



# Stan

Software Ecosystem for Modern Bayesian Inference

**Course materials:**

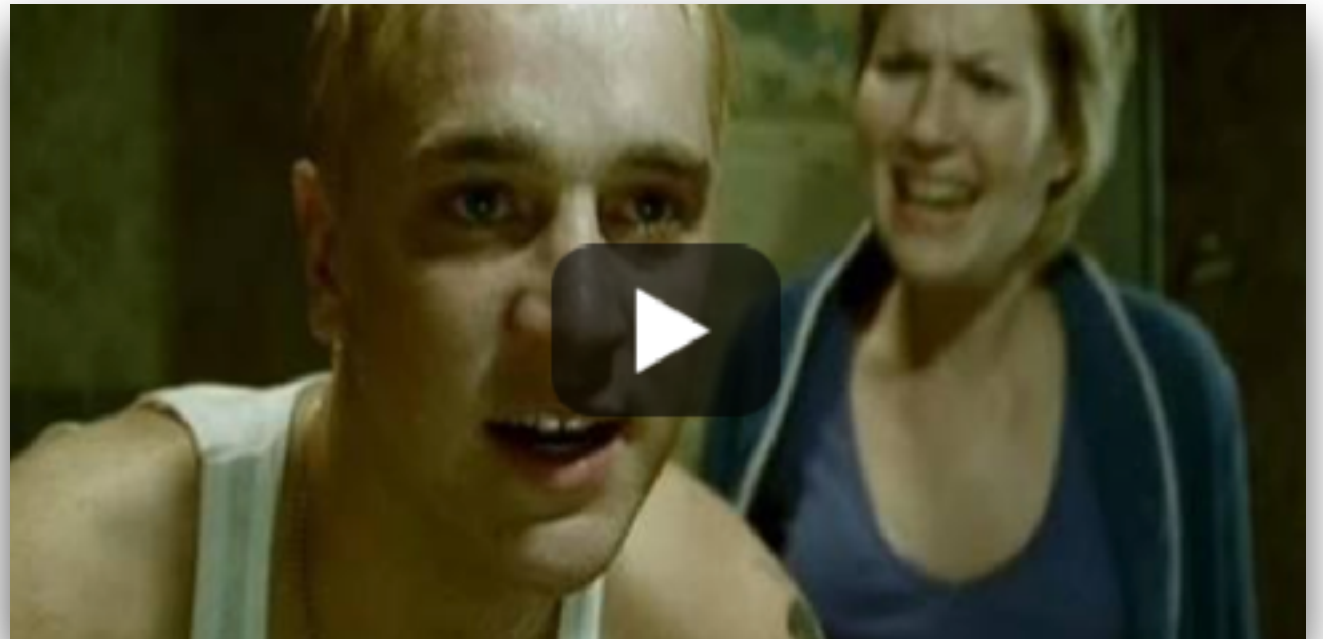
**<https://github.com/rpruim/StanWorkshop>**

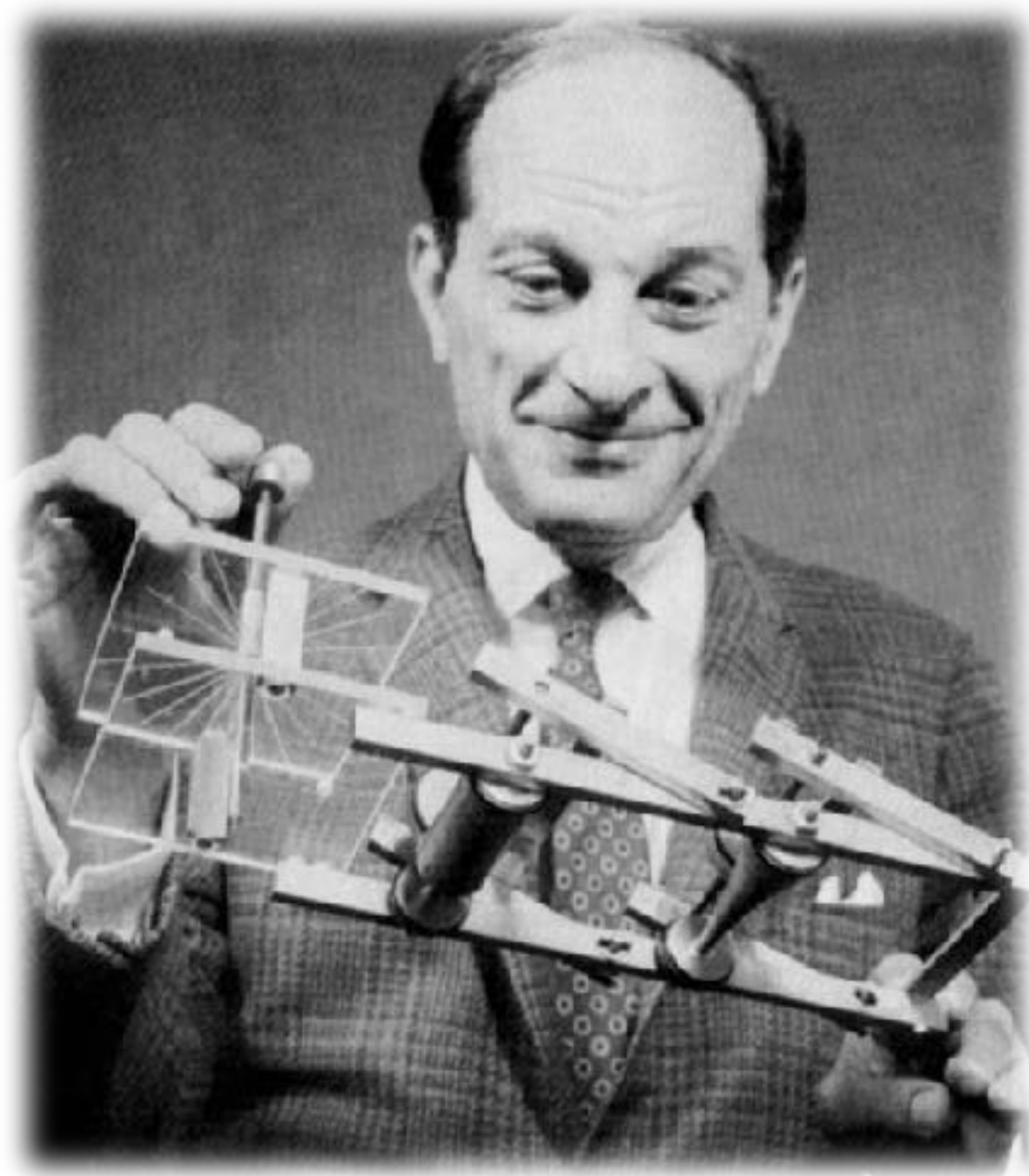
**Jonah Gabry**  
**Columbia University**

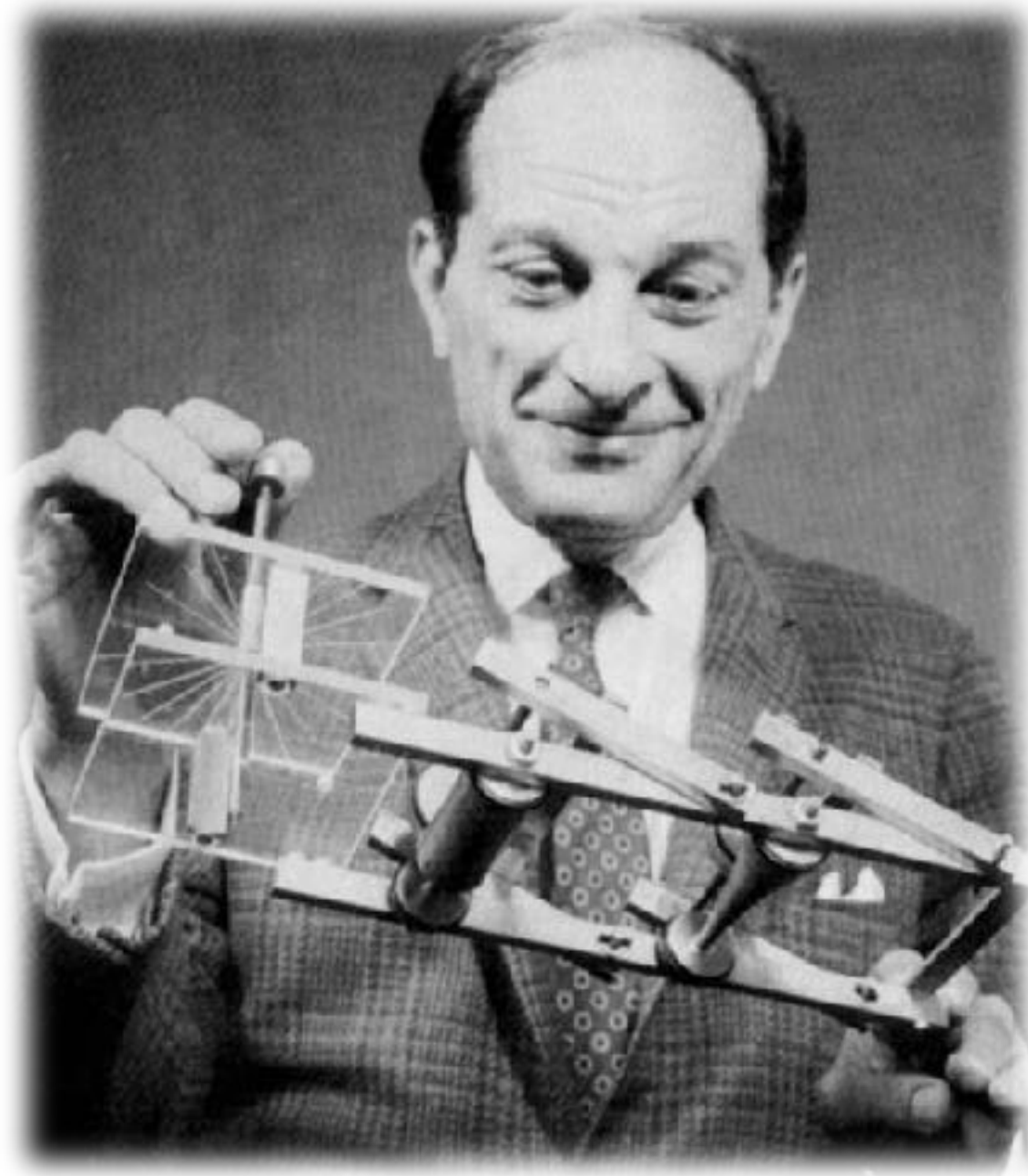
**Vianey Leos Barajas**  
**Iowa State University**

