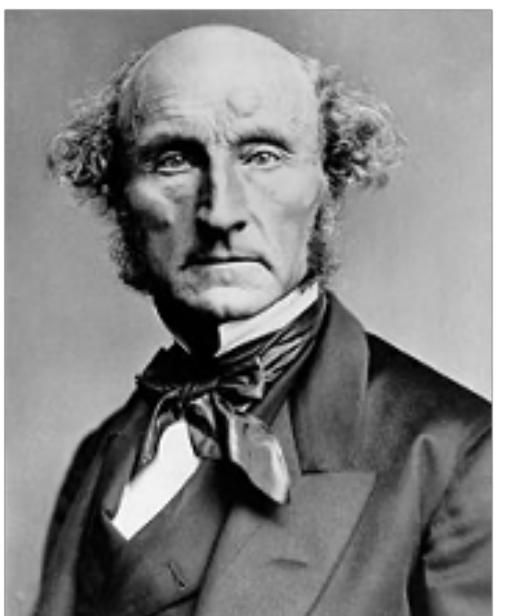


We are halfway done!

**Mar 14**

# How can we give scientific explanations for observed data?

We need **a model** to explain what we observed and tested.

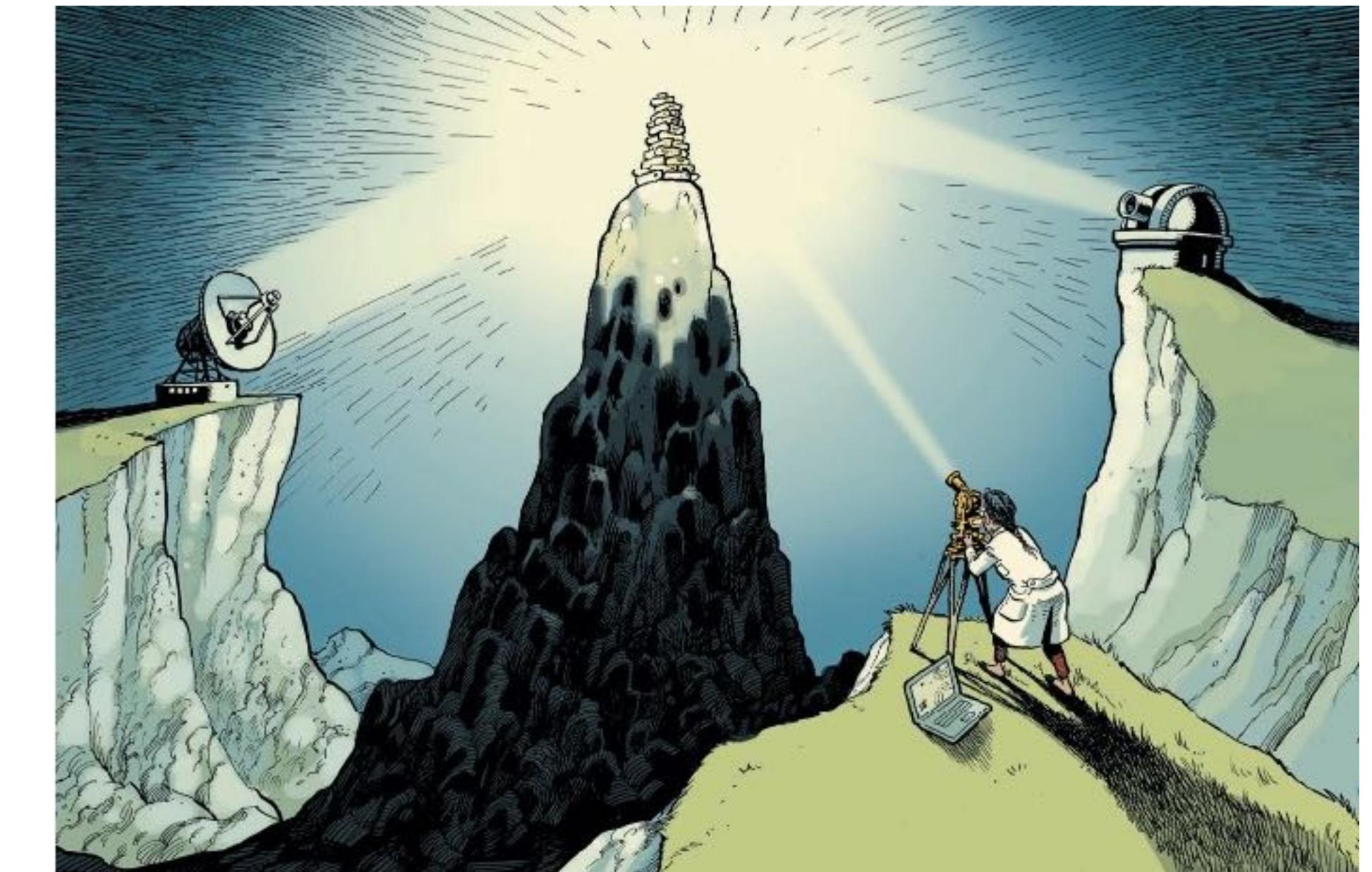


John Stuart Mill

## SECOND CANON.

*If an instance in which the phenomenon under investigation occurs, and an instance in which it does not occur, have every circumstance save one in common, that one occurring only in the former; the circumstance in which alone the two instances differ, is the effect, or cause, or a necessary part of the cause, of the phenomenon.*

JS Mill, A system of Logic (1843)



Munafo & Davey Smith, "Repeating Experiments is not enough" Nature (2018)



Peter Lipton

Contrastive Explanation & causal triangulation, Philosophy of Science (1991)



George Davey Smith

# Today's lecture: Bayesian, PGM, Causality

- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

# **Statistical Methods for High-dimensional Biology**



## **Bayesian Inference & Probabilistic Graphical Models, Causal Inference Tools**

Yongjin Park, UBC Path&Lab, STAT, BC Cancer

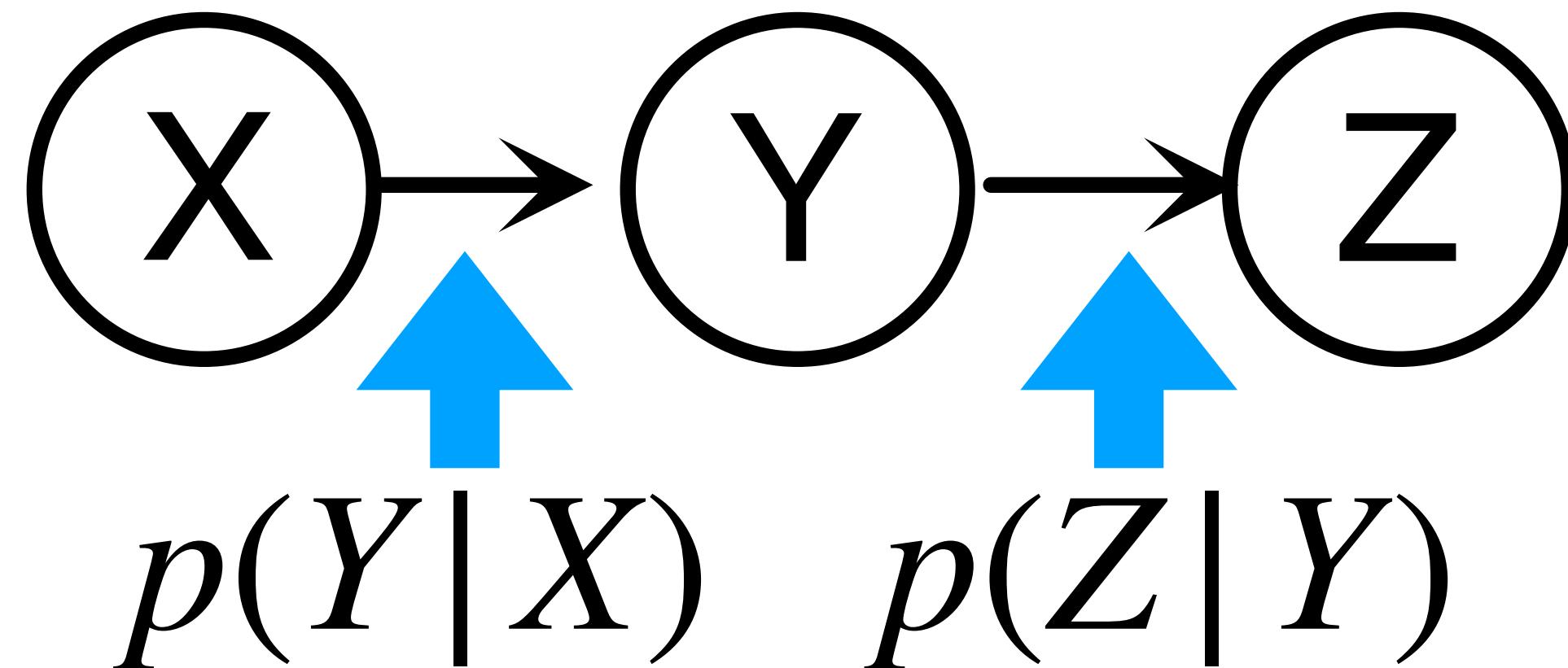
# **What will be a good language to describe our model, hypothesis, belief, or assumptions?**

- It must be easy to understand (non-mathematical audience should be able to grasp the meaning intuitively)
- It must generalize well across many different scientific fields
- It can represent the notion of uncertainty
- It must allow a general-purpose computer algorithm to simulate and compute.

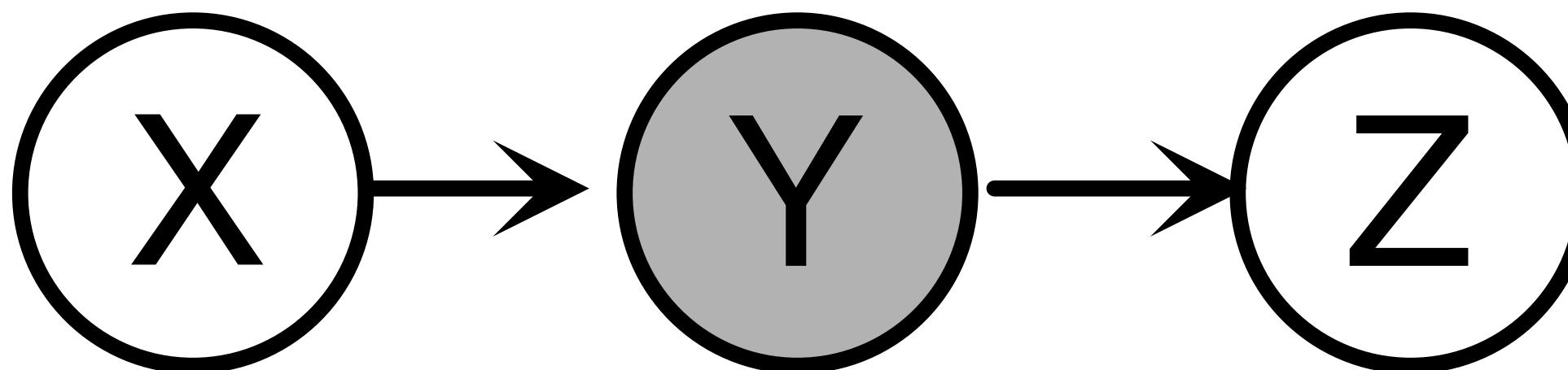
# Today's lecture: Bayesian, PGM, Causality

- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

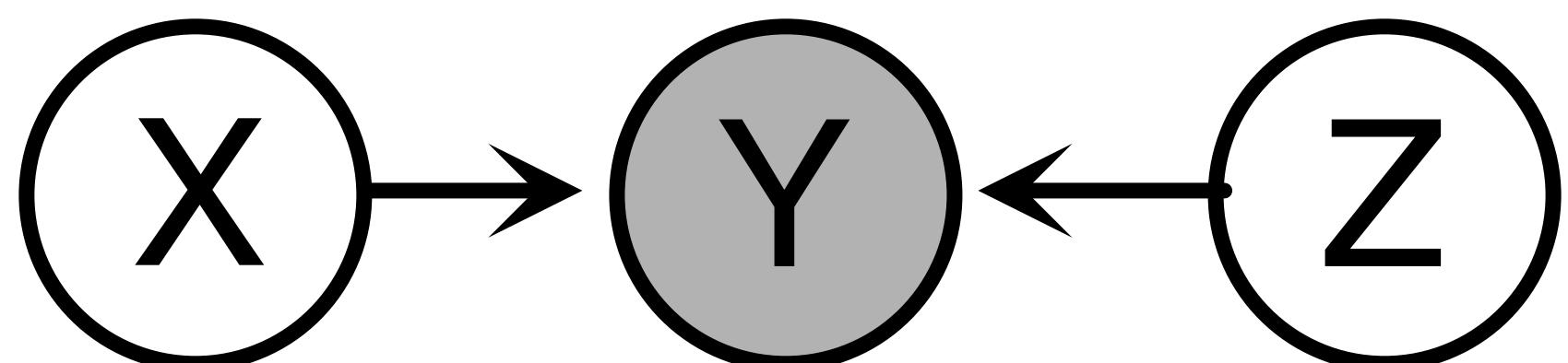
# A graphical language for a probabilistic model



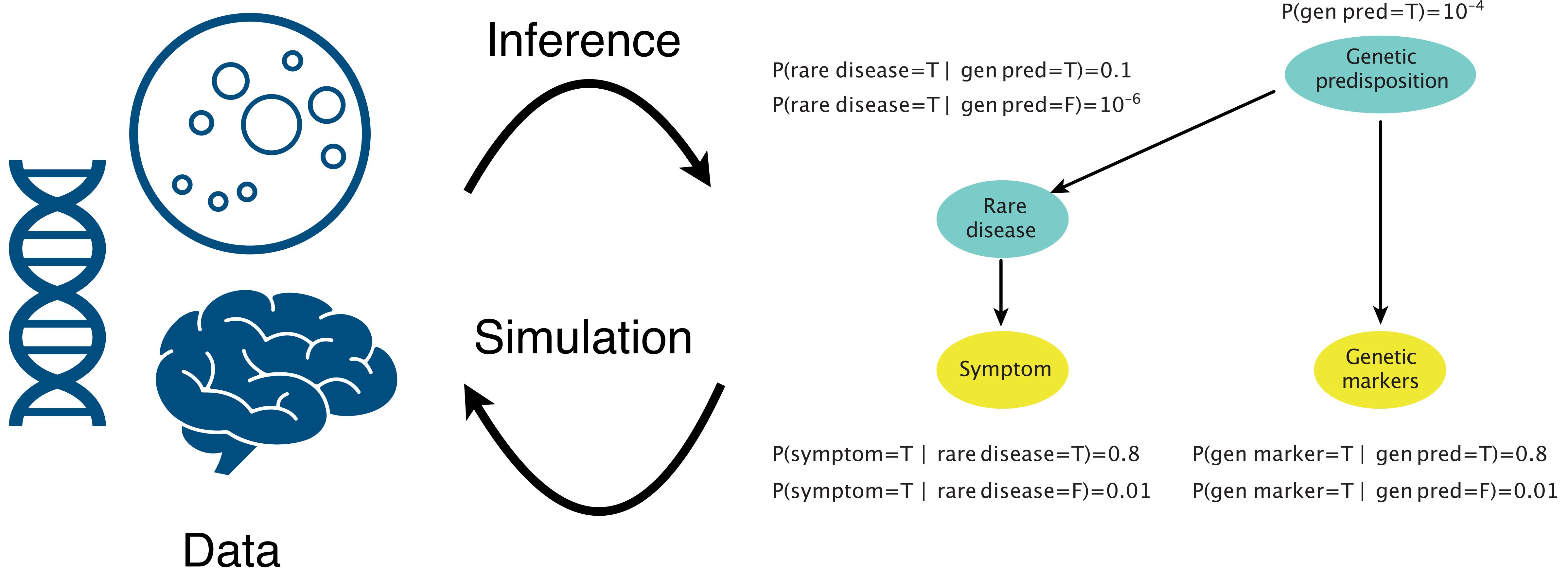
- Each node: a random variable
- Arrow/edge: dependency, conditional probability
- Shaded: conditioned/observed
- Open: unobserved/unknown (yet)



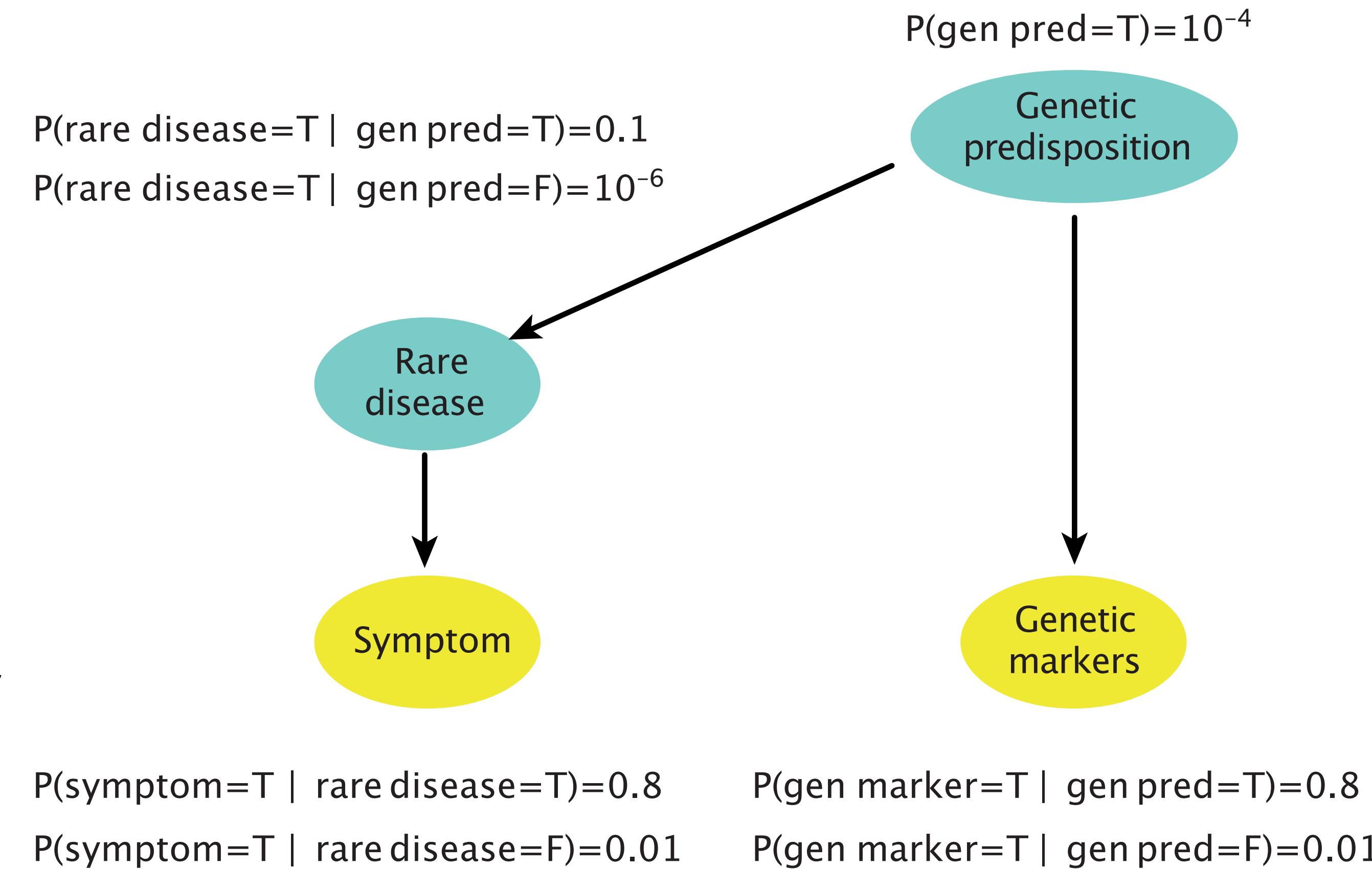
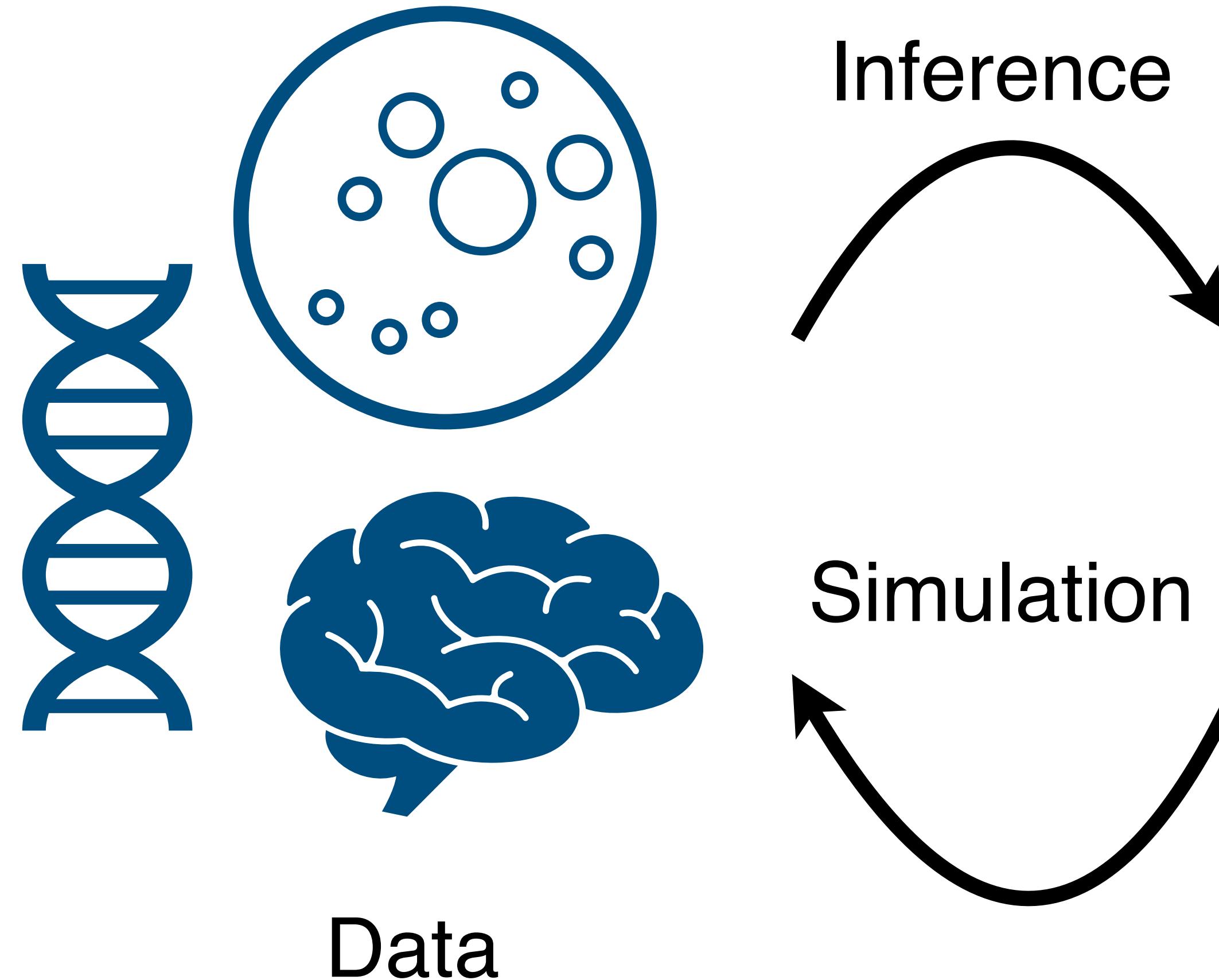
# Almost all the joint probabilities (parametric models) can be described by a graph language



# Modelling = synthesizing conditional probabilities!



# Bayesian inference is a model-based approach



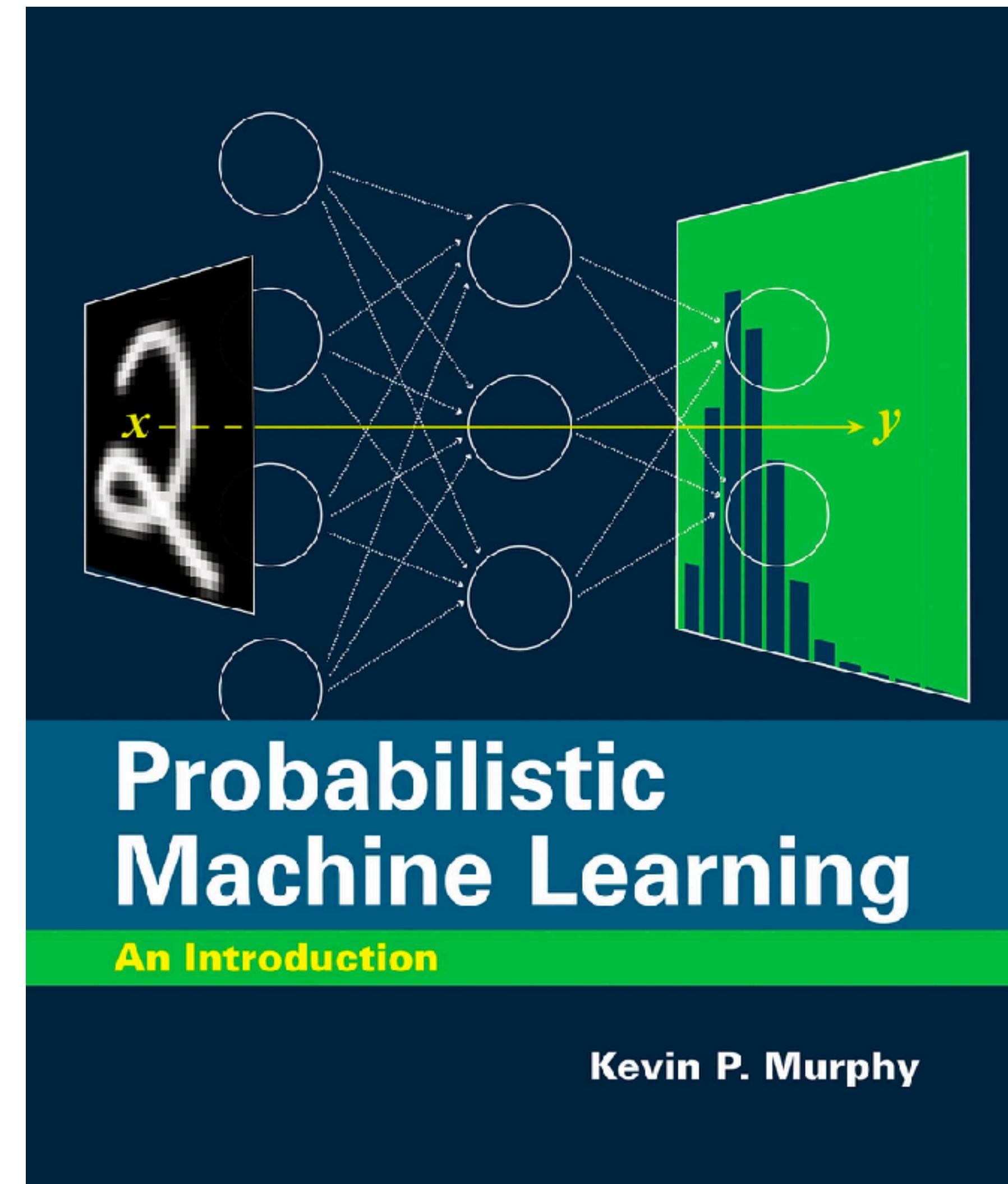
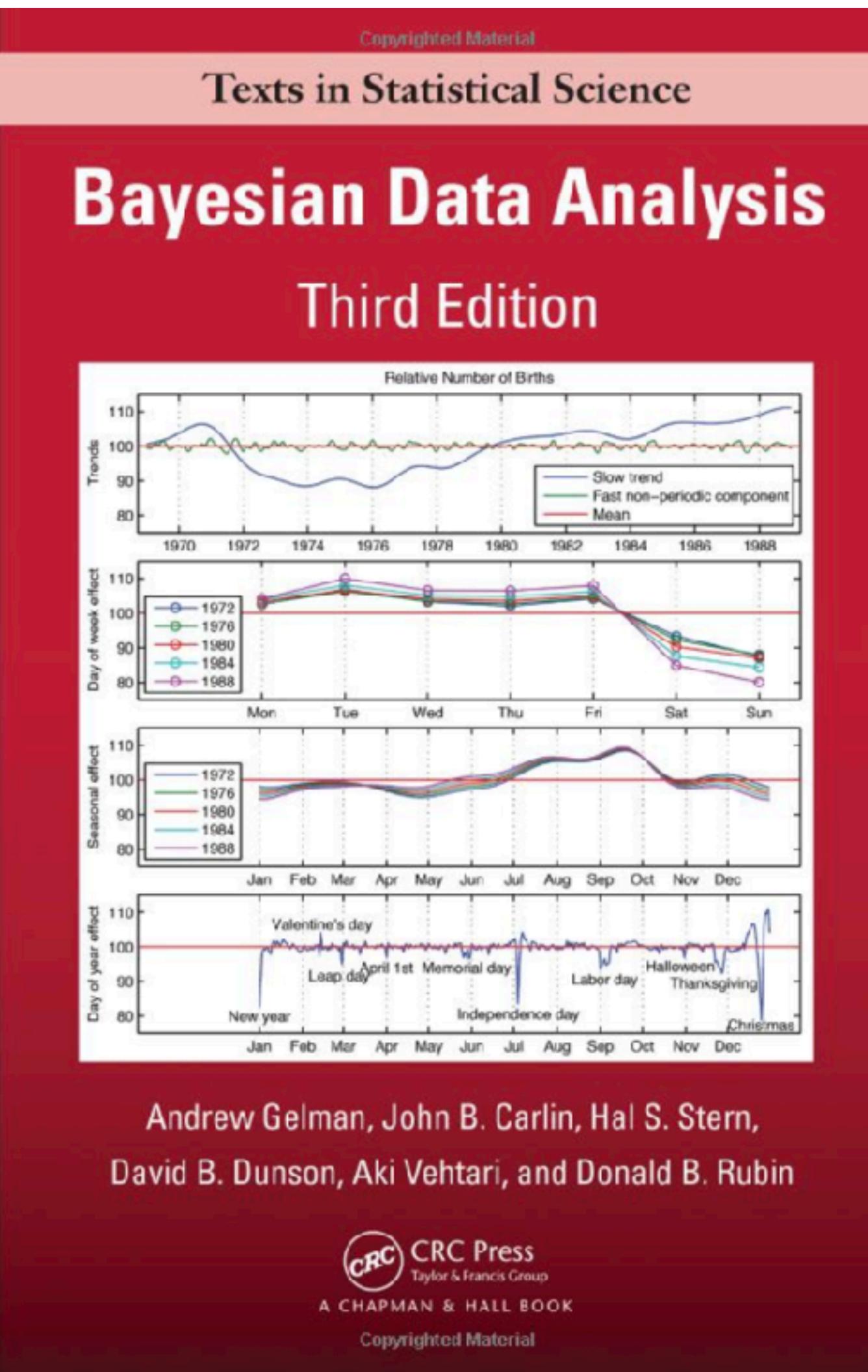
# Source code available!

The screenshot shows a GitHub repository page for 'STAT540-UBC / lectures'. The repository is public, has 1 fork, 2 stars, and 0 contributors. The 'Code' tab is selected. A specific file, 'bayesian\_stan.Rmd', is shown. It was last updated 2 days ago by Yongjin Park. The file contains 660 lines (455 sloc) and is 17 KB. The code content is as follows:

```
1 ---  
2 title: "Bayesian Inference"  
3 author: "Yongjin Park"  
4 classoption: "aspectratio=169"  
5 ---
```

[https://github.com/STAT540-UBC/lectures/blob/main/lect13-causality\\_bayesian/bayesian\\_stan.Rmd](https://github.com/STAT540-UBC/lectures/blob/main/lect13-causality_bayesian/bayesian_stan.Rmd)

# If you're more interested in Bayesian statistics and probabilistic graphical models



# Today's lecture: Bayesian, PGM, Causality

- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

## A warm-up example: the mean and variance of log-Normal

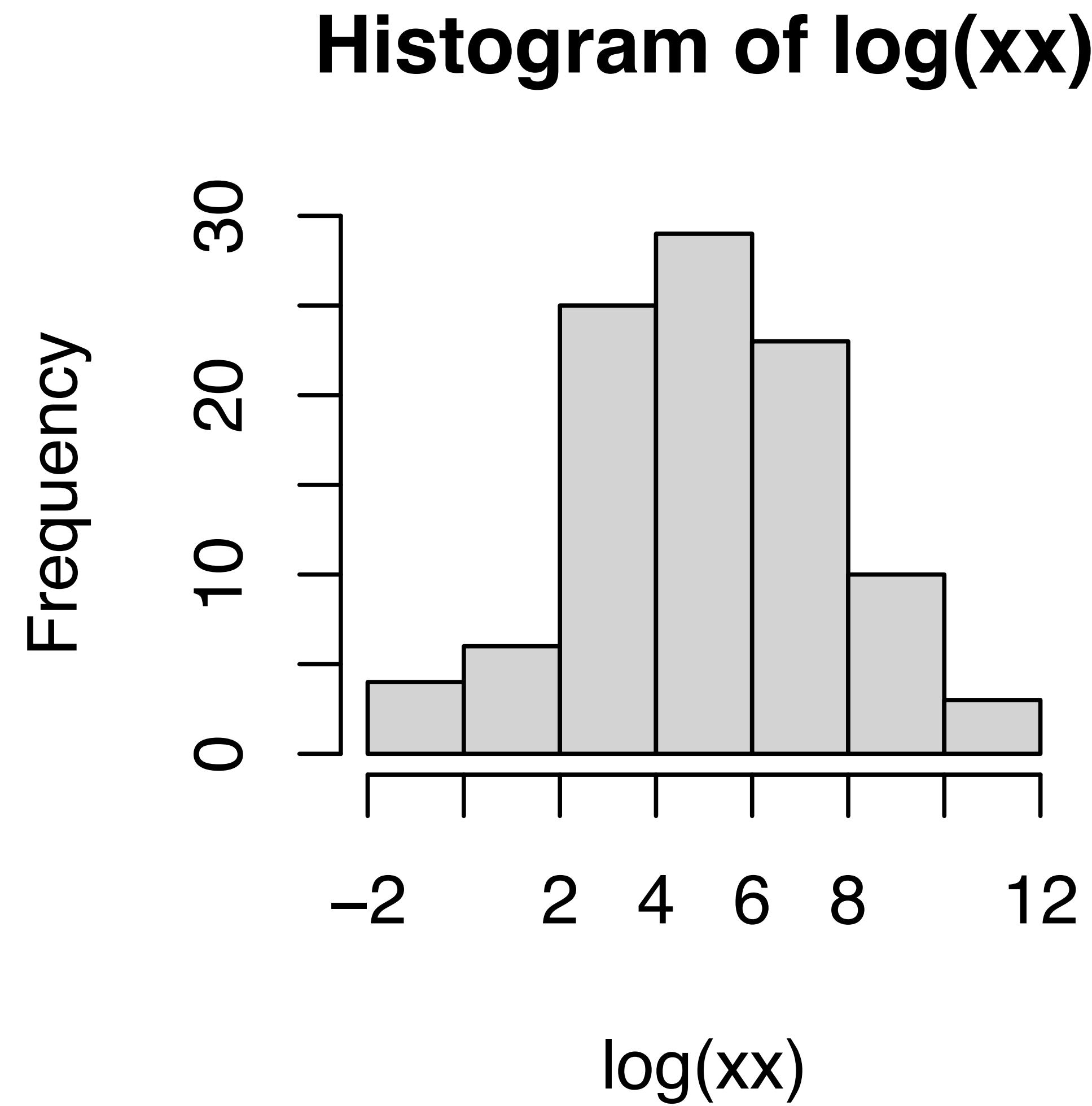
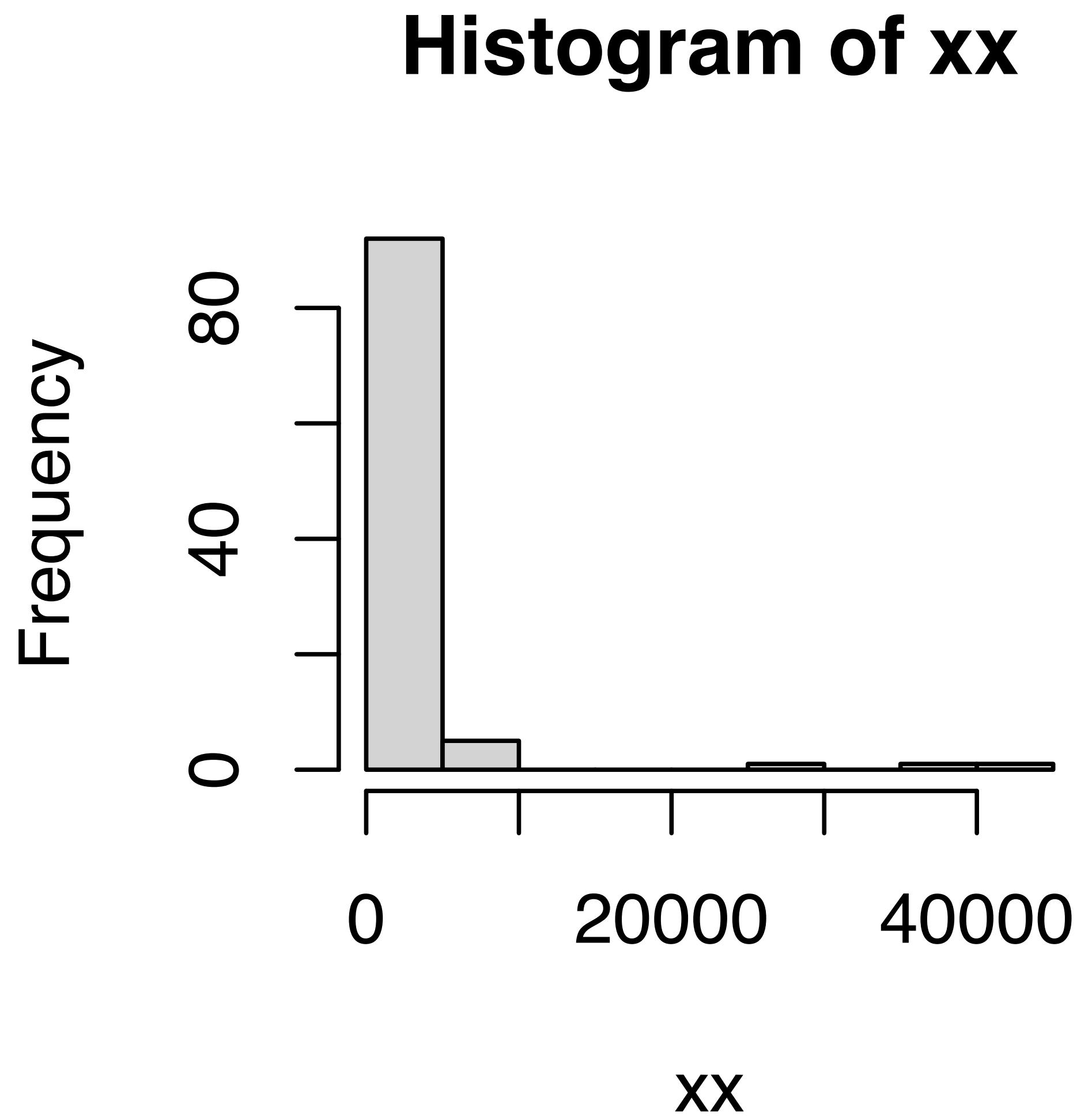
Suppose we have a sequence of real-numbered experimental data observed:  $x_1, x_2, \dots$

- ▶ A previous graduate student told me that she calibrated the device, so that the numbers generally follow log-Normal distribution.
- ▶ However, she rushed to join some cool biotech start-up and did not inform us what were the desired mean and error parameters.
- ▶ Your PI desperately wants to know about the parameters to complete the Aim 3 in recent grant applications.

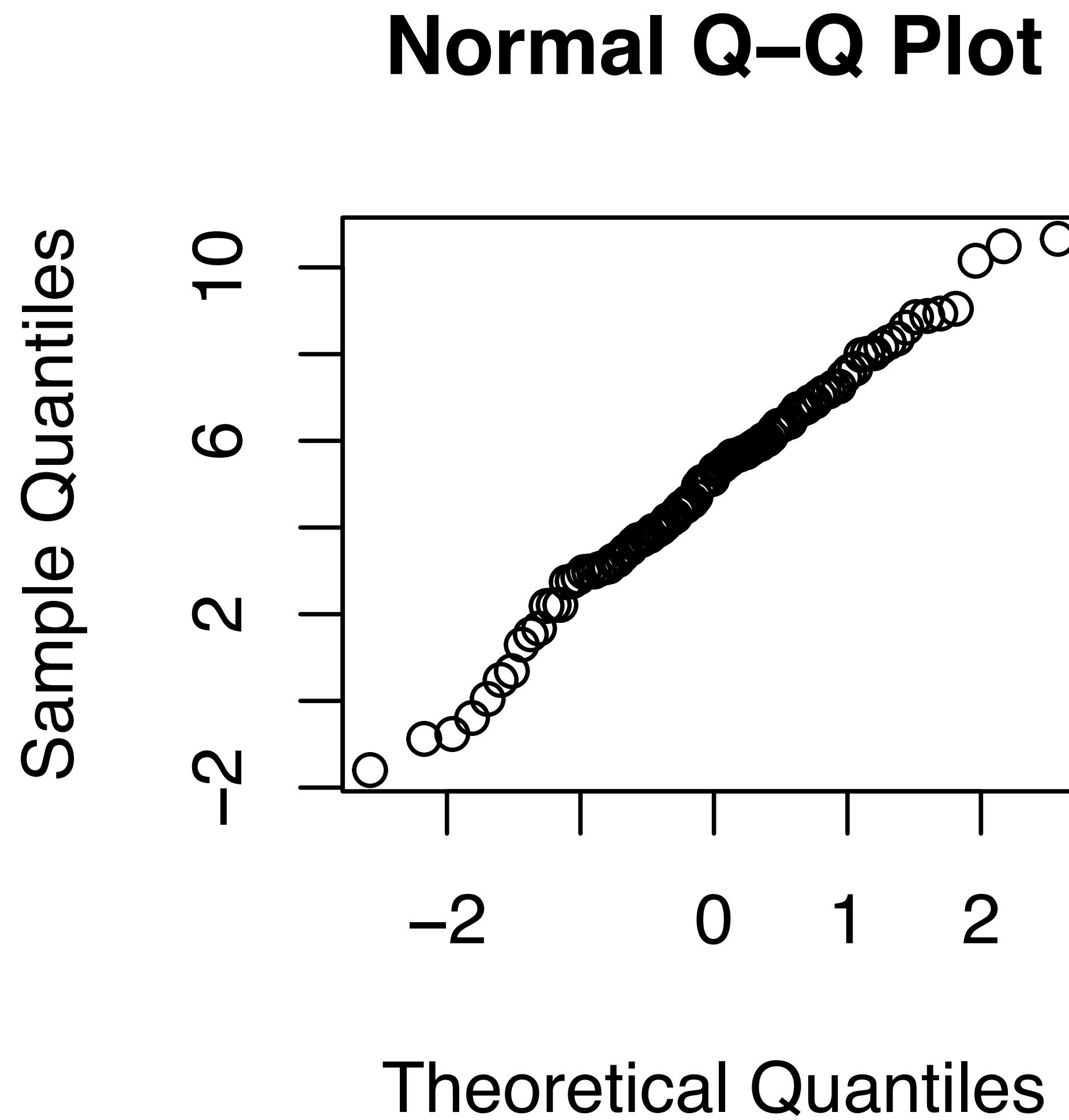
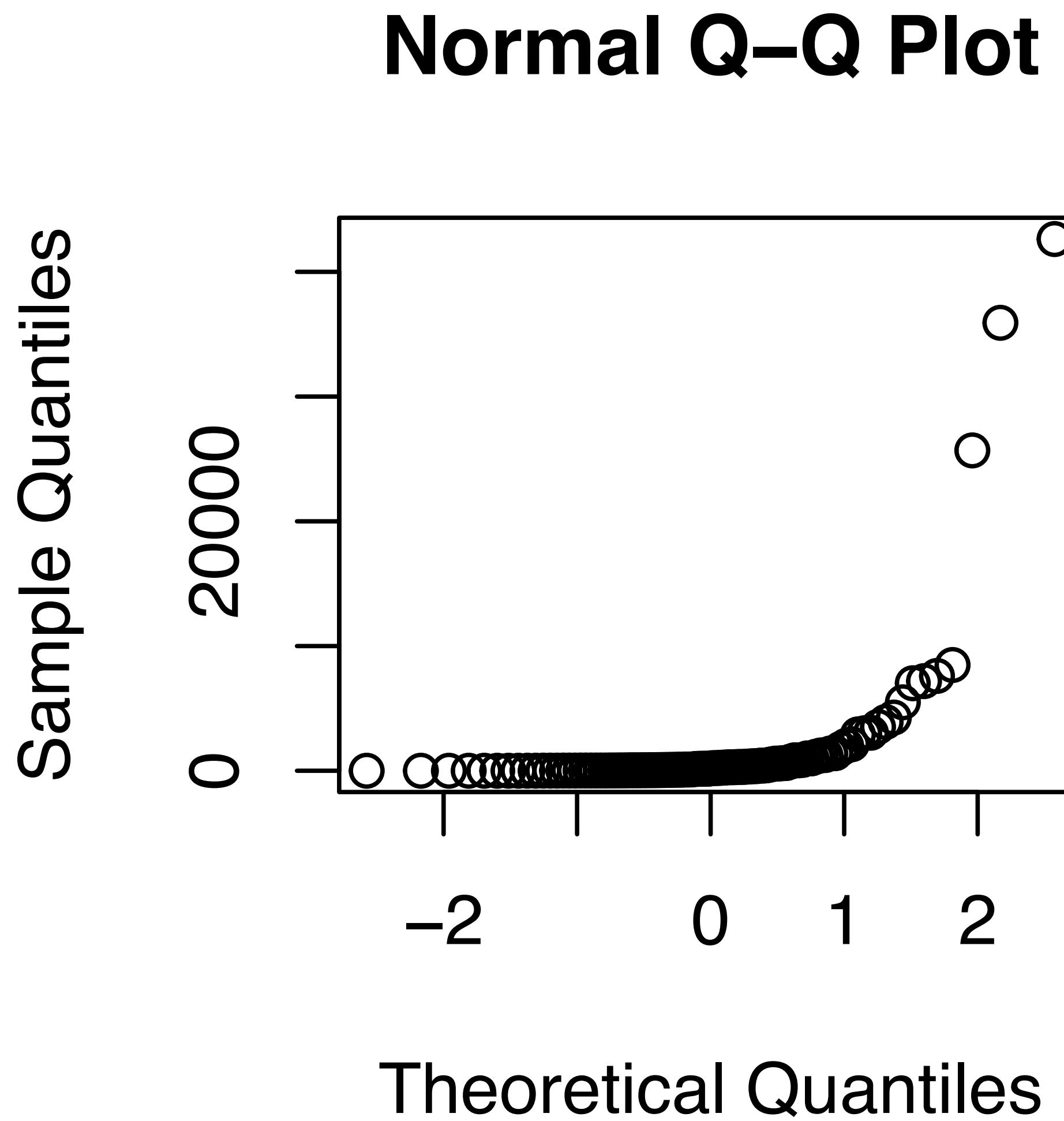
## The actual generative scheme

```
set.seed(1331)
nn <- 100
xx <- rlnorm(nn, 5, sqrt(7))
```

First, what do they look like?



