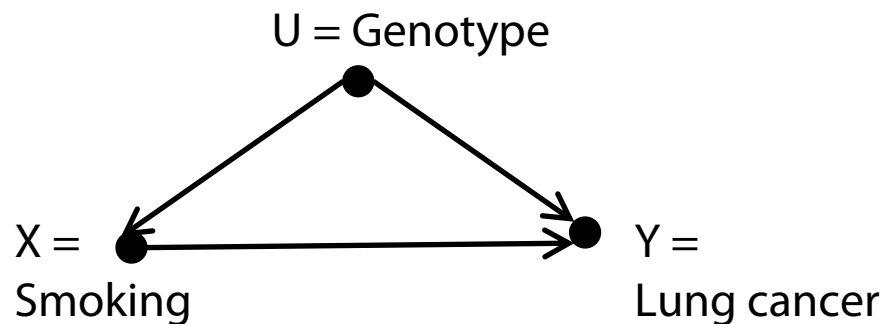
- Sometimes also need to condition on colliders
- There are four backdoor paths from  $X$  to  $Y$ 
  1.  $X \leftarrow E \rightarrow R \rightarrow Y$
  2.  $X \leftarrow E \rightarrow R \leftarrow A \rightarrow Y$
  3.  $X \leftarrow R \rightarrow Y$
  4.  $X \leftarrow R \leftarrow A \rightarrow Y$
- $R$  needed to block 3. path
- But  $R$  collider on 2. path, hence need further blocking variable
- Can use as blocking set  $Z$   
 $\{E, R\}$ ,  $\{R, A\}$  or  $\{E, R, A\}$



# Front-door Criterion (Motivating Example)

## Example

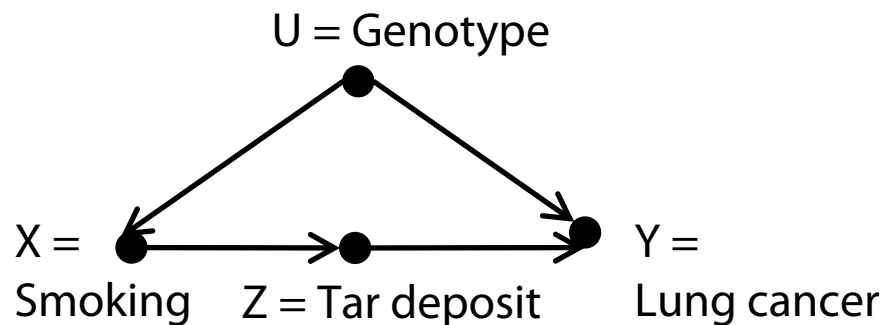
- Sometimes backdoor criterion not applicable
  - $P(y \mid \text{do}(x)) = ?$
  - Genotype  $U$  not observed in data
  - Hence conditioning on  $U$  does not help



# Front-door Criterion (Motivating Example)

## Example

- Sometimes backdoor criterion not applicable
  - $P(y \mid \text{do}(x)) = ?$
  - Genotype  $U$  not observed in data
  - Hence conditioning on  $U$  does not help
  - But sometimes a mediating variable helps



# Front-door Criterion (Motivating Example)

	Tar (400)		No tar (400)		All subjects (800)	
	Smokers (380)	Nonsmokers (20)	Smokers (20)	Nonsmokers (380)	Smokers (400)	Nonsmokers (400)
<b>No cancer</b>	323 (85%)	1 (5%)	18 (90%)	38 (10%)	341 (85%)	39 (9.75%)
<b>Cancer</b>	57 (15%)	19 (95%)	2 (10%)	342 (90%)	59 (15%)	361 (92.25%)

Tobacco industry argues:

- 15% of smoker w/ cancer < 92.25% nonsmoker w/ cancer
- Tar: 15% smoker w/ cancer < 95% nonsmoker w/ cancer
- Non tar: 10% smoker w/ cancer < 90% nonsmoker w/ cancer

# Front-door Criterion (Motivating Example)

	Smokers (400)		Nonsmokers (400)		All subjects (800)	
	Tar (380)	No tar (20)	Tar (20)	No tar (380)	Tar (400)	No tar (400)
<b>No cancer</b>	323 (85%)	18 (90%)	1 (5%)	38 (10%)	324 (81%)	56 (19%)
<b>Cancer</b>	57 (15%)	2 (10%)	19 (95%)	342 (90%)	76 (9%)	344 (81%)

Who is right?