# Why “Stan”?

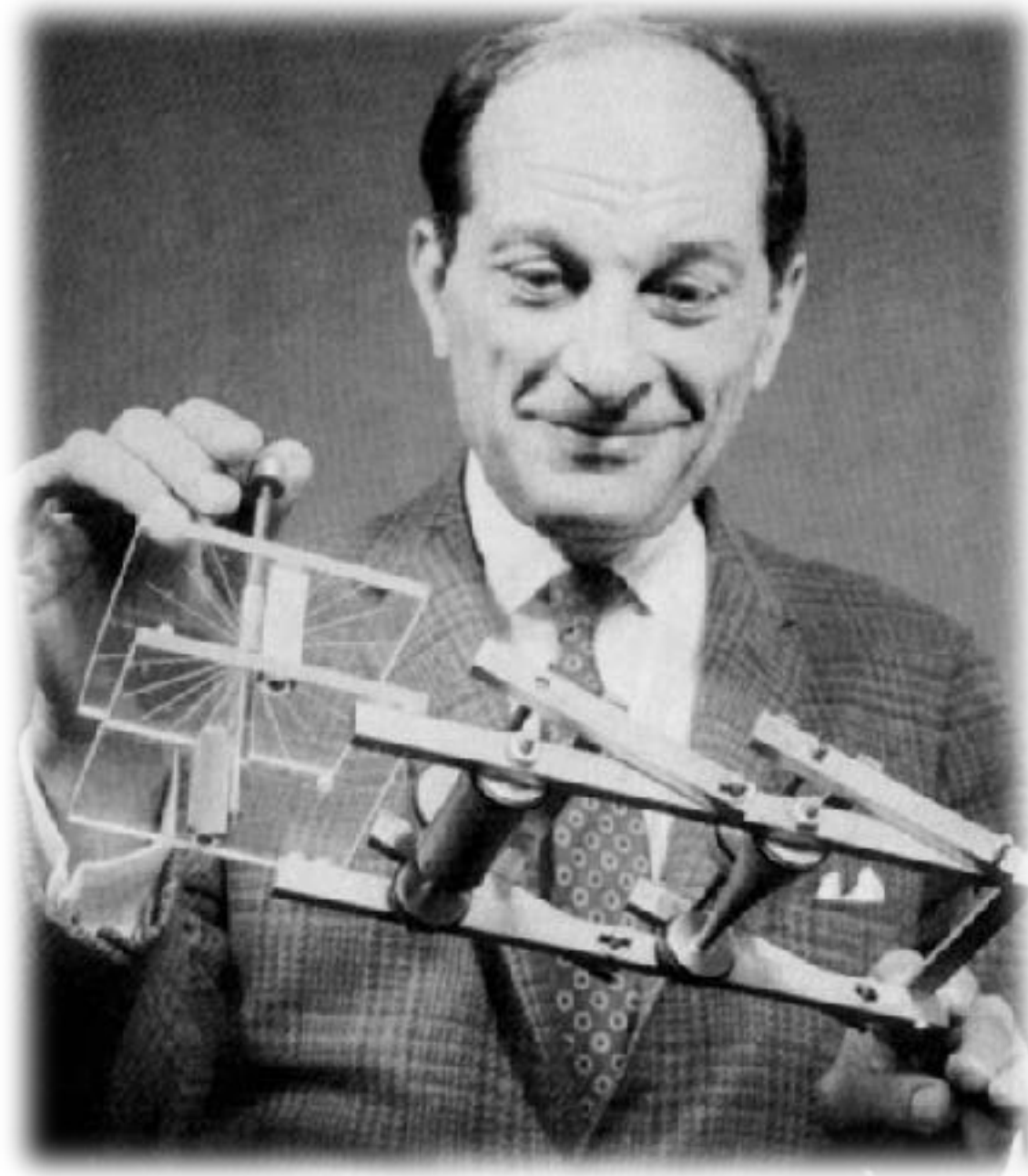
suboptimal SEO





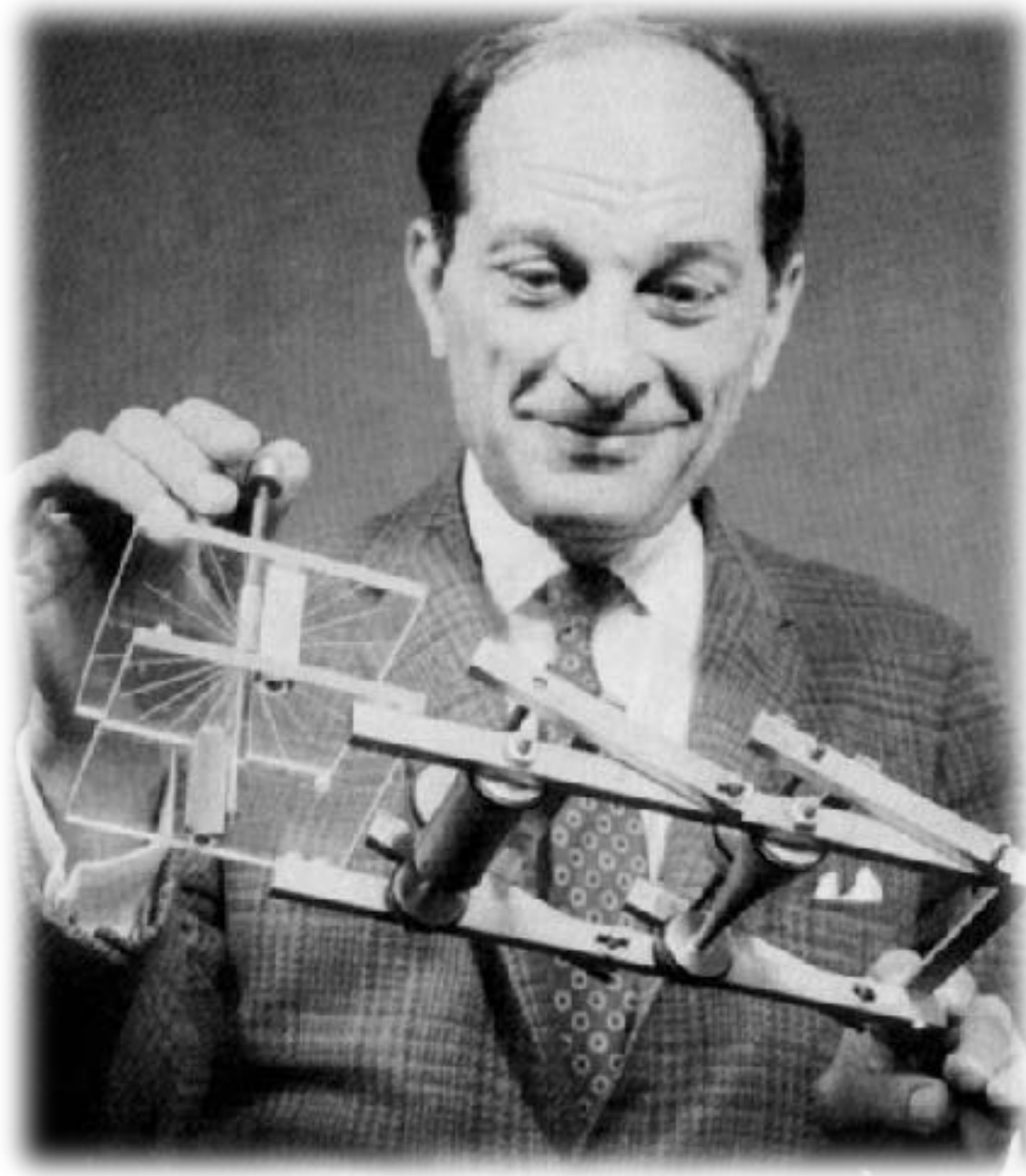


Stanislaw Ulam  
(1909–1984)



Stanislaw Ulam  
(1909–1984)

Monte Carlo  
Method



**Stanislaw Ulam  
(1909–1984)**

**H-Bomb**

**Monte Carlo  
Method**

# What is Stan?



# What is Stan?

- Open source probabilistic **programming language, inference algorithms**

# What is Stan?

- Open source probabilistic **programming language, inference algorithms**
- Stan **program**
  - declares data and (constrained) parameter variables
  - defines log posterior (or penalized likelihood)

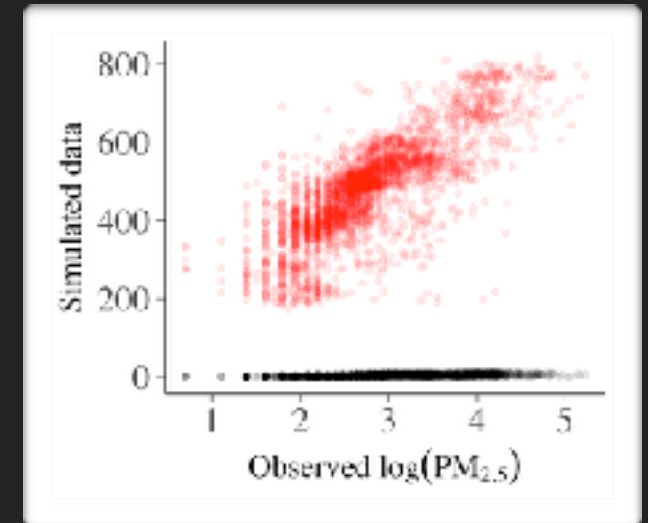
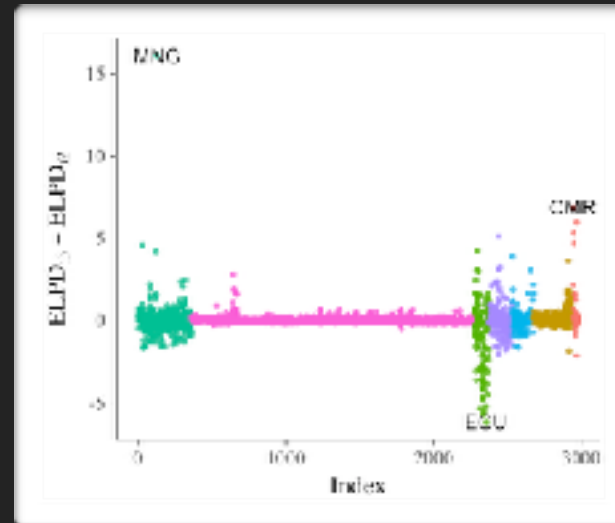
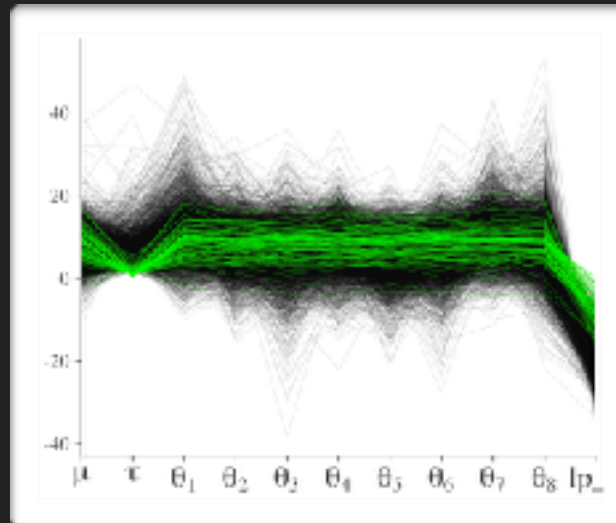
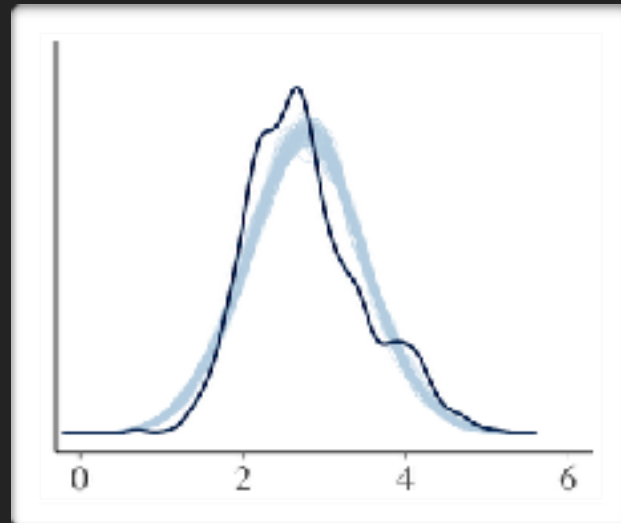
# What is Stan?

- Open source probabilistic **programming language, inference algorithms**
- Stan **program**
  - declares data and (constrained) parameter variables
  - defines log posterior (or penalized likelihood)
- Stan **inference**
  - MCMC for full Bayes
  - VB for approximate Bayes
  - Optimization for (penalized) MLE

# What is Stan?

- Open source probabilistic **programming language, inference algorithms**
- Stan **program**
  - declares data and (constrained) parameter variables
  - defines log posterior (or penalized likelihood)
- Stan **inference**
  - MCMC for full Bayes
  - VB for approximate Bayes
  - Optimization for (penalized) MLE
- Stan **ecosystem**
  - lang, math library (C++)
  - interfaces and tools (R, Python, many more)
  - documentation ([example model repo](#), [user guide](#) & [reference manual](#), [case studies](#), R package vignettes)
  - online community ([Stan Forums](#) on Discourse)

# Visualization in Bayesian workflow



**Jonah Gabry**

Columbia University  
Stan Development Team

# Workflow

Bayesian data analysis

# **Workflow**

## Bayesian data analysis

- Exploratory data analysis

# Workflow

## Bayesian data analysis

- Exploratory data analysis
- *Prior* predictive checking



# Workflow

## Bayesian data analysis

- Exploratory data analysis
- *Prior* predictive checking
- Model fitting and algorithm diagnostics

# Workflow

## Bayesian data analysis

- Exploratory data analysis
- *Prior* predictive checking
- Model fitting and algorithm diagnostics
- *Posterior* predictive checking

# Workflow

## Bayesian data analysis

- Exploratory data analysis
- *Prior* predictive checking
- Model fitting and algorithm diagnostics
- *Posterior* predictive checking
- Model comparison (e.g., via cross-validation)

# Workflow

## Bayesian data analysis

- Exploratory data analysis
- *Prior* predictive checking
- Model fitting and algorithm diagnostics
- *Posterior* predictive checking
- Model comparison (e.g., via cross-validation)

Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., and Gelman, A. (2019).