First, what do they look like?



# Can we infer the mean and variance of the observed data?

Now we are wondering what will be the suitable parameters...

To answer this question, we can turn to statistics.

Let's define data likelihood:

$$P(\{x_i\} | \mu, \gamma) = \prod_{i=1}^n \log \mathcal{N}(x_i | \mu, \gamma^{-1})$$

A **Frequentist** Question/Approach:

- ▶ What are the **most likely**  $\mu$  and  $\gamma$  values?
- ▶ **How confident** are we about the estimated  $\hat{\mu}$  and  $\hat{\gamma}$ ?
  - ▶ Standard error (error bars)
  - ▶ p-value of some hypothesis test (rejecting the null)

# Slightly different world views: Frequentist vs. Bayesian

## Frequentist inference

- ▶ If  $n \rightarrow \infty$ , how confident are you about the estimated parameter  $\hat{\theta}$ ?
- ▶ Will a true parameter  $\theta$  (theoretical quantity) be consistently captured around a confidence interval around the estimated parameter  $\hat{\theta}$ ?
- ▶ If so, how “frequently” the constructed CI will include the true parameter as  $n \rightarrow \infty$ ?

$$p(\theta \in (\hat{\theta}_n - 2\hat{s}\hat{e}_n, \hat{\theta}_n + 2\hat{s}\hat{e}_n))$$

## Bayesian inference

- ▶ Given the data observed ( $n < \infty$ ), what is the distribution of the unknown parameter?
- ▶ What was the underlying generative model?
- ▶ Can we take into account uncertainty across multiple types of models?
- ▶ How likely is this new observation  $x^*$  given the other observed data points?

$$p(\theta | \{i \in [n] : x_i\})$$

# Frequentist Inference

What can I say about  
this parameter (model)  $\theta$ ,  
given plenty of  
**unseen data?**

Can I give the same,  
consistent  $\hat{\theta}$ ,  
given plenty of  
**unseen data?**



# Bayesian Inference

<https://unsplash.com/photos/77AW8rM9KGg>

**Based on what I've seen so far,  
the parameter  $\theta$  is about here...**

# A Bayesian way to know about the “how confident” question (given data)

- ▶ Given data  $\{x_i\}$ , what is the probability of the parameters (posterior probability)?

$$p(\underbrace{\mu, \gamma}_{\text{unknown parameter}} \mid \underbrace{\{x_i\}}_{\text{data}})$$

- ▶ Eventually we also want to predict future outcomes (posterior prediction),

$$p(\underbrace{x^*}_{\text{future data}} \mid \underbrace{\{x_i\}}_{\text{data}}) = \int d\mu d\gamma \overbrace{p(\underbrace{x^*}_{\text{new observation}} \mid \mu, \gamma)}^{\text{data generating}} \overbrace{p(\mu, \gamma \mid \underbrace{\{x_i\}}_{\text{training data}})}^{\text{posterior probability}}$$

- ▶ Moreover, we want to find a better model comparing the model evidence:

$$p(\text{data} \mid \text{model 1}) \text{ vs. } p(\text{data} \mid \text{model 2})$$

# A Bayesian way to know about the “how confident” question (given data)

- ▶ Given data  $\{x_i\}$ , what is the probability of the parameters (posterior probability)?

$$p(\underbrace{\mu, \gamma}_{\text{unknown parameter}} \mid \underbrace{\{x_i\}}_{\text{data}})$$

- ▶ Eventually we also want to predict future outcomes (posterior prediction),

$$p(\underbrace{x^*}_{\text{future}} \mid \underbrace{\{x_i\}}_{\text{data}}) = \int d\mu d\gamma p(x^* \mid \underbrace{\mu, \gamma}_{\text{averaging over uncertainty}}) p(\mu, \gamma \mid \{x_i\})$$

- ▶ Moreover, we want to find a better model comparing the model evidence:

$$p(\text{data} \mid \text{model 1}) \text{ vs. } p(\text{data} \mid \text{model 2})$$

*Side note:* Bayesian inference often involves heavy computational work

In the order of “easiness” in statistical learning (generally):

- ▶ Model parameter estimation/optimization
- ▶ Classification/categorical value prediction
- ▶ Real-valued prediction
- ▶ Model averaging ← Bayesian inference
- ▶ Probability/density estimation ← Bayesian inference
- ▶ Evidence computation (aka, partition function in physics) ← Bayesian inference

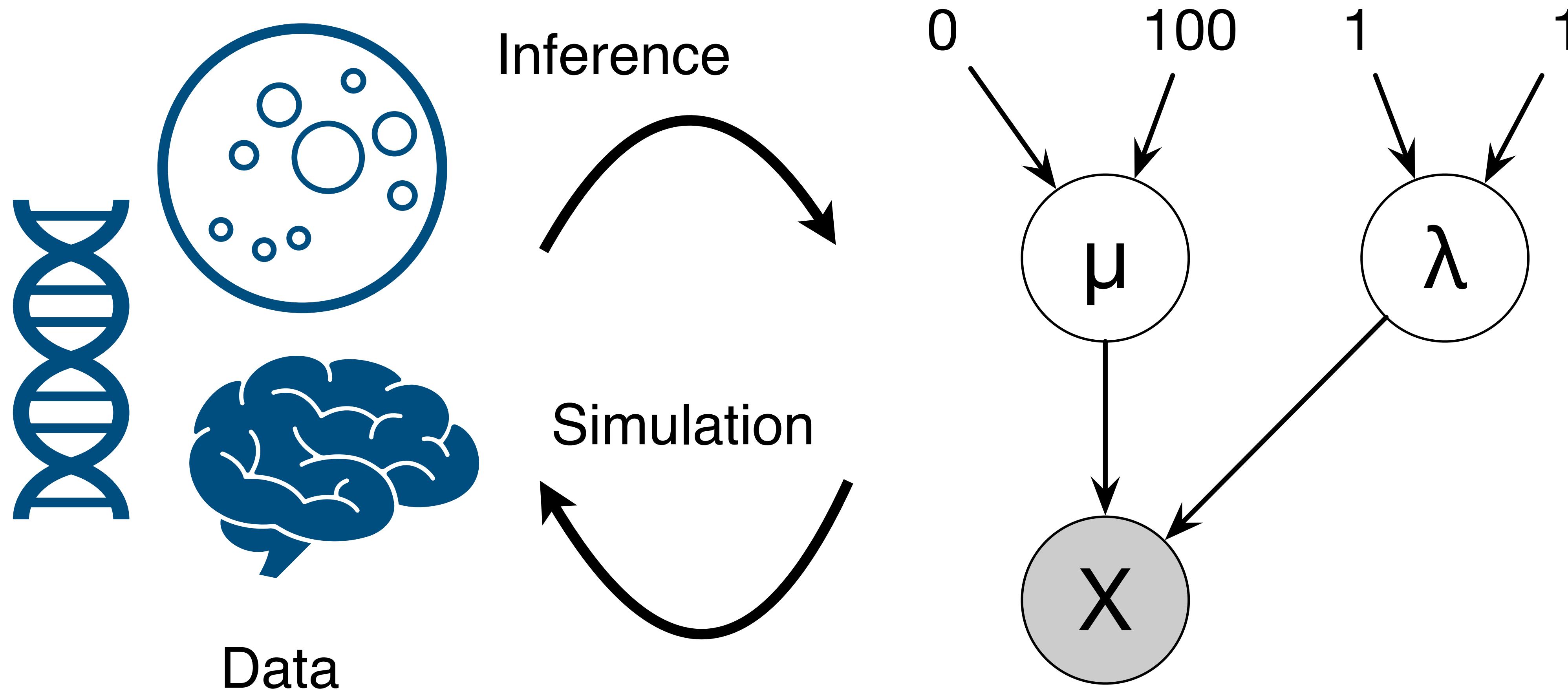
# The first step in Bayesian inference is to build a generative model

In our example:

1. Sample  $\mu \sim \mathcal{N}(0, 10^2)$
  2. Sample  $\gamma \sim \text{Gamma}(1, 1)$
  3. Sample  $X \sim \log\mathcal{N}(\mu, \gamma^{-1})$
- ▶ A prior generative scheme of  $\mu$
  - ▶ A prior generative scheme of  $\gamma$
  - ▶ A data-generating of  $x$  given  $\mu, \gamma$

- ▶ “What I cannot create, what I cannot understand” - Richard Feynman
- ▶ How do we recover (reverse-engineer) these unknown parameters?

# The goal is to infer $\mu, \lambda$

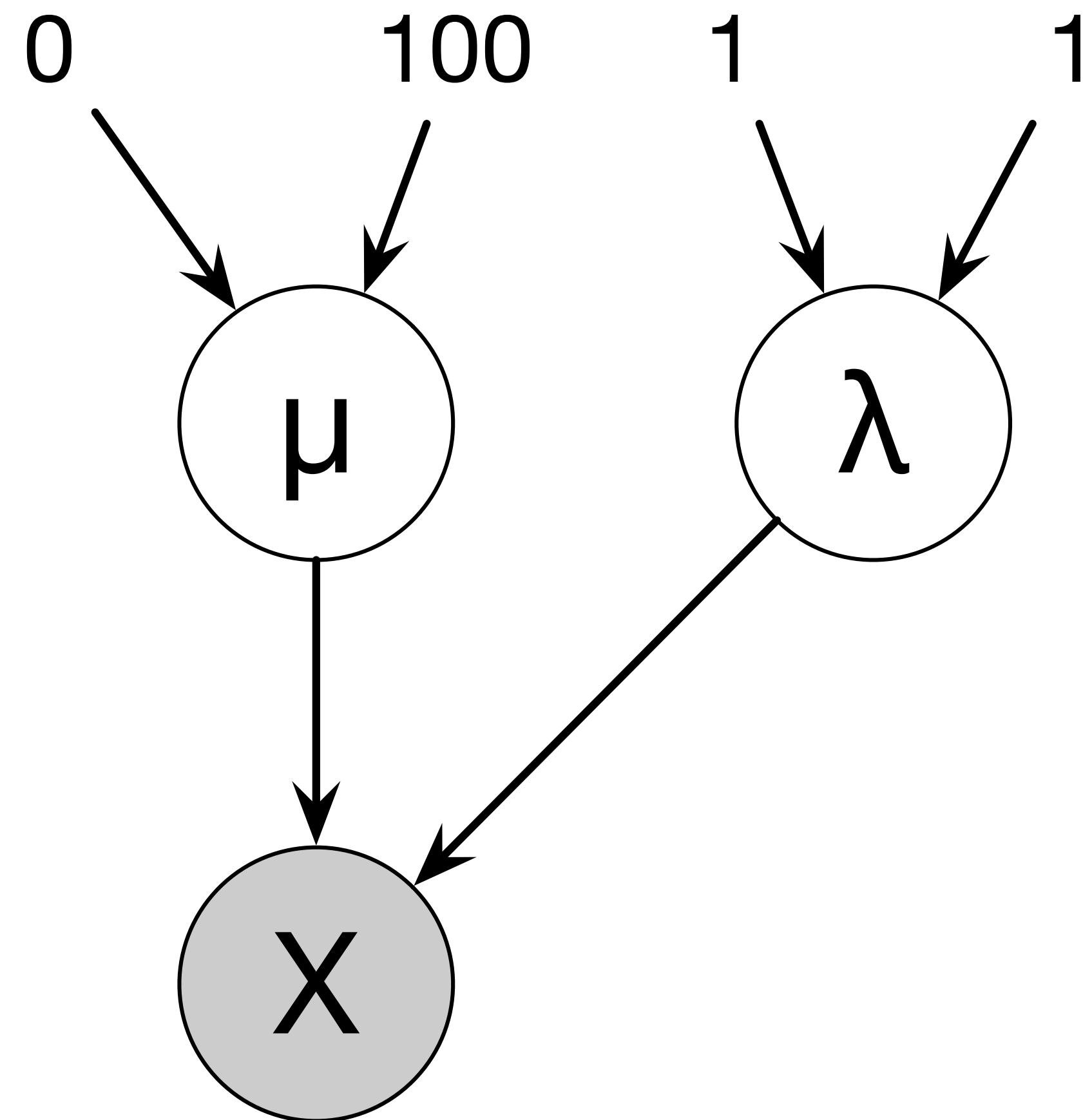


# The goal is to infer $\mu, \lambda$

$$\mu \sim \mathcal{N}(0, 10^2)$$

$$\lambda \sim \text{Gamma}(1, 1)$$

$$X \sim \text{log-}\mathcal{N}(\mu, \lambda)$$



We will use *stan* to infer/simulate the "posterior" distribution of  $\mu$  and  $\lambda$



<https://mc-stan.org/>

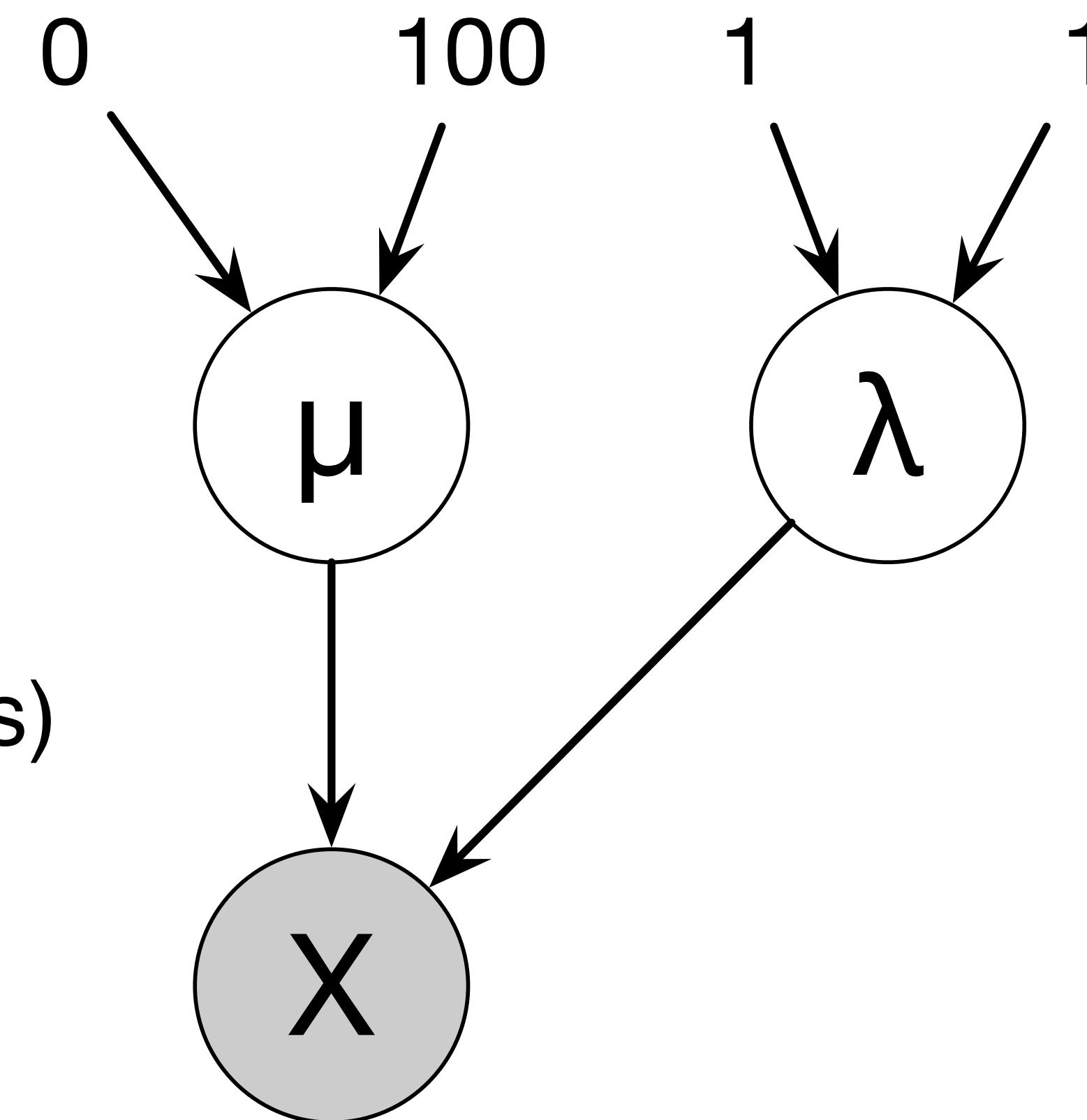


# Simulation of posterior sampling: What we need and what *Stan* does for us



For each Markov Chain (simulation process)

1. Simulate  $\mu$  given current  $\lambda$  and  $X$
2. Simulate  $\lambda$  given current  $\mu$  and  $X$
3. Repeat 1-2



# stan: A good news! All you need is to complete the model description!

```
data {  
    int n;                      // sample size  
    real x[n];                  // n data points  
}  
parameters {  
    real mu;                    // mean parameter  
    real<lower=1e-8> lambda;   // precision  
}  
model {  
    mu ~ normal(0, 1e2);        // prior  
    lambda ~ gamma(1e-2, 1e-2);  // prior  
    for(i in 1:n){              // data generating  
        x[i] ~ lognormal(mu, sqrt(1/lambda));  
    }  
}
```

We can have a separate file with model description in stan language:

1. data block to describe expected data dimensions and types
2. parameters block to describe (unknown) model parameters to infer
3. model block to describe data-generating schemes (a probabilistic graph model)

# rstan incorporates existing stan codes into R scripts to run Bayesian inference algorithm

In R, it can be a simple string

```
.code <- "
data {
  int n;
  real x[n];
}
parameters {
  real mu;
  real<lower=1e-8> lambda;
}
model {
  mu ~ normal(0, 1e2);
  lambda ~ gamma(1e-2, 1e-2);
  for(i in 1:n){
    x[i] ~ lognormal(mu, sqrt(1/lambda));
  }
}"
```

We can copy/paste stan's model code in some R string.

1. data block to describe expected data dimensions and types
2. parameters block to describe (unknown) model parameters to infer
3. model block to describe data-generating schemes (a probabilistic graph model)

# Running stan is as easy as calling other R functions

```
library(rstan)
options(mc.cores = parallel::detectCores())

if.needed("example_stan_lgaussian.rds", {          ## don't re-run

    ## Run MCMC inference algorithm
    .fit <- stan(model_code = .code,
                  data = list(n = nn, x = xx),      ## code
                  pars = c("mu", "lambda"),       ## a list of data
                  chains = 5,                   ## parameters of interest
                  iter=1000)                    ## number of parallel MCMC chains
                                         ## how many iterations?

    saveRDS(.fit, "example_stan_lgaussian.rds") ## save the results
})

## Extract the sampled parameters
.fit <- readRDS("example_stan_lgaussian.rds")
.mu <- rstan::extract(.fit, pars="mu", inc_warmup=TRUE, permuted=FALSE)
.lambda <- rstan::extract(.fit, pars="lambda", inc_warmup=TRUE, permuted=FALSE)
```

As a result of MCMC...

```
dim(.mu)
```

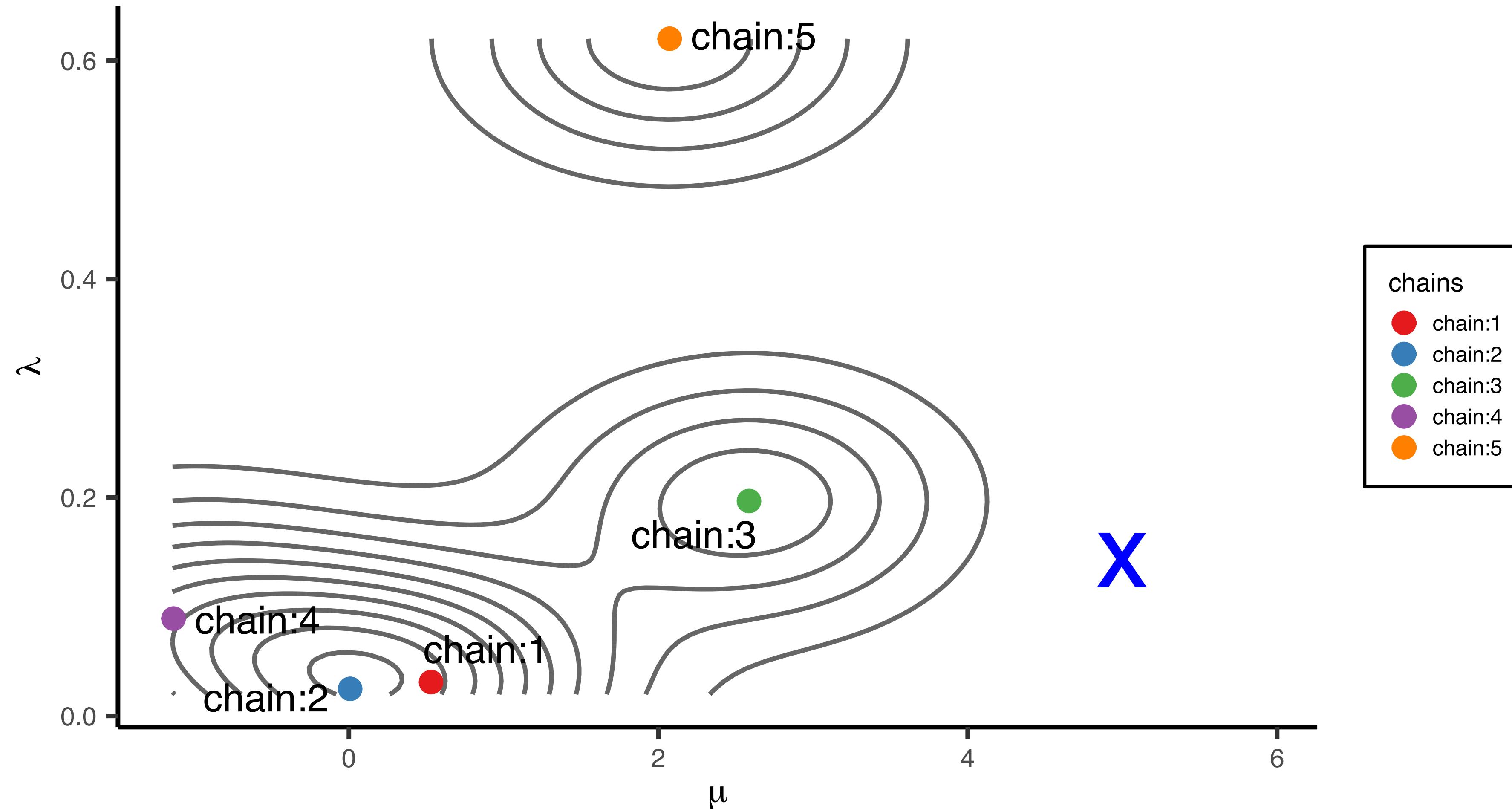
```
[1] 1000 5 1
```

```
dim(.lambda)
```

```
[1] 1000 5 1
```

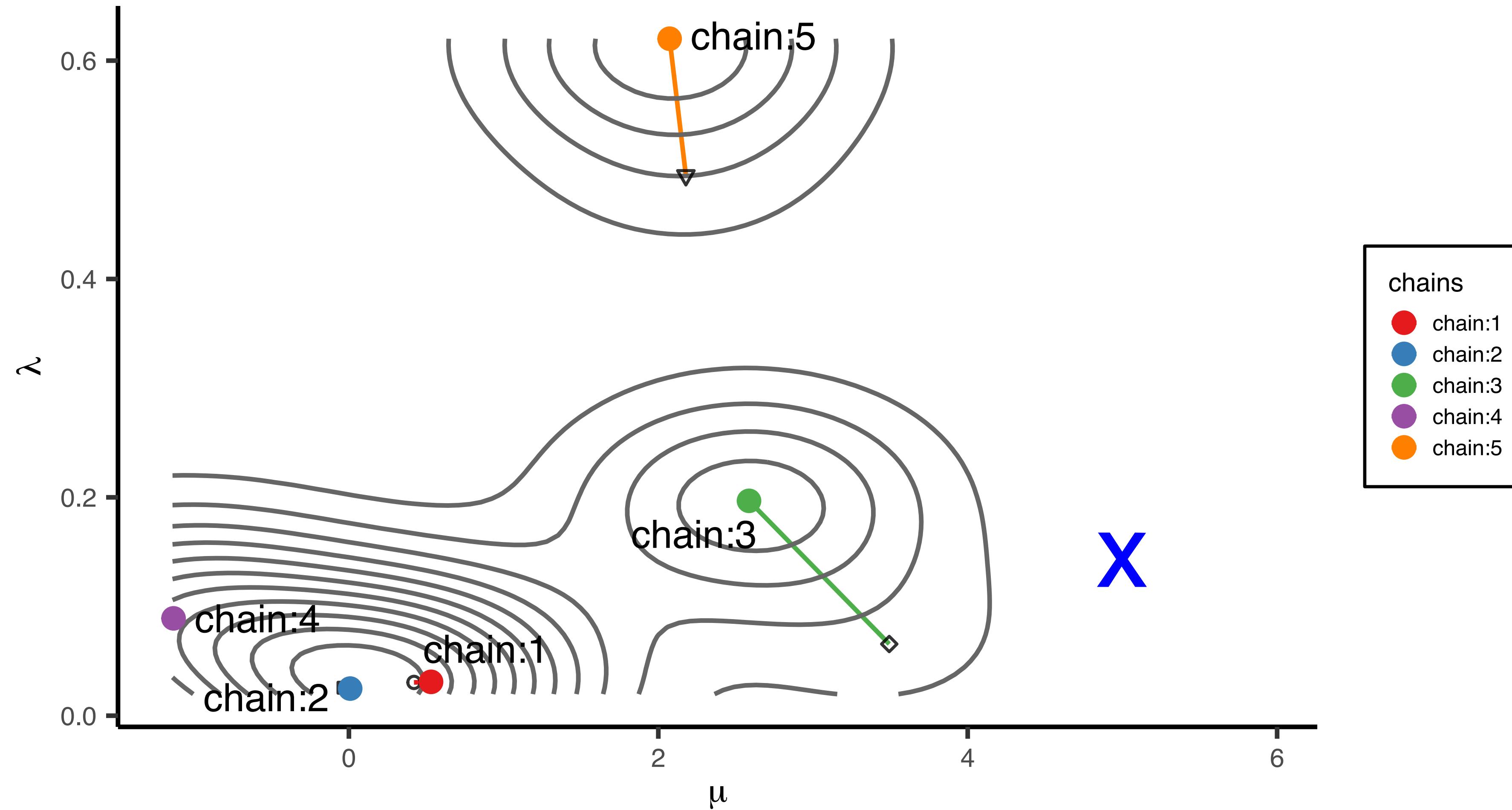
As a result of four independent MCMC runs

first 3 iterations



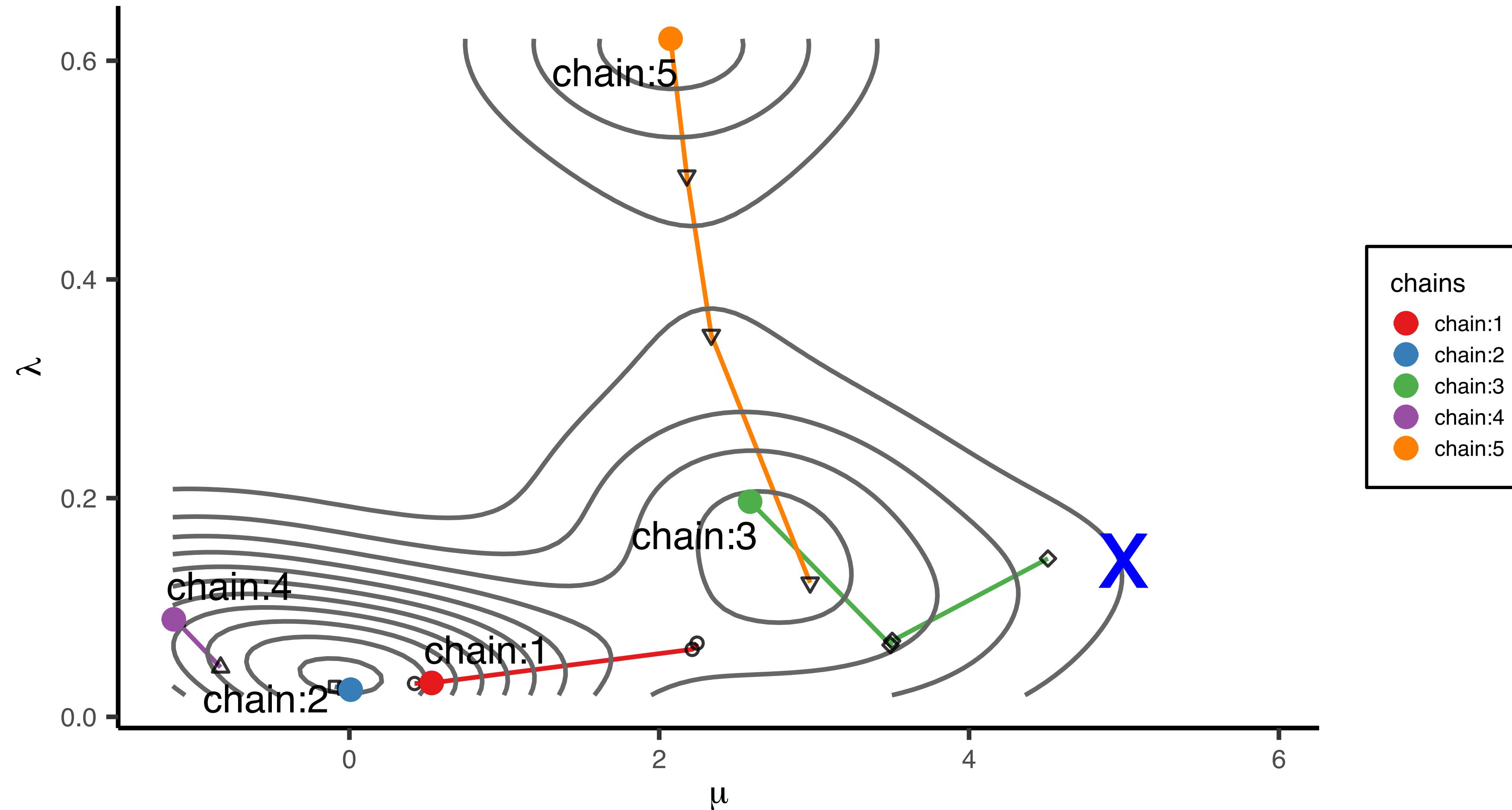
As a result of four independent MCMC runs

first 5 iterations



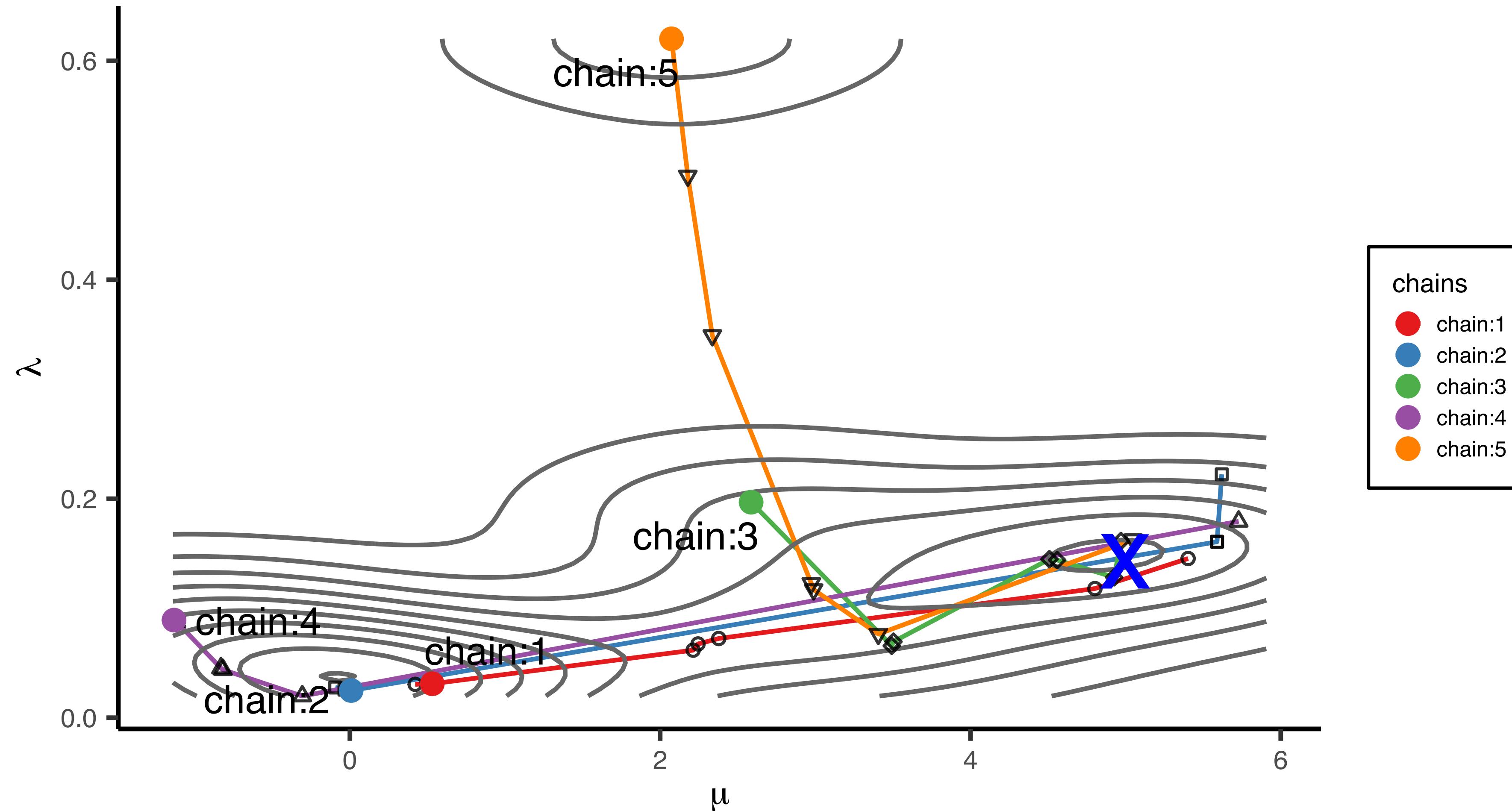
As a result of four independent MCMC runs

first 7 iterations



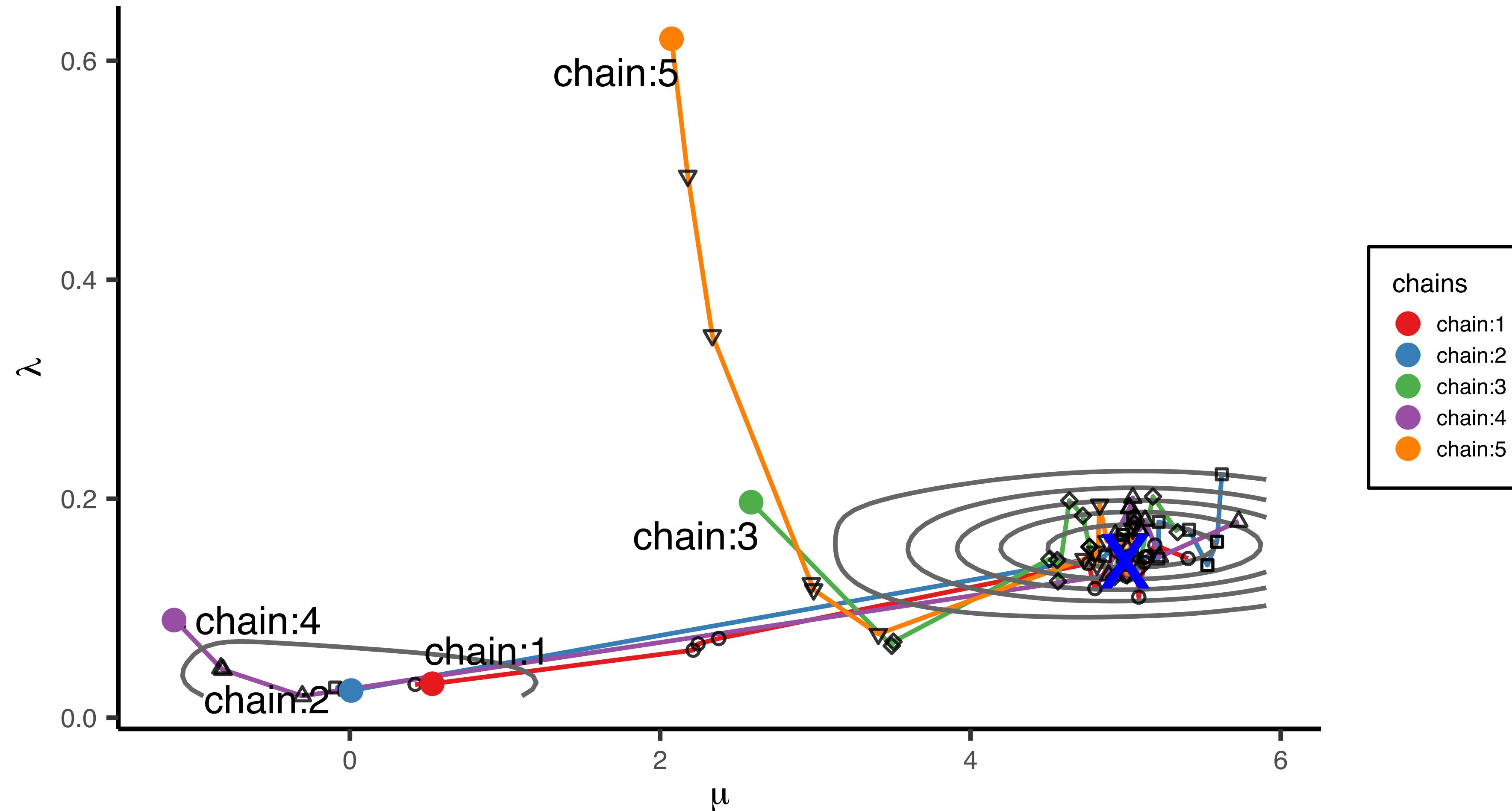
As a result of four independent MCMC runs

first 10 iterations



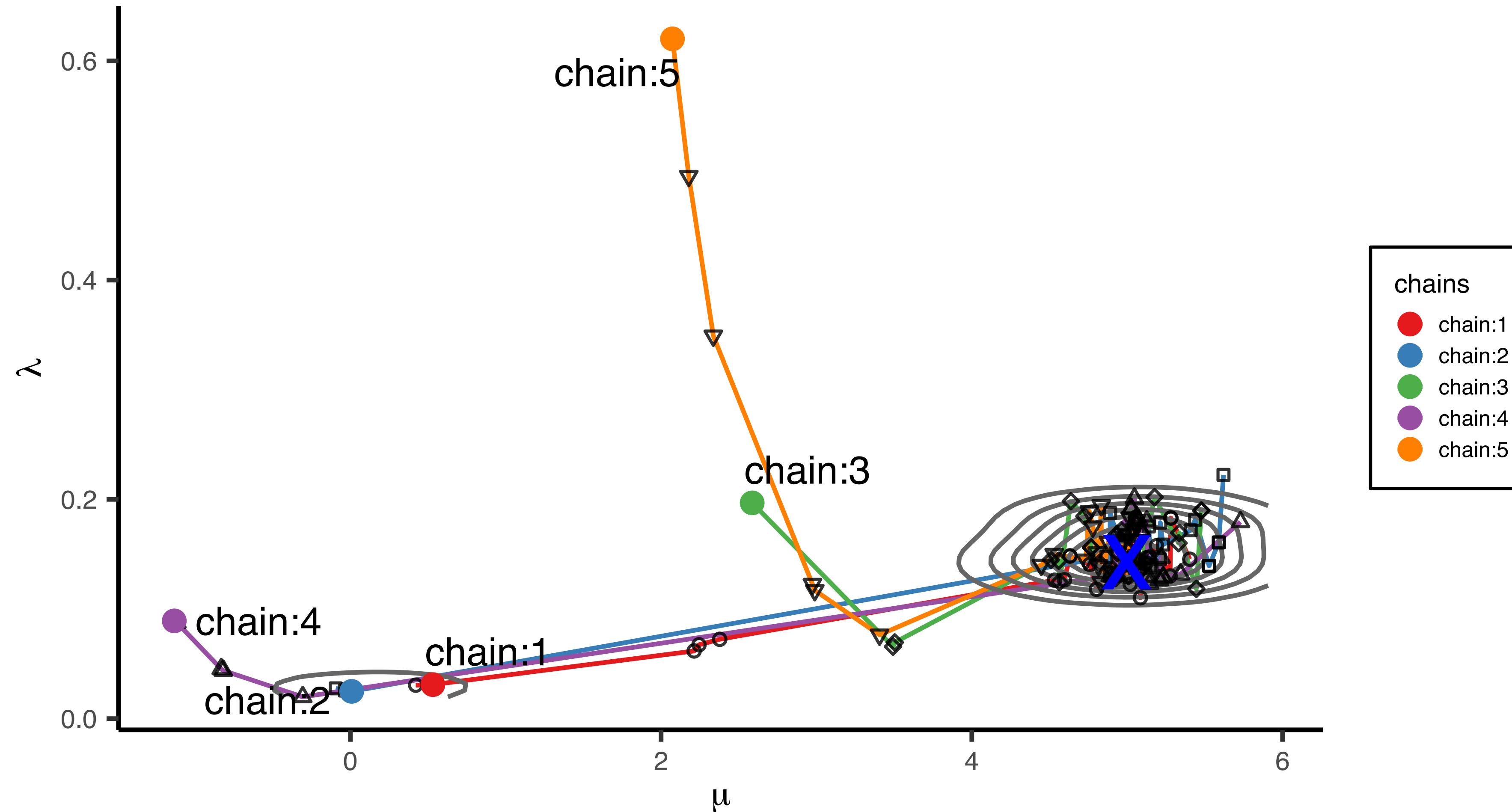
As a result of four independent MCMC runs