Antismoking lobby argues:

- Choosing to smoke increases chances of tar deposit ( $95\% = 380/400$ )
- Effect of tar deposit: look separately at smokers vs. Non-smokers

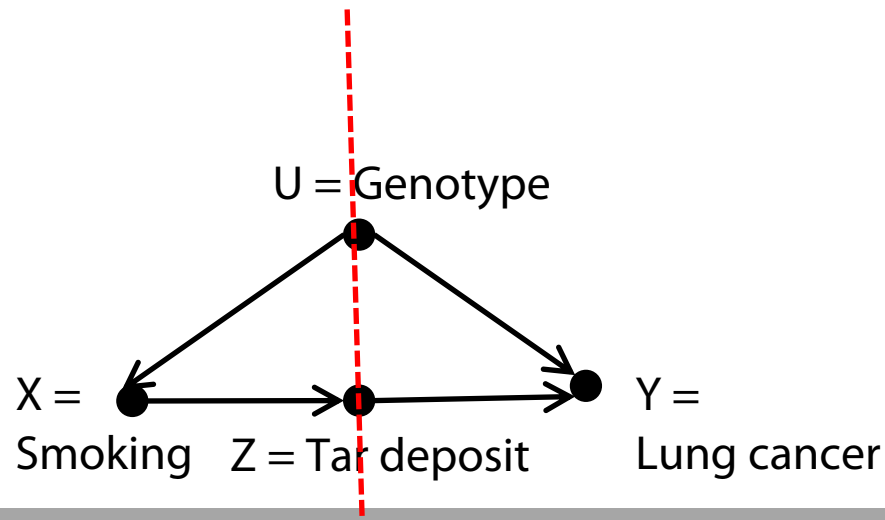
• Smokers: 10 % cancer  $\xrightarrow{+tar}$  15 % cancer

• Nonsmokers: 90 % cancer  $\xrightarrow{+tar}$  95 % cancer

# Front-door Criterion (Intuition)

- Separate effect of  $X$  on  $Y$ :

Effect of  $X$  on  $Y$  = effect of  $X$  on  $Z$  + effect of  $Z$  on  $Y$





# Front-door Criterion (Intuition)

- Effect of  $X$  on  $Z$ :

$$P(Z = z \mid \text{do}(X = x)) = P(Z = z \mid X = x)$$

(No unblocked  
 $X$ - $Z$  backdoor path)

- Effect of  $Z$  on  $Y$ :

$$P(Y = y \mid \text{do}(Z = z)) = \sum_x P(Y = y \mid Z = z, X = x)P(X=x)$$

( $X$  blocks  $Z$ - $Y$ -backdoorpath)

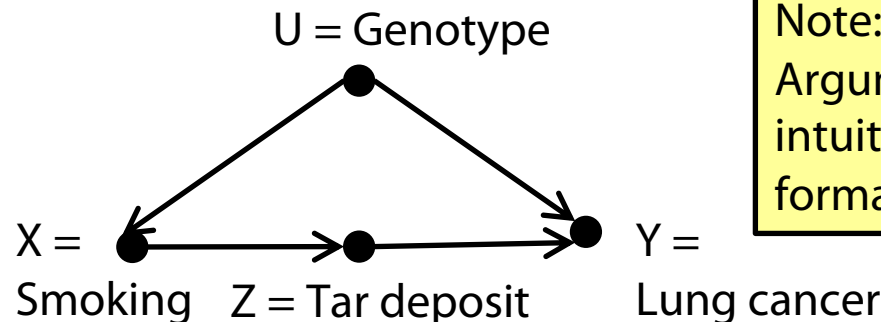
- Effect of  $X$  on  $Y$ :

$$P(Y = y \mid \text{do}(X=x))$$

$$= \sum_z P(Y=y \mid \text{do}(Z=z))P(Z=z \mid \text{do}(X=x))$$

$$= \sum_z \sum_{x'} P(Y=y \mid Z=z, X=x')P(X=x')P(Z=z \mid X=x')$$

(Chaining and summing out)



Note:

Argument in last step rather intuitive. See next slide for formal derivation

# More detailed derivation

$$P(y|\text{do}(X=x))$$

$$= \sum_u P(Y=y|x,u)P(u)$$

(adjustment on  $U$ )

$$= \sum_u \sum_z P(Y=y|z,x,u)P(z|x,u)P(u)$$

(conditioning on  $Z$ )

$$= \sum_u \sum_z P(Y=y|z,x,u)P(z|x)P(u)$$

( $Z$  independent of  $U$

given  $X$  by (d-separation))

$$= \sum_z P(z|x) \sum_u P(Y=y|z,x,u) P(u)$$

(factoring out)

$$= \sum_z P(z|x) \sum_u P(Y=y|z,u) P(u)$$

( $Y$  independent of  $X$  given  $Z,U$ )