**Visualization in Bayesian workflow.**

*Journal of the Royal Statistical Society Series A*

Journal version: [rss.onlinelibrary.wiley.com/doi/full/10.1111/rssa.12378](https://rss.onlinelibrary.wiley.com/doi/full/10.1111/rssa.12378)

arXiv preprint: [arxiv.org/abs/1709.01449](https://arxiv.org/abs/1709.01449)

Code: [github.com/jgabry/bayes-vis-paper](https://github.com/jgabry/bayes-vis-paper)

# Example

# Example

## **Goal**

Estimate global PM<sub>2.5</sub> concentration

# Example

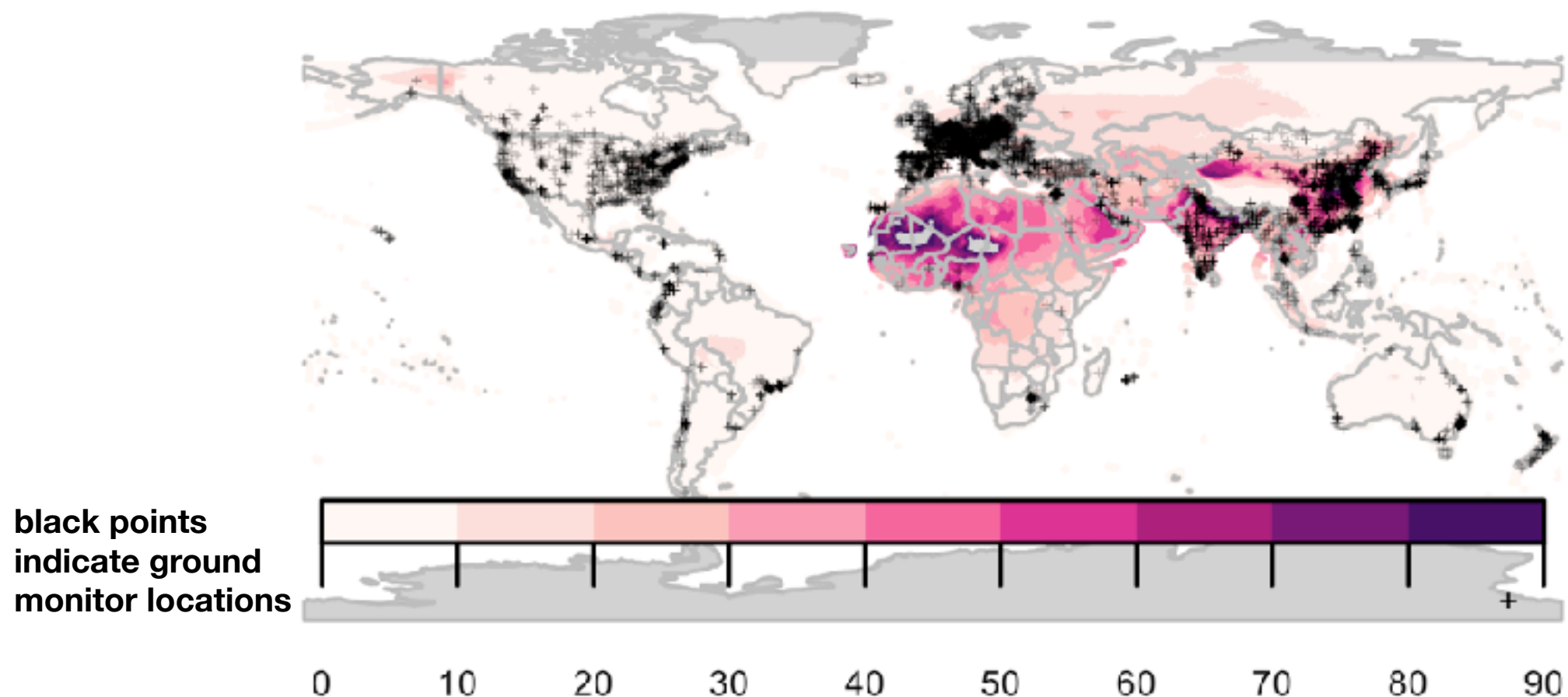
**Goal**      Estimate global PM2.5 concentration

**Problem**    Most data from noisy satellite measurements (ground monitor network provides sparse, heterogeneous coverage)

# Example

**Goal** Estimate global PM2.5 concentration

**Problem** Most data from noisy satellite measurements (ground monitor network provides sparse, heterogeneous coverage)



**Satellite estimates of PM2.5 and ground monitor locations**

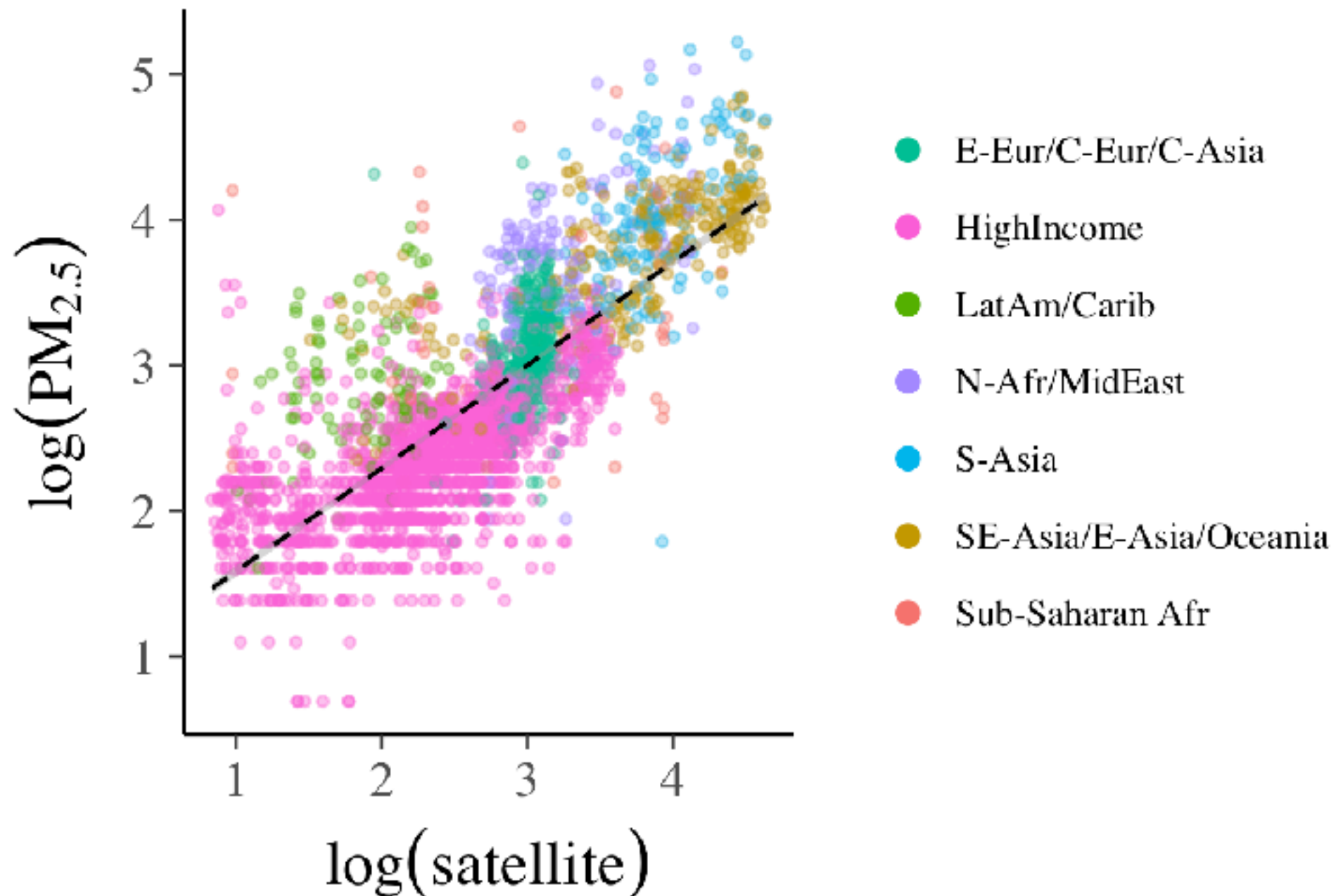


# Exploratory Data Analysis

*Building a network of models*

# Exploratory data analysis

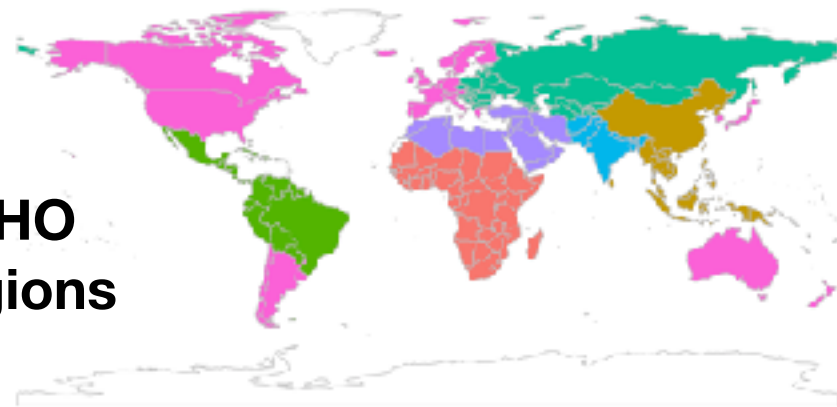
building a network of models



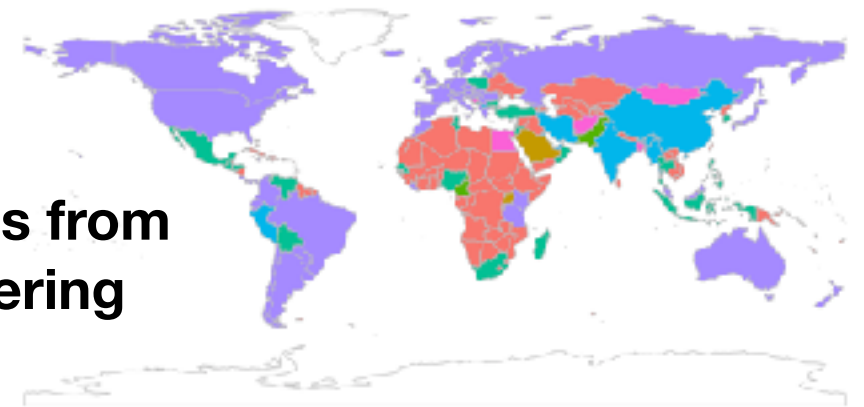
# Exploratory data analysis

building a network of models

**WHO  
Regions**



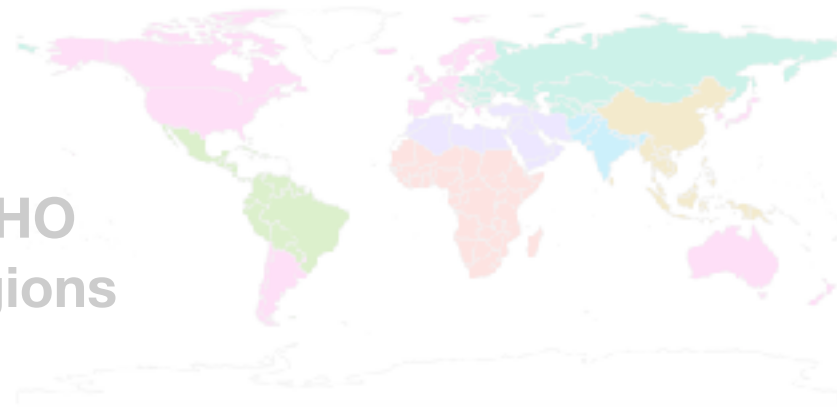
**Regions from  
clustering**



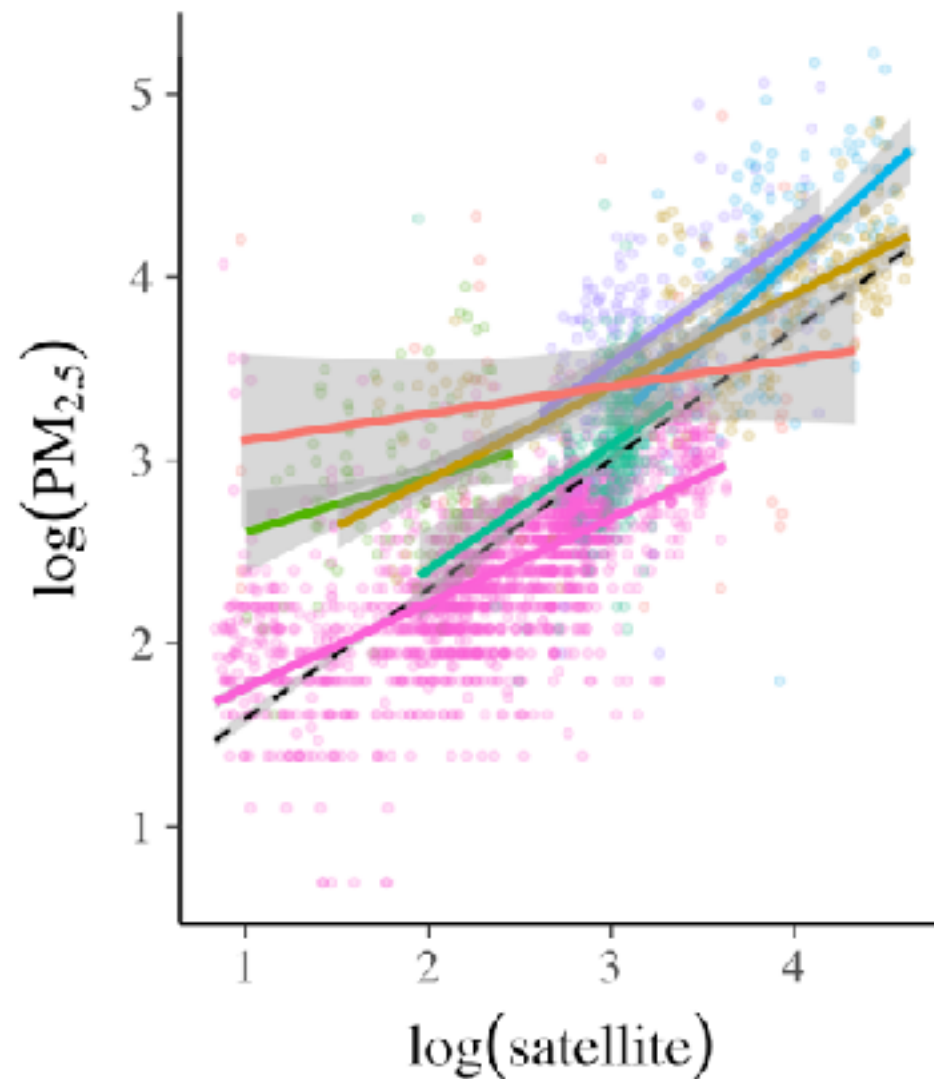
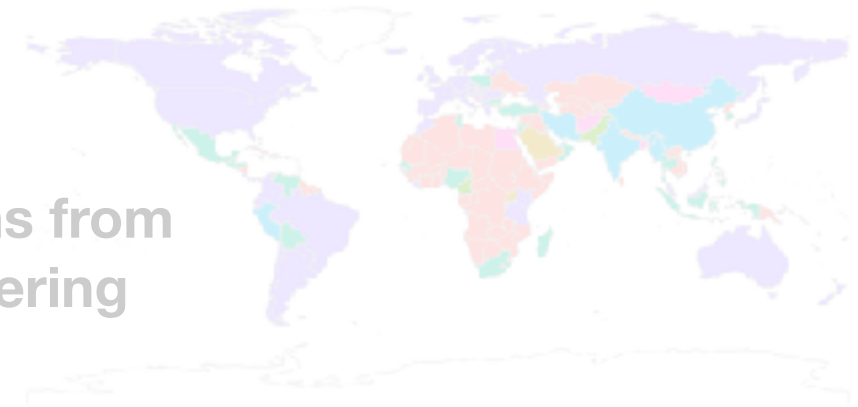
# Exploratory data analysis