first 20 iterations



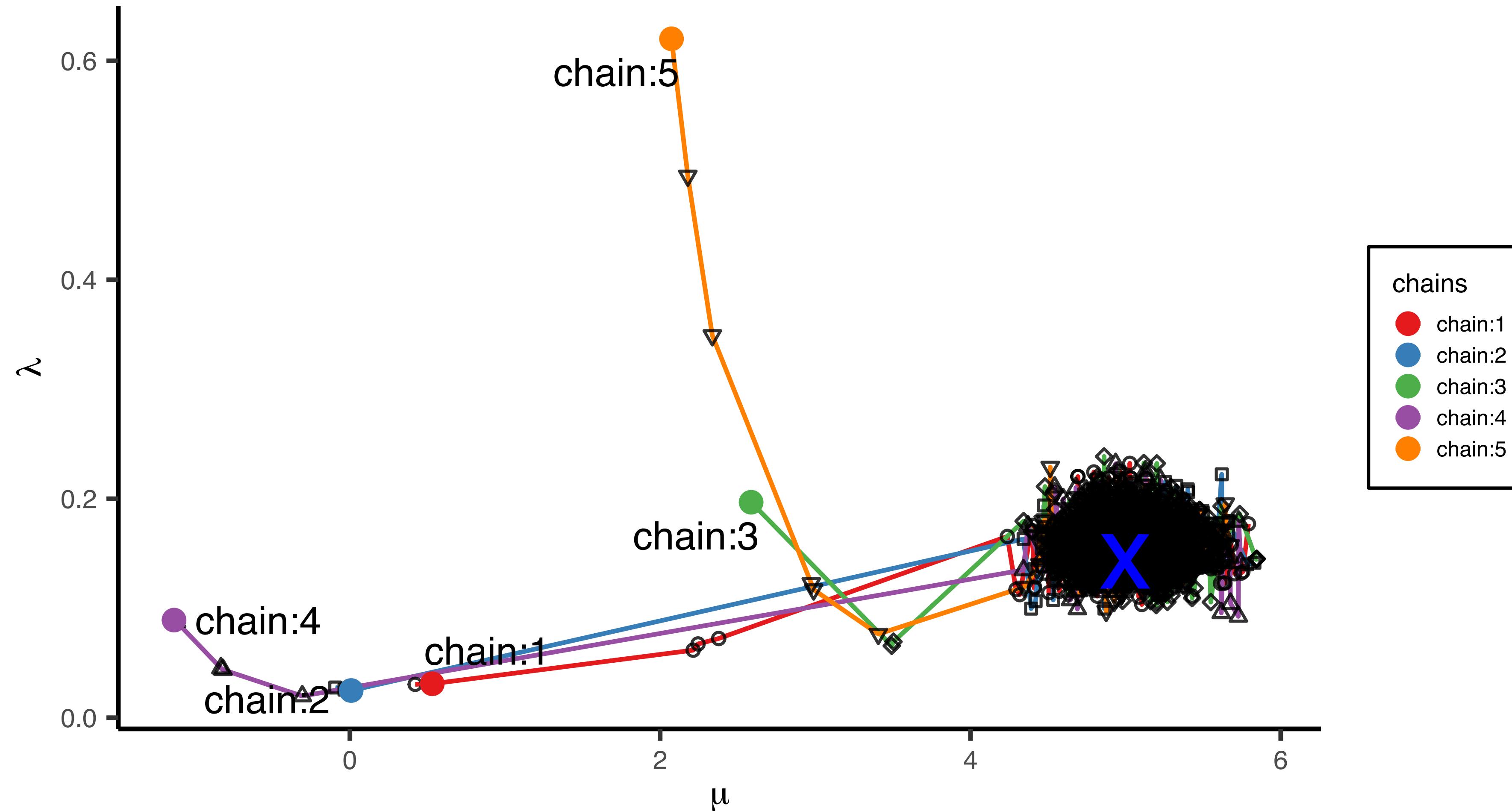
As a result of four independent MCMC runs

first 30 iterations

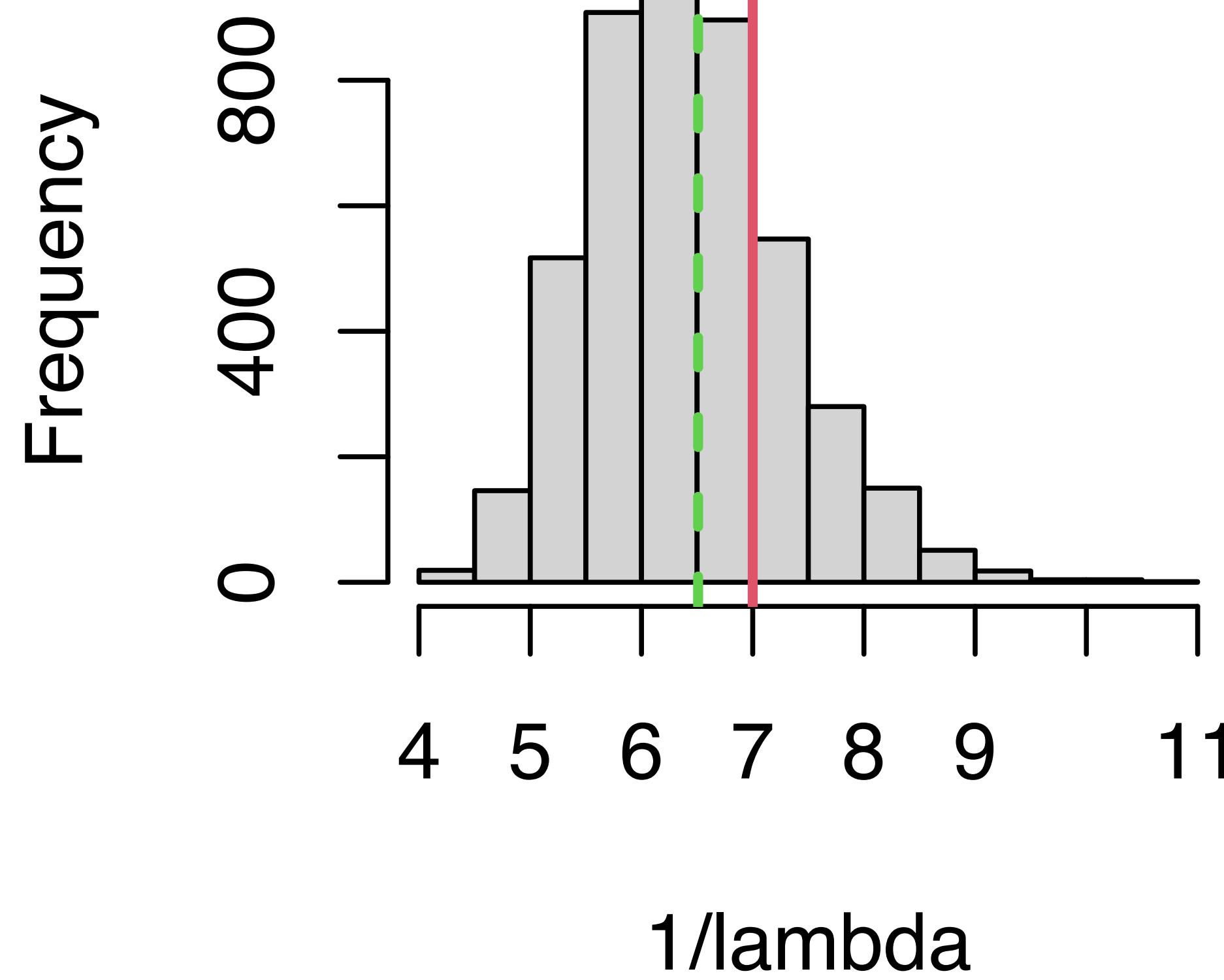
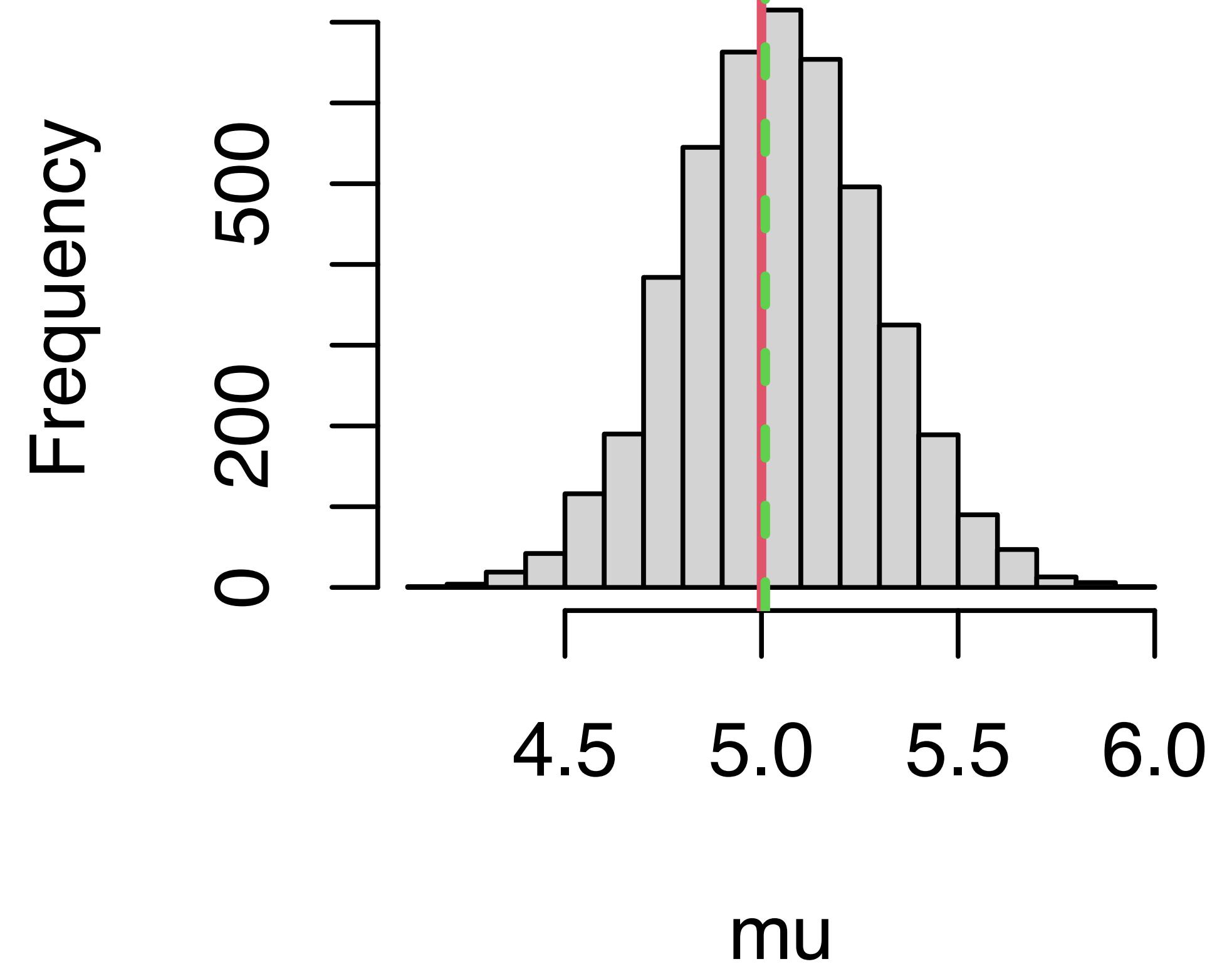


As a result of four independent MCMC runs

first 500 iterations



## Posterior probability



# Bayesian linear regression

Let's consider another toy example, Poisson regression

```
n <- 100      ## sample
p.tot <- 10    ## total no. of predictors
p <- 2        ## first `p` are true predictors
sigma <- 0.1   ## error StdDev

set.seed(1331)
X <- matrix(rnorm(n * p.tot), nrow=n, ncol=p.tot)
theta <- matrix(2 * rnorm(p), nrow=p, ncol=1)
y <- round(exp(X[, 1:p] %*% theta + rnorm(n) * sigma))
```

Can you draw the corresponding probabilistic graphical model?

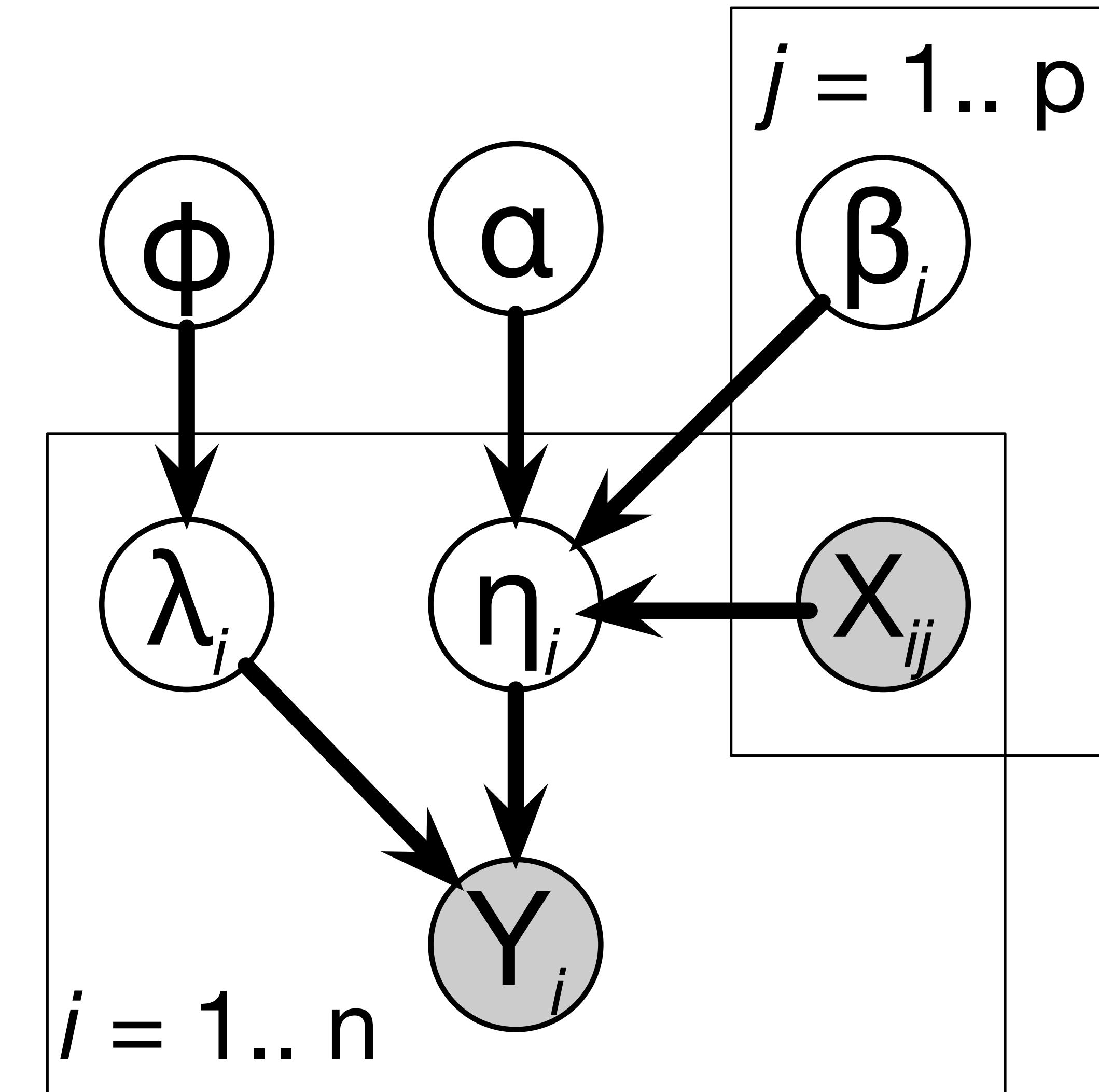
$$\alpha, \beta_j \sim \mathcal{N}(0, 1)$$

$$\phi \sim \text{Gamma}(1, 1)$$

$$\lambda_i | \phi \sim \text{Gamma}(\phi, 1)$$

$$\eta_i \leftarrow \mathbf{x}_i \boldsymbol{\beta} + \alpha$$

$$Y_i | \eta_i, \lambda_i \sim \text{Poisson}(\lambda_i \exp(\eta_i))$$



# Let's write out a Poisson regression generative model

```
data {  
    int n; int p;                                // number of observations & predictors  
    matrix[n, p] X;                            // design matrix  
    int y[n];                                    // response vector  
}  
parameters {  
    vector[p] beta;                             // regression coefficients  
    real alpha;                                 // intercept  
    vector<lower=1e-8>[n] lambda;                // overdispersion  
    real phi;                                   // hyperparameter  
}  
transformed parameters {  
    vector[n] eta;  
    eta = X * beta + alpha;  
}  
model {  
    alpha ~ normal(0, 1); phi ~ gamma(1e4, 1e4); //  
    for (j in 1:p) { beta[j] ~ normal(0, 1); }      // coefficients  
    for (i in 1:n) {  
        lambda[i] ~ gamma(phi, 1);                  // random overdispersion  
        y[i] ~ poisson(lambda[i] * exp(eta[i]));    // poisson rate  
    }  
}
```

```
.code <- "
data {
  int n; int p;                                // number of observations & predictors
  matrix[n, p] X;                             // design matrix
  int y[n];                                    // response vector
}
parameters {
  vector[p] beta;                            // regression coefficients
  real alpha;                               // intercept
  vector<lower=1e-8>[n] lambda;           // overdispersion
  real phi;                                 // hyperparameter
}
transformed parameters {
  vector[n] eta;
  eta = X * beta + alpha;
}
model {
  alpha ~ normal(0, 1);   phi ~ gamma(1e4, 1e4); //
  for (j in 1:p) { beta[j] ~ normal(0, 1); }      //
  for (i in 1:n) {
    lambda[i] ~ gamma(phi, 1);                  // random overdispersion
    y[i] ~ poisson(lambda[i] * exp(eta[i]));    // poisson rate
  }
}
"
```

## Run stan to sample the parameters of interest

```
if.needed("example_stan_regression.rds", {          ## don't re-run

  lm.fit <- stan(model_code=.code,
                  data = list(n=nrow(X),
                               p=ncol(X),
                               y=y[,1],
                               X=X),
                  pars = c("beta","alpha"),
                  iter=1000,
                  chains=5)                         ## model code
                                              ## samples
                                              ## predictors
                                              ## response
                                              ## design matrix
                                              ## parameters
                                              ## MCMC
                                              ## chains

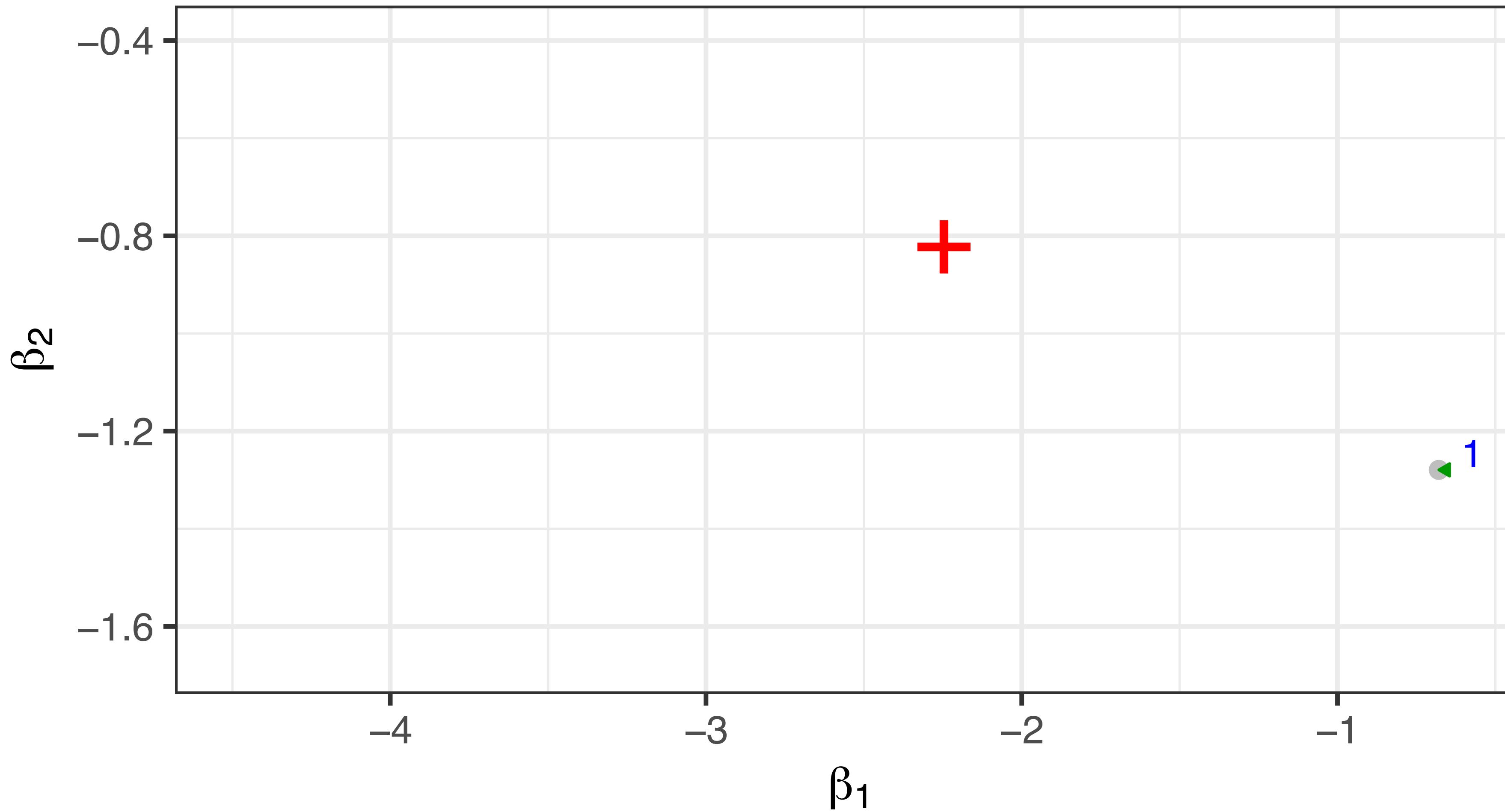
  saveRDS(lm.fit, "example_stan_regression.rds")
})

lm.fit <- readRDS("example_stan_regression.rds") ## save the results

.beta <- extract(lm.fit, pars="beta", inc_warmup=TRUE, permuted=FALSE)
```

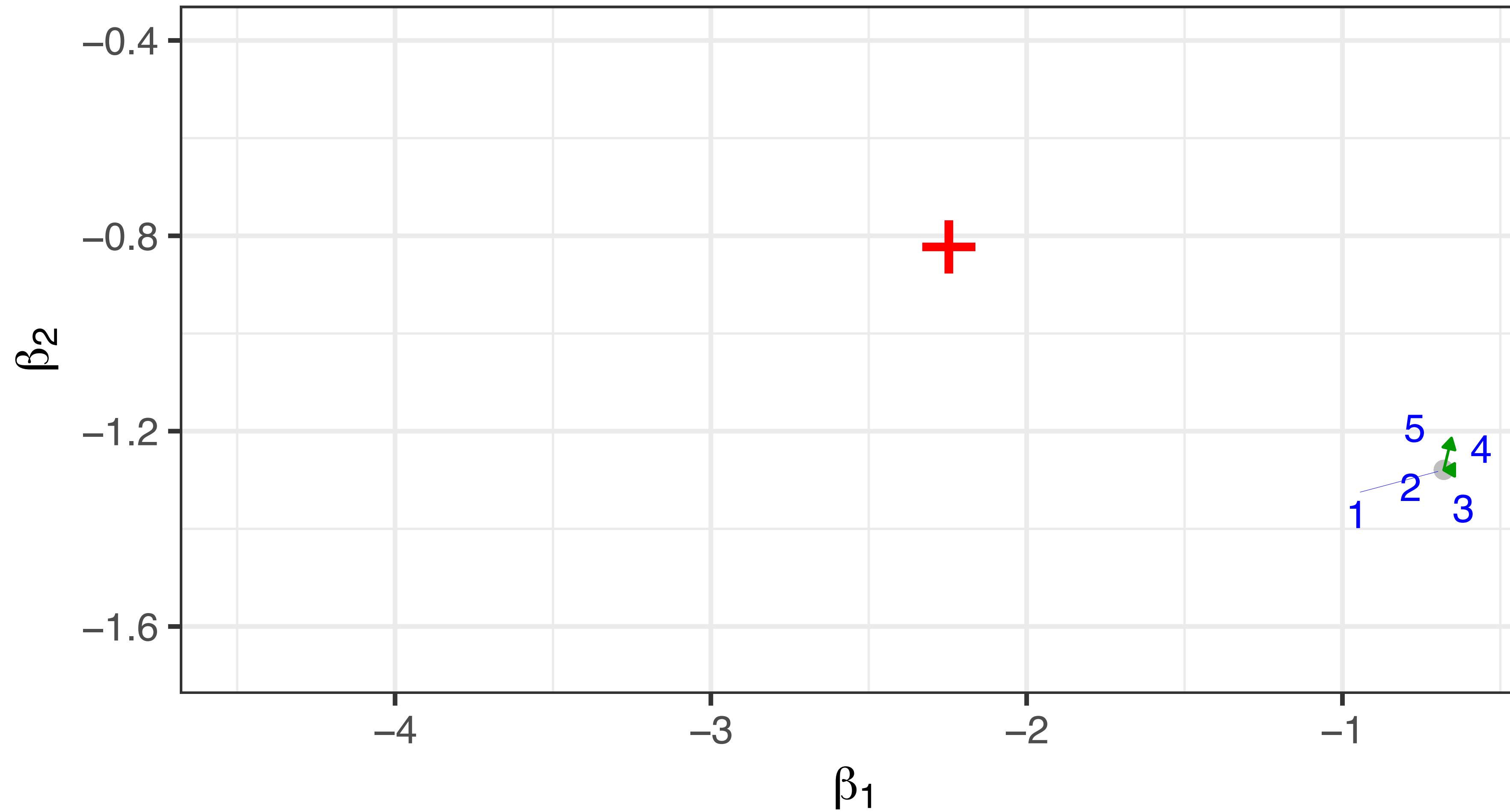
Let's take a look at truly non-zero coefficients

1 MCMC iterations



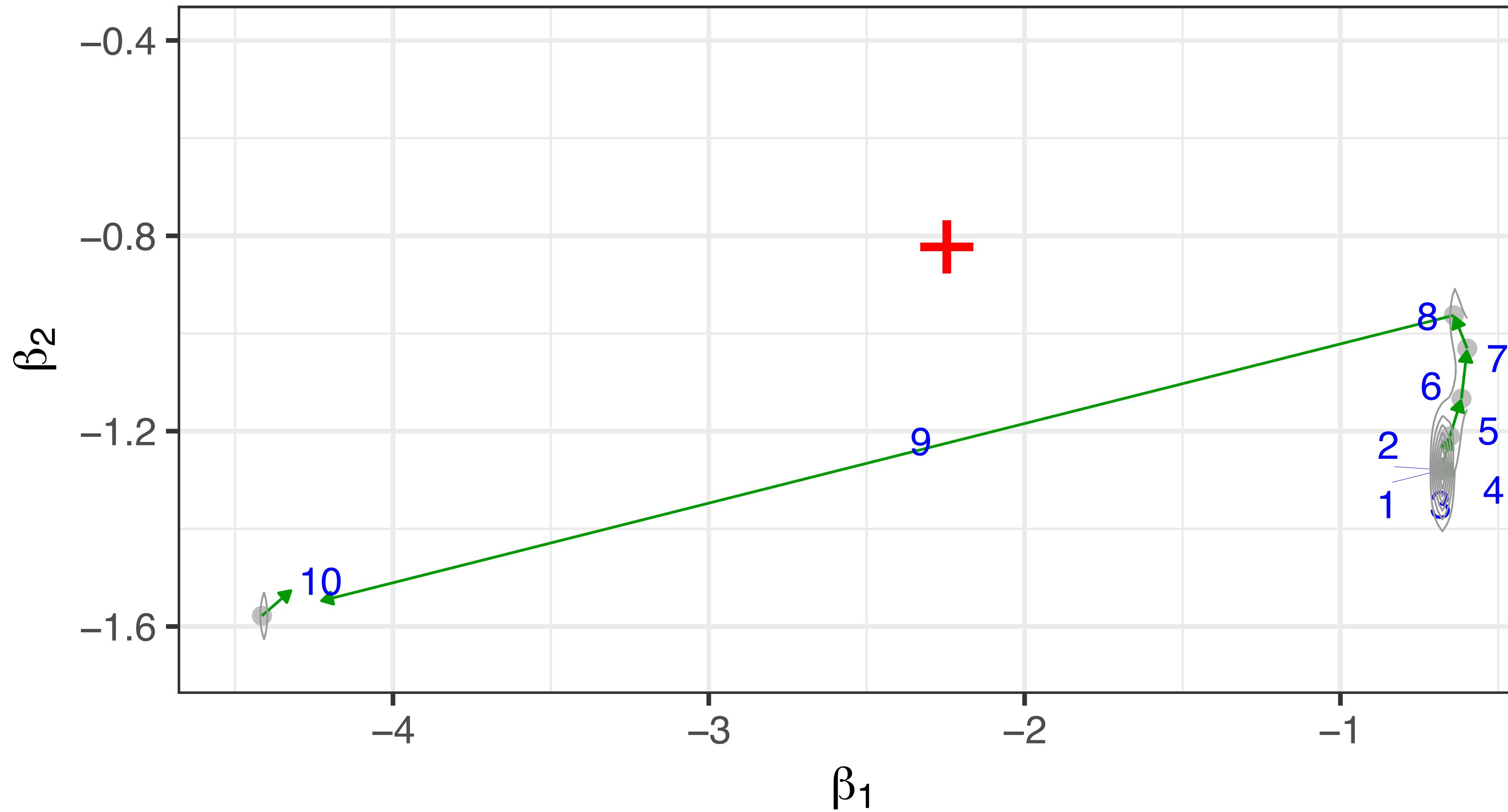
Let's take a look at truly non-zero coefficients

5 MCMC iterations



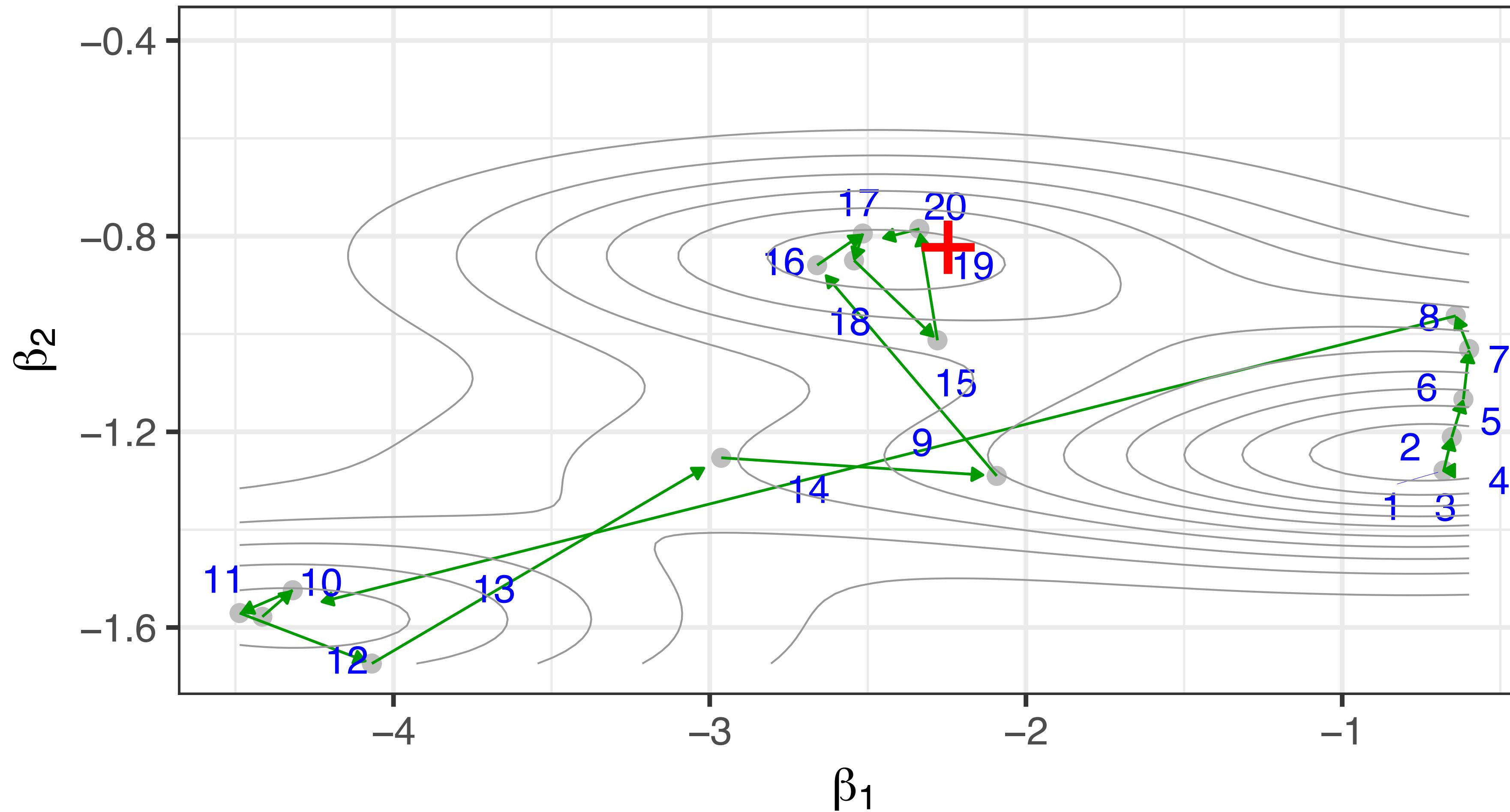
Let's take a look at truly non-zero coefficients

10 MCMC iterations



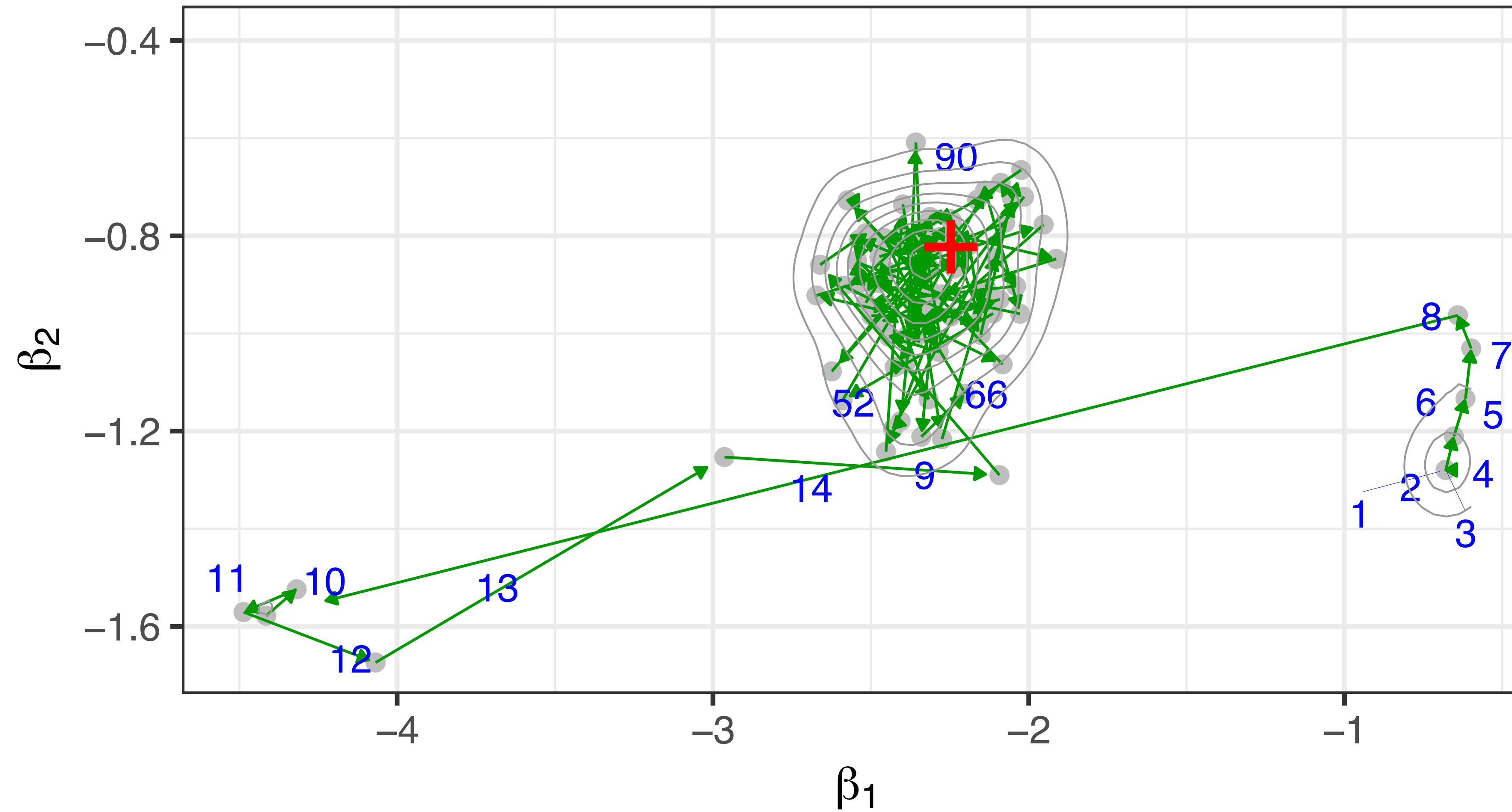
Let's take a look at truly non-zero coefficients

20 MCMC iterations



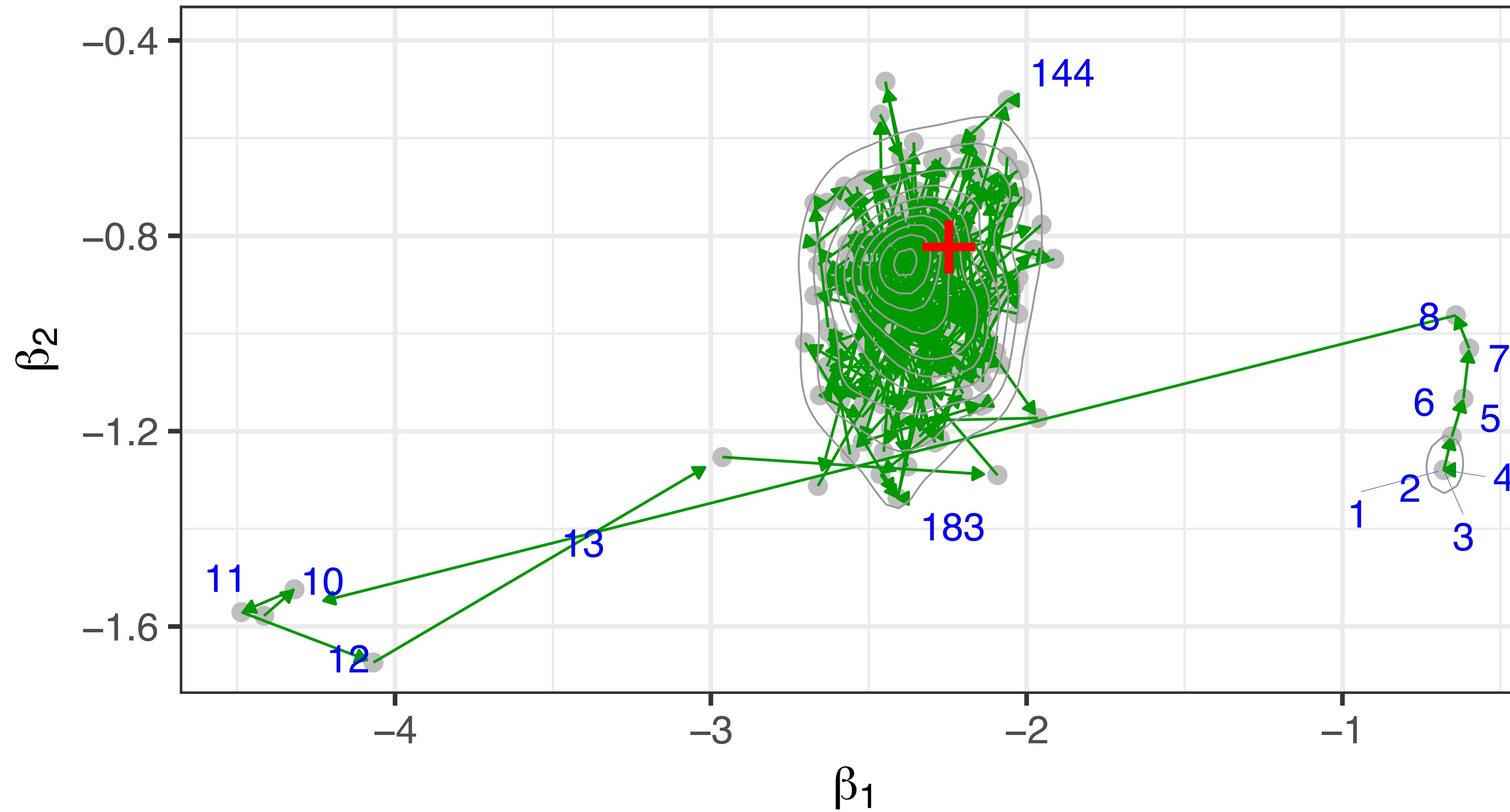
Let's take a look at truly non-zero coefficients