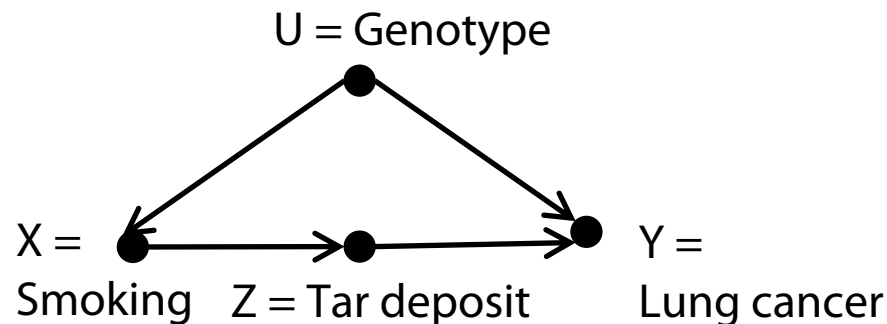
$$= \sum_z P(z|x)P(Y|\text{do}(z))$$

(definition of  $\text{do}()$ )

$$= \sum_z P(z|x) \sum_{x'} P(Y|x',z) P(x')$$

(adjustment via  $X$ )

$$= \sum_z \sum_{x'} P(z|x) P(Y|x',z) P(x')$$



# Front-door Criterion (Formulation & Theorem)

## Definition

Set of variables  $Z$  satisfies front-door criterion w.r.t. pair of variables  $(X,Y)$  iff

1.  $Z$  intercepts all directed paths from  $X$  to  $Y$
2. Every backdoor path from  $X$  to  $Z$  is blocked (by collider)
3. All  $Z$ - $Y$  backdoor paths are blocked by  $X$

## Theorem (Front-door adjustment)

If  $Z$  fulfills front-door criterion w.r.t.  $(X,Y)$  and  $P(x,z) > 0$   
then  $P(y|\text{do}(x)) = \sum_z P(z|x) \sum_{x'} P(y|z, x')P(x')$

# Conditional Interventions (Example)

## **Example** (conditioned drug administering)

- Administer drug ( $X = 1$ ) if fever  $Z > z$

- Formally:

$$P(Y = y \mid \text{do}(X = g(Z)))$$

where  $g(Z) = 1$  if  $Z > z$  and  $g(Z) = 0$  otherwise

- Can be reduced to calculating **z-specific effect**

$$P(Y = y \mid \text{do}(X = x), Z = z)$$

# Conditional Interventions (Rule)

## **Rule** (z-specific effect)

If there is set  $S$  of variables s.t.  $S \cup Z$  satisfies  
backdoor criterion

then the z-specific effect is given by

$$P(y \mid \text{do}(x), z) = \sum_s P(y \mid x, s, z) P(s \mid z)$$

Reduction of conditional intervention to z-specific effect:

$$P(Y = y \mid \text{do}(X = g(Z))) =$$

$$= \sum_z P(Y = y \mid \text{do}(X = g(Z), Z = z) P(Z = z \mid \text{do}(X = g(Z)))$$

(conditioning on  $Z$ )

$$= \sum_z P(Y = y \mid \text{do}(X = g(Z), Z = z) P(Z = z)$$

( $Z$  before  $X$ )

$$= \sum_z P(Y = y \mid \text{do}(X = x), z)_{|x=g(z)} P(Z = z)$$

# Intervention Calculation in Practice?

## JMHUEBNER'S BLOG

Just another WordPress.com

(GCE) calculation by intervention useful as long as (domains of) conditioned variable set  $Z$  and values small (i.e., few summations)

## Theory VS Practice



"In theory, there is no difference between theory and practice.

But in practice, there is." Jan L.A. van de Snepschaut

# Inverse Probability Weighting

---

- Inverse probability weighting gives estimation of GCE on small sample size  $\ll |z|$
- Estimation with propensity score  $P(X=x|Z=z)$ 
  - Propensity score can be estimated similarly as in linear regression
  - Weight **small** sample set with propensity
  - Estimation of  $P(y|do(x))$   
by counting all events for  $y$  for each stratum  $X = x$   
(No summation over all instances of  $Z$  required)