building a network of models

WHO  
Regions



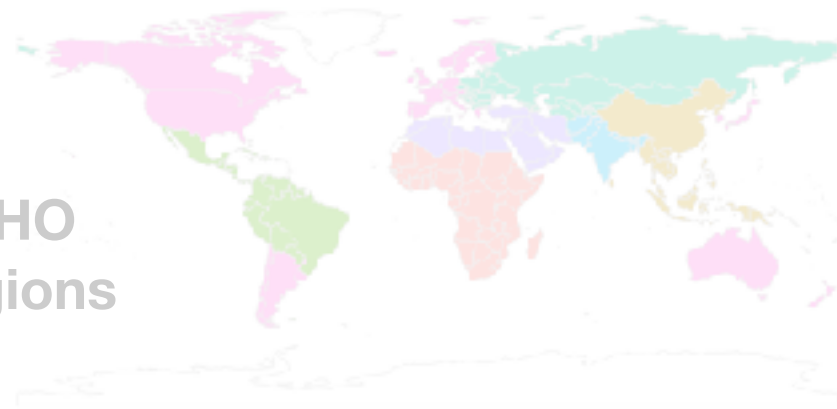
Regions from  
clustering



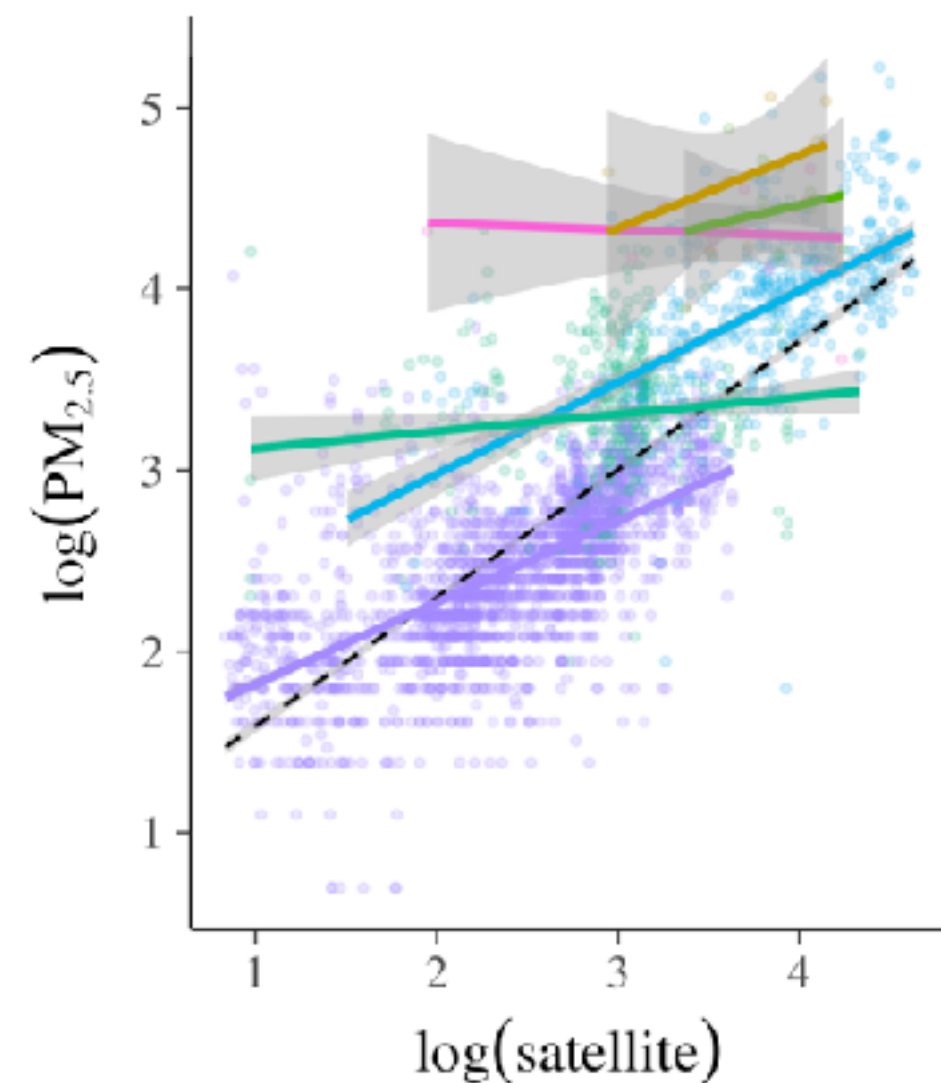
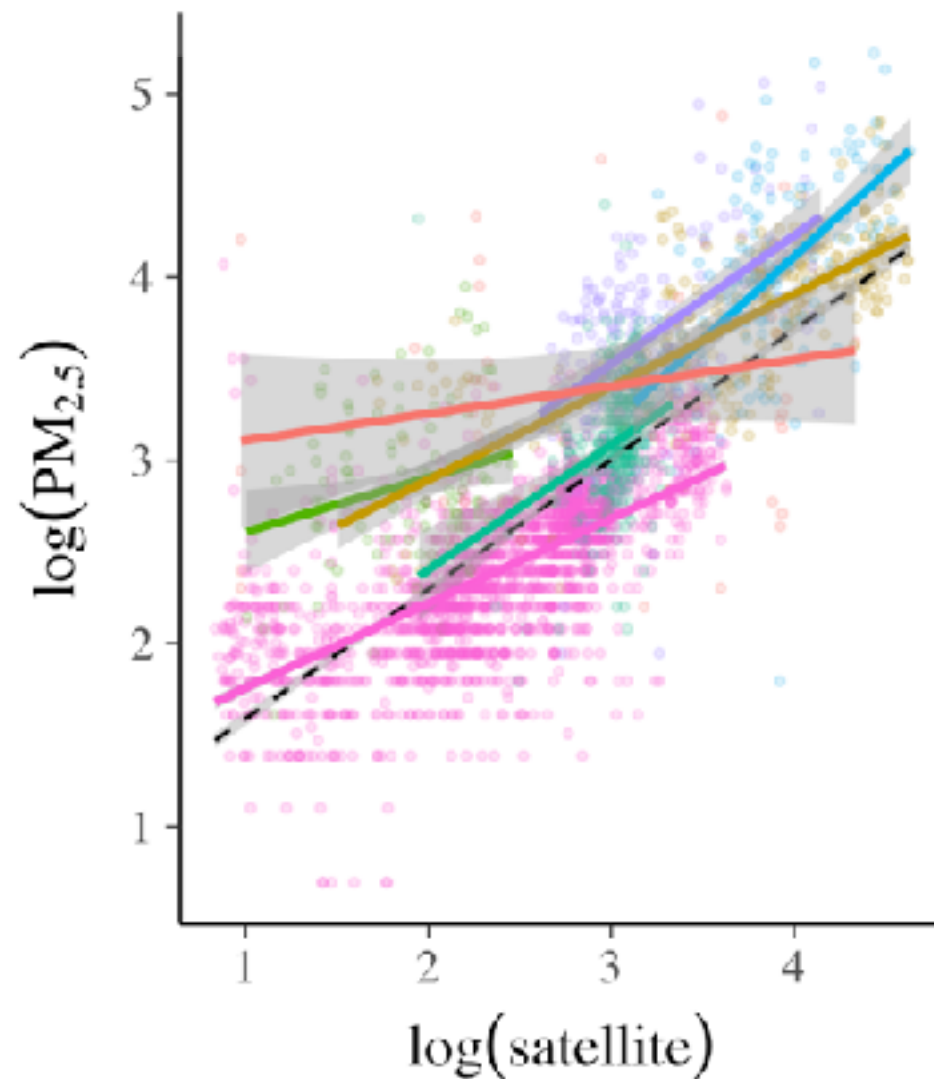
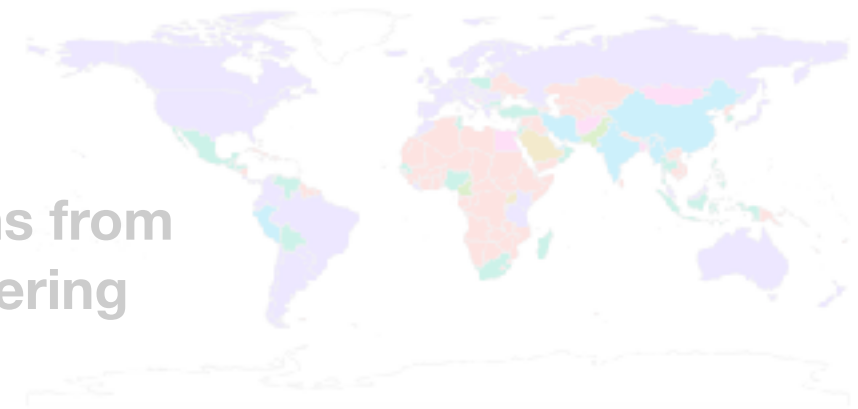
# Exploratory data analysis

building a network of models

WHO  
Regions



Regions from  
clustering



# Exploratory data analysis

building a network of models

For measurements  $n = 1, \dots, N$   
and regions  $j = 1, \dots, J$

## Model 1

# Exploratory data analysis

building a network of models

For measurements  $n = 1, \dots, N$   
and regions  $j = 1, \dots, J$

## Model 1

$$\log(\text{PM}_{2.5, n_j}) \sim N(\alpha + \beta \log(\text{sat}_{n_j}), \sigma)$$

# Exploratory data analysis

building a network of models

For measurements  $n = 1, \dots, N$   
and regions  $j = 1, \dots, J$

## Models 2 and 3



# Exploratory data analysis

building a network of models

For measurements  $n = 1, \dots, N$   
and regions  $j = 1, \dots, J$

## Models 2 and 3

$$\log(\text{PM}_{2.5, nj}) \sim N(\mu_{nj}, \sigma)$$

# Exploratory data analysis

building a network of models

For measurements  $n = 1, \dots, N$   
and regions  $j = 1, \dots, J$

## Models 2 and 3

$$\log(\text{PM}_{2.5, nj}) \sim N(\mu_{nj}, \sigma)$$

$$\mu_{nj} = \alpha_0 + \alpha_j + (\beta_0 + \beta_j) \log(\text{sat}_{nj})$$

# Exploratory data analysis

building a network of models

For measurements  $n = 1, \dots, N$   
and regions  $j = 1, \dots, J$

## Models 2 and 3

$$\log(\text{PM}_{2.5, n_j}) \sim N(\mu_{n_j}, \sigma)$$

$$\mu_{n_j} = \boxed{\alpha_0 + \alpha_j} + \boxed{(\beta_0 + \beta_j)} \log(\text{sat}_{n_j})$$

# Exploratory data analysis

building a network of models

For measurements  $n = 1, \dots, N$   
and regions  $j = 1, \dots, J$

## Models 2 and 3

$$\log(\text{PM}_{2.5, n_j}) \sim N(\mu_{n_j}, \sigma)$$

$$\mu_{n_j} = \boxed{\alpha_0 + \alpha_j} + \boxed{(\beta_0 + \beta_j)} \log(\text{sat}_{n_j})$$

$$\alpha_j \sim N(0, \tau_\alpha) \quad \beta_j \sim N(0, \tau_\beta)$$

# Prior predictive checks

*Fake data can be almost as valuable as real data*

# A Bayesian modeler commits to an a priori *joint distribution*

*Likelihood x Prior*

$$p(\mathbf{y}, \boldsymbol{\theta}) = p(\mathbf{y} \mid \boldsymbol{\theta})p(\boldsymbol{\theta}) = p(\boldsymbol{\theta} \mid \mathbf{y})p(\mathbf{y})$$

*Posterior x  
Marginal Likelihood*

**Data**  
(observed)

**Parameters**  
(unobserved)

# Generative models

# Generative models

- If we disallow improper priors, then Bayesian modeling is generative



# Generative models

- If we disallow improper priors, then Bayesian modeling is generative
- In particular, we have a simple way to simulate from  $p(y)$ :