100 MCMC iterations

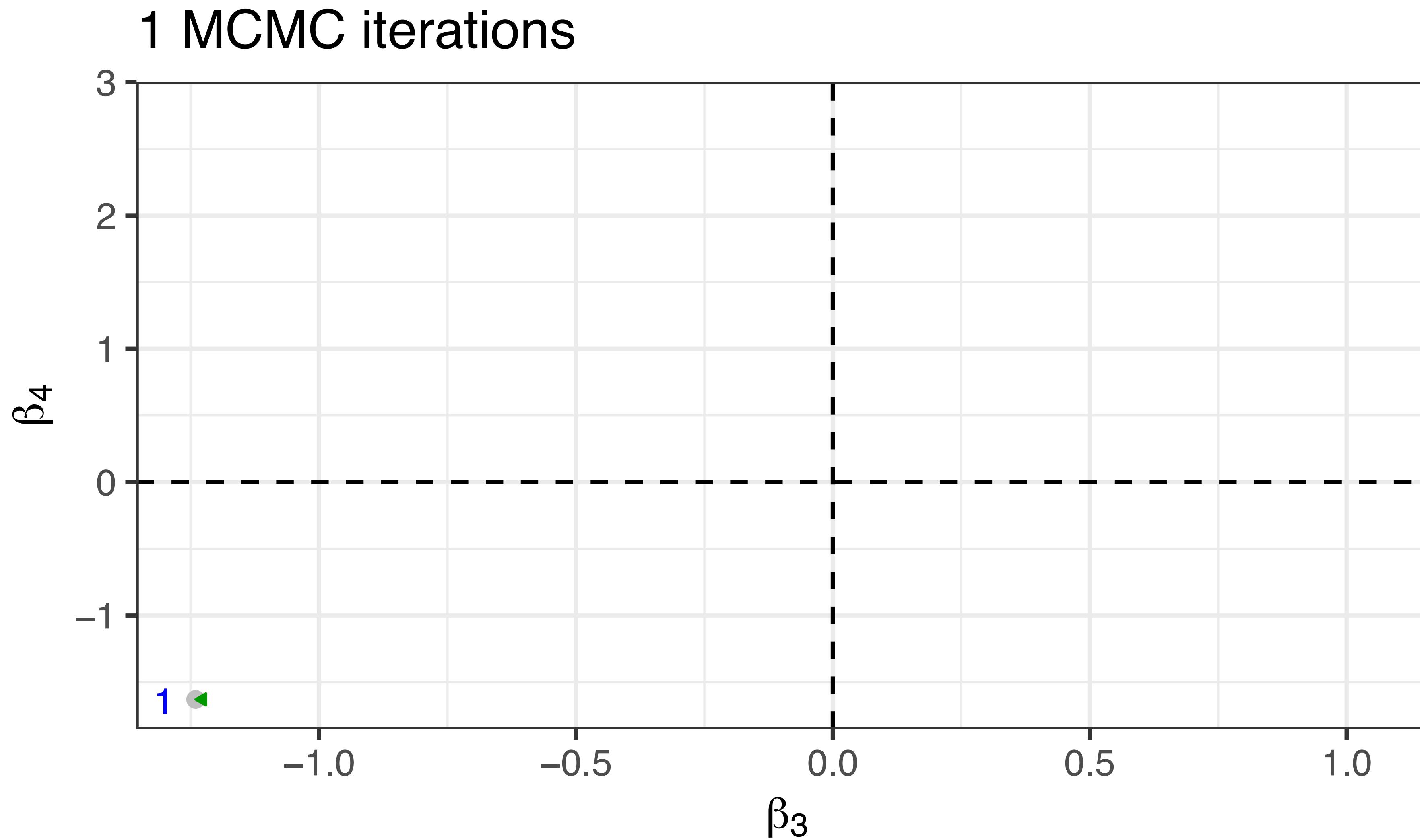


Let's take a look at truly non-zero coefficients

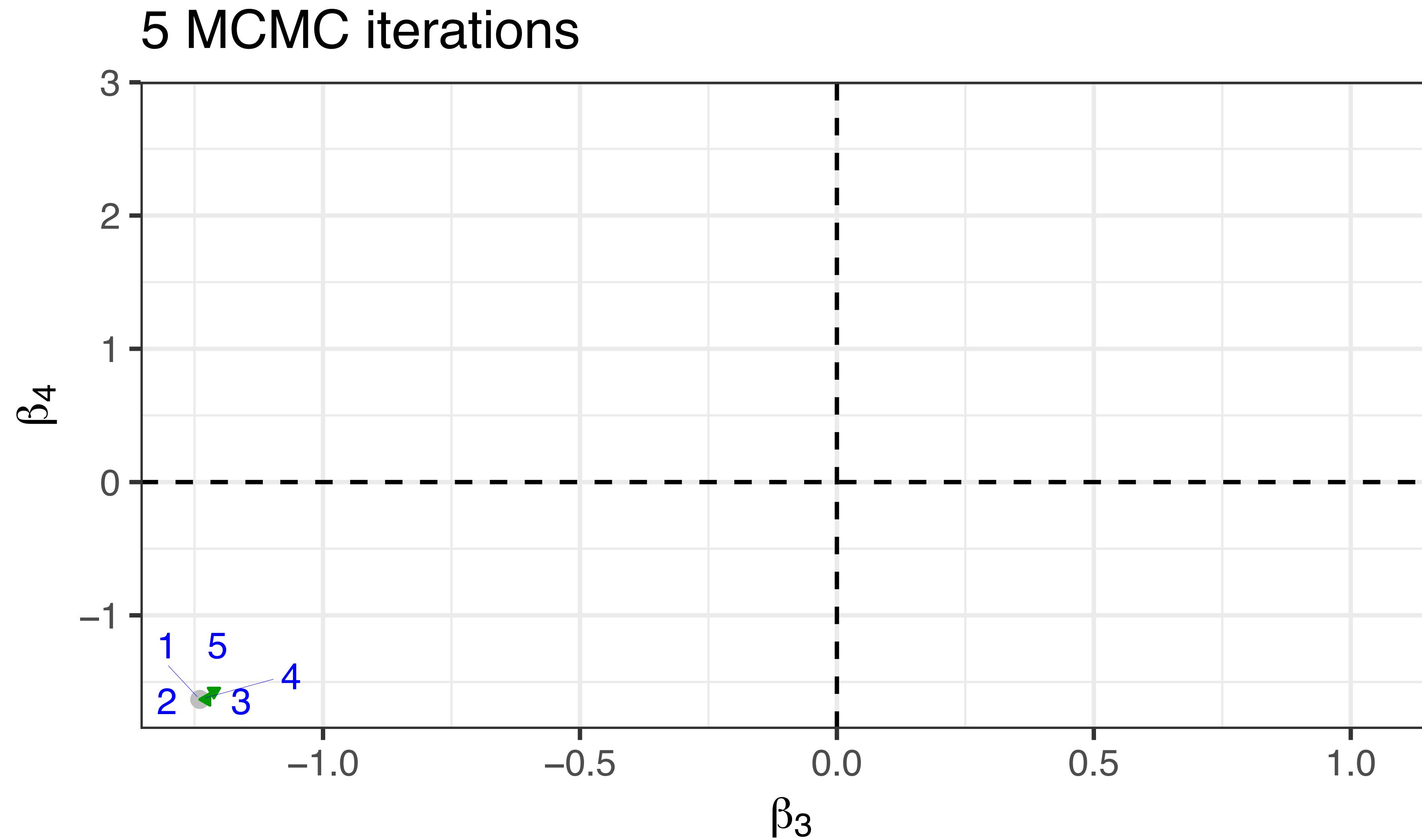
300 MCMC iterations



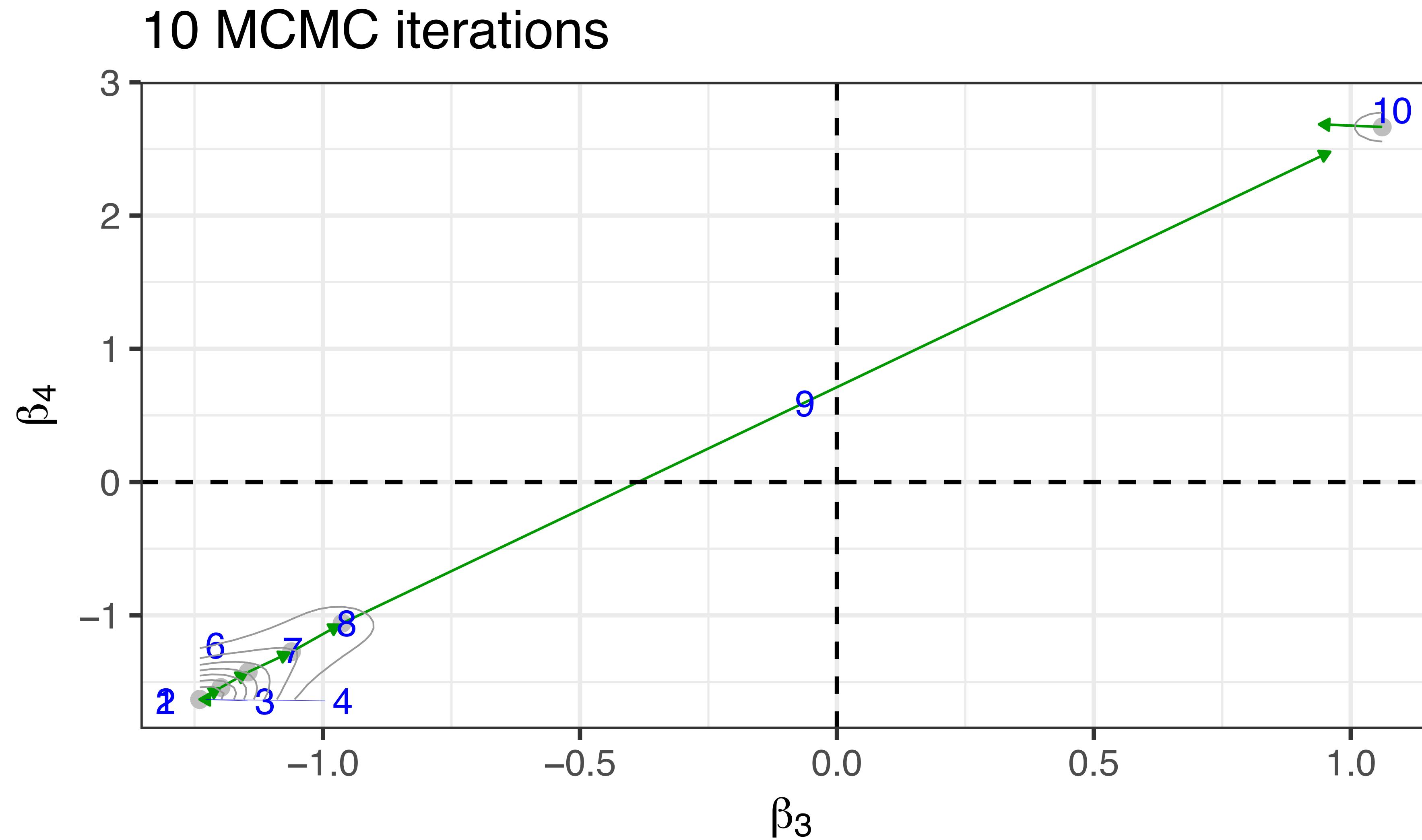
Let's take a look at other “null” coefficients



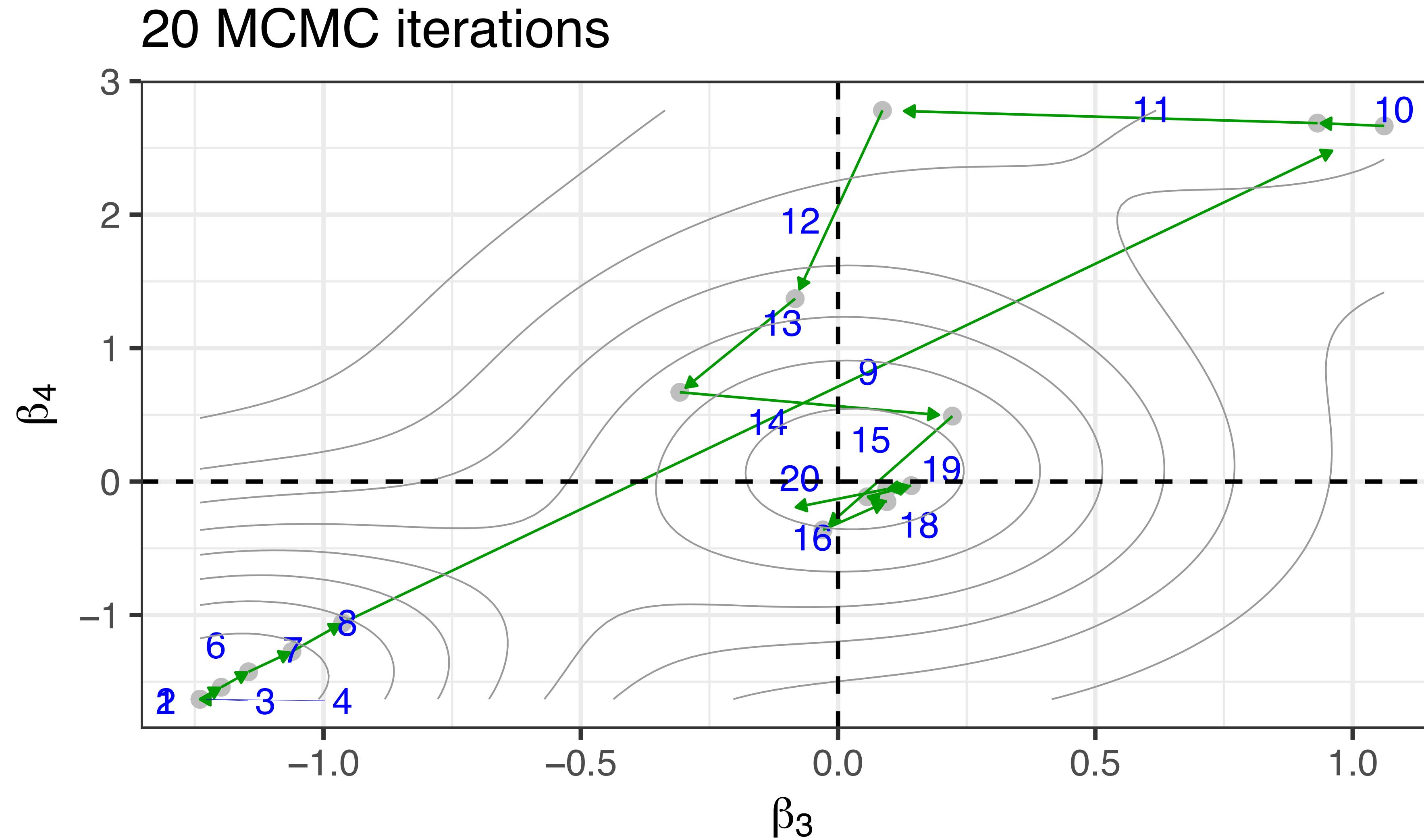
Let's take a look at other “null” coefficients



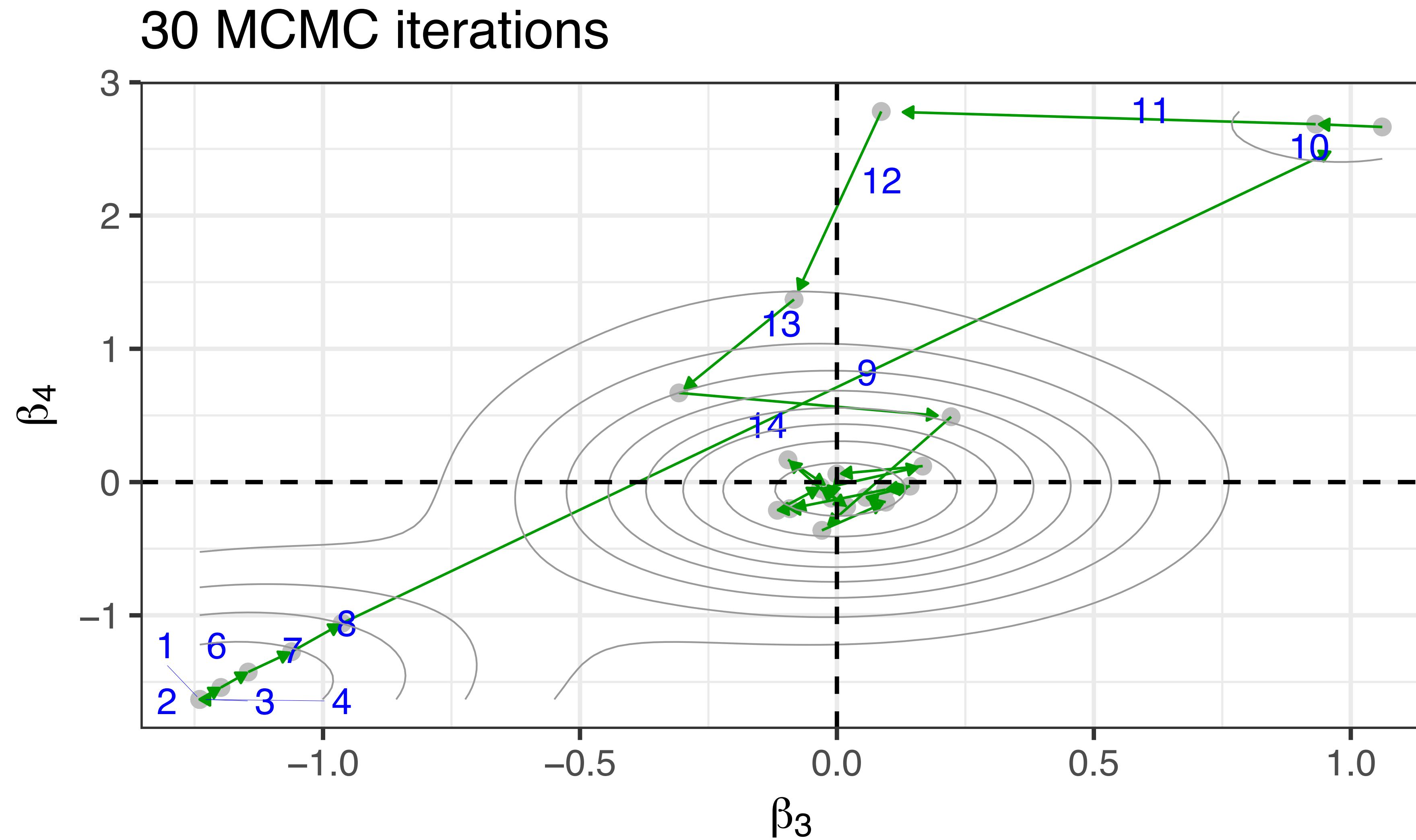
Let's take a look at other “null” coefficients



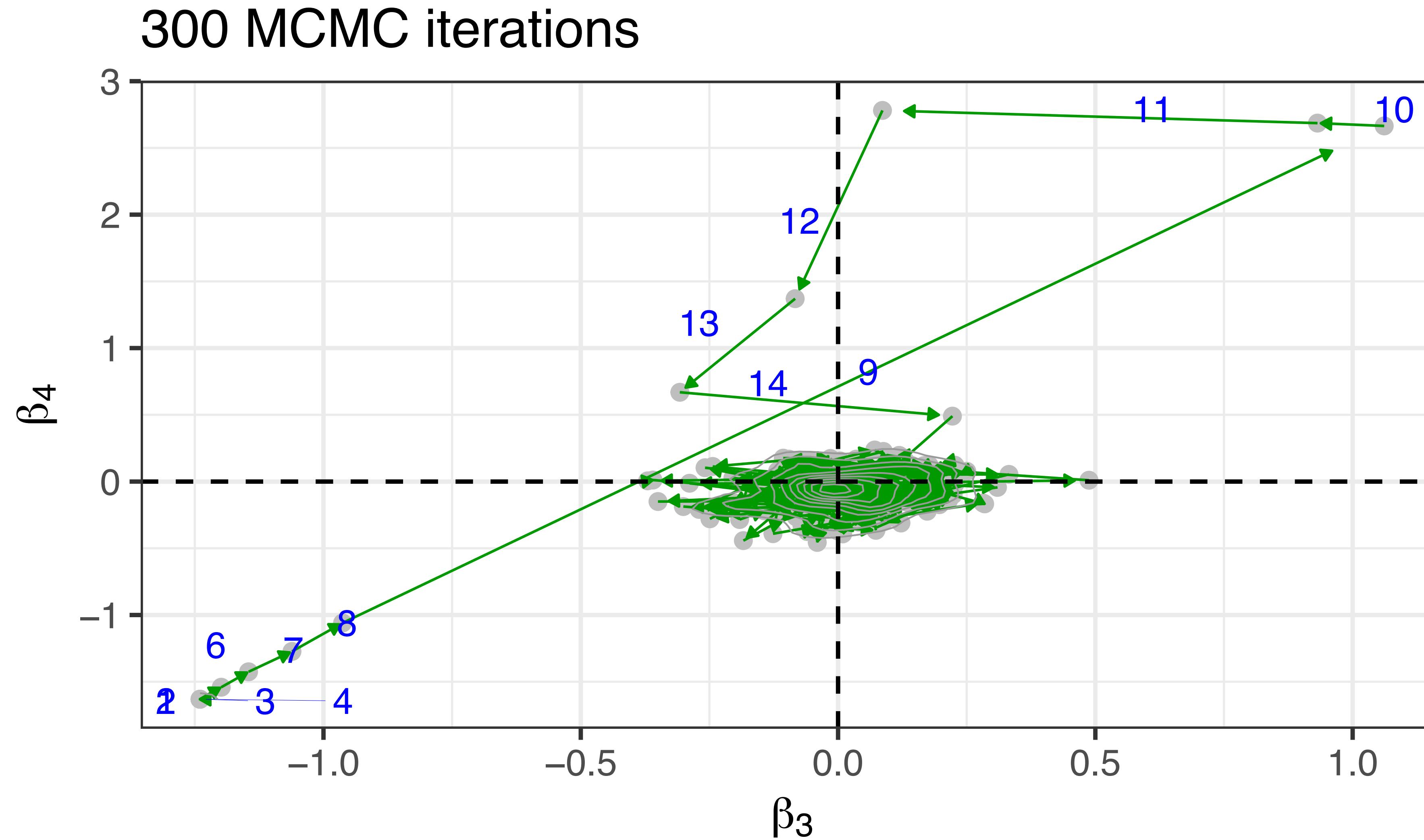
Let's take a look at other “null” coefficients



Let's take a look at other “null” coefficients

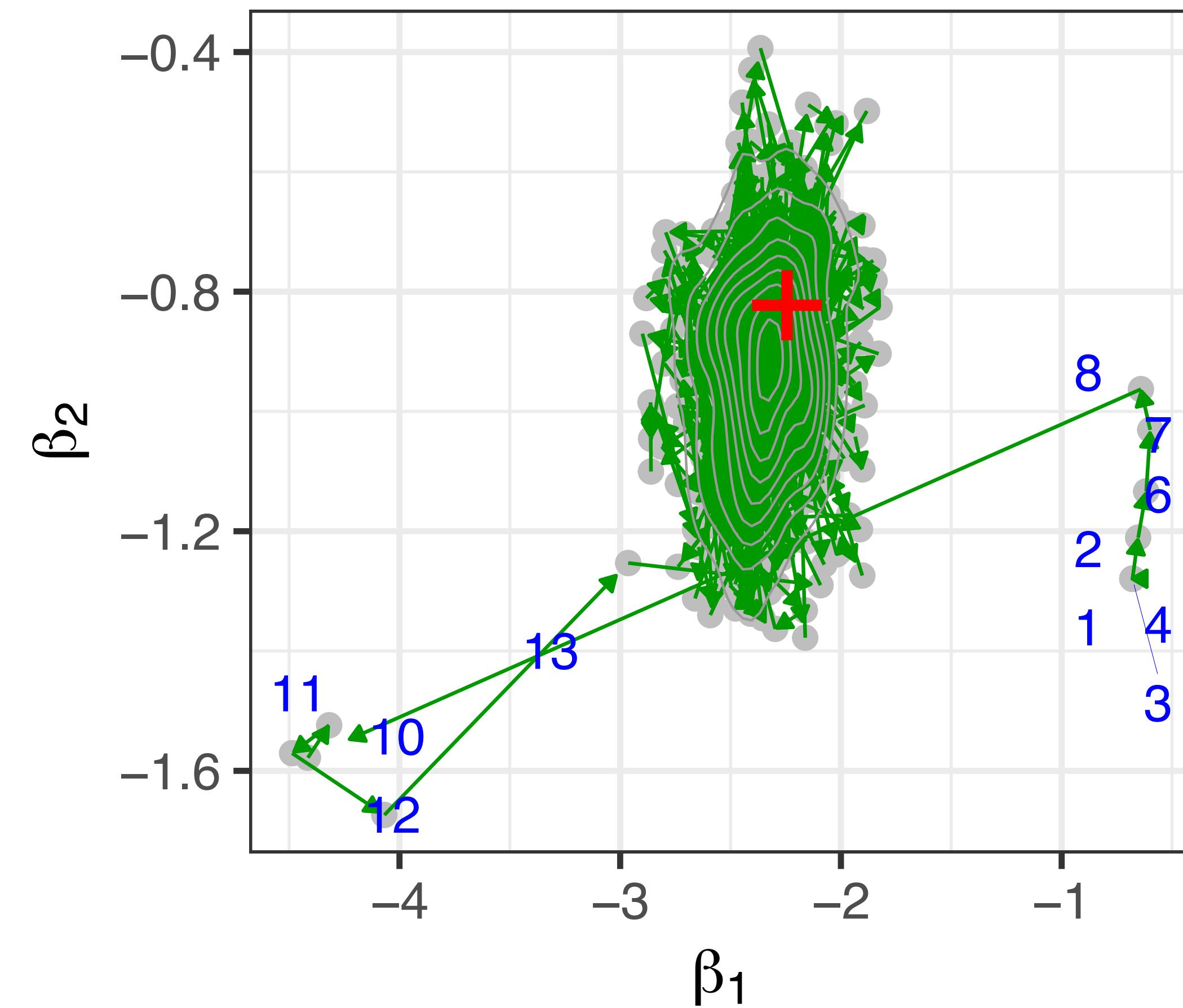


Let's take a look at other “null” coefficients

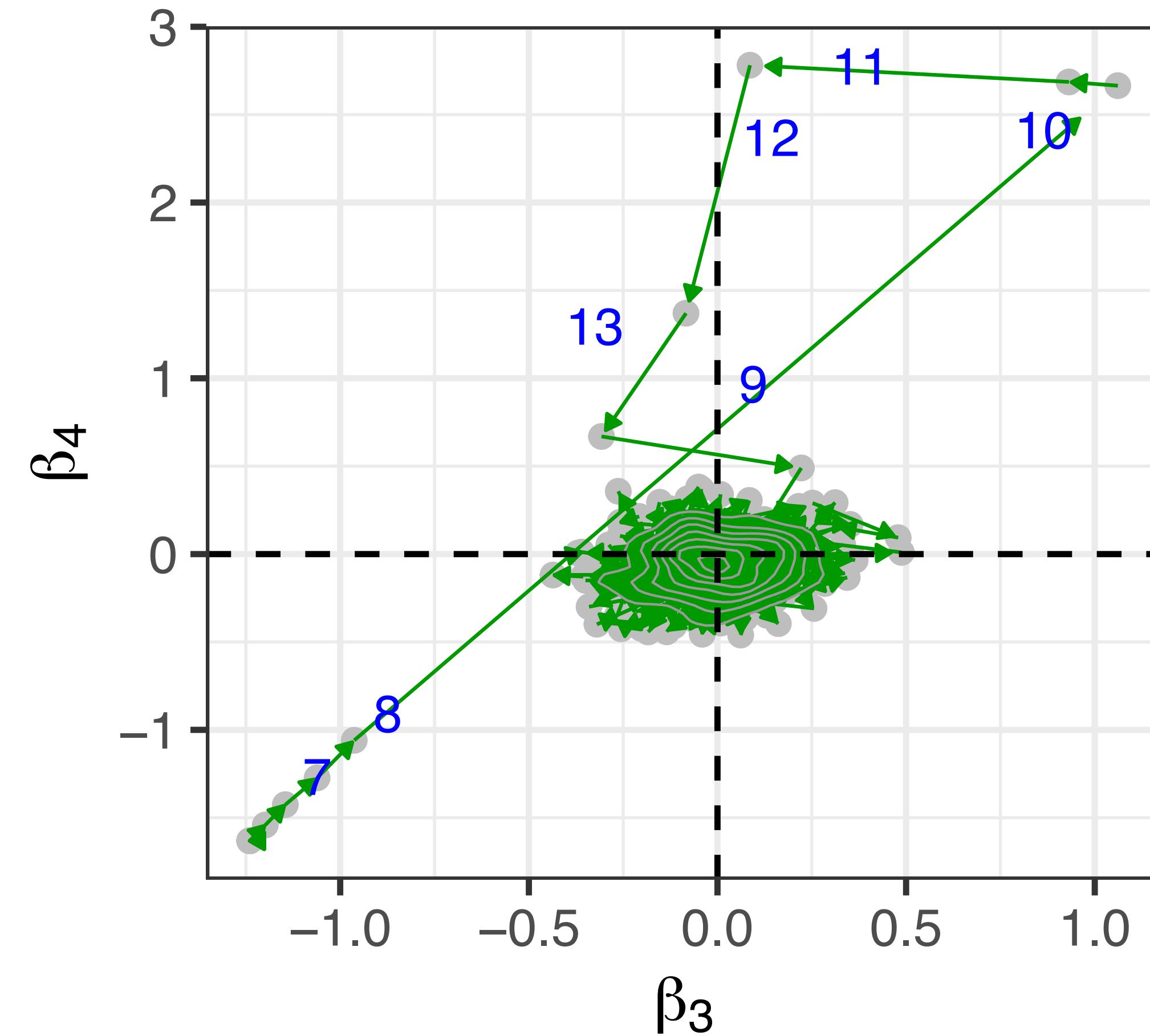


# Compare true non-zero vs. zero coefficients

800 MCMC iterations



800 MCMC iterations



# Today's lecture: Bayesian, PGM, Causality

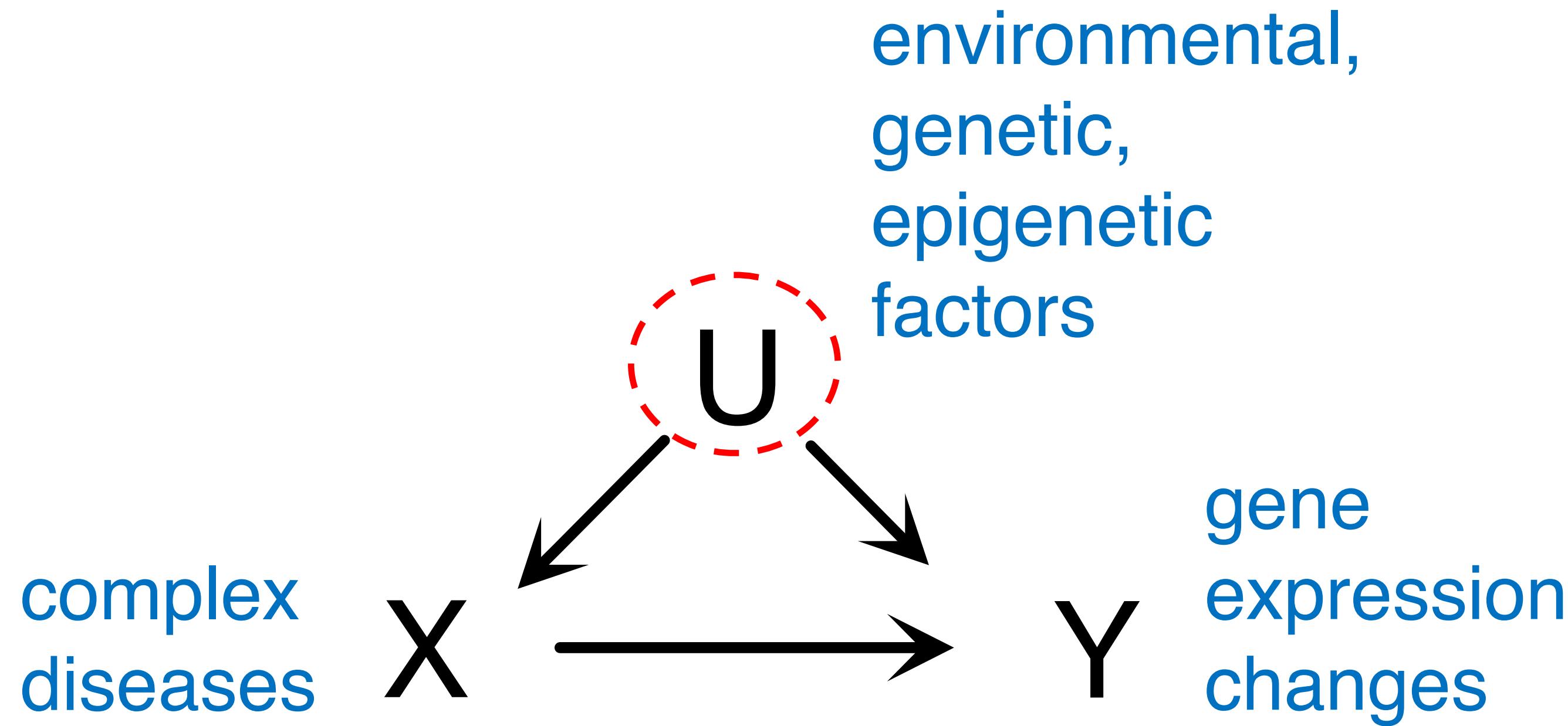
- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

# Causal inference came back!

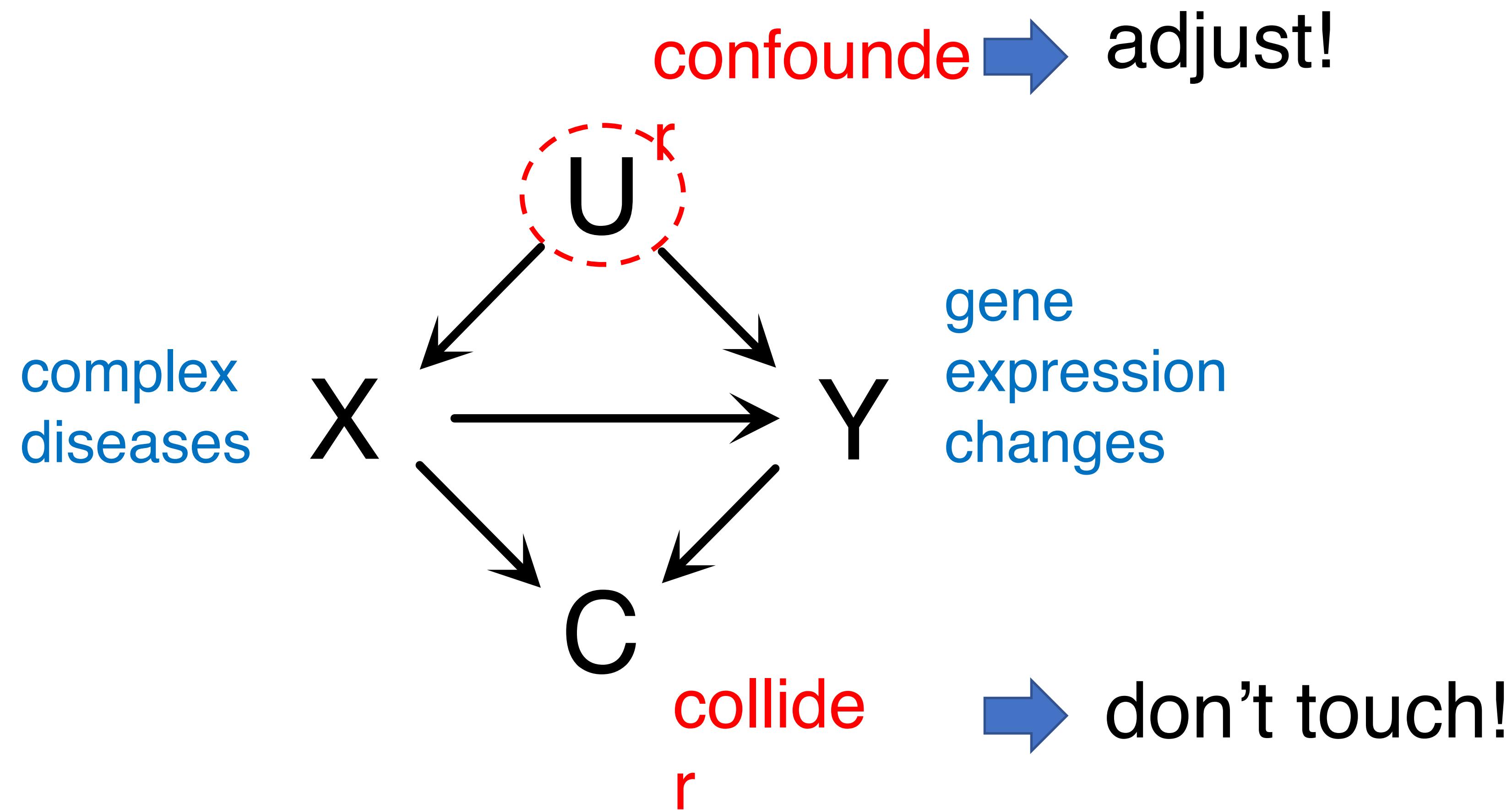


In the experimental design part of Lecture #3, we wanted to discuss causality.

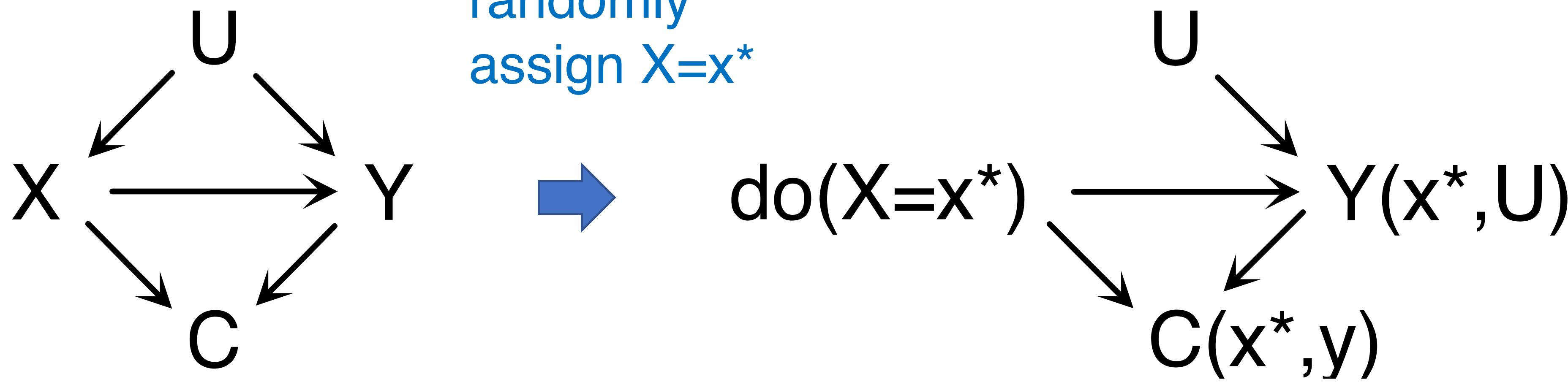
# What if we have very little knowledge of confounding effects?



# What if we don't know how to distinguish confounder vs. collider?



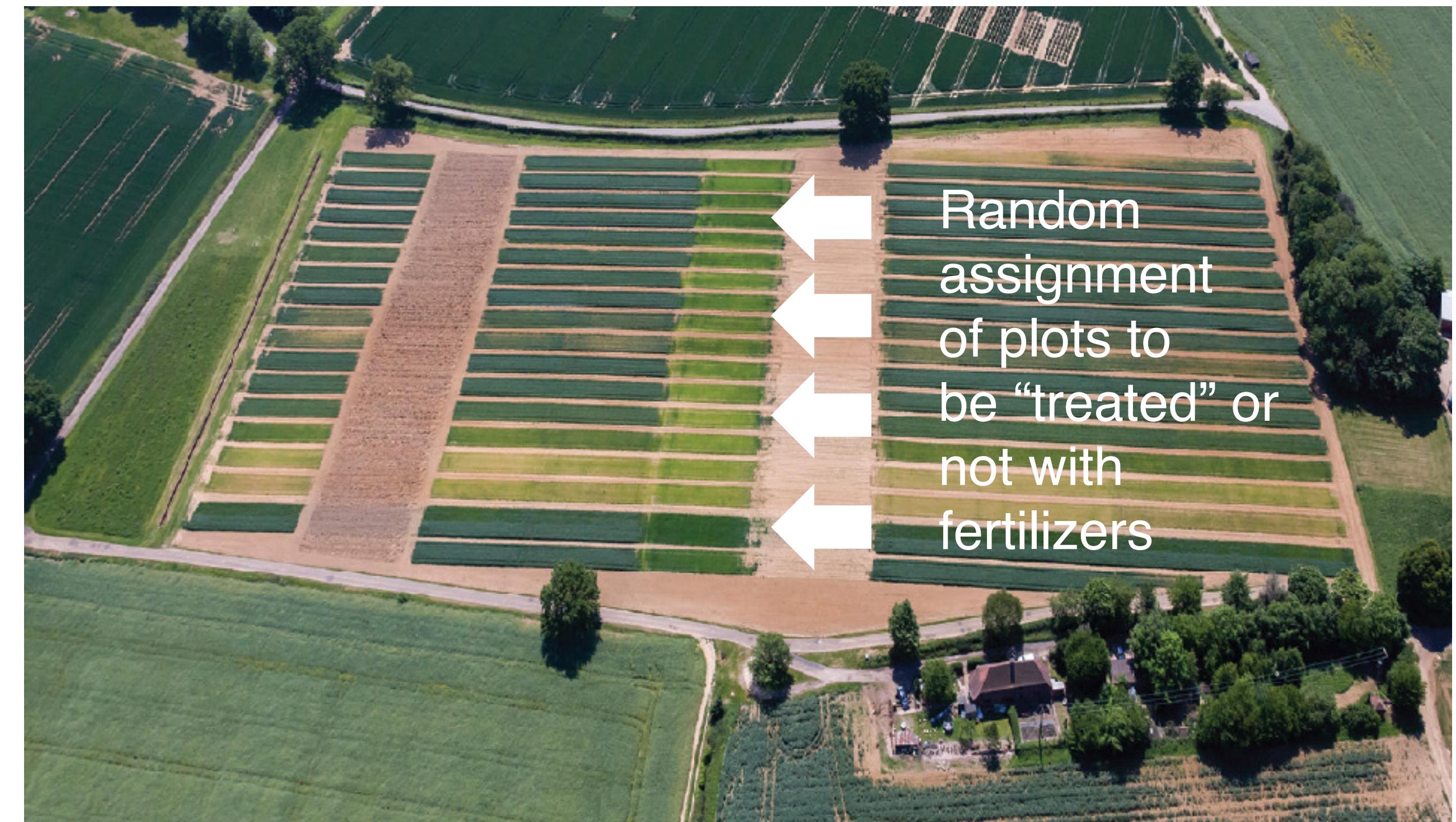
# Randomized Control Trial: the gold standard experiment for causal discovery



experiment/  
intervention to  
change the  
underlying causal  
structure

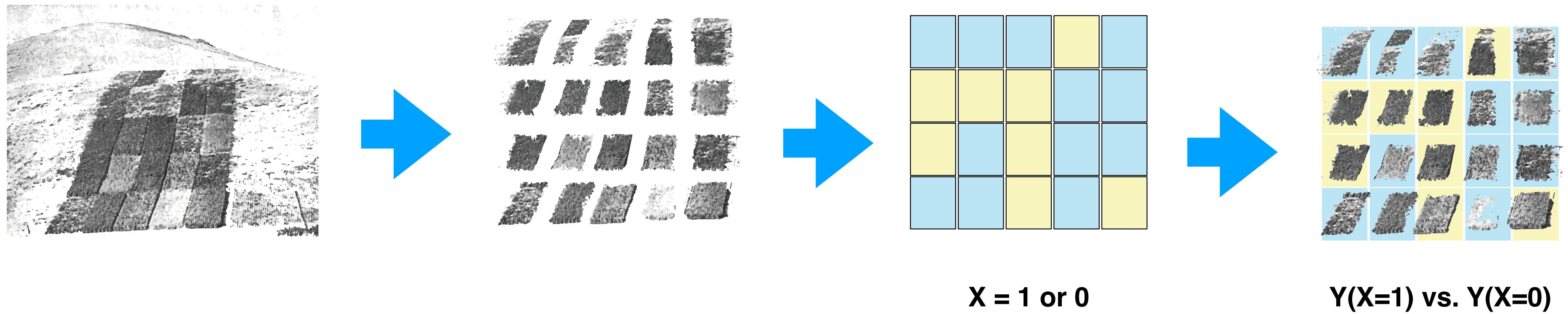
# The arrangement of Field Experiments (RA Fisher, 1922-1926)

- Randomly select plots to treat manure ( $X=1$ ) or not ( $X=0$ ).
- Measure the yield of crops ( $Y$ ).