# Inverse Probability Weighting

- Filtering-Case  $P(Y=y, Z=z | X=x)$ : Evidence leads to re-normalization of full joint probability
  - $P(Y=y, Z=z | X=x) = P(Y=y, Z=z, X=x) / P(X=x)$
  - Have to weight  $(Y, Z, X)$  samples by  $1/P(X=x)$
- Intervention-Case  $P(y | \text{do}(x))$ : Weighting by propensity
  - $P(y | \text{do}(x))$ 
    - $= \sum_z P(Y=y | X=x, Z=z) P(Z=z)$
    - $= \sum_z P(Y=y | X=x, Z=z) P(Z=z) P(X=x | Z=z) / P(X=x | Z=z)$
    - $= \sum_z P(X=x, Y=y, Z=z) / P(X=x | Z=z)$

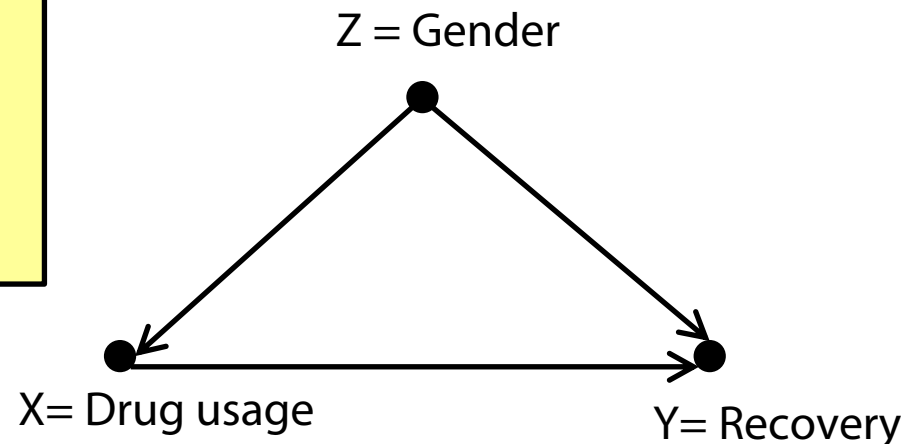
Weighting joint distribution by inverse propensity



# Inverse Probability Weighting (Example)

	Recovery rate <b>with</b> drug	Recovery rate <b>without</b> drug
Men	81/87 (93%)	234/270 (87%)
Women	192/263 (73%)	55/80 (69%)
Combined	273/350 (78%)	289/350 (83%)

- Rewrite table to get % of population for each (X,Y,Z) instance
- Example:  
 $\%(yes, yes, male) = 81/700 = 0.116$



# Sample percentages

	Recovery rate <b>with drug</b>	Recovery rate <b>without drug</b>
Men	81/87 (93%)	234/270 (87%)
Women	192/263 (73%)	55/80 (69%)
Combined	273/350 (78%)	289/350 (83%)

X	Y	Z	% of population
yes	yes	male	0.116
yes	yes	female	0.274
yes	no	male	0.01
yes	no	female	0.101
no	yes	male	0.334
no	yes	female	0.079
no	no	male	0.051
no	no	female	0.036

# Weighting when Filtering for X=yes

X	Y	Z	% of population
yes	yes	male	0.116
yes	yes	female	0.274
yes	no	male	0.01
yes	no	female	0.101
no	yes	male	0.334
no	yes	female	0.079
no	no	male	0.051
no	no	female	0.036

Consider  $X = \text{yes}$  & weight  $(X,Y,Z)$  with  $1/P(X=\text{yes}) = 1/(0.116+0.274+0.01+0.101)$

X	Y	Z	% of population
yes	yes	male	0.232
yes	yes	female	0.547
yes	no	male	0.02
yes	no	female	0.202

# Weighting when Intervening $\text{do}(X=\text{yes})$

X	Y	Z	% of population
yes	yes	male	0.116
yes	yes	female	0.274
yes	no	male	0.01
yes	no	female	0.101
no	yes	male	0.334
no	yes	female	0.079
no	no	male	0.051
no	no	female	0.036

Consider  $X = \text{yes}$  & weight  $(X,Y,Z)$  with  $1/P(X=\text{yes}|Z=z)$   
 $P(X=\text{yes}|Z=\text{male}) = (0.116 + 0.01)/(0.116+0.01 + 0.334 + 0.051)$   
 $P(X=\text{yes}|Z=\text{female}) = (0.274 + 0.101)/(0.274+0.101 + 0.079 + 0.036)$

In this example no real savings!  
 These come into play when  
 $\text{dom}(Z) \gg \text{sample size}$

X	Y	Z	% of population
yes	yes	male	0.476
yes	yes	female	0.357
yes	no	male	0.042
yes	no	female	0.132