# Generative models

- If we disallow improper priors, then Bayesian modeling is generative
- In particular, we have a simple way to simulate from  $p(y)$ :

$$\theta^* \sim p(\theta)$$

# Generative models

- If we disallow improper priors, then Bayesian modeling is generative
- In particular, we have a simple way to simulate from  $p(y)$ :

$$\begin{array}{c} \theta^* \sim p(\theta) \\ \downarrow \\ y^* \sim p(y|\theta^*) \end{array}$$

# Generative models

- If we disallow improper priors, then Bayesian modeling is generative
- In particular, we have a simple way to simulate from  $p(y)$ :

$$\begin{array}{ccc} \theta^* \sim p(\theta) & & \\ \downarrow & \longleftrightarrow & y^* \sim p(y) \\ y^* \sim p(y|\theta^*) & & \end{array}$$

# **Prior predictive checking:**

fake data is almost as useful as real data

*What do vague/non-informative priors imply about the data our model can generate?*

# Prior predictive checking:

fake data is almost as useful as real data

*What do vague/non-informative priors imply about the data our model can generate?*

$$\alpha_0 \sim N(0, 100)$$

$$\beta_0 \sim N(0, 100)$$

$$\tau_\alpha^2 \sim \text{InvGamma}(1, 100)$$

$$\tau_\beta^2 \sim \text{InvGamma}(1, 100)$$

# Prior predictive checking:

fake data is almost as useful as real data

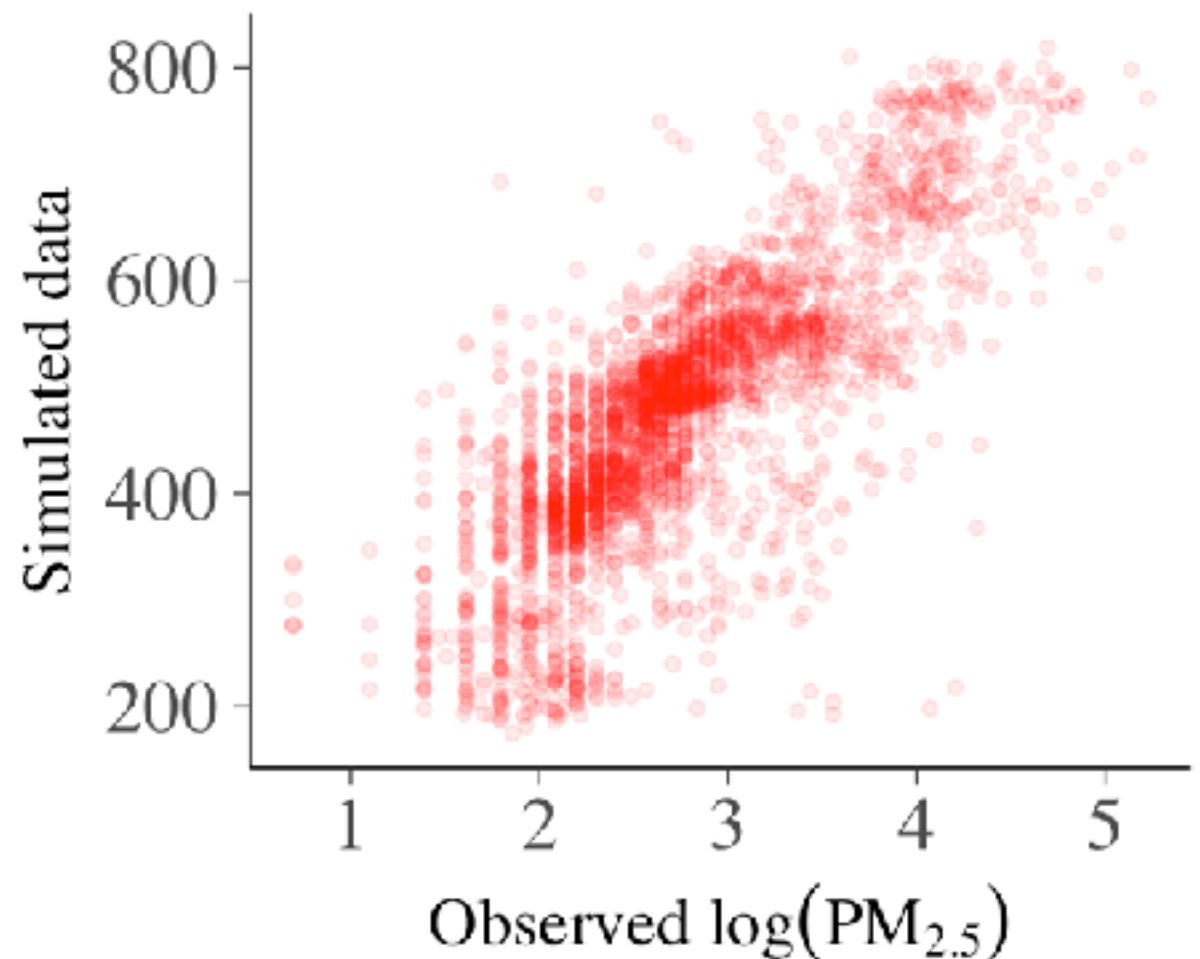
*What do vague/non-informative priors imply about the data our model can generate?*

$$\alpha_0 \sim N(0, 100)$$

$$\beta_0 \sim N(0, 100)$$

$$\tau_\alpha^2 \sim \text{InvGamma}(1, 100)$$

$$\tau_\beta^2 \sim \text{InvGamma}(1, 100)$$



# Prior predictive checking:

fake data is almost as useful as real data

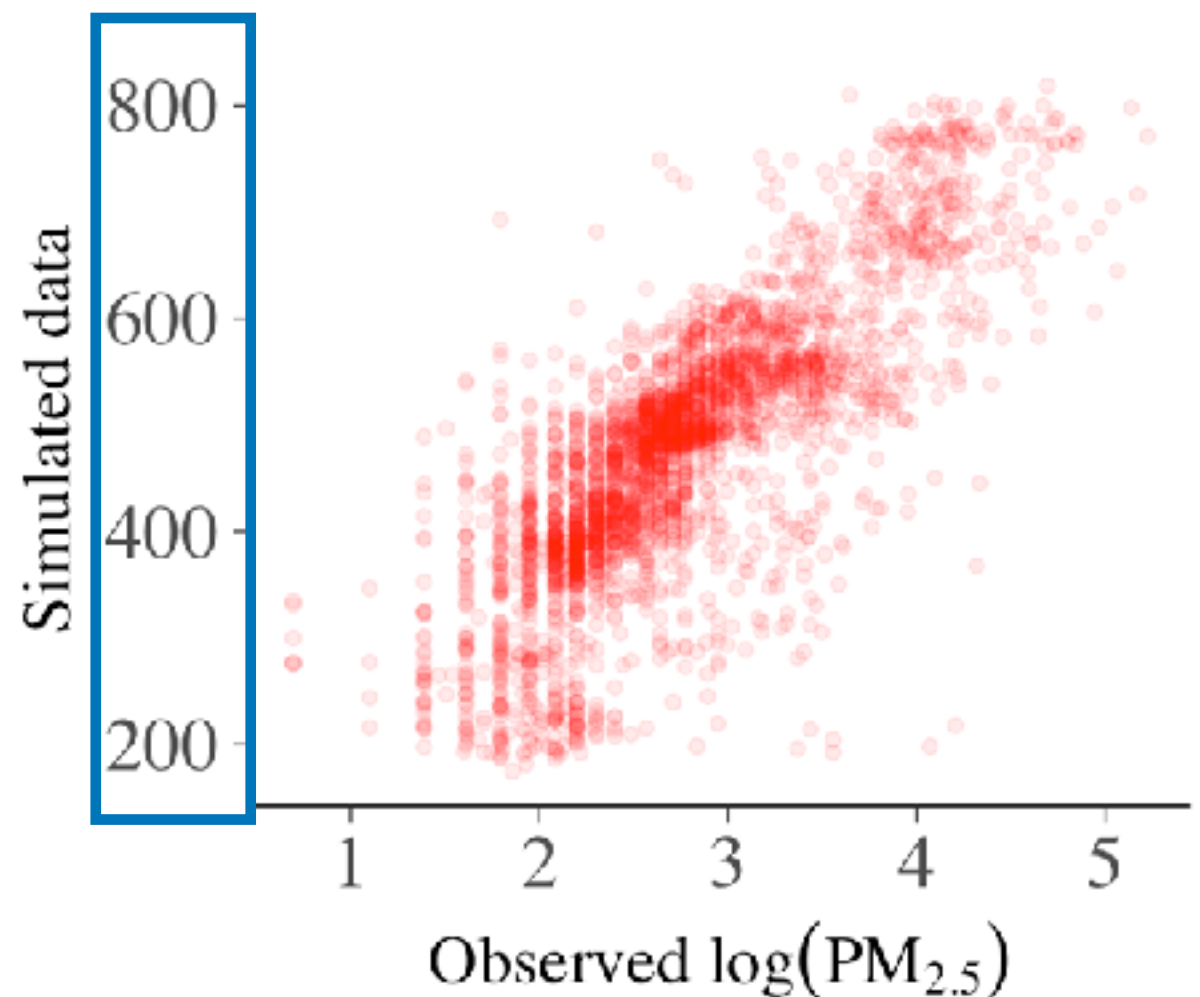
*What do vague/non-informative priors imply about the data our model can generate?*