Langkjær-Bain, “Where the seeds of modern statistics were sown” (2

# The first Randomized Control Trial experiments



# Today's lecture: Bayesian, PGM, Causality

- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

**Again, the R codes for the examples are all available online**

**[https://github.com/STAT540-UBC/lectures/blob/main/lect13-causality\\_bayesian/causality.Rmd](https://github.com/STAT540-UBC/lectures/blob/main/lect13-causality_bayesian/causality.Rmd)**

# A graphical language for causal inference (used throughout this lecture)

## Causal relationship



$\Leftrightarrow$

“X causes Y”

$\Rightarrow$

?

$P(Y|X)$

Very frequently not  
true

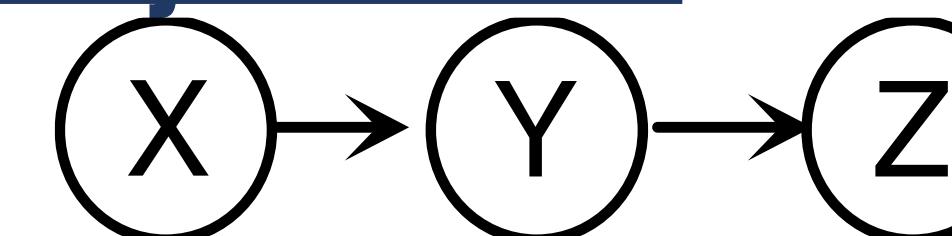
Probabilistic  
graphical  
models care  
about  
dependency

## Causal path & reachability

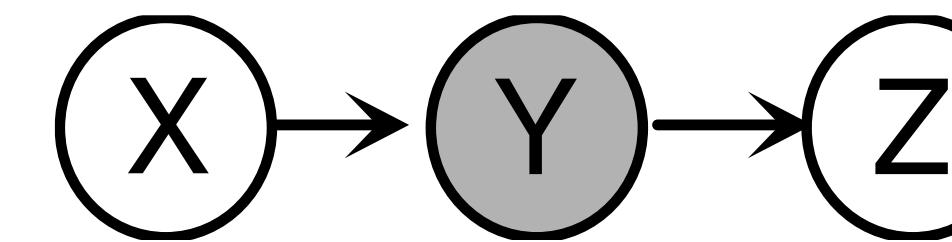


The effect of X can influence Z (flow)

## Conditioning/adjustment

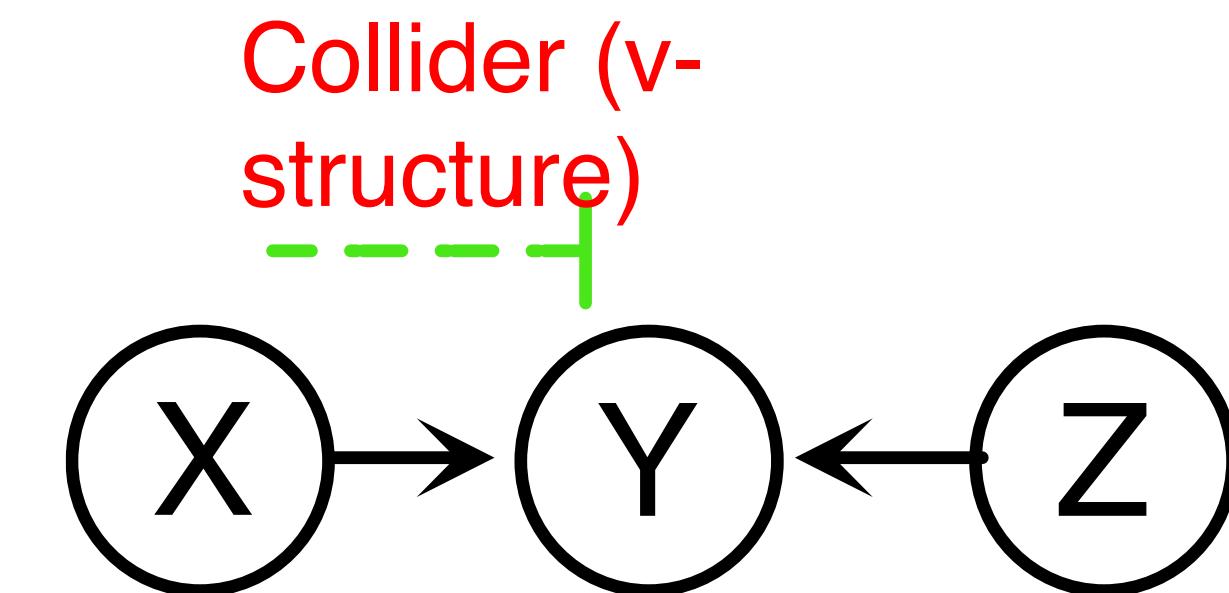
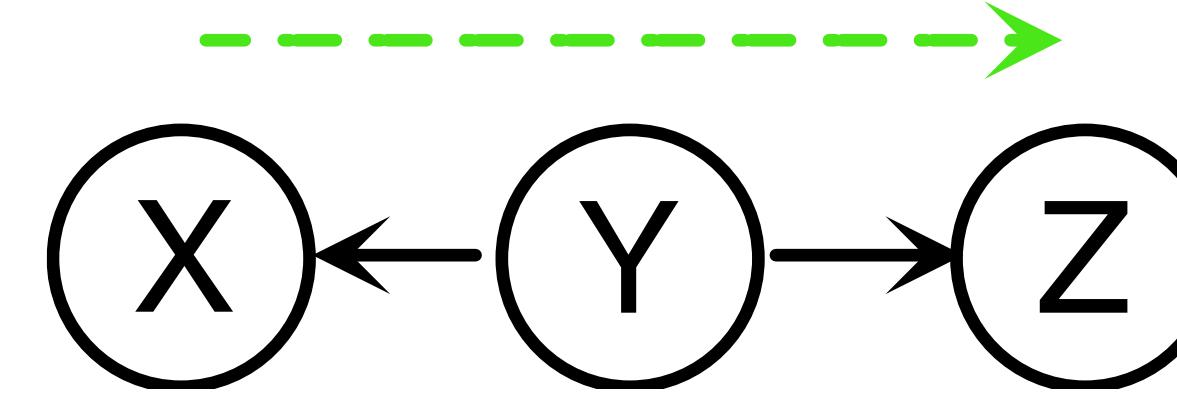
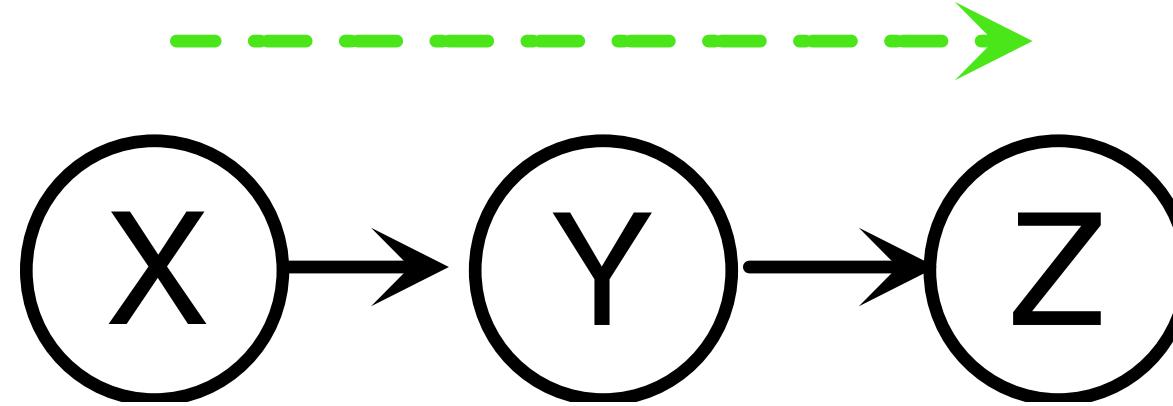


Open or no circle: not conditioned



Closed circle: conditioned  
(e.g., setting/given  $Y = y^*$ )

# d-separation: testing conditional independency (flow vs. no flow)

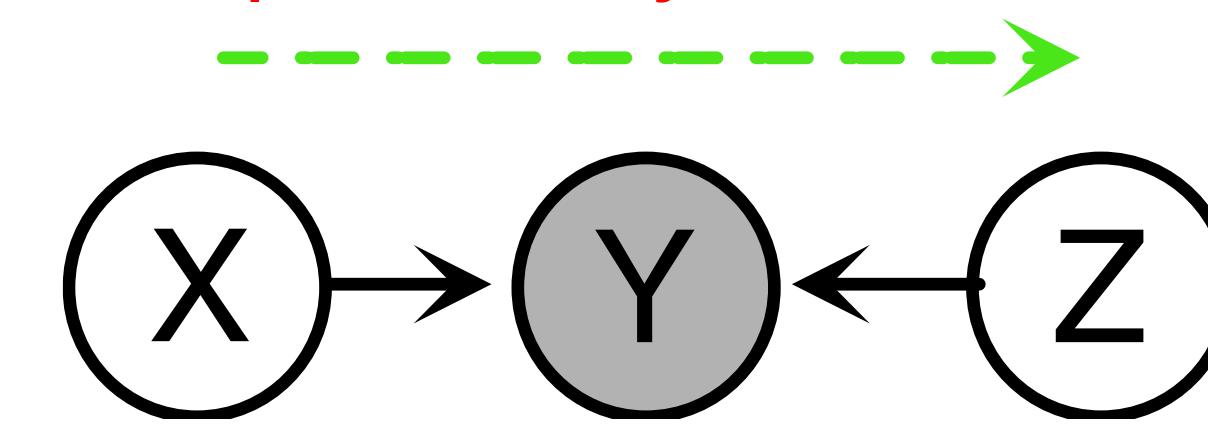
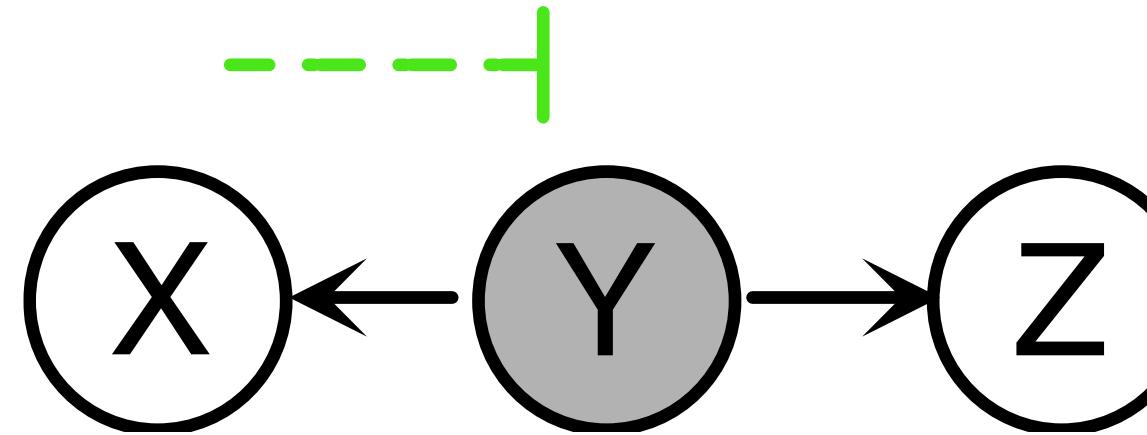
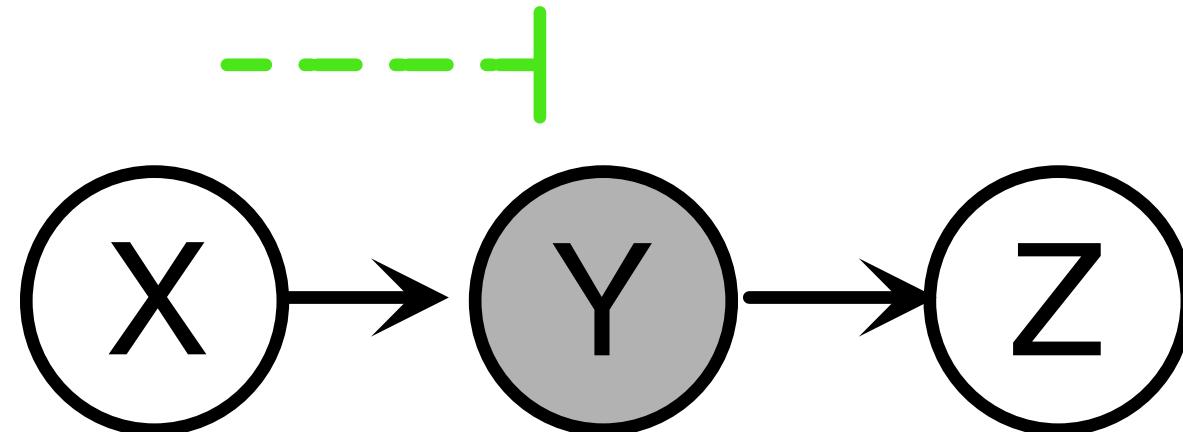
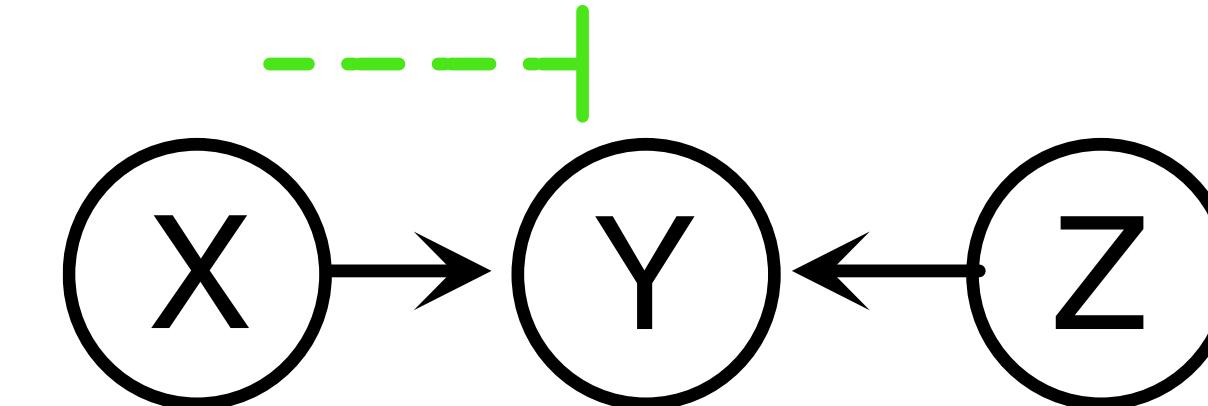
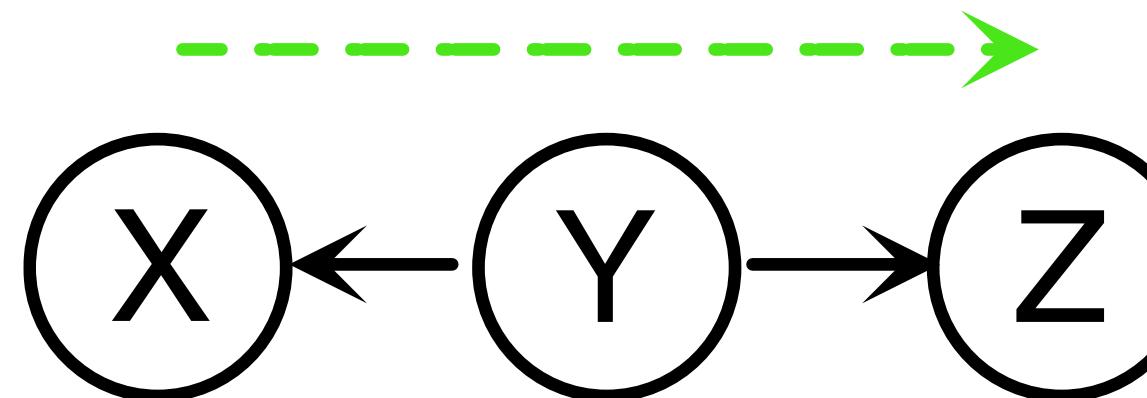
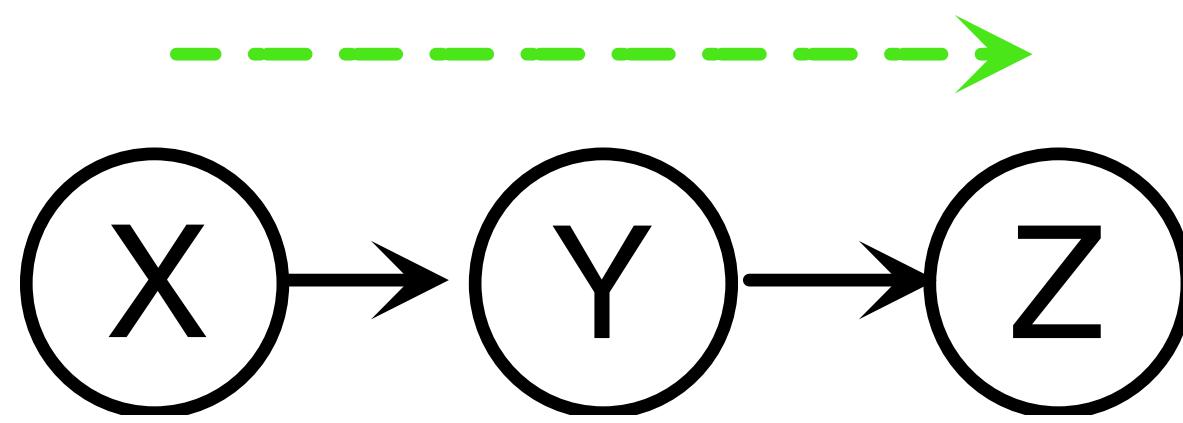


- “The outcome of Z depends on the effect of X”
- Because the outcome of Z depends on the outcome of Y
- And the outcome of Y depends on X

- “The outcome of Z depends on the effect of Y”
- Since “the outcome of X depends on Y”
- The outcomes of X and Z are dependent on each other

- The outcome of Z doesn’t depend on Y or X
- (Neither X nor Y causes Z)

# d-separation: testing conditional independency (flow vs. no flow)



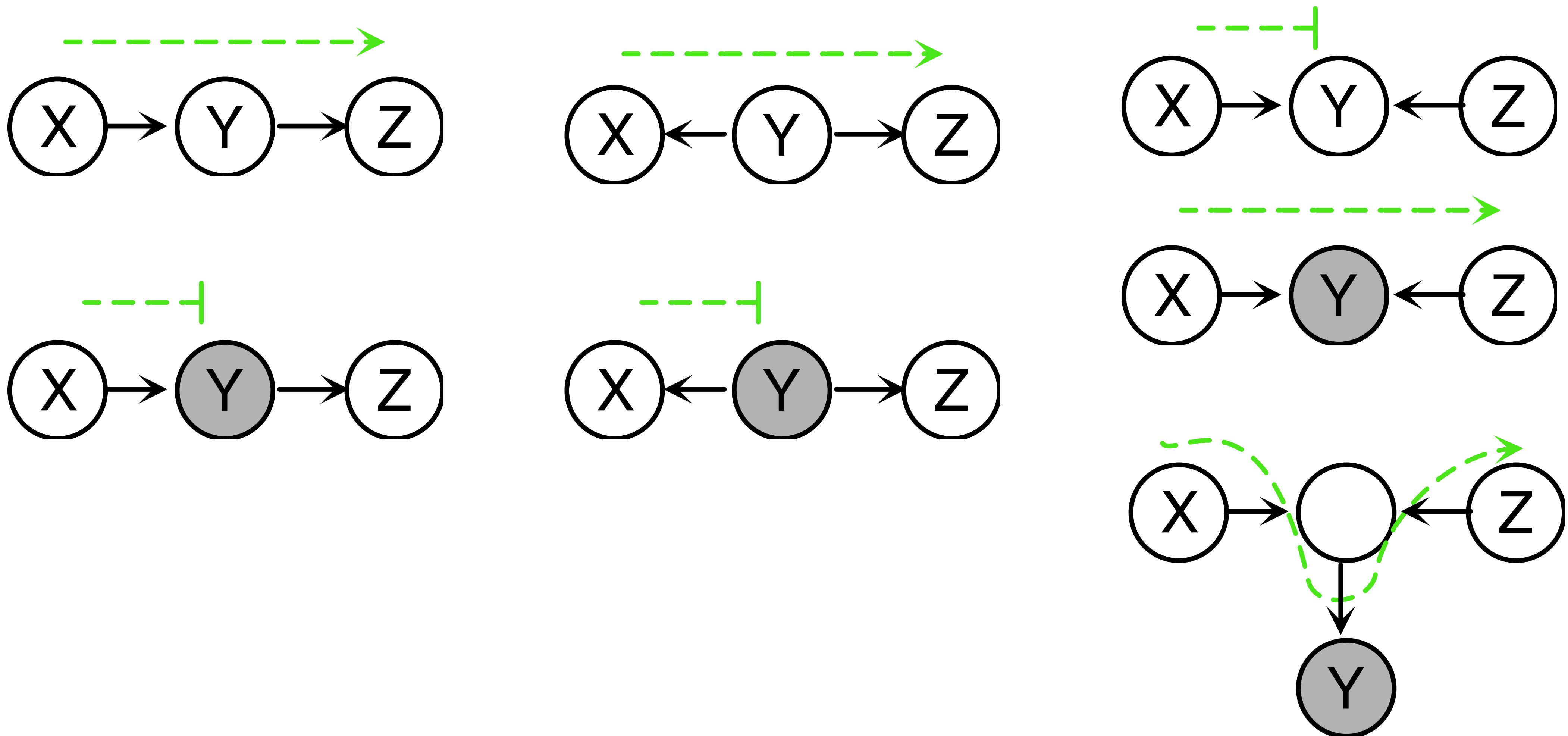
- Given  $Y=y^*$ ,  $Z$  doesn't depend on  $X$

- Given  $Y=y^*$ ,  $Z$  only depends on  $Y$

"Explain away"

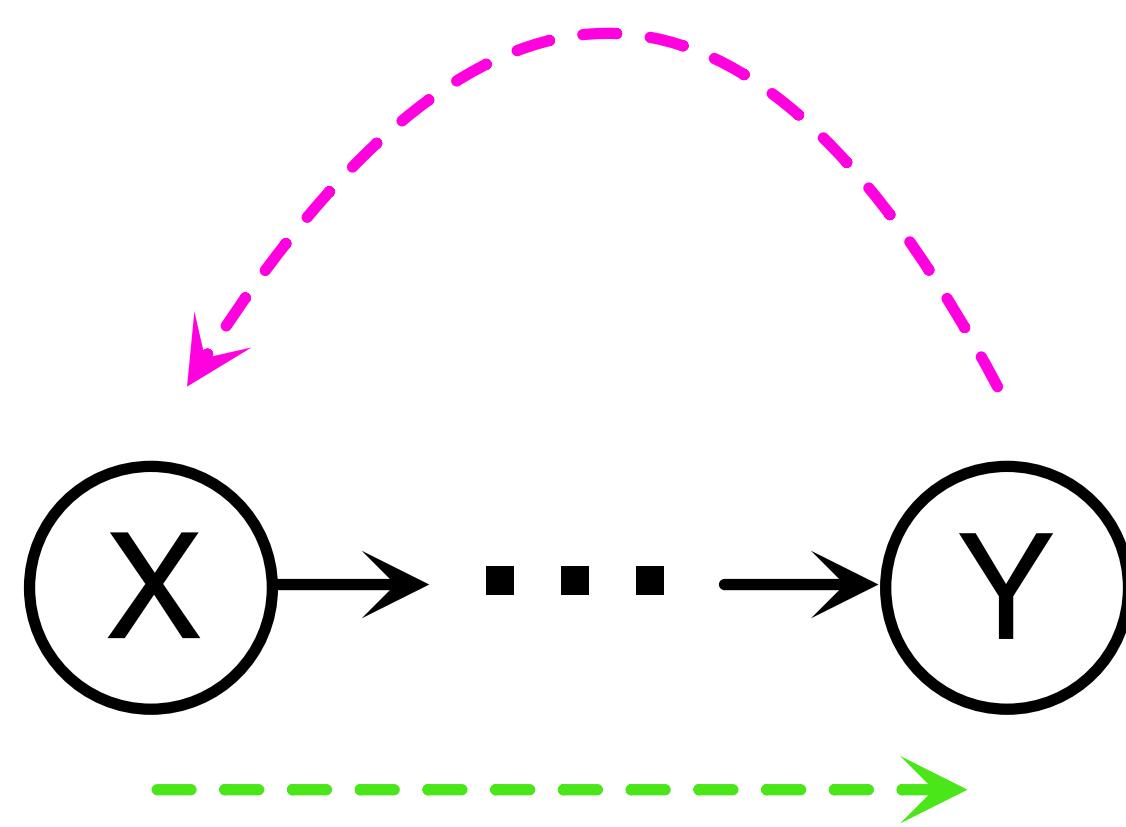
- Either  $X$  or  $Z$  can explain  $\{Y = y^*\}$
- If  $\{X = x^*\}$  can explain  $\{Y = y^*\}$ ,  $Z$  doesn't explain  $Y$
- vice versa

# d-separation: testing conditional independency (flow vs. no flow)



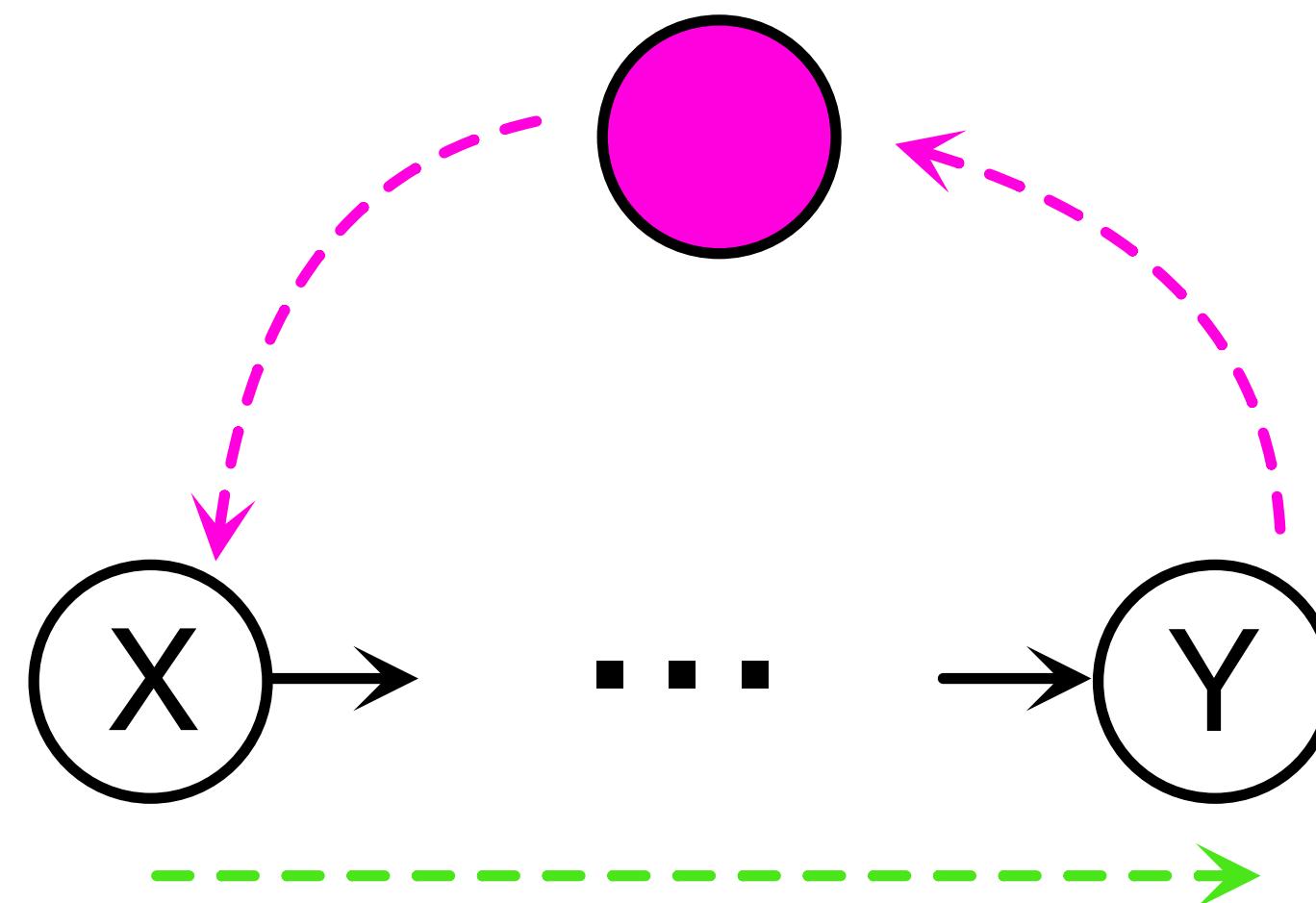
## “Backdoor” exercise

**Question:** Which nodes should be “conditioned” and/or “adjusted” to block a reverse path from  $Y$  to  $X$ ?

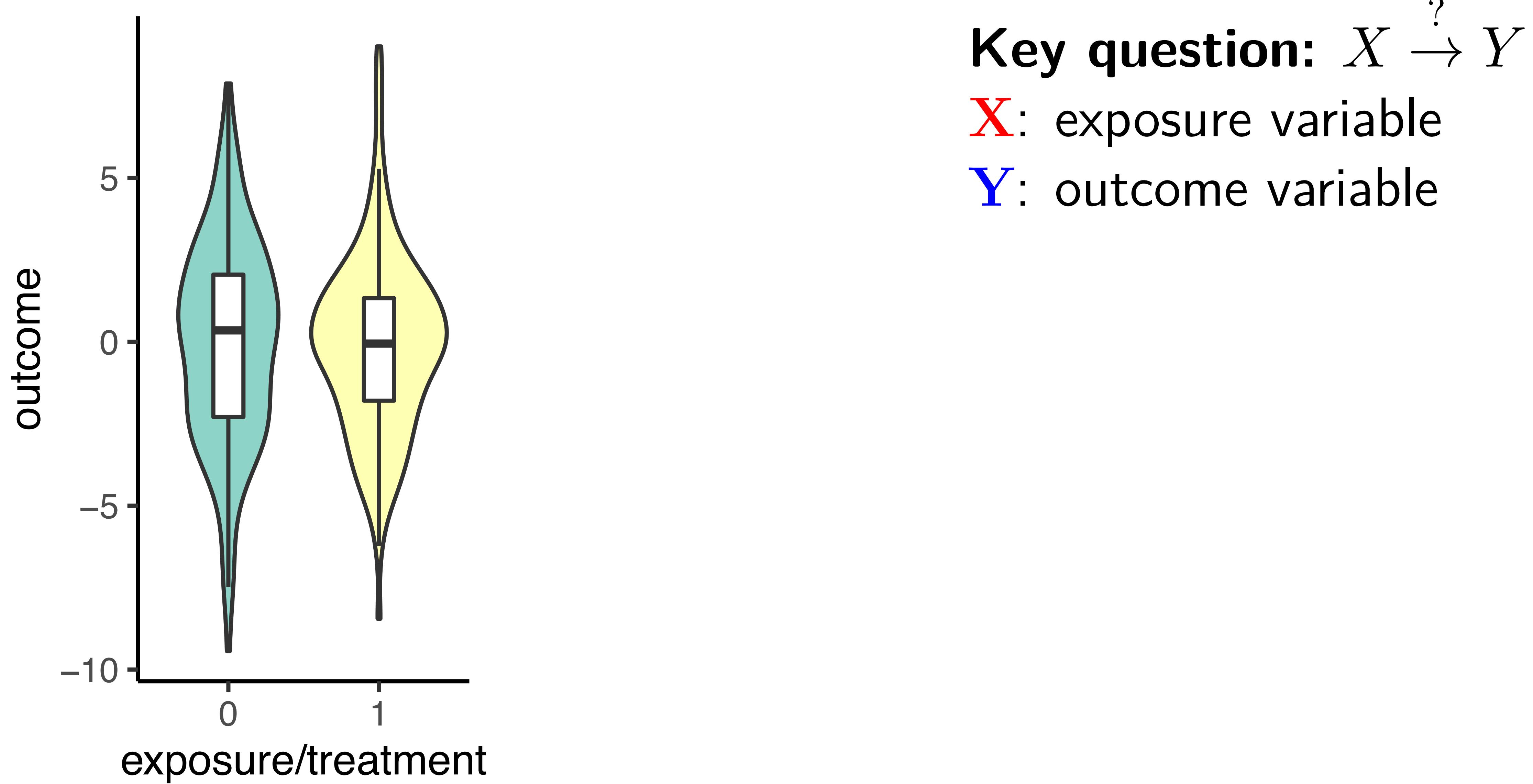


## “Backdoor” exercise

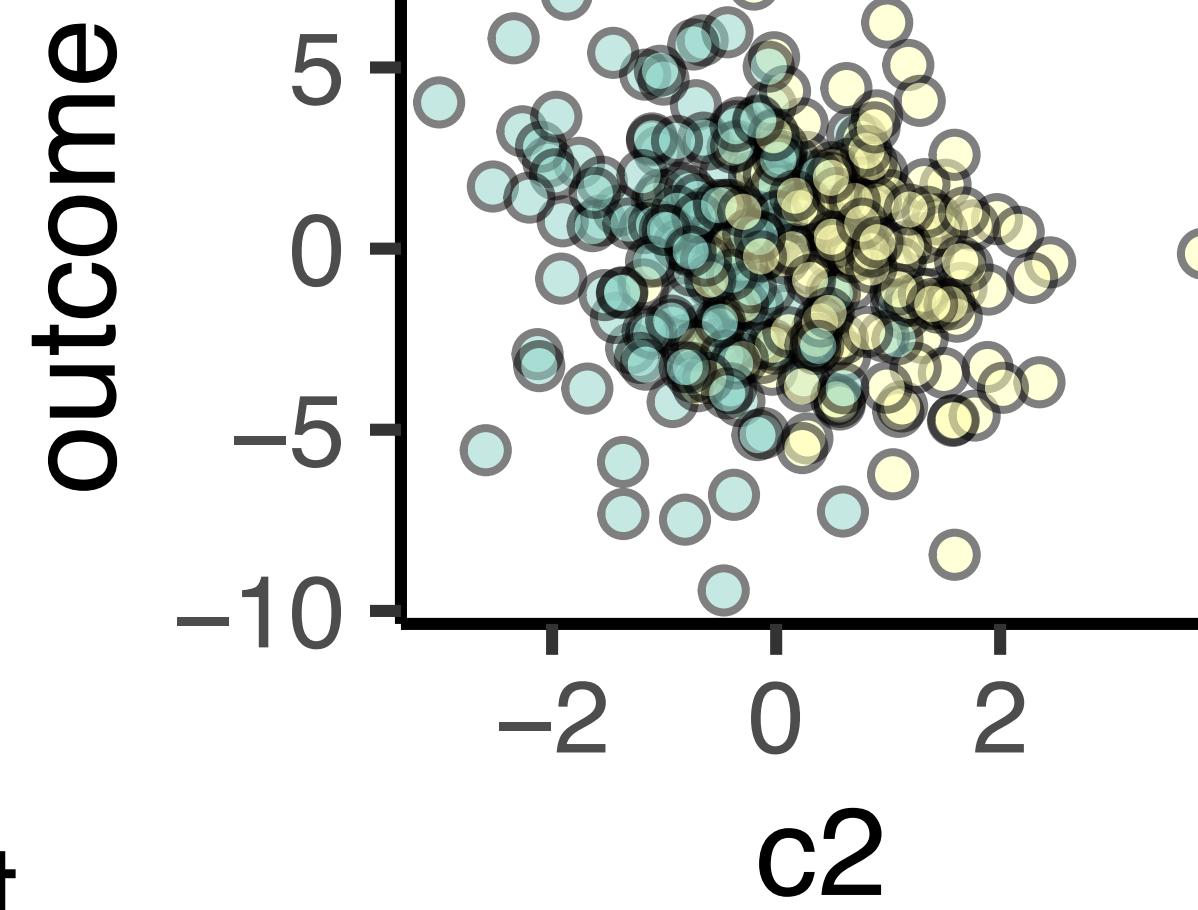
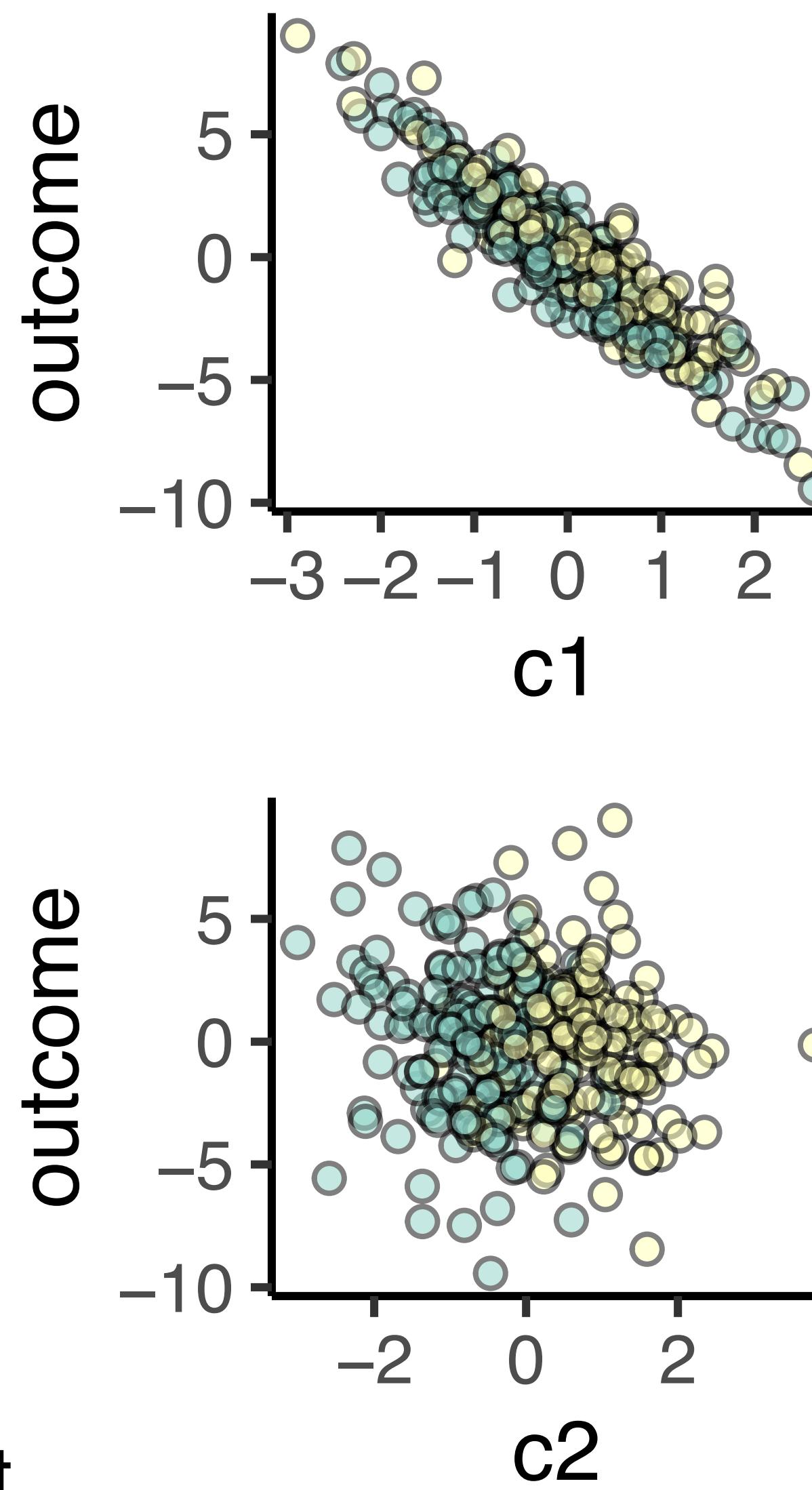
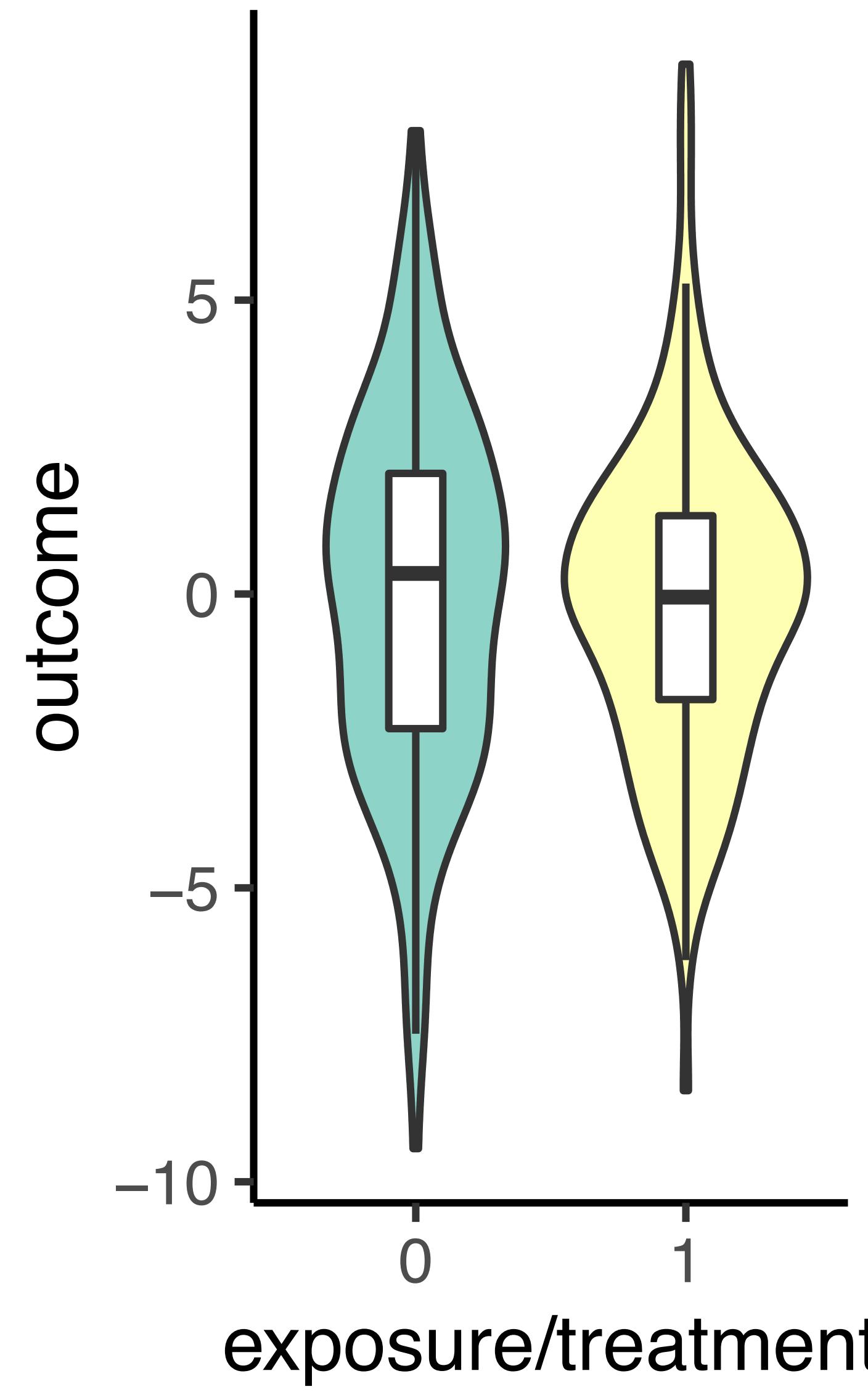
**Question:** Which nodes should be “conditioned” and/or “adjusted” to block a reverse path from  $Y$  to  $X$ ?



# A working example: confounder adjustment in case-control study



# A working example: confounder adjustment in case-control study



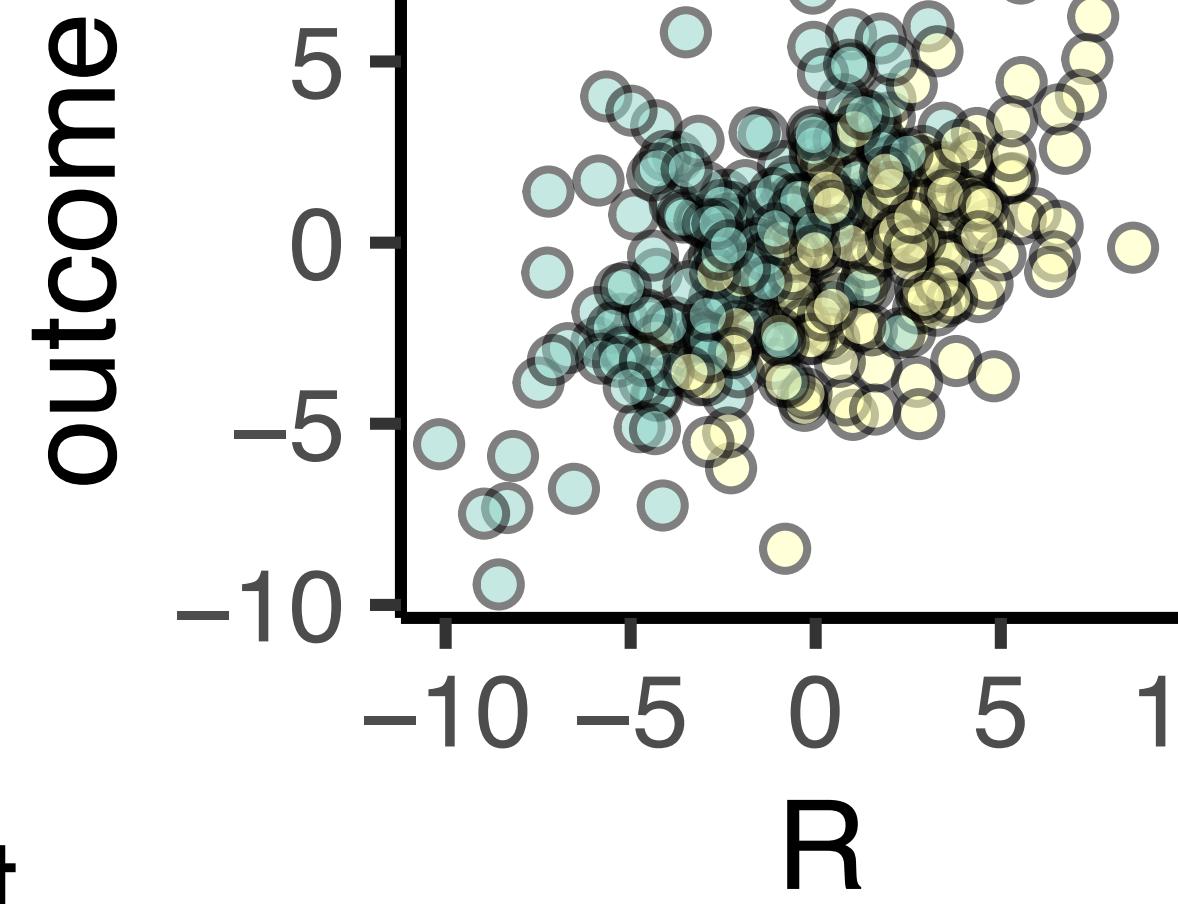
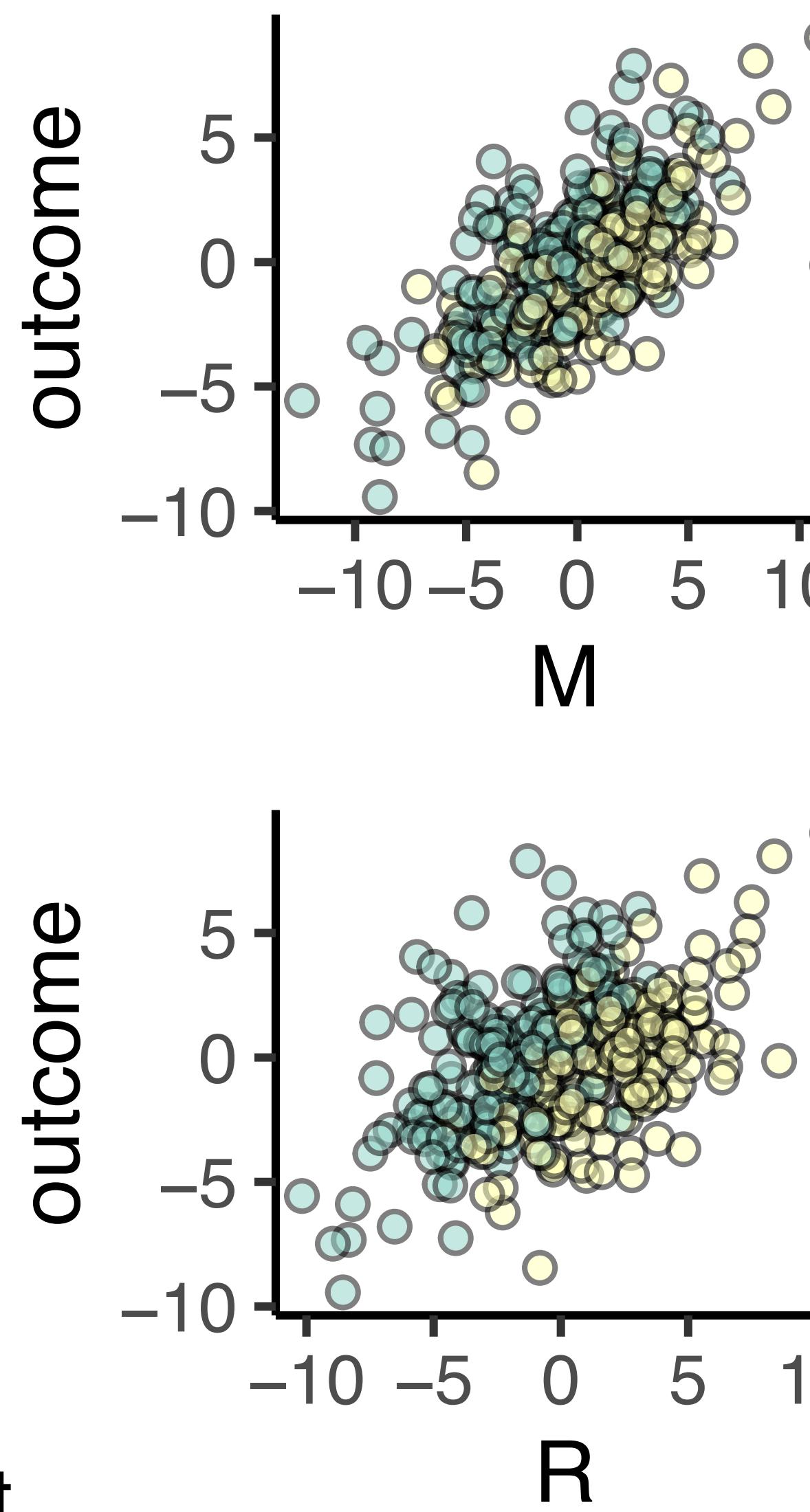
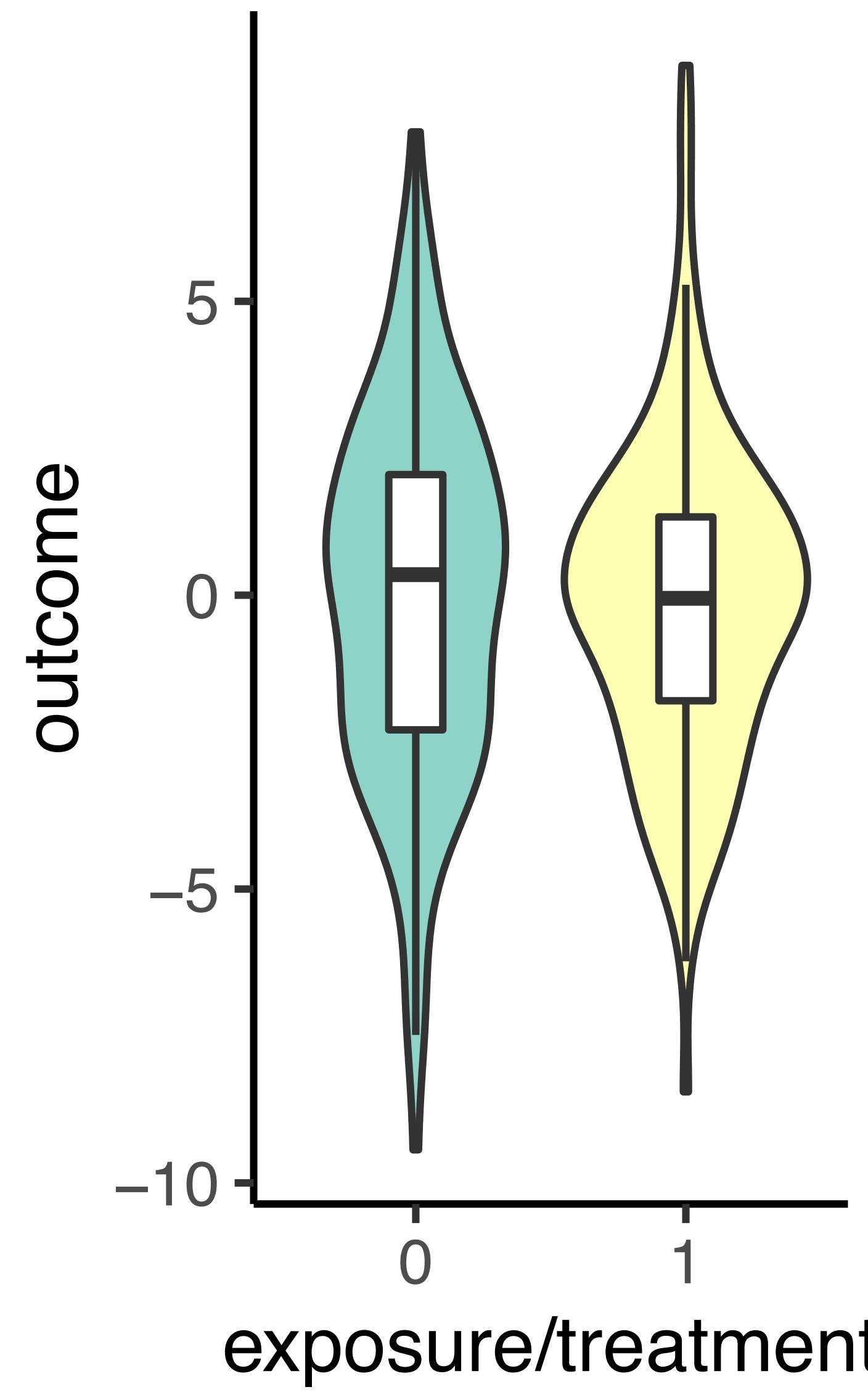
**Key question:**  $X \xrightarrow{?} Y$

**X:** exposure variable

**Y:** outcome variable

**C<sub>1</sub> and C<sub>2</sub>:** covariates

# A working example: confounder adjustment in case-control study



**Key question:**  $X \xrightarrow{?} Y$

**X:** exposure variable

**Y:** outcome variable

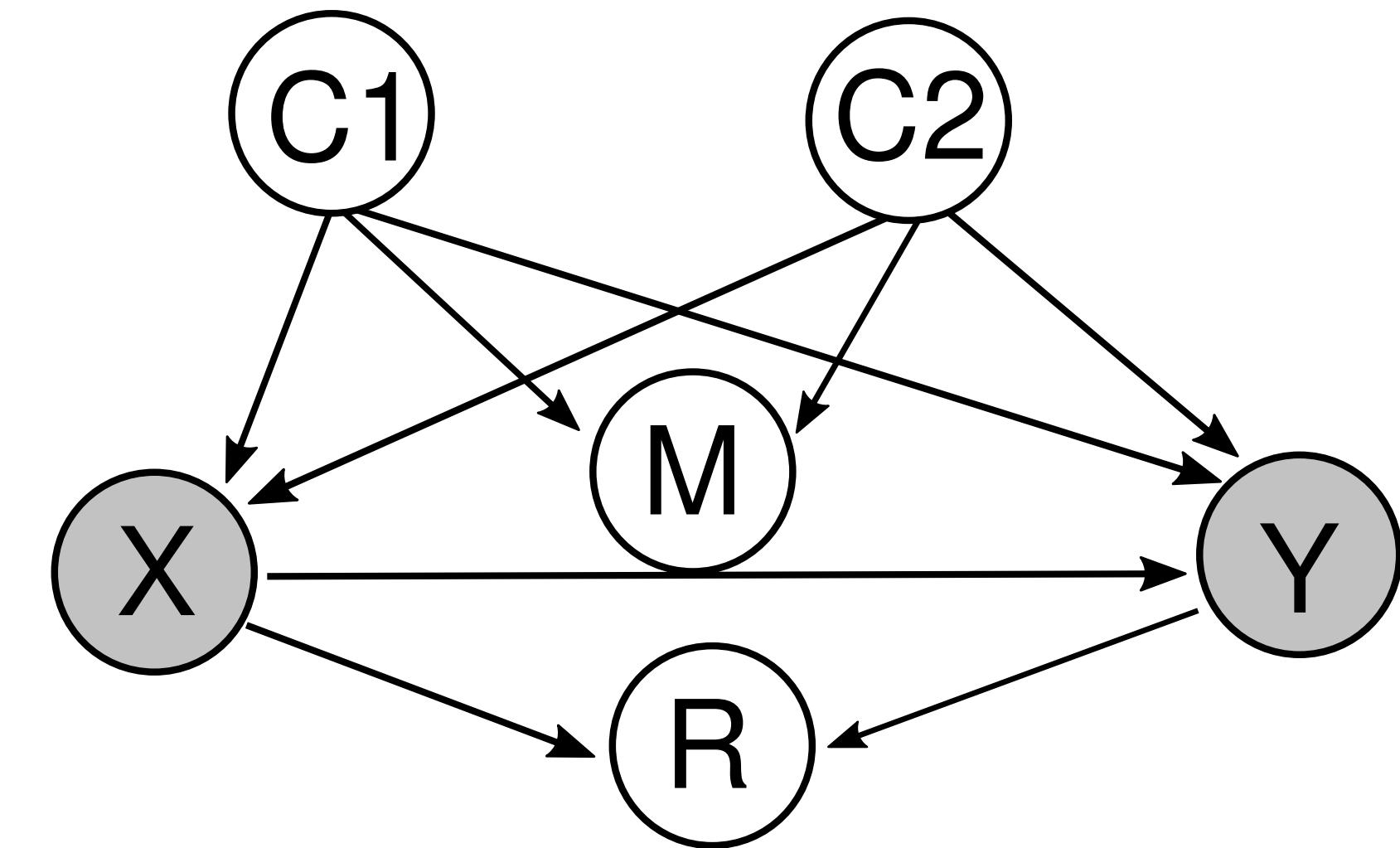
**C<sub>1</sub>** and **C<sub>2</sub>**: covariates

**M:** other covariate

**R:** other covariate

# Causal inference with a graphical model

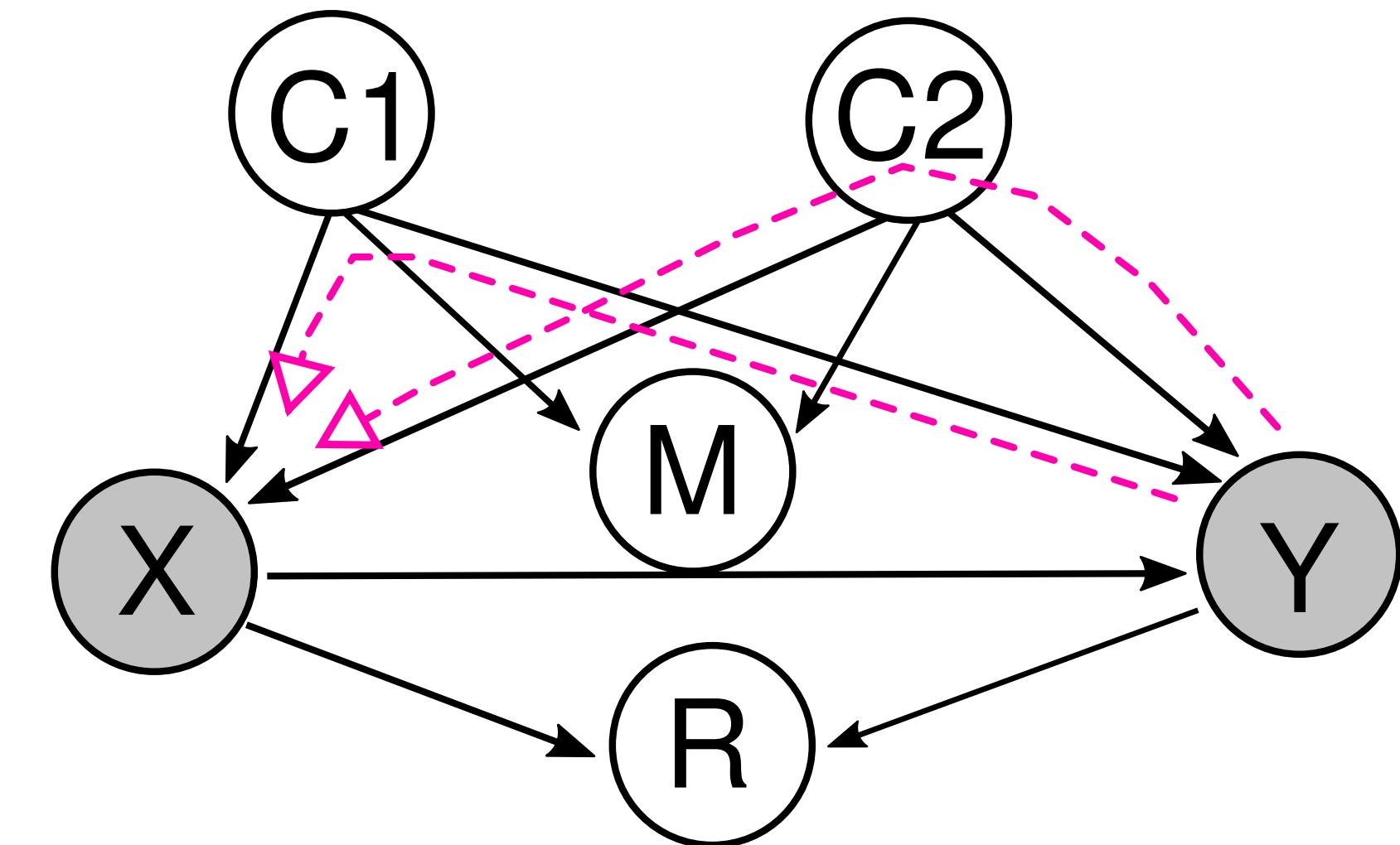
1. Build a causal structural model



What are potential backdoors?

# Causal inference with a graphical model

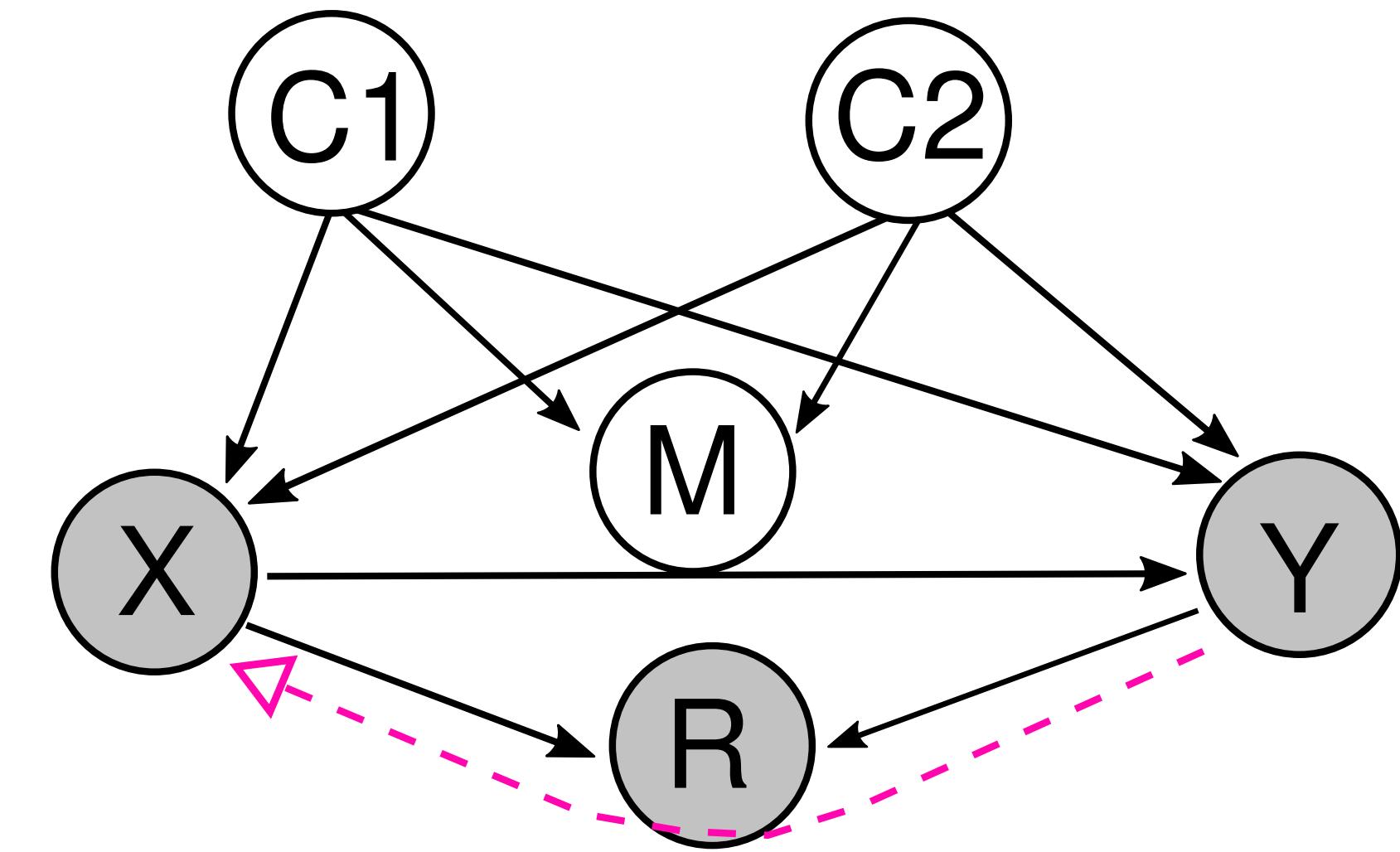
1. Build a causal structural model
2. Identify "back-door" paths/variables (*closing*  $Y \rightarrow X$ , *while opening*  $X \rightarrow Y$ )



How do we close them?

# Causal inference with a graphical model

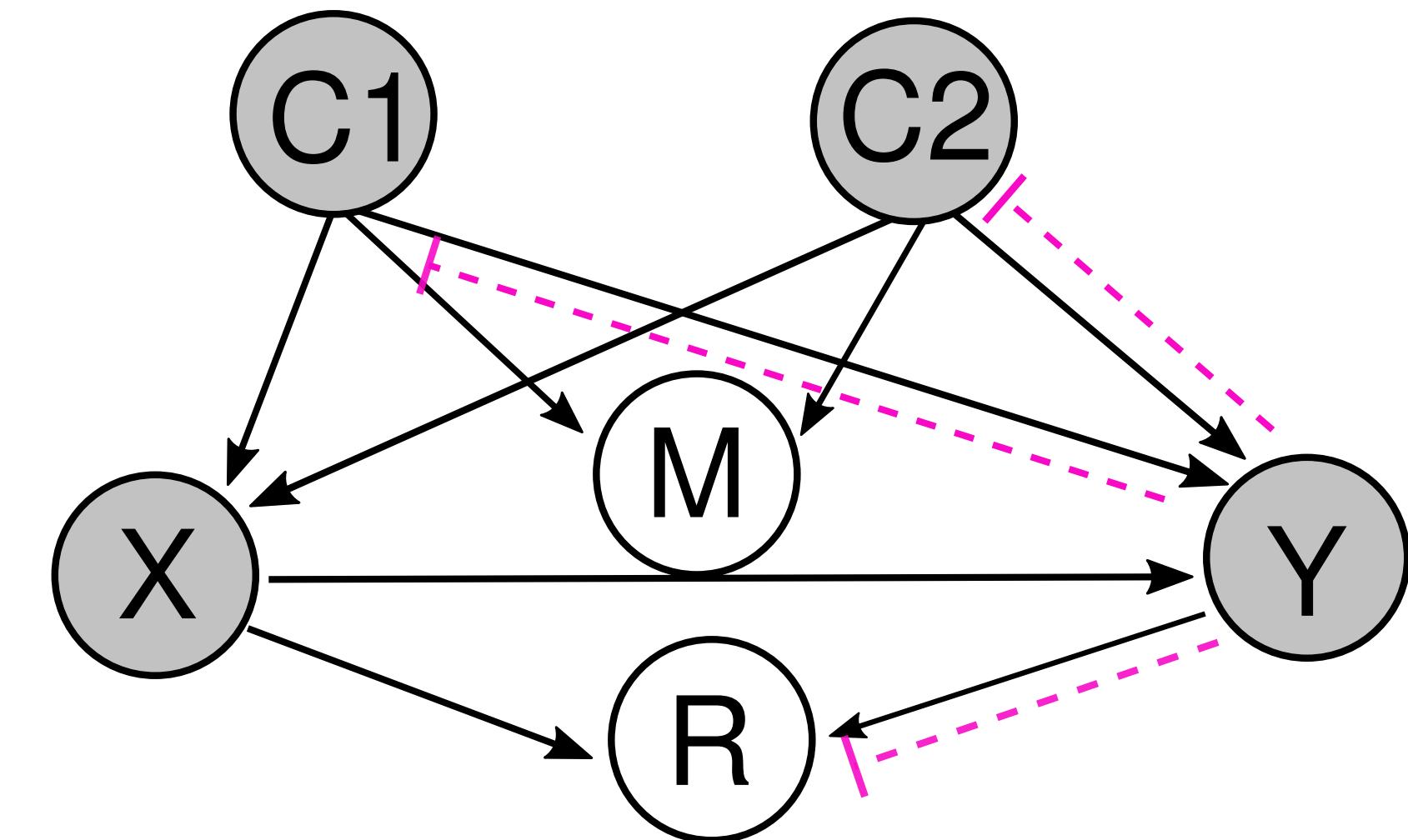
1. Build a causal structural model
2. Identify "back-door" paths/variables (*closing*  $Y \rightarrow X$ , *while opening*  $X \rightarrow Y$ )



What about this?

# Causal inference with a graphical model

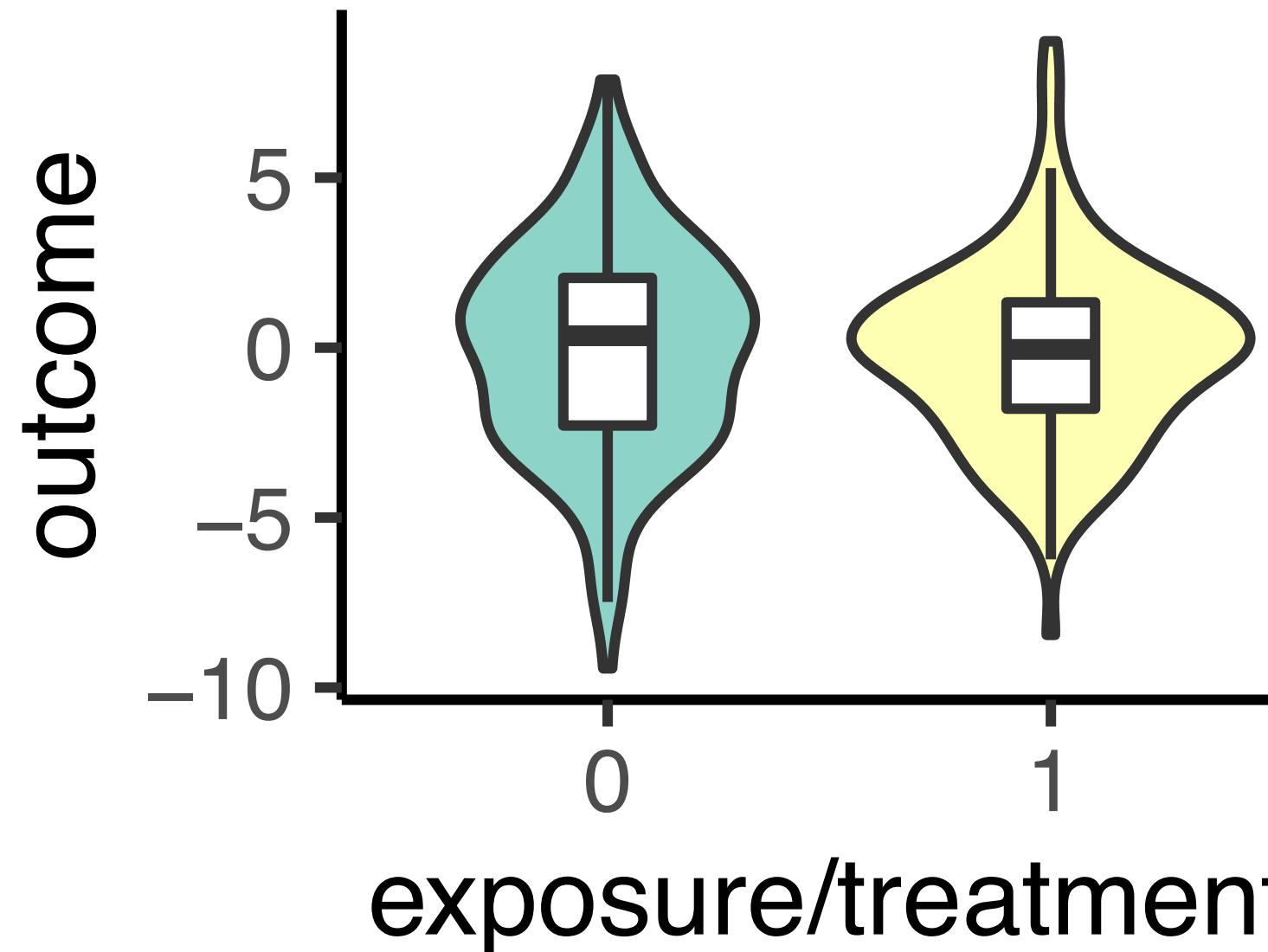
1. Build a causal structural model
2. Identify "back-door" paths/variables (*closing*  $Y \rightarrow X$ , *while opening*  $X \rightarrow Y$ )



Is this enough?

# Causal inference with a graphical model

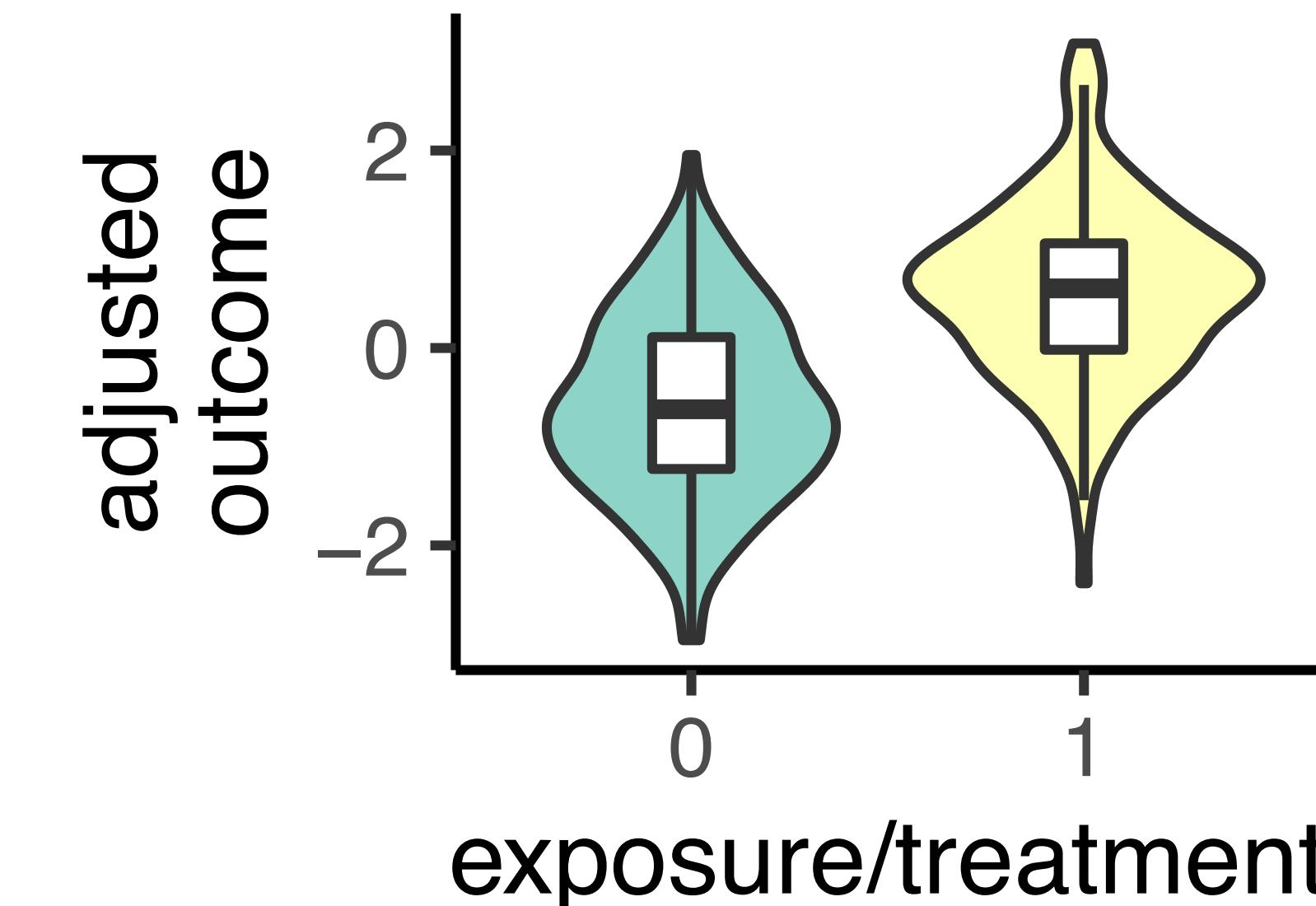
1. Build a causal structural model
2. Identify "back-door" paths/variables (*closing*  $Y \rightarrow X$ , *while opening*  $X \rightarrow Y$ )
3. Adjust "back-door" variables
4. Estimate causal effects



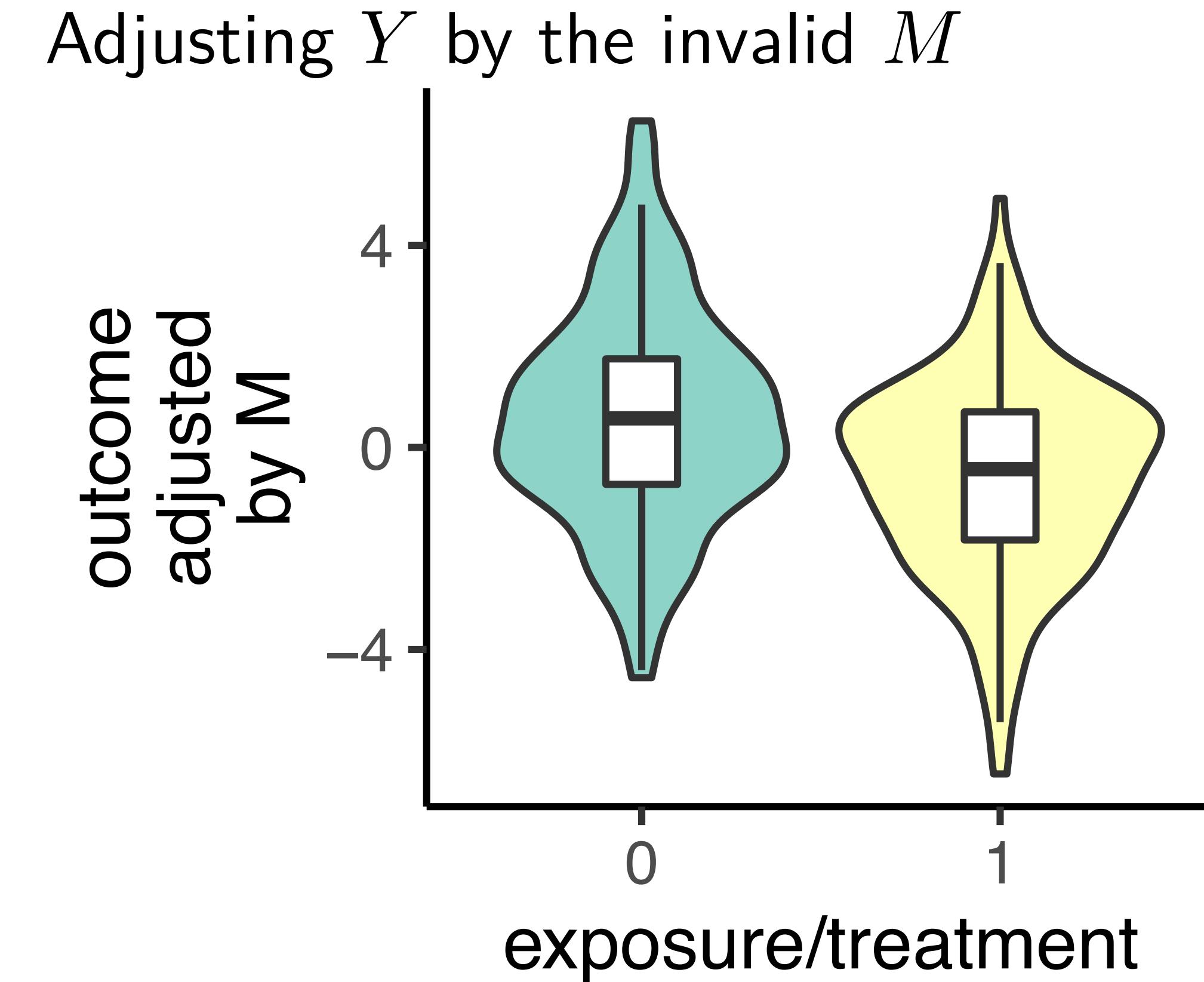
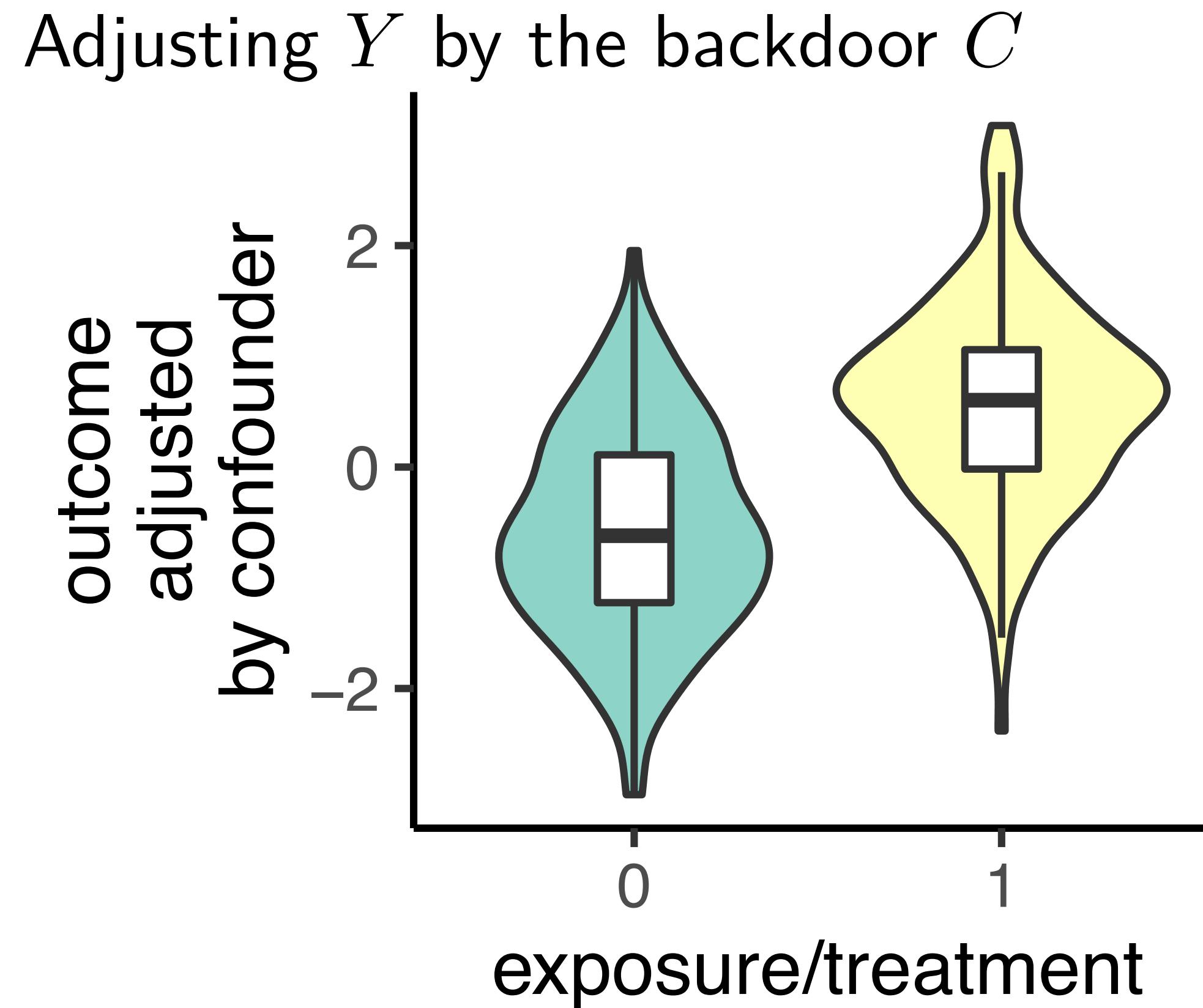
$$Y \leftarrow Y - \sum_{k=1}^2 C_k \hat{\beta}_k$$