$$\alpha_0 \sim N(0, 100)$$

$$\beta_0 \sim N(0, 100)$$

$$\tau_\alpha^2 \sim \text{InvGamma}(1, 100)$$

$$\tau_\beta^2 \sim \text{InvGamma}(1, 100)$$

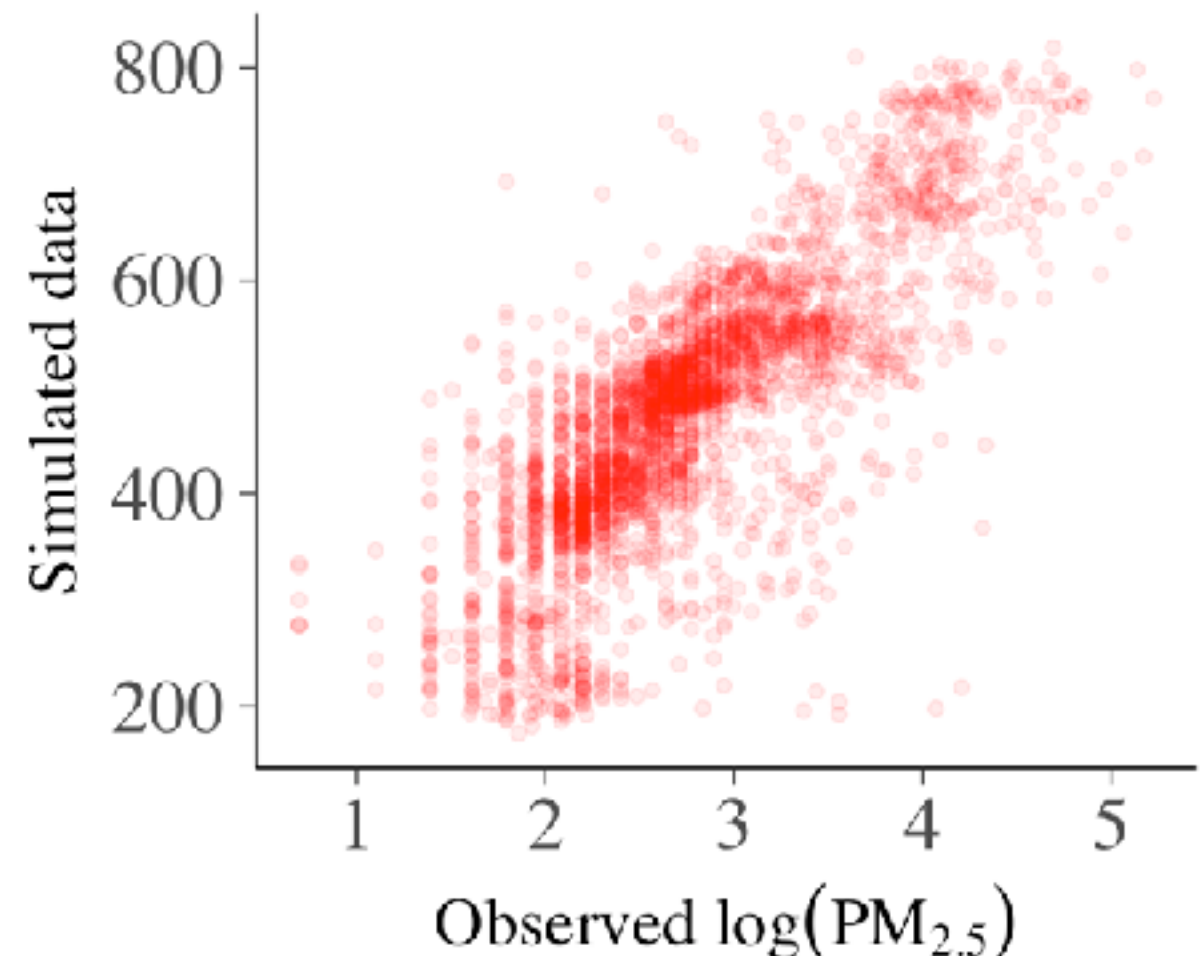




# Prior predictive checking:

fake data is almost as useful as real data

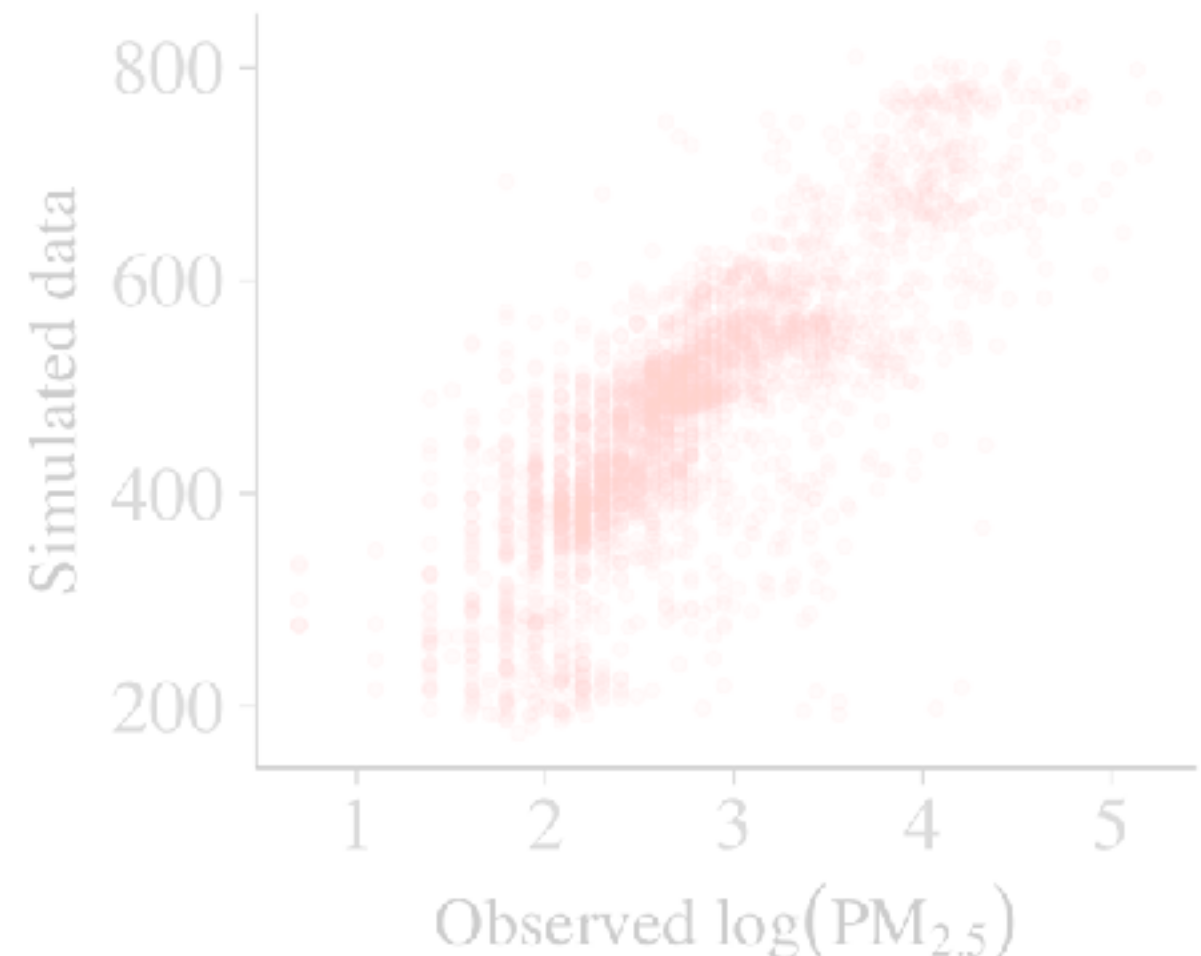
- The prior model is **two orders of magnitude** off the real data
- Two orders of magnitude **on the log scale!**



# Prior predictive checking:

fake data is almost as useful as real data

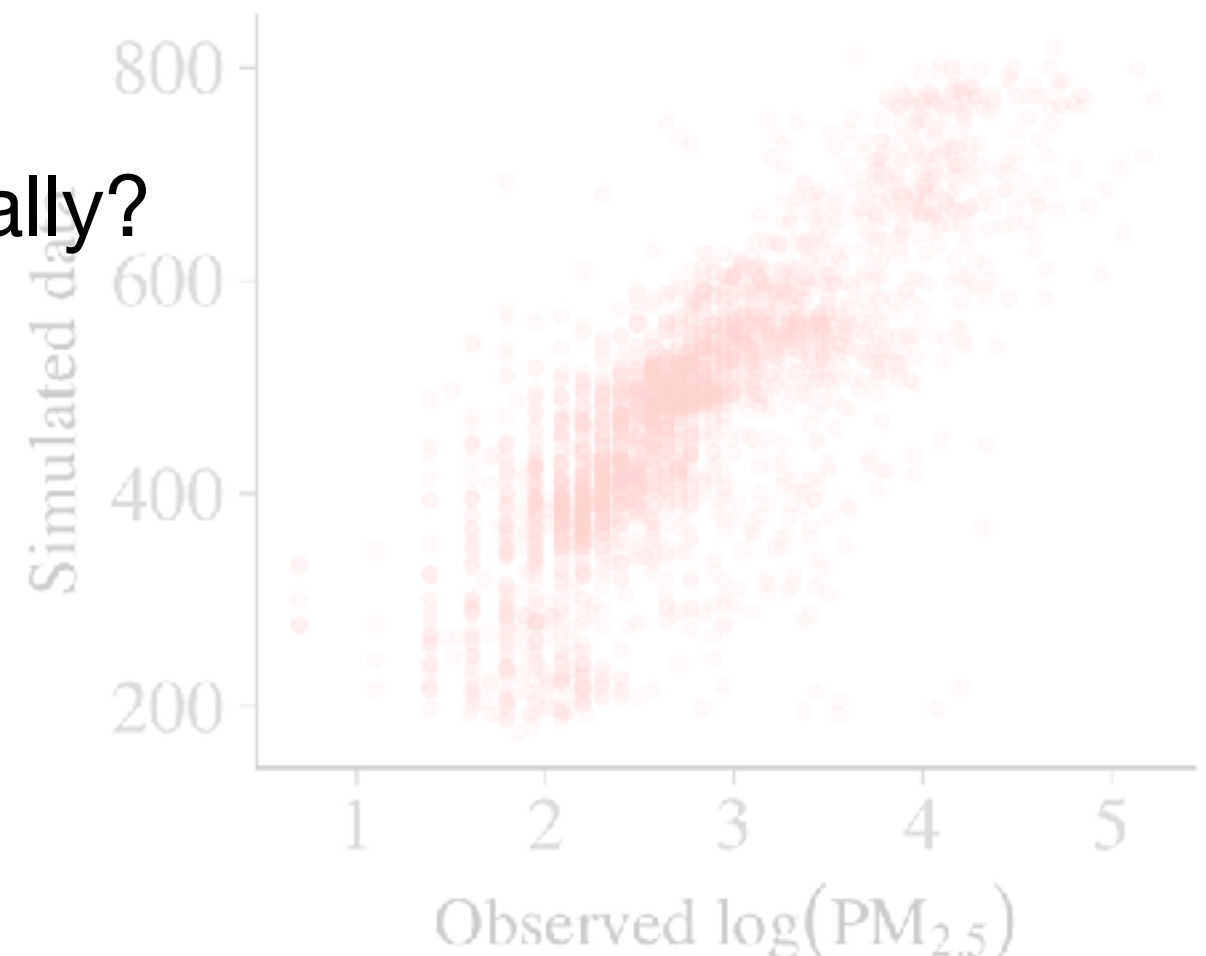
- The prior model is **two orders of magnitude** off the real data
- Two orders of magnitude **on the log scale!**



# Prior predictive checking:

fake data is almost as useful as real data

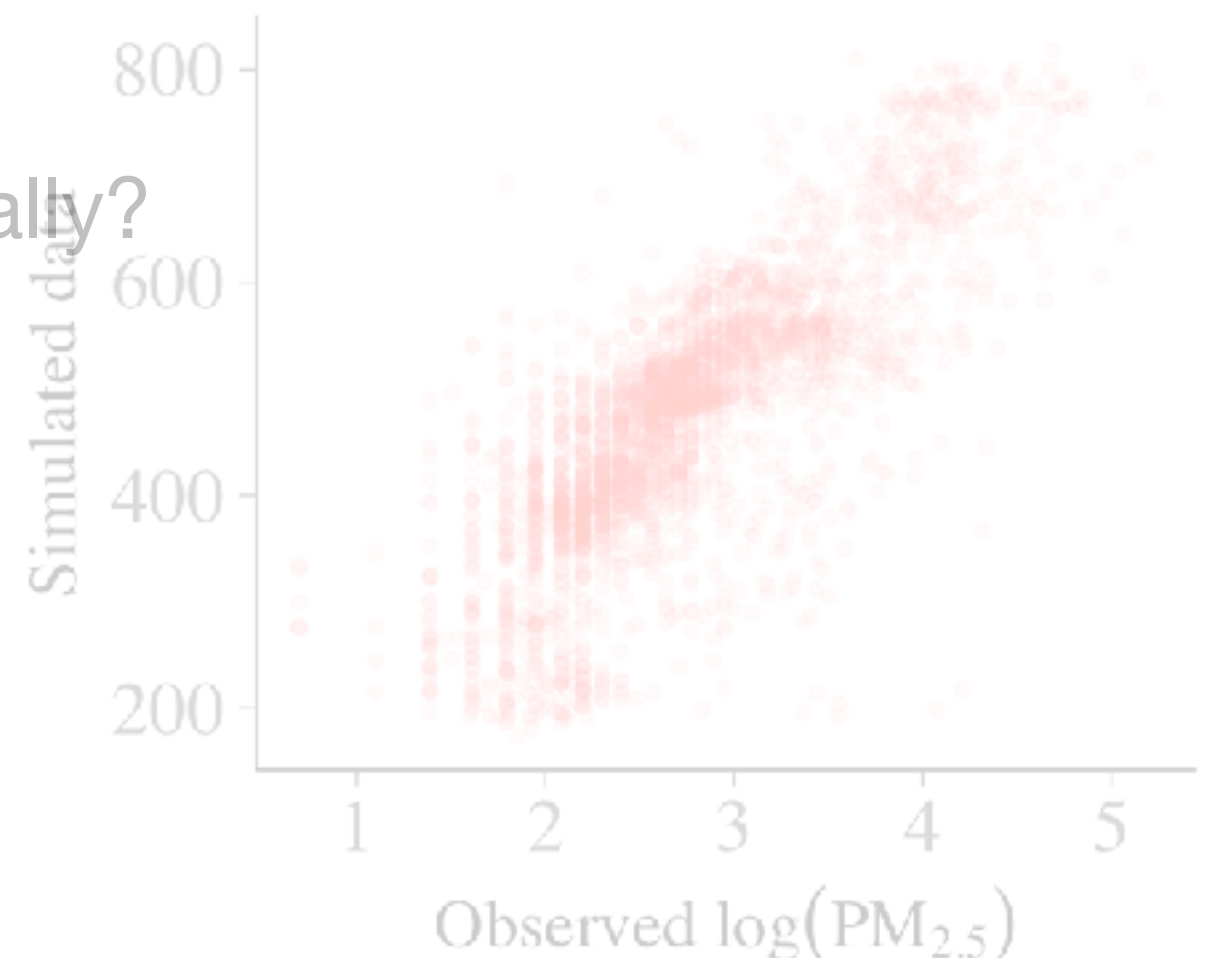
- The prior model is **two orders of magnitude** off the real data
- Two orders of magnitude **on the log scale!**
- What does this mean practically?



# Prior predictive checking:

fake data is almost as useful as real data

- The prior model is **two orders of magnitude** off the real data
- Two orders of magnitude **on the log scale!**
- What does this mean practically?
- The data will have to overcome the prior...



## **Prior predictive checking:**

fake data is almost as useful as real data

*What are better priors for the global intercept and slope  
and the hierarchical scale parameters?*

# Prior predictive checking:

fake data is almost as useful as real data

*What are better priors for the global intercept and slope  
and the hierarchical scale parameters?*

$$\alpha_0 \sim N(0, 1)$$

$$\beta_0 \sim N(1, 1)$$

$$\tau_\alpha \sim N_+(0, 1)$$

$$\tau_\beta \sim N_+(0, 1)$$

# Prior predictive checking:

fake data is almost as useful as real data

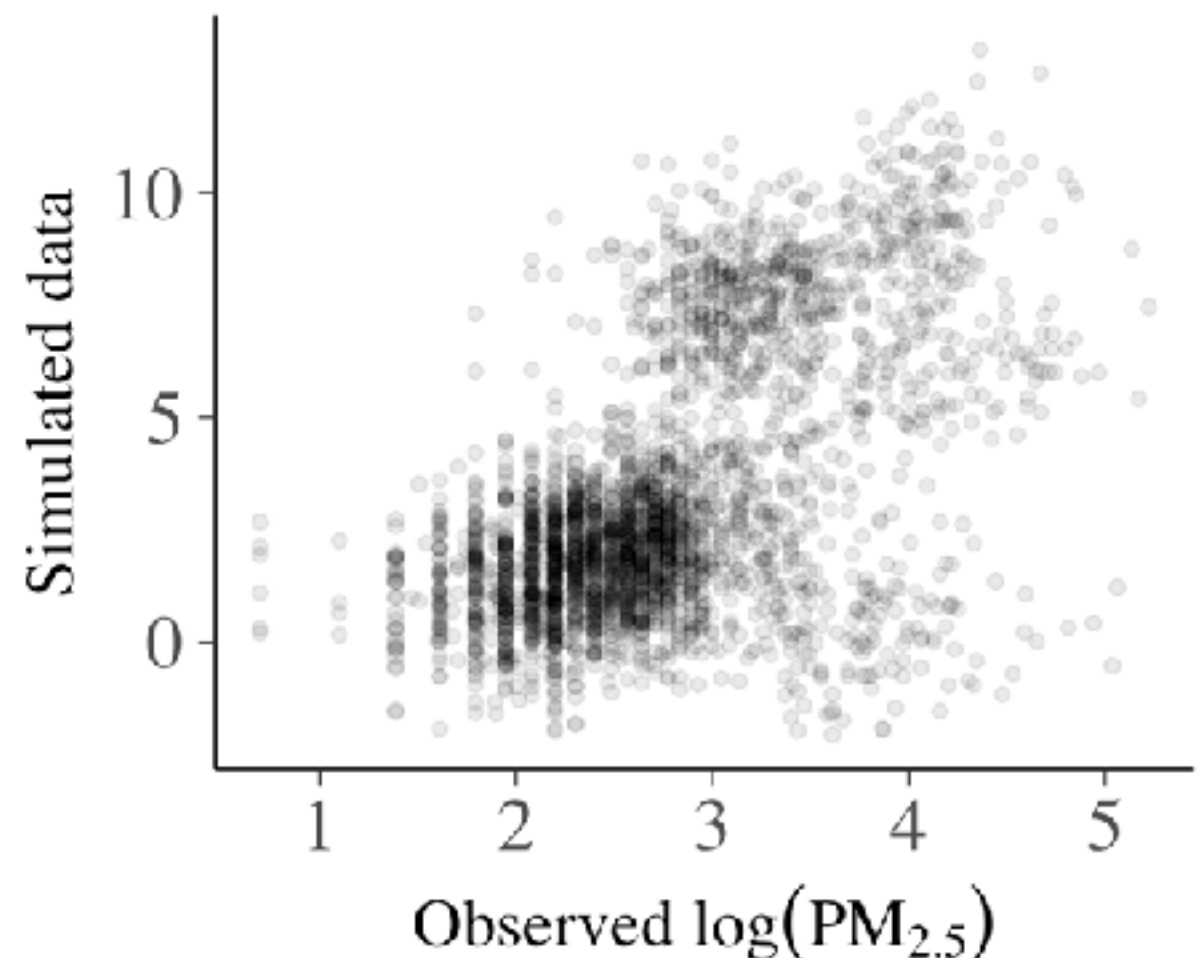
*What are better priors for the global intercept and slope and the hierarchical scale parameters?*