which approximates

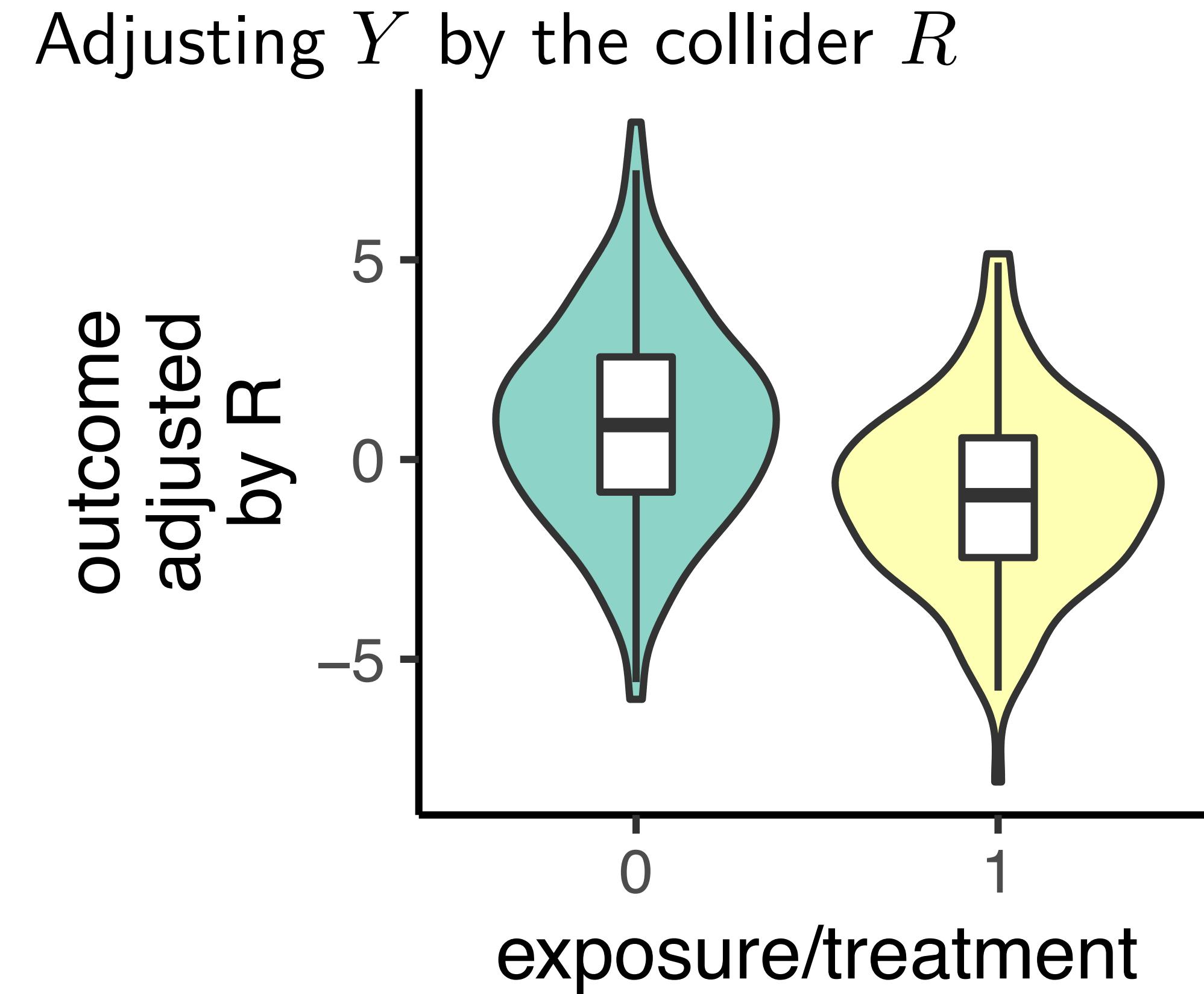
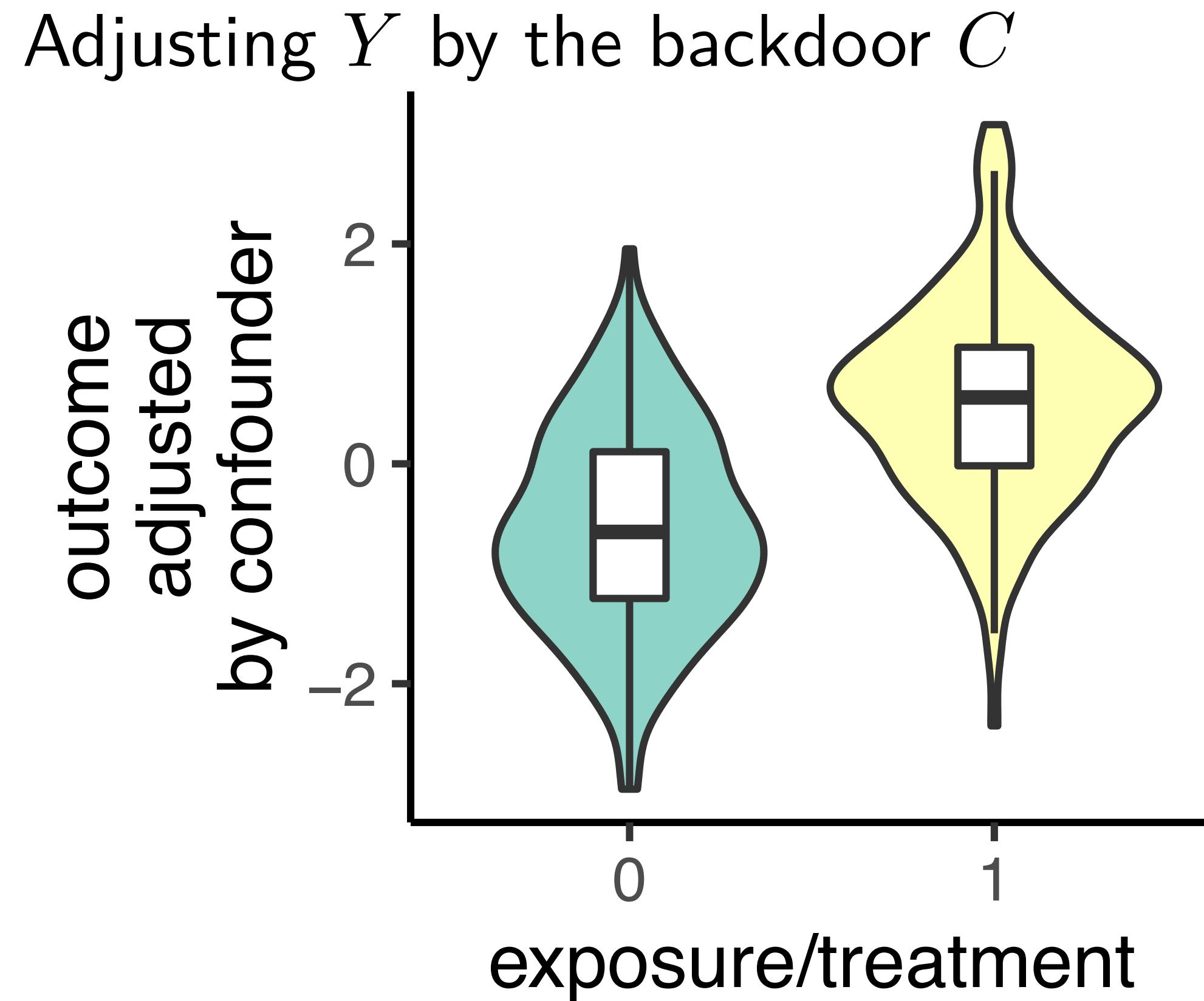
$$p(Y|X) = \int_C p(Y|X, C)p(C)dC$$



# What would happen if we close a wrong “backdoor” variable?



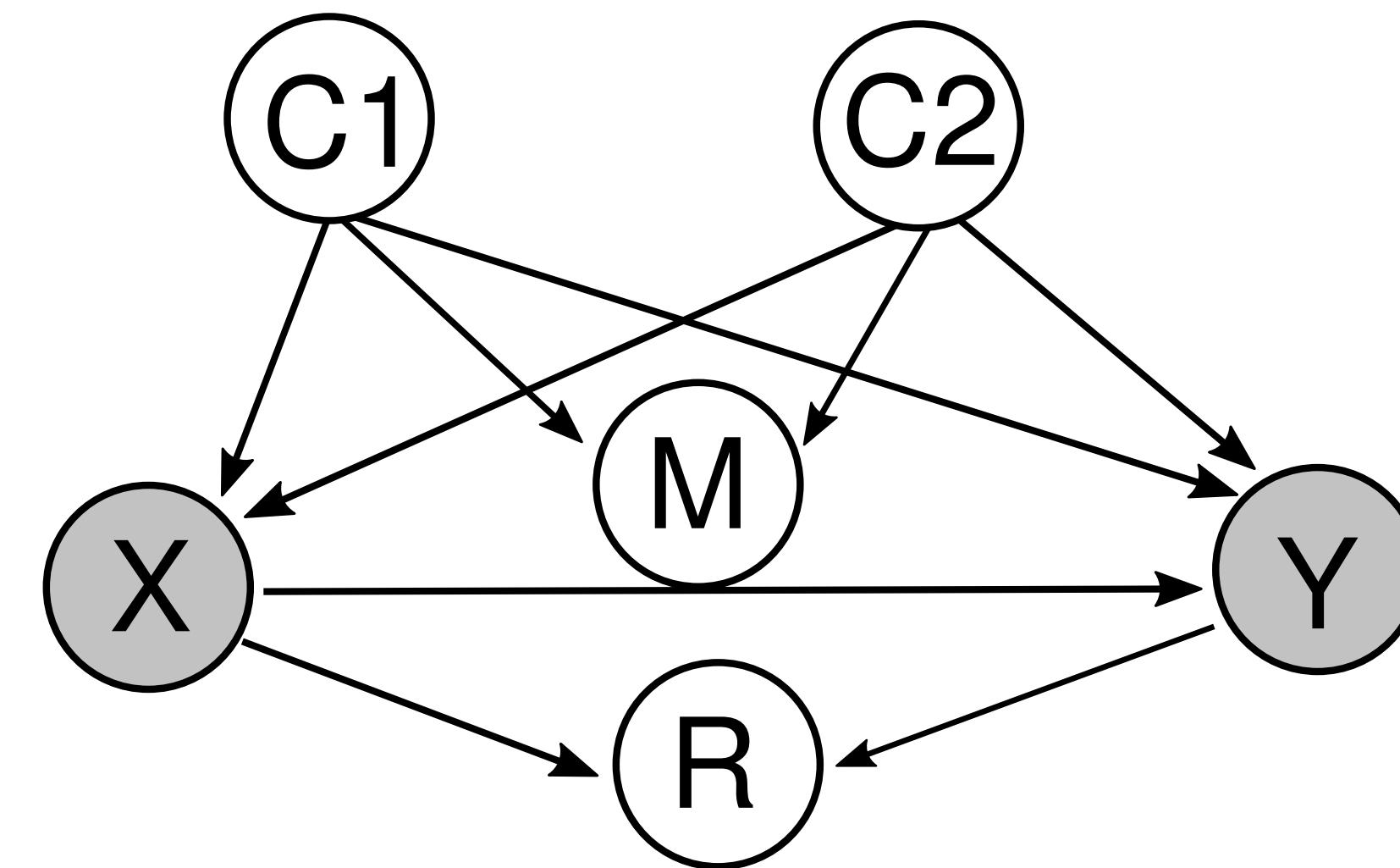
# What would happen if we close a wrong “backdoor” variable?



# Today's lecture: Bayesian, PGM, Causality

- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

The same example with causal relationship from  $X$  to  $Y$

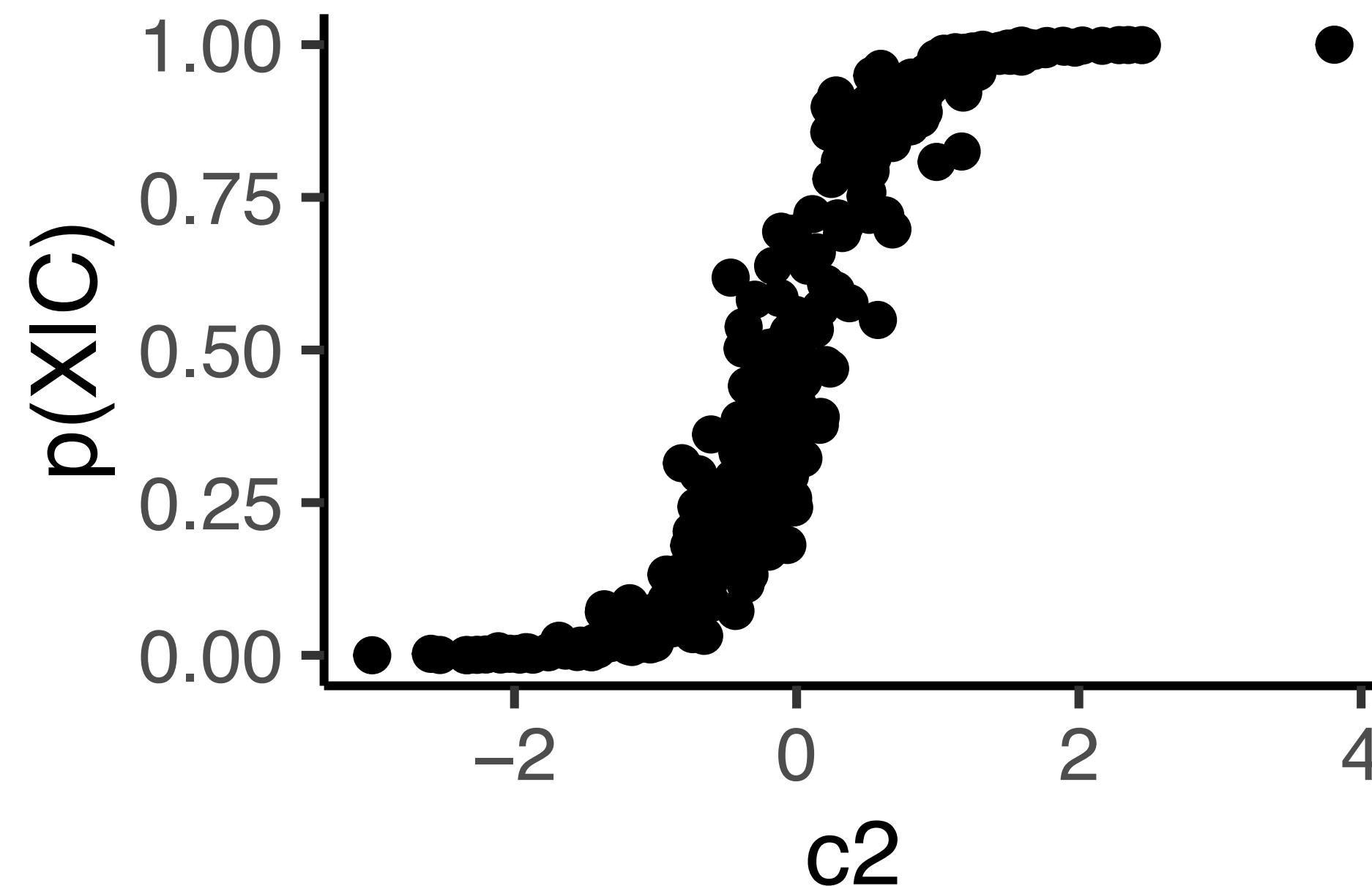
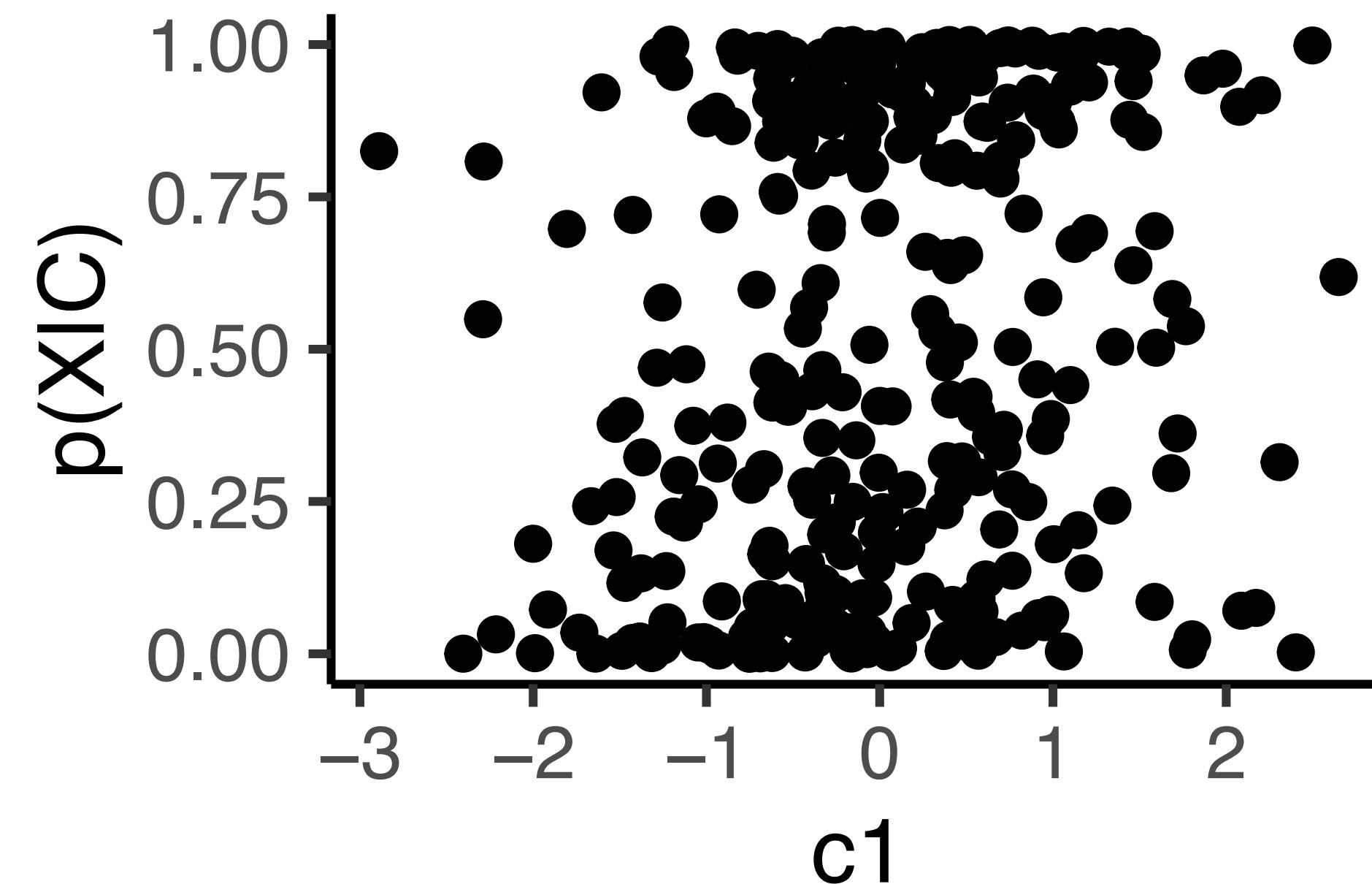


# Inverse Propensity Weighting

What is the model for  $X$  given confounders (the backdoor variables)?

Propensity = probability of assignment  $X = 1$ :

$$p(X|C_1, C_2) \approx \frac{1}{1 + \exp(-\beta_0 - \beta_1 C_1 - \beta_2 C_2)}$$

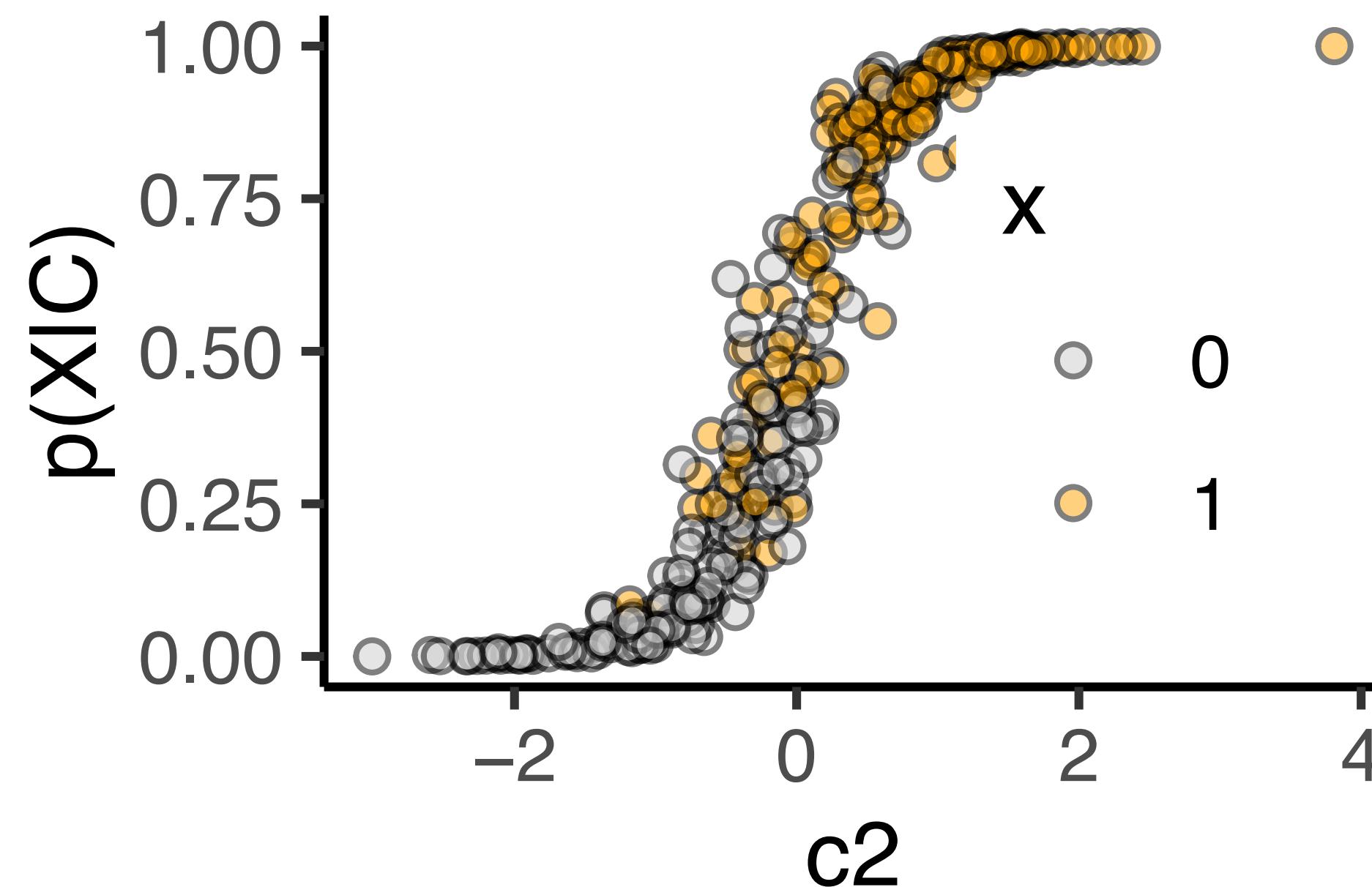
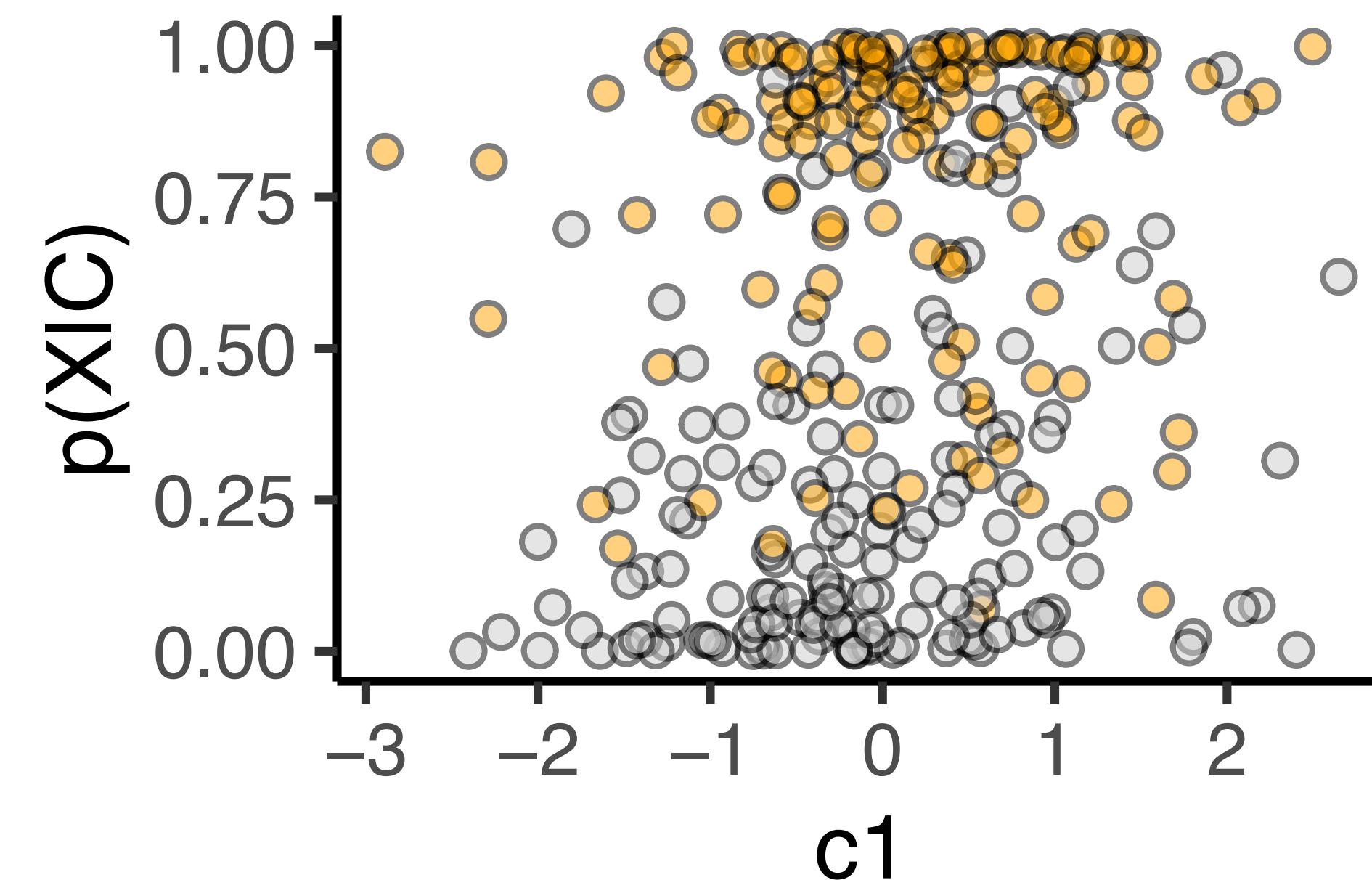


# Inverse Propensity Weighting

What is the model for  $X$  given confounders (the backdoor variables)?

Propensity = probability of assignment  $X = 1$ :

$$p(X|C_1, C_2) \approx \frac{1}{1 + \exp(-\beta_0 - \beta_1 C_1 - \beta_2 C_2)}$$



## Intuition behind IPW

What if we have assigned  $X = 1$  unconfounded by  $C$ ?

## Intuition behind IPW

What if we have assigned  $X = 1$  unconfounded by  $C$ ?

In other words, what if samples could be drawn more from the underrepresented group?

## Intuition behind IPW

What if we have assigned  $X = 1$  unconfounded by  $C$ ?

In other words, what if samples could be drawn more from the underrepresented group?

Likewise, what if samples could be dropped in the overrepresented group?

# Inverse Propensity Weighting “inverse” confounded assignment

$$\hat{Y}_i^{(1)} = \frac{X_i Y_i}{\hat{p}(X_i = 1 | C_i)}$$

$$\hat{Y}_i^{(0)} = \frac{(1 - X_i) Y_i}{1 - \hat{p}(X_i = 1 | C_i)}$$

```
p.xc <-
  glm(x~cc, family="binomial") %>%
  predict() %>%
  sigmoid() %>%
  clamp() # avoid 0 or 1

ww <- x / p.xc + (1-x) / (1-p.xc)
```

# Inverse Propensity Weighting “inverse” confounded assignment

$$\hat{Y}_i^{(1)} = \frac{X_i Y_i}{\hat{p}(X_i = 1 | C_i)}$$

$$\hat{Y}_i^{(0)} = \frac{(1 - X_i) Y_i}{1 - \hat{p}(X_i = 1 | C_i)}$$

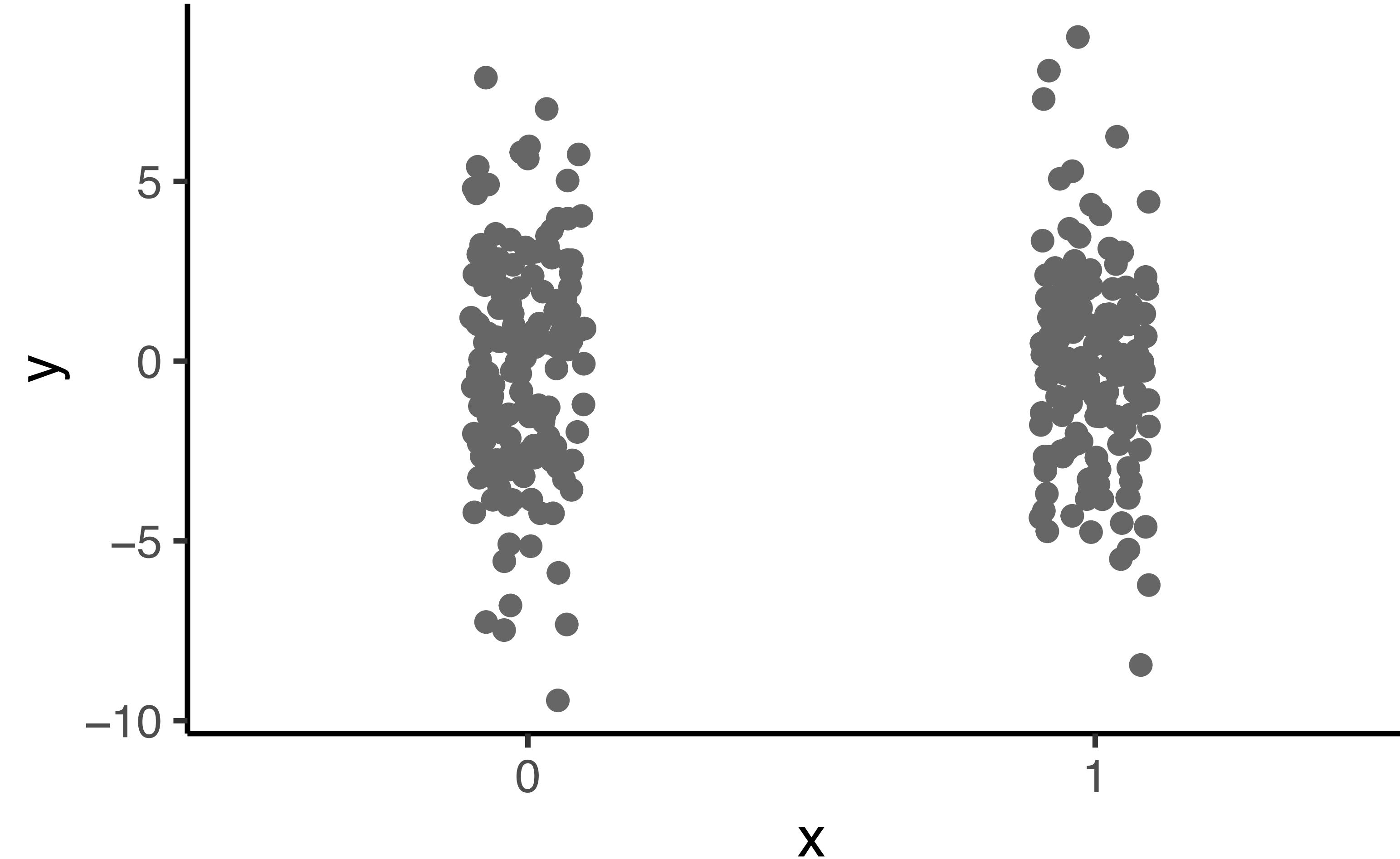
equivalently give weights for  $\forall i$

$$W_i \propto \begin{cases} 1/p(X_i = 1 | C_i) & X_i = 1 \\ 1/p(X_i = 0 | C_i) & X_i = 0 \end{cases}$$

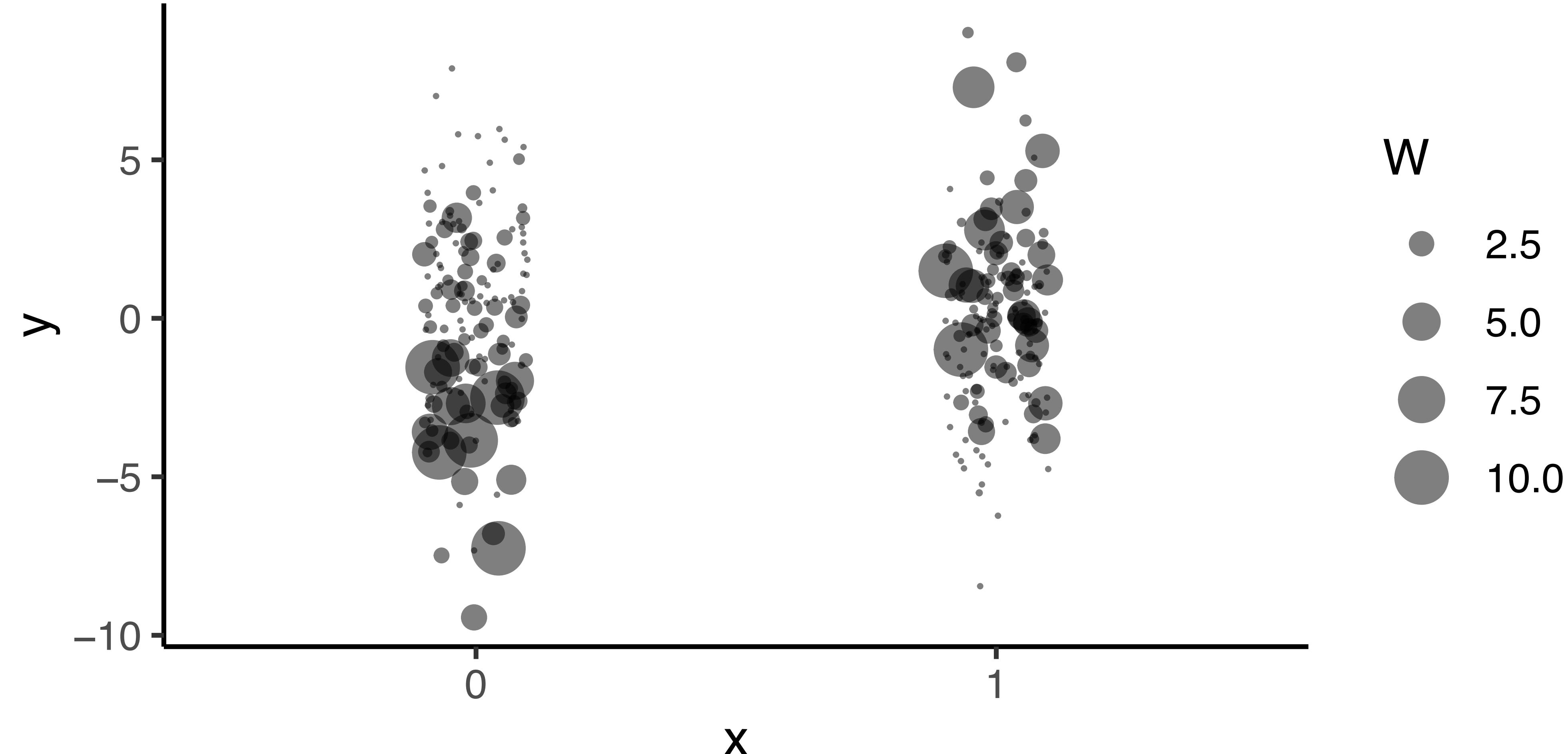
```
p.xc <-  
  glm(x~cc, family="binomial") %>%  
  predict() %>%  
  sigmoid() %>%  
  clamp() # avoid 0 or 1
```

```
ww <- x / p.xc + (1-x) / (1-p.xc)
```

Take samples inversely proportional to propensity

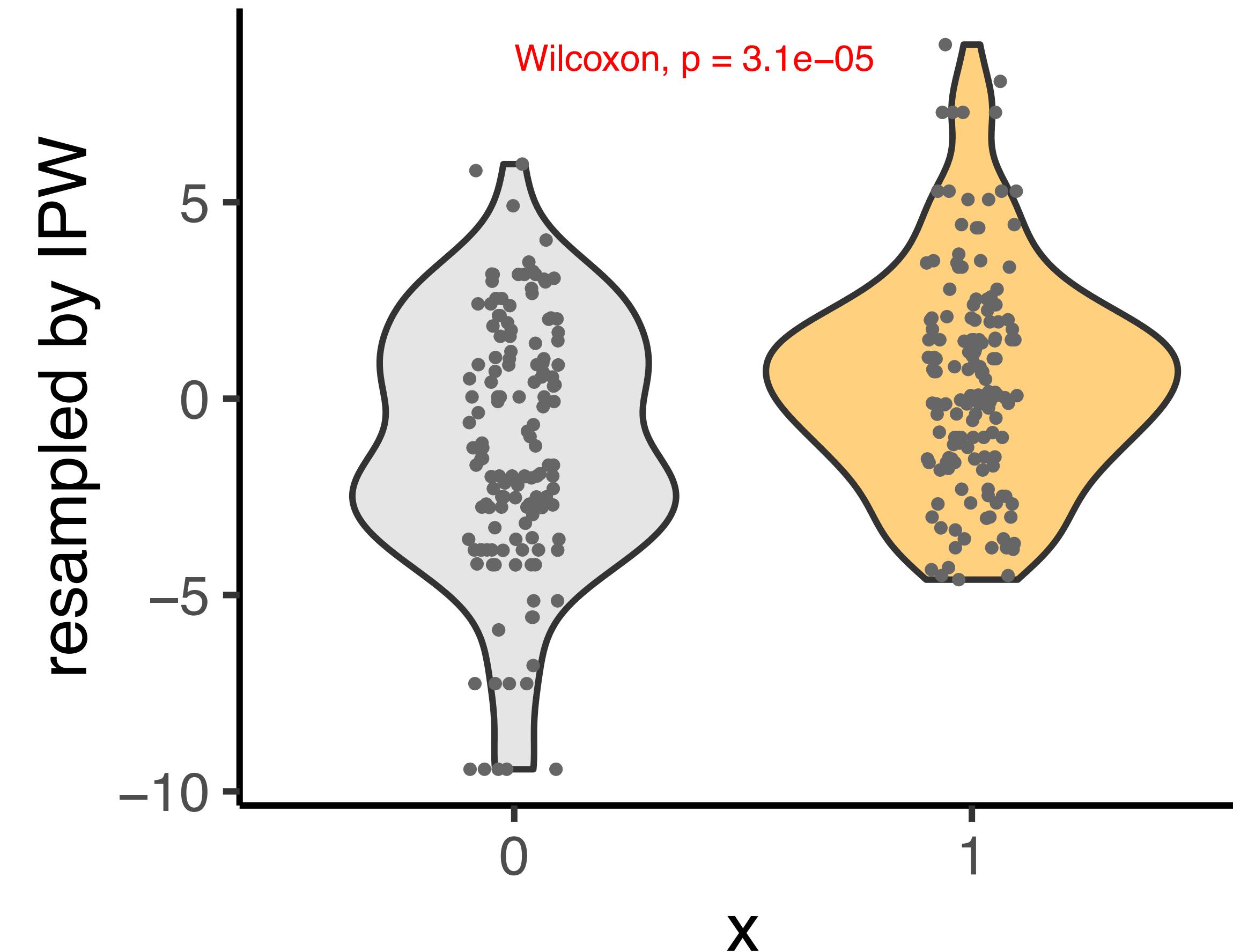
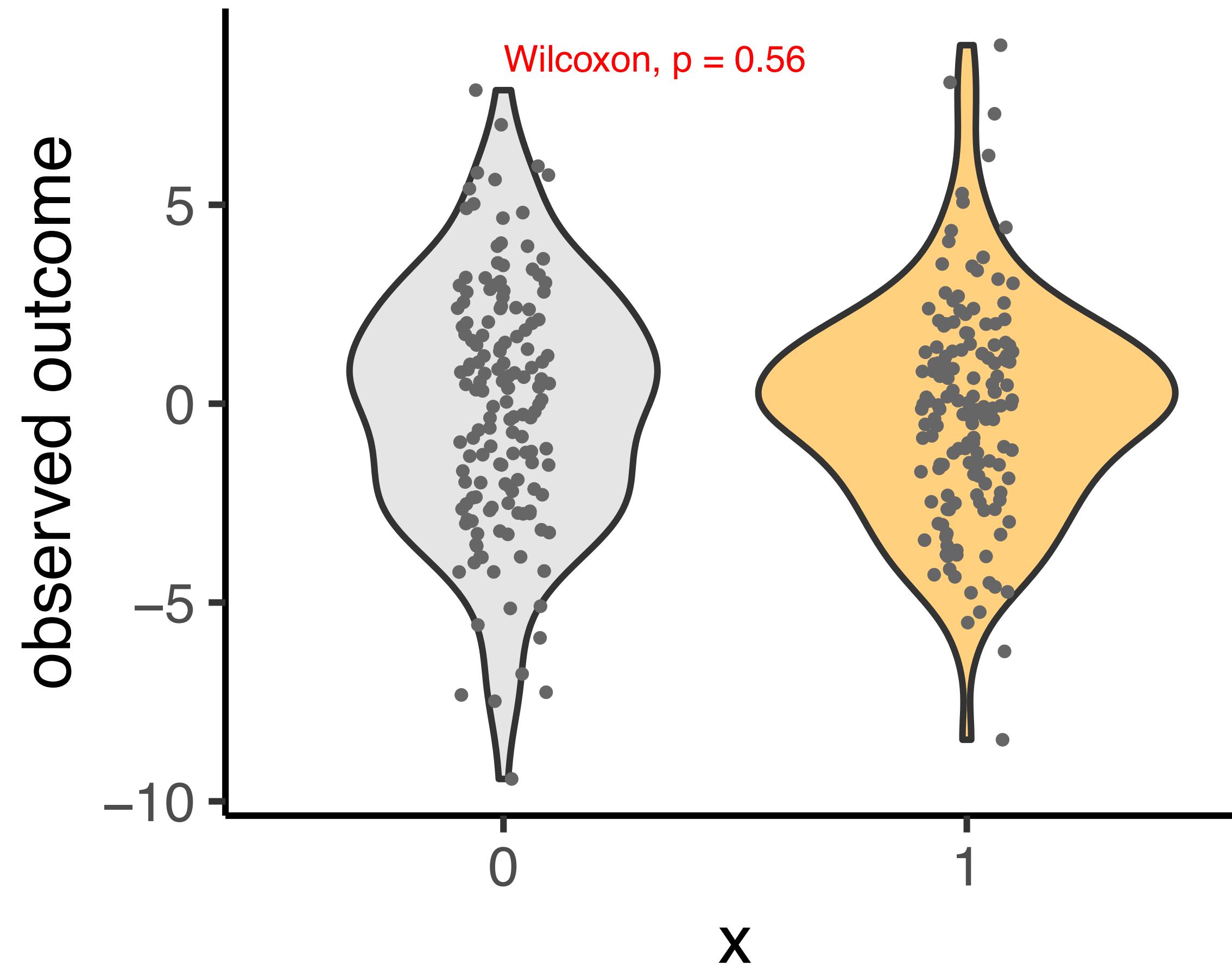


Take samples inversely proportional to propensity



$$W_i \propto \begin{cases} 1/p(X_i = 1|C_i) & X_i = 1 \\ 1/p(X_i = 0|C_i) & X_i = 0 \end{cases}$$

# Take samples inversely proportional to propensity



$$W_i \propto \begin{cases} 1/p(X_i = 1|C_i) & X_i = 1 \\ 1/p(X_i = 0|C_i) & X_i = 0 \end{cases}$$

## Why IPW works? Unbiased estimate potential outcome

Letting  $e(z) = \hat{p}(X = 1|C = z)$ ,

we can prove  $\mathbb{E}[XY/e(X)] \rightarrow \mathbb{E}[Y^{(1)}]$

using

- ▶ Strong ignorability
- ▶ Smoothness
- ▶ Stable Unit Treatment (exposure) Variable

# Why IPW works? Unbiased estimate potential outcome

Letting  $e(z) = \hat{p}(X = 1|C = z)$ ,

$$\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)}\right] =$$

# Why IPW works? Unbiased estimate potential outcome

Letting  $e(z) = \hat{p}(X = 1|C = z)$ ,

$$\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)} \middle| C_i\right]\right] \quad (\text{law of total expectation})$$

# Why IPW works? Unbiased estimate potential outcome

Letting  $e(z) = \hat{p}(X = 1|C = z)$ ,

$$\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)} \middle| C_i\right]\right] \quad (\text{law of total expectation})$$

(C is sufficient backdoor)  $= \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i^{(1)}}{e(C_i)} \middle| C_i\right]\right]$

# Why IPW works? Unbiased estimate potential outcome

Letting  $e(z) = \hat{p}(X = 1|C = z)$ ,

$$\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)} \middle| C_i\right]\right] \quad (\text{law of total expectation})$$

(C is sufficient backdoor)  $= \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i^{(1)}}{e(C_i)} \middle| C_i\right]\right]$

(strong ignorability)  $= \mathbb{E}\left[Y_i^{(1)} \mathbb{E}\left[\frac{X_i}{e(C_i)} \middle| C_i\right]\right] \quad Y^{(1)} \perp\!\!\!\perp X|C$

# Why IPW works? Unbiased estimate potential outcome

Letting  $e(z) = \hat{p}(X = 1|C = z)$ ,

$$\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)} \middle| C_i\right]\right] \quad (\text{law of total expectation})$$

(C is sufficient backdoor)  $= \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i^{(1)}}{e(C_i)} \middle| C_i\right]\right]$

(strong ignorability)  $= \mathbb{E}\left[Y_i^{(1)} \mathbb{E}\left[\frac{X_i}{e(C_i)} \middle| C_i\right]\right] \quad Y^{(1)} \perp\!\!\!\perp X|C$

(smoothness)  $= \mathbb{E}\left[Y_i^{(1)} \frac{1}{e(C_i)} \mathbb{E}[X_i | C_i]\right] \quad 0 < e(C) < 1$

# Why IPW works? Unbiased estimate potential outcome

Letting  $e(z) = \hat{p}(X = 1|C = z)$ ,

$$\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i}{e(C_i)} \middle| C_i\right]\right] \quad (\text{law of total expectation})$$

(C is sufficient backdoor)  $= \mathbb{E}\left[\mathbb{E}\left[\frac{X_i Y_i^{(1)}}{e(C_i)} \middle| C_i\right]\right]$

(strong ignorability)  $= \mathbb{E}\left[Y_i^{(1)} \mathbb{E}\left[\frac{X_i}{e(C_i)} \middle| C_i\right]\right] \quad Y^{(1)} \perp\!\!\!\perp X|C$

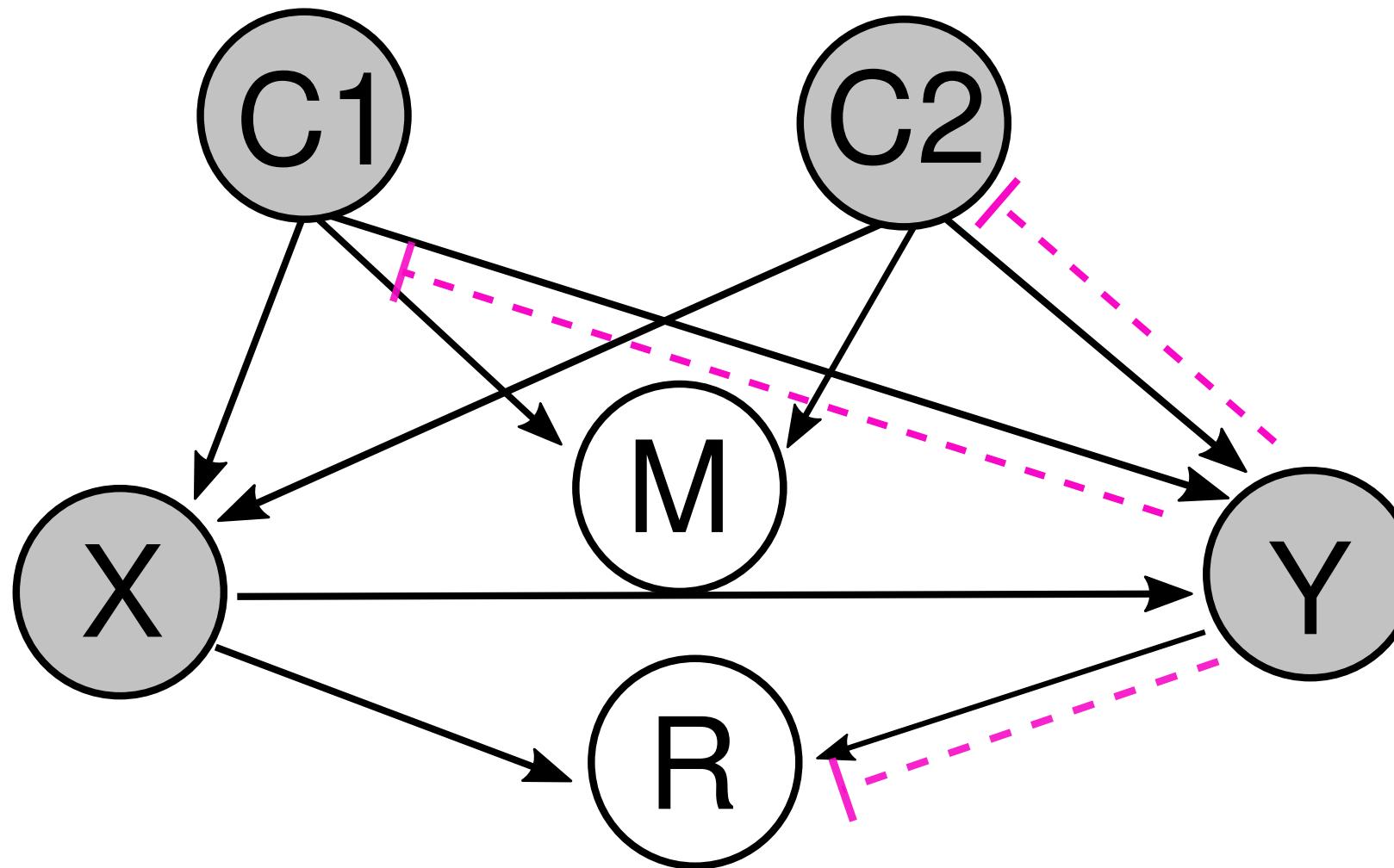
(smoothness)  $= \mathbb{E}\left[Y_i^{(1)} \frac{1}{e(C_i)} \mathbb{E}[X_i | C_i]\right] \quad 0 < e(C) < 1$

$$= \mathbb{E}[Y^{(1)}]$$

# Today's lecture: Bayesian, PGM, Causality

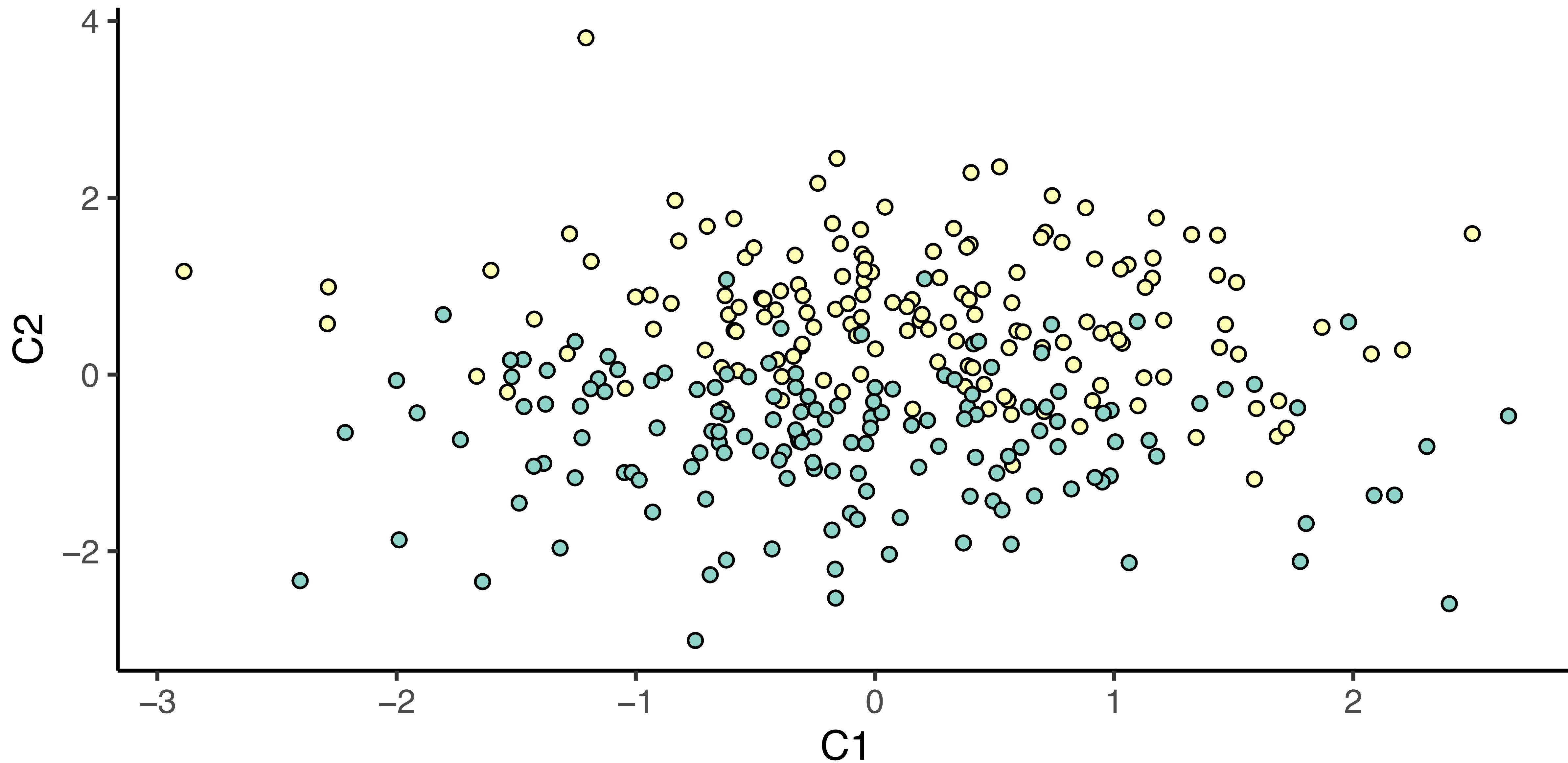
- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

# Estimating potential outcome by matching

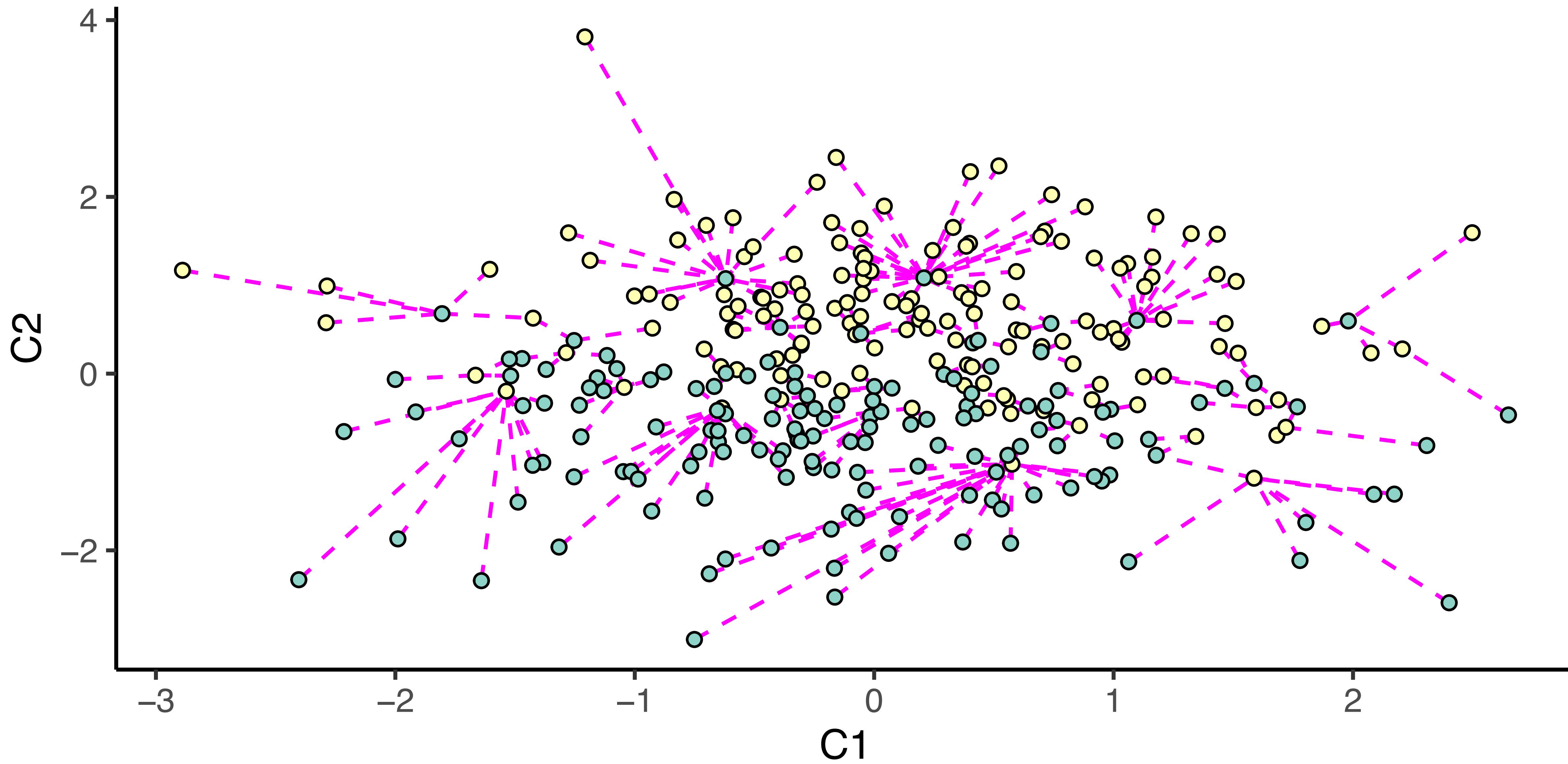


- ▶ Estimate  $\mathbb{E}[Y_i^{(0)}|C_{i1}, C_{i2}]$  for  $X_i = 1$  to compare with  $\mathbb{E}[Y_i|X_i = 1]$
- ▶ Estimate  $\mathbb{E}[Y_i^{(1)}|C_{i1}, C_{i2}]$  for  $X_i = 0$  to compare with  $\mathbb{E}[Y_i|X_i = 0]$

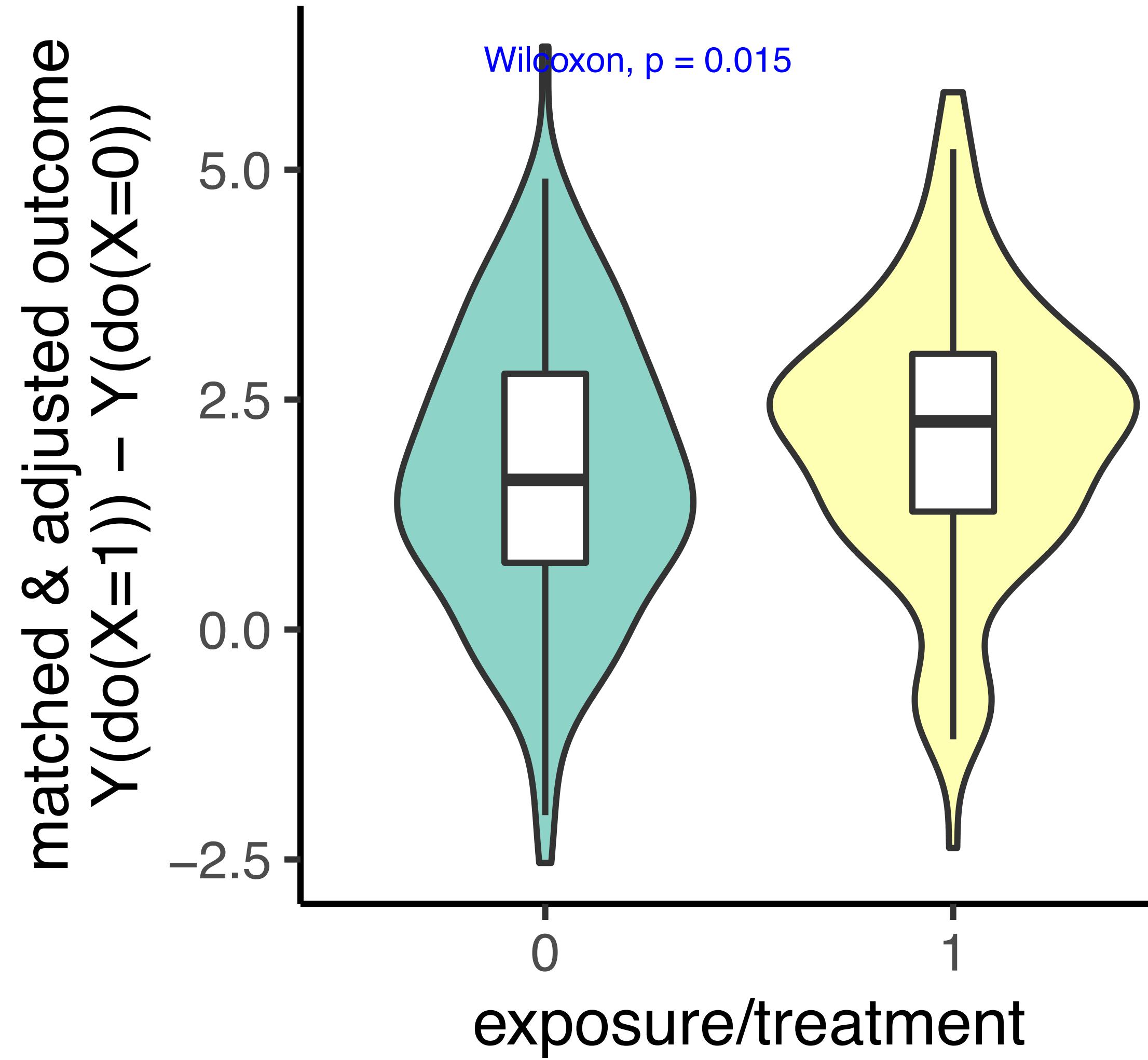
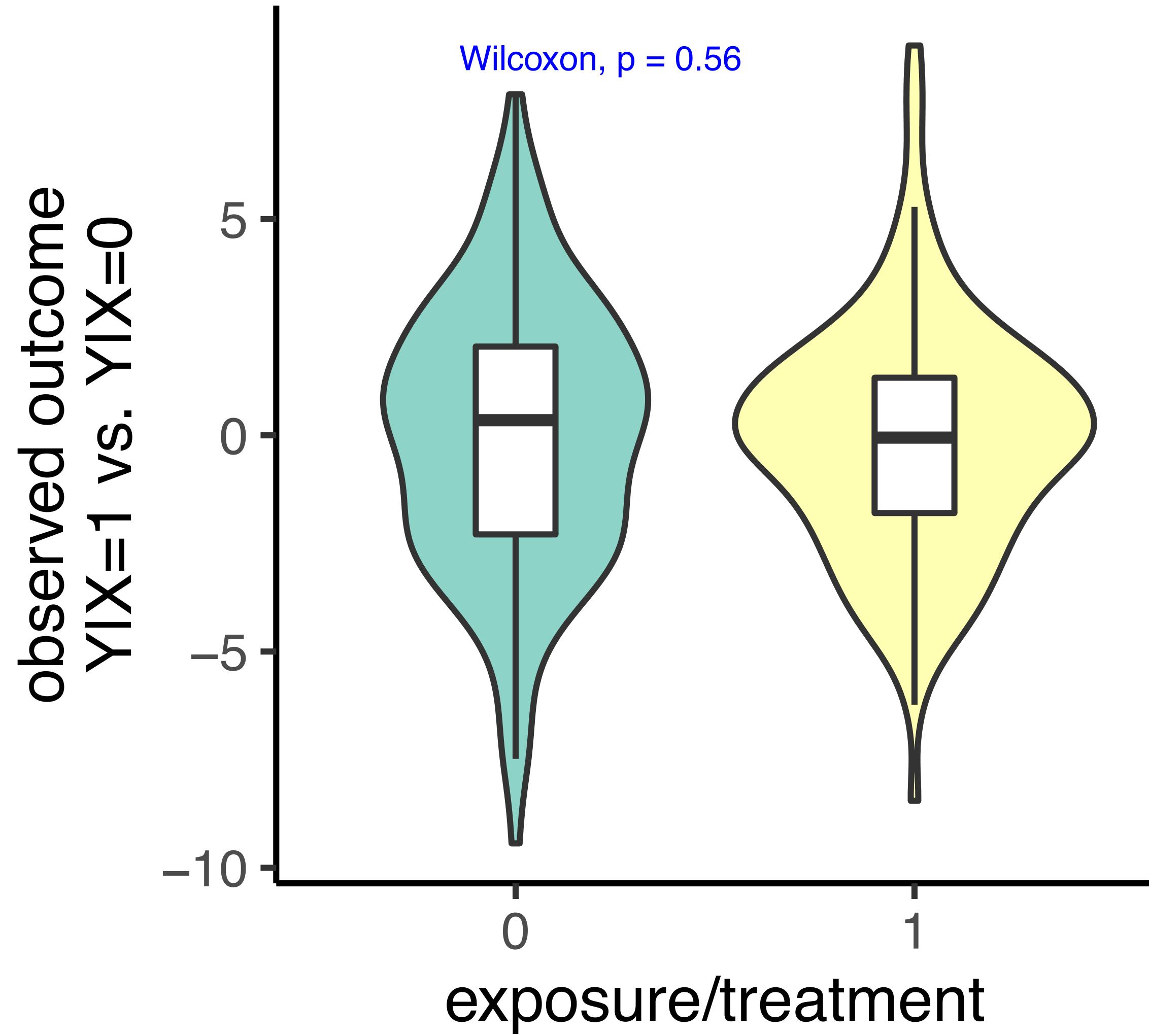
# Estimating potential outcome by matching



# Estimating potential outcome by matching



# Estimating potential outcome by matching



# Today's lecture: Bayesian, PGM, Causality

- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

# Bayesian Additive Regression Tree (BART) approach

- ▶ G-formula → outcome regression models
- ▶ Regression model for  $Y \sim S$  for each  $X = 1$  and  $X = 0$  using BART

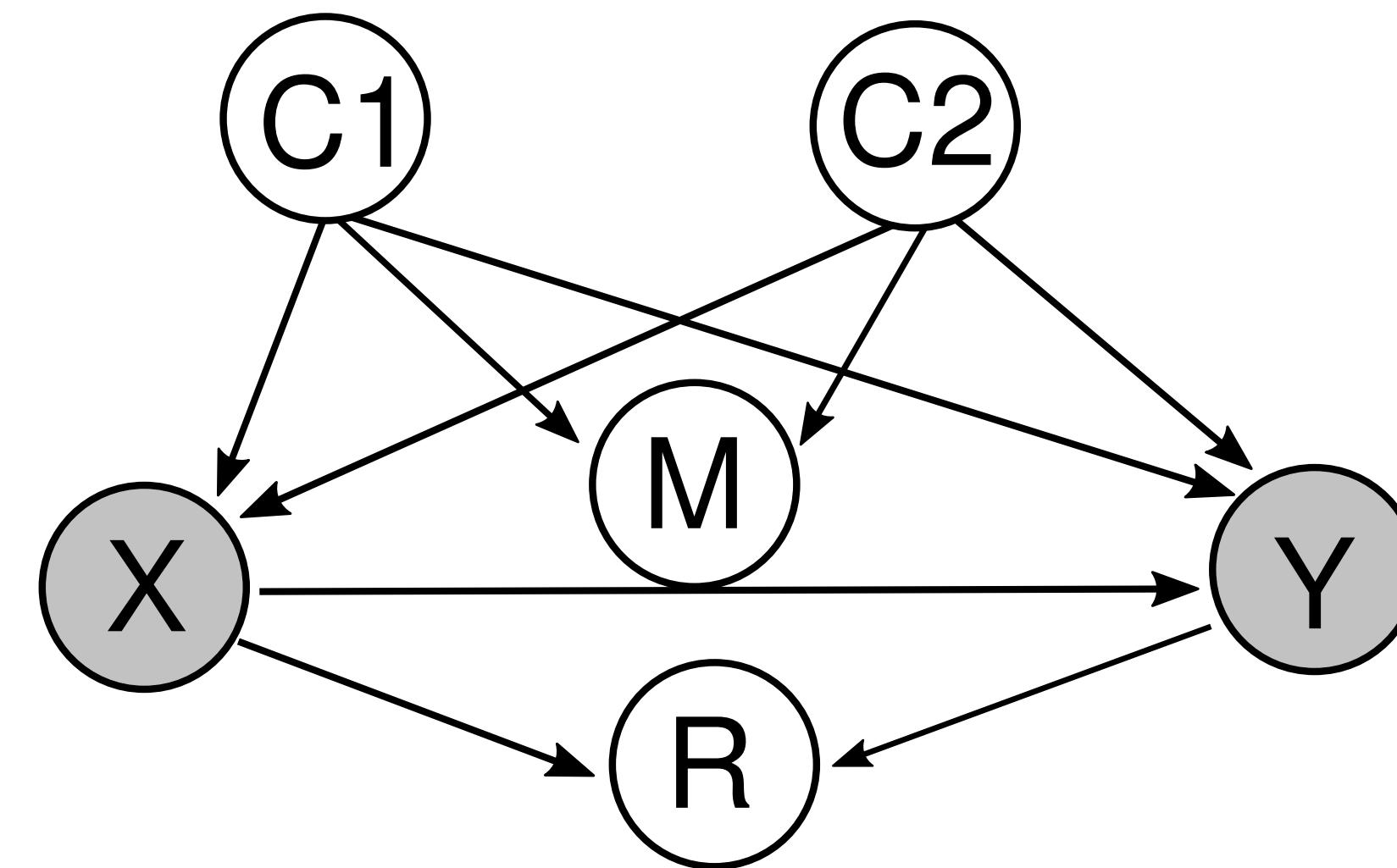
$$\mathbb{E}[Y^{(1)}] = \int_S \mathbb{E}[Y|X=1, S] dS$$

$$\mathbb{E}[Y^{(0)}] = \int_S \mathbb{E}[Y|X=0, S] dS$$

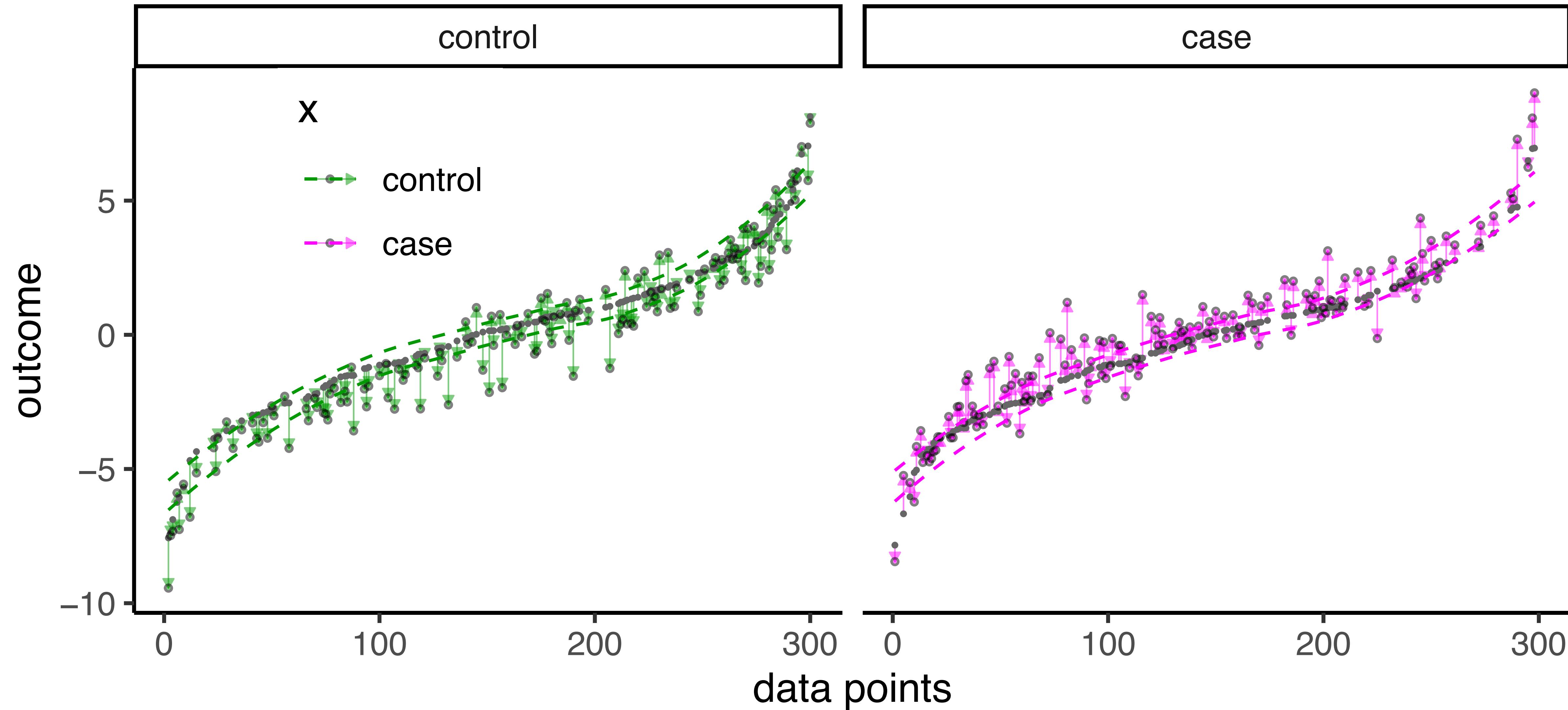
- ▶ Estimate causal effect:  $\mathbb{E}[Y^{(1)}] - \mathbb{E}[Y^{(0)}]$

Hill, *Bayesian Nonparametric Modeling for Causal Inference* (2011)

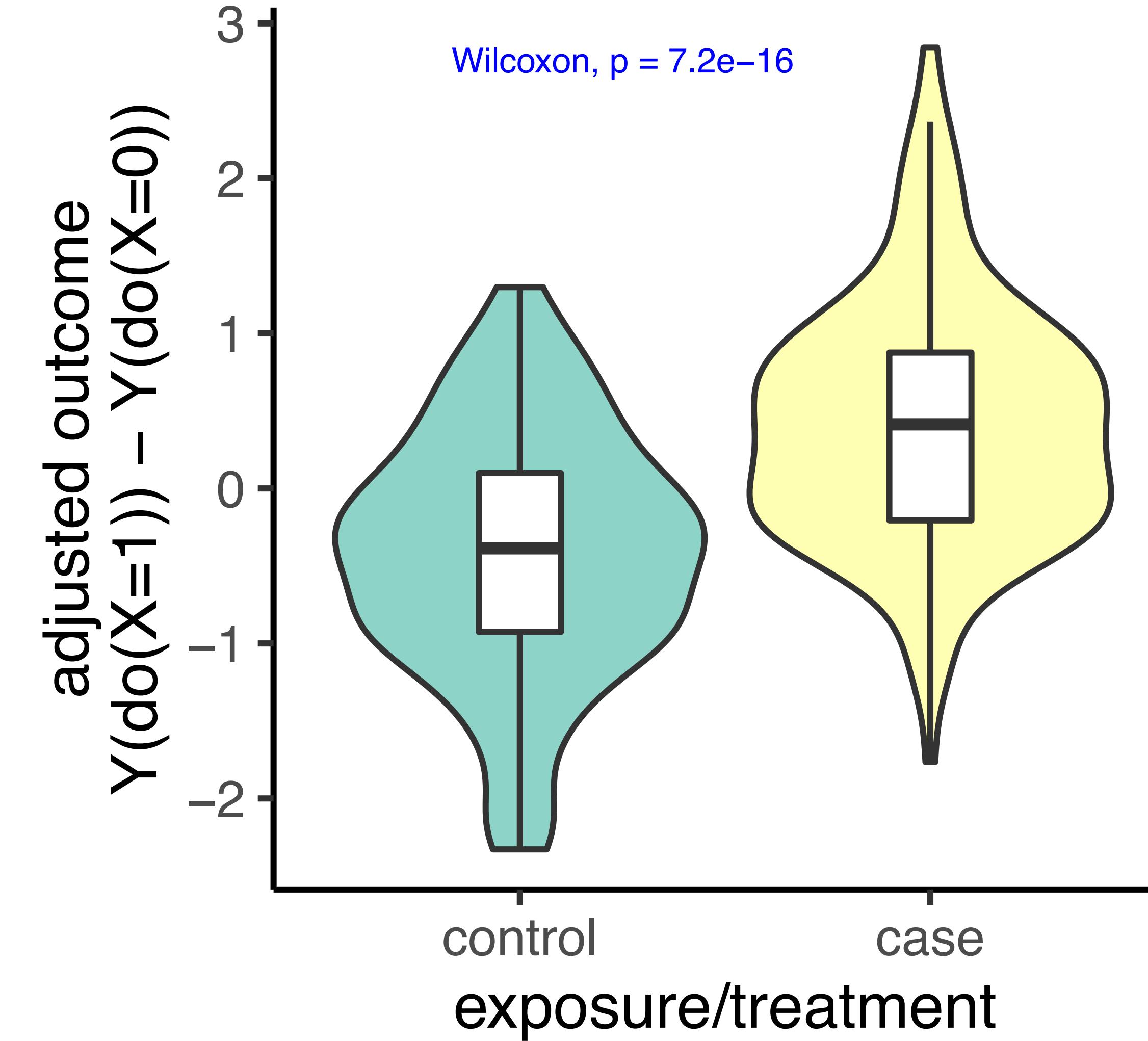
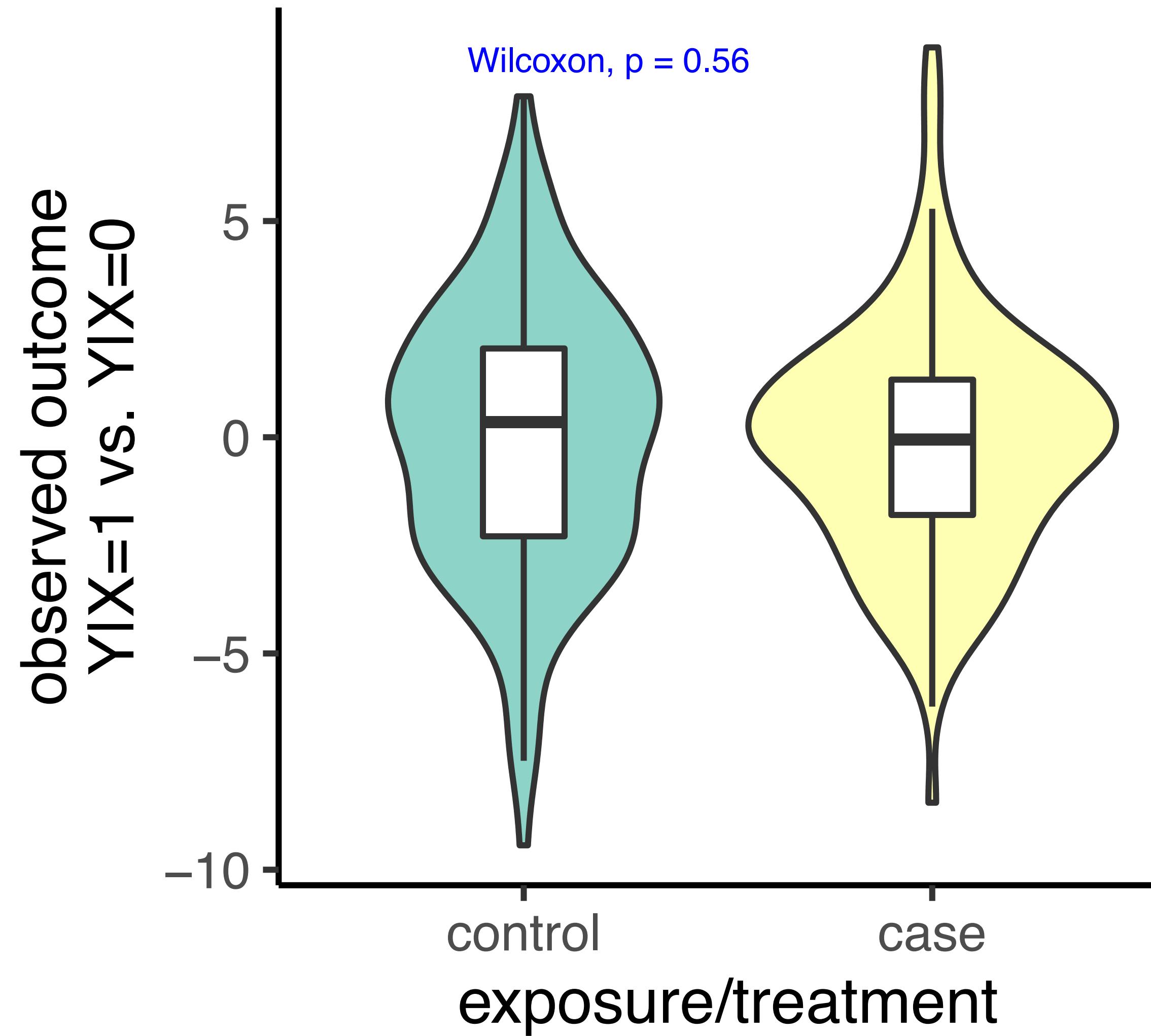
The same example with causal relationship from  $X$  to  $Y$



# BART: a regression model to impute potential outcomes



# BART: a regression model to impute potential outcomes

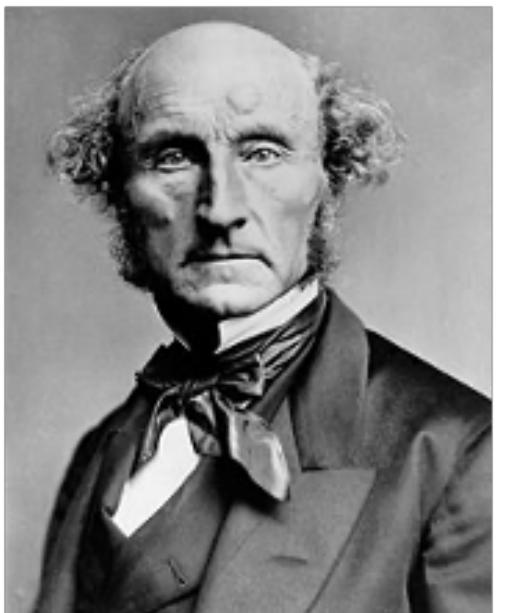


# Today's lecture: Bayesian, PGM, Causality

- **Bayesian Inference**
  - Why is it worth knowing about Bayesian inference?
  - Graphical language in probabilistic modelling
  - Examples of (practical) Bayesian inference
- **Causal inference**
  - Observation vs. Experimentation
  - Identification of unwanted bias/variance
  - More general causal inference approaches

# How can we learn "true" scientific mechanisms from data collected in observational studies?

A causal model!

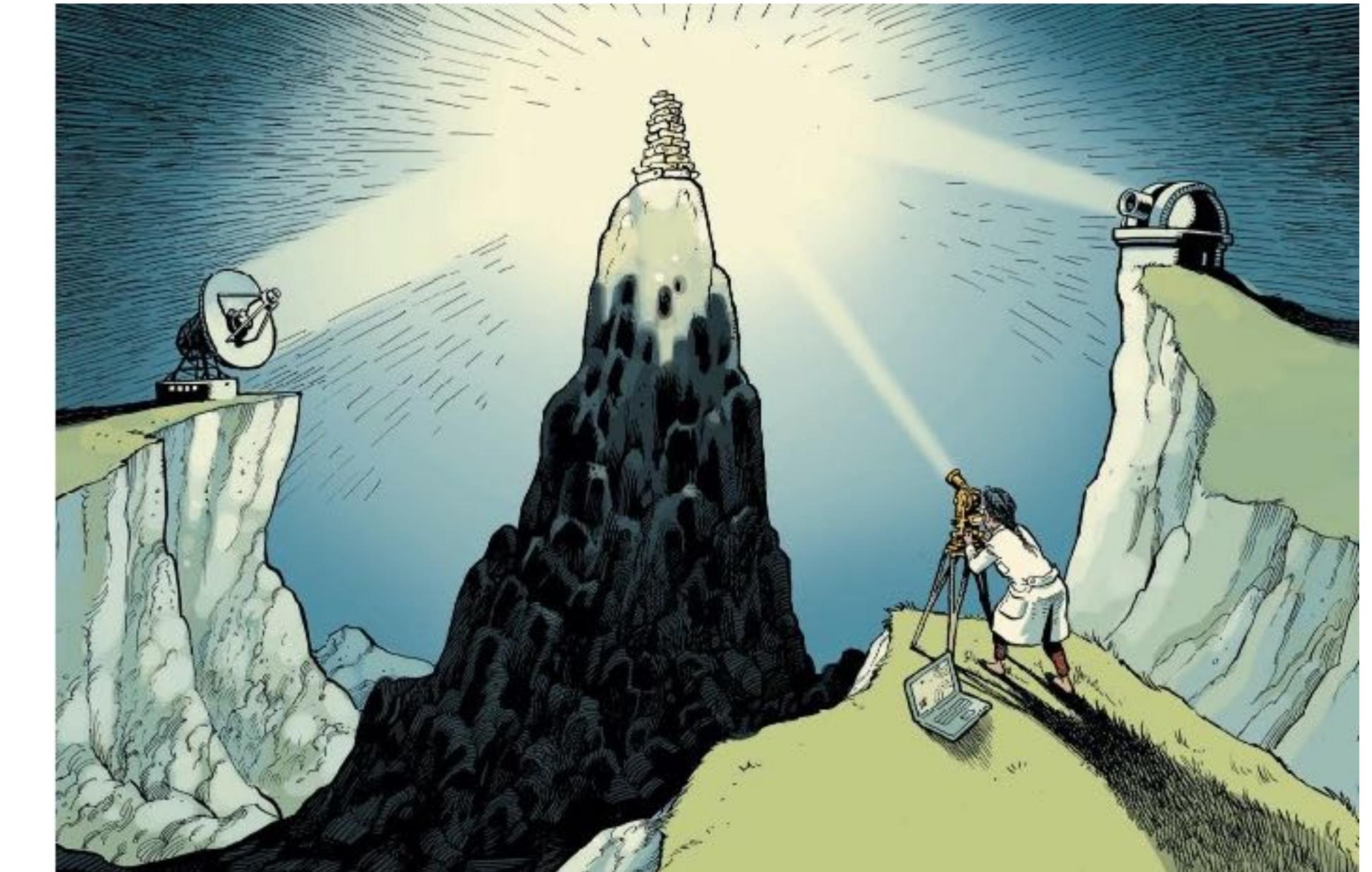


John Stuart Mill

## SECOND CANON.

*If an instance in which the phenomenon under investigation occurs, and an instance in which it does not occur, have every circumstance save one in common, that one occurring only in the former; the circumstance in which alone the two instances differ, is the effect, or cause, or a necessary part of the cause, of the phenomenon.*

JS Mill, A system of Logic (1843)



Munafo & Davey Smith, "Repeating Experiments is not enough" Nature (2018)



Peter Lipton

Contrastive Explanation & causal triangulation, Philosophy of Science (1991)



George Davey Smith