$$\alpha_0 \sim N(0, 1)$$

$$\beta_0 \sim N(1, 1)$$

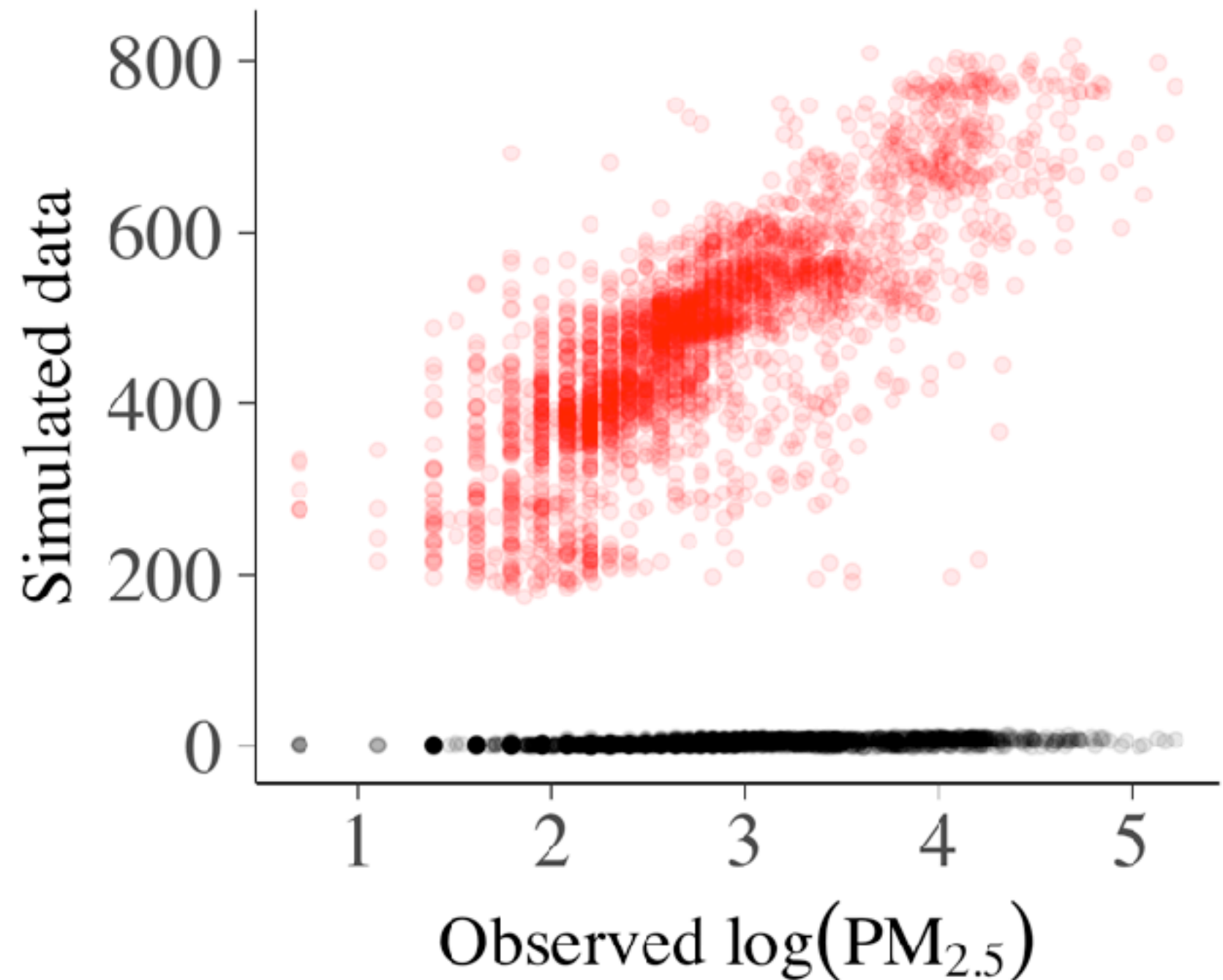
$$\tau_\alpha \sim N_+(0, 1)$$

$$\tau_\beta \sim N_+(0, 1)$$



# Prior predictive checking:

fake data is almost as useful as real data

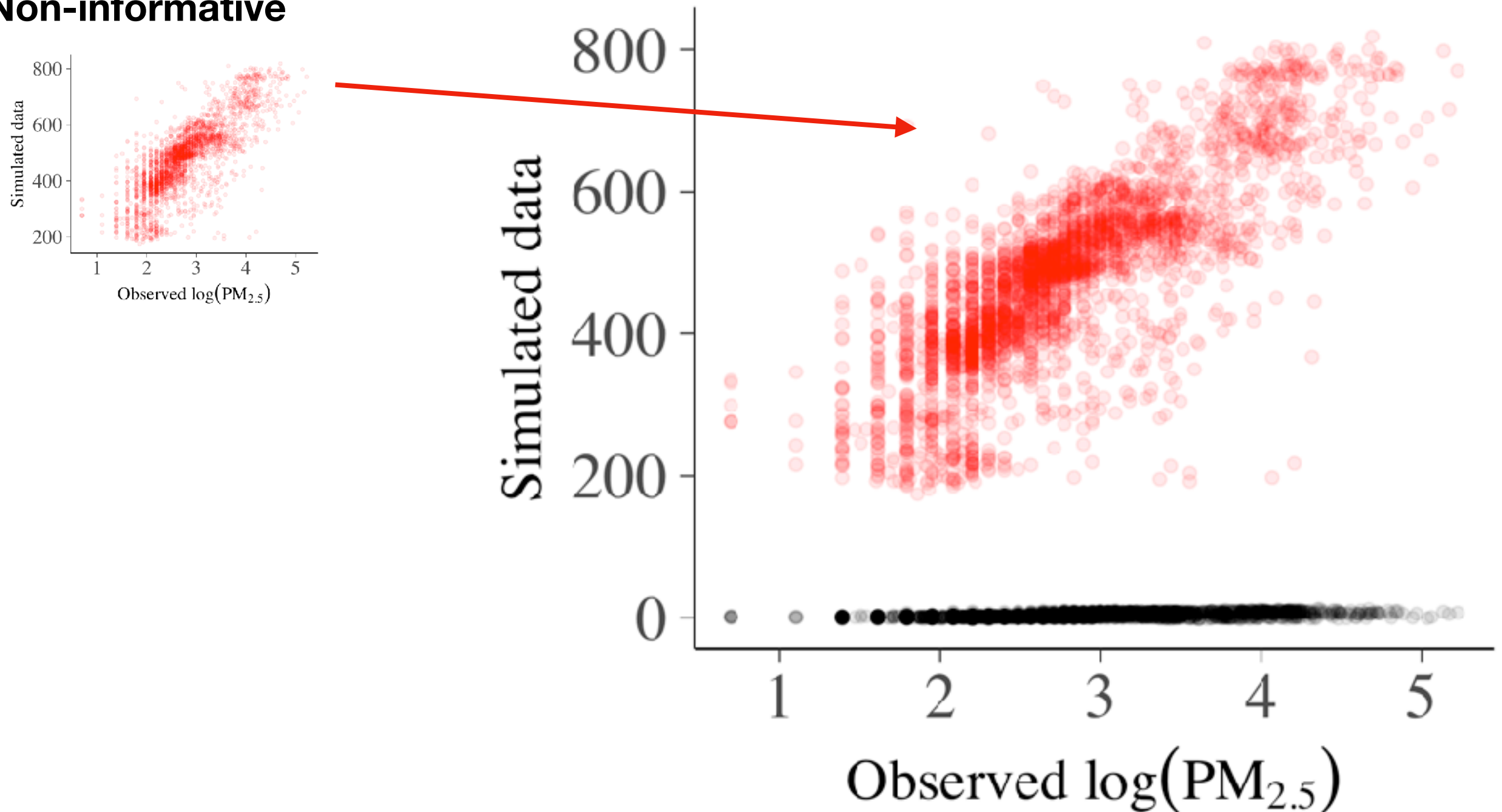




# Prior predictive checking:

fake data is almost as useful as real data

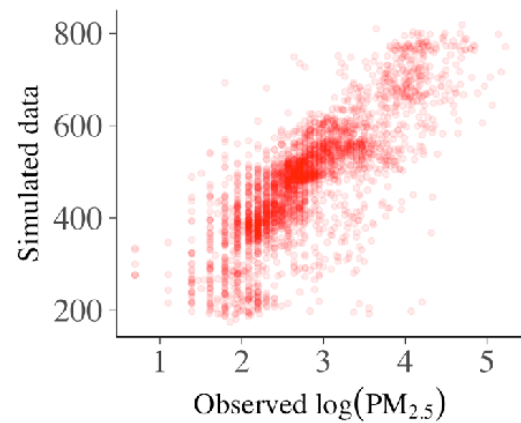
## Non-informative



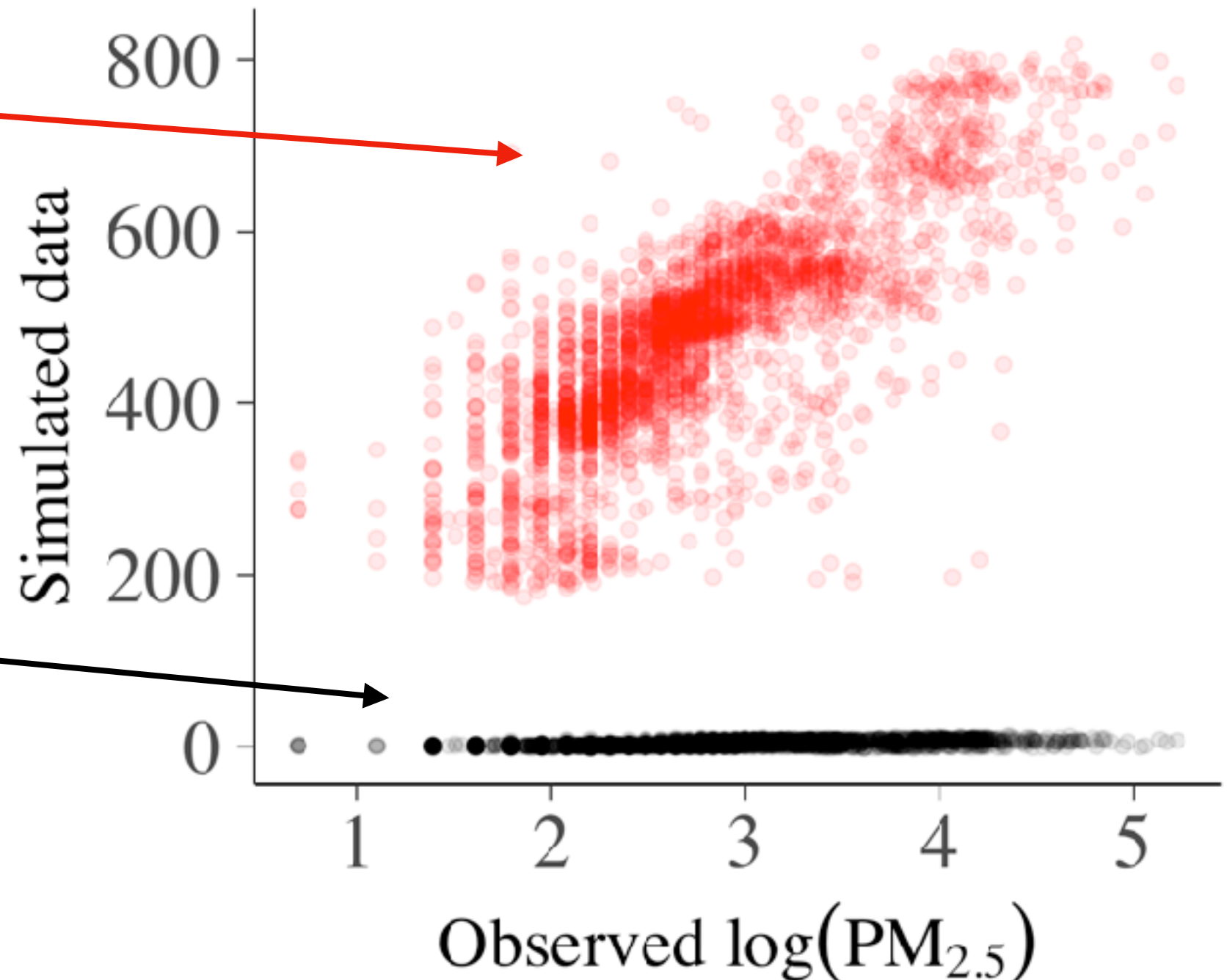
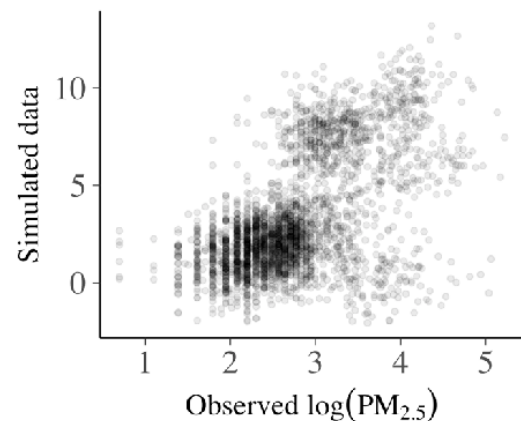
# Prior predictive checking:

fake data is almost as useful as real data

## Non-informative



## Weakly informative

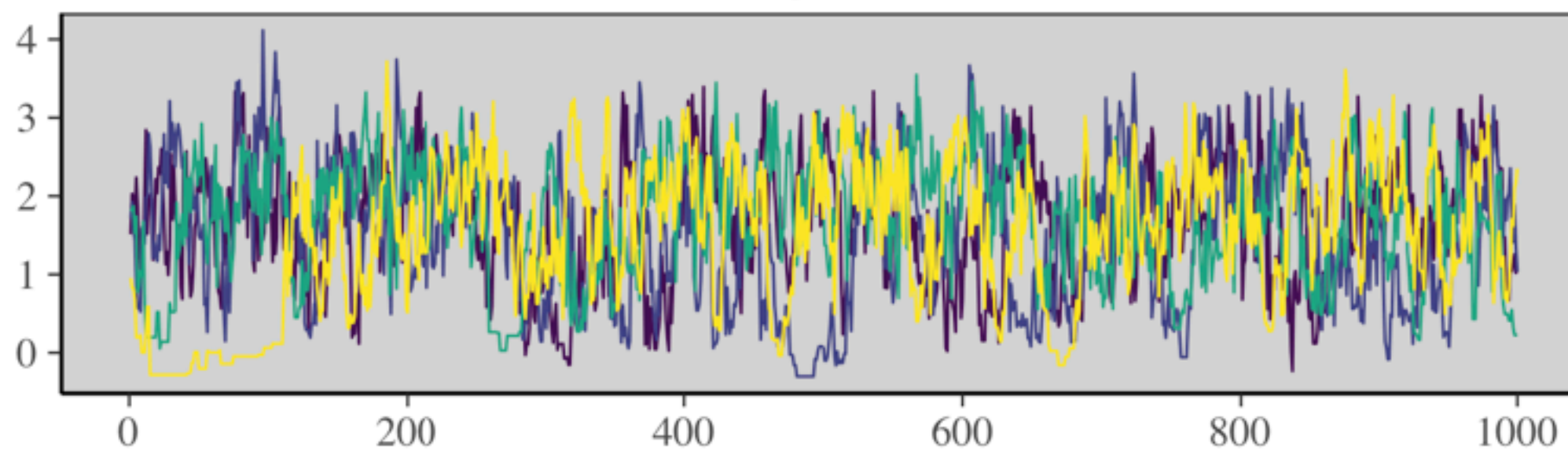


# MCMC diagnostics

*Beyond trace plots*

<https://chi-feng.github.io/mcmc-demo/>

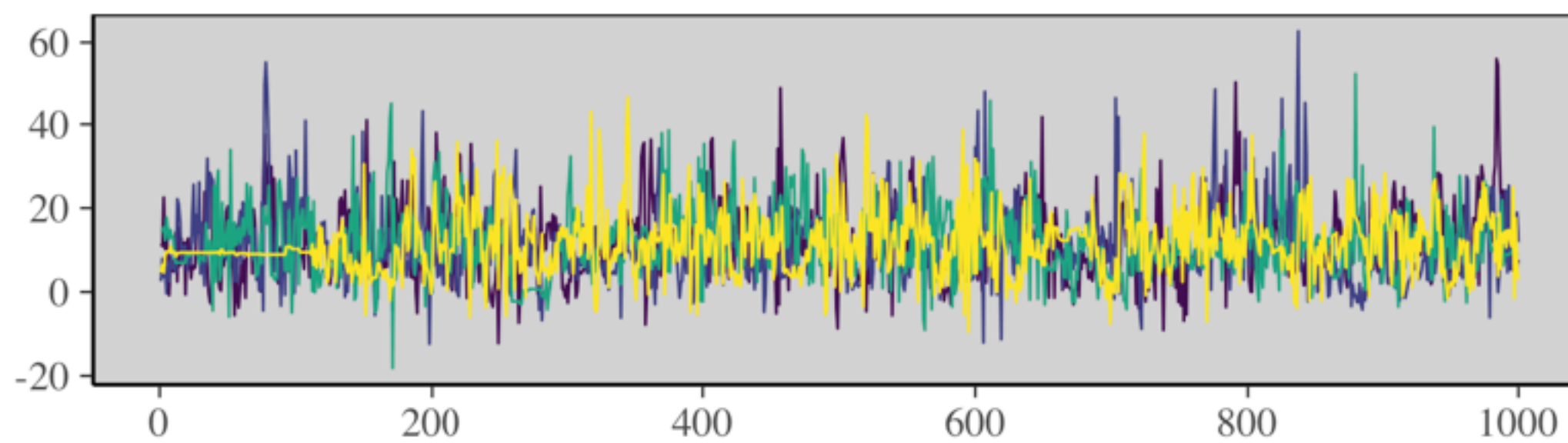
$\log(\tau)$



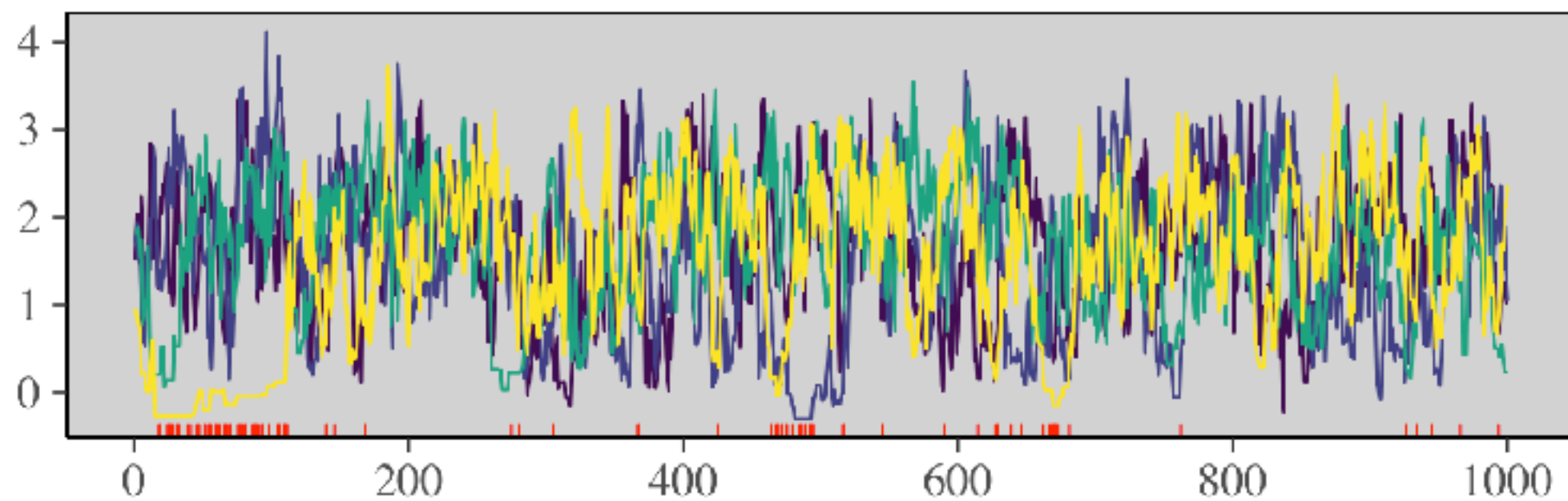
Chain

- 1
- 2
- 3
- 4

$\theta[1]$



$\log(\tau)$



Chain

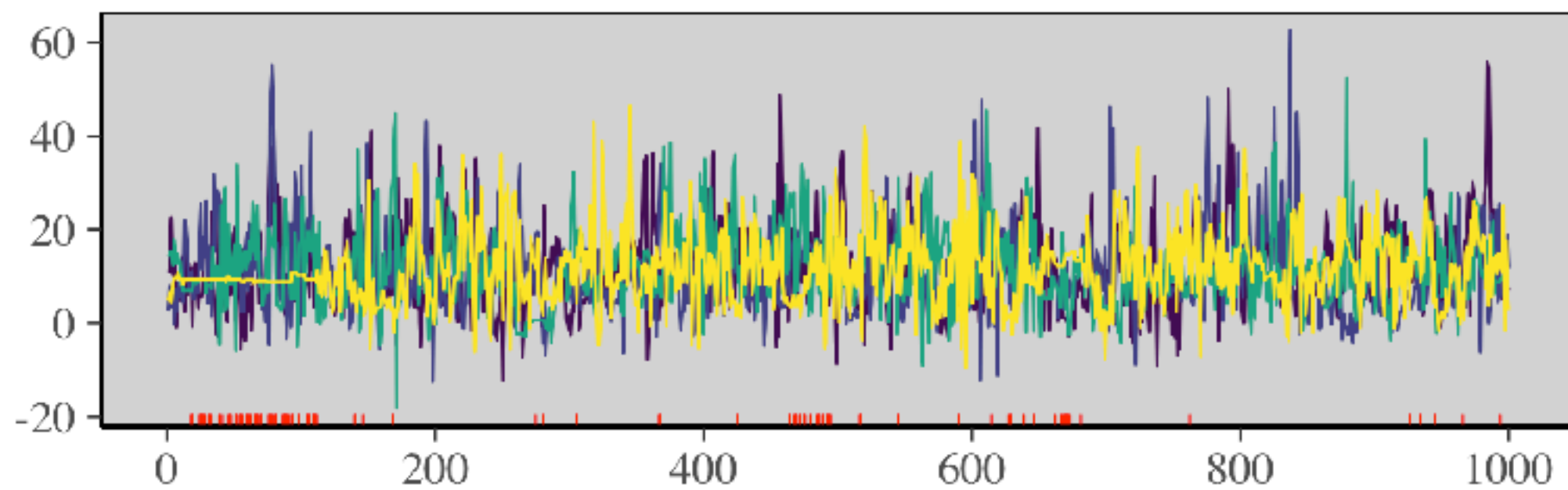
1

2

3

4

$\theta[1]$



Divergence

# MCMC diagnostics

## beyond trace plots

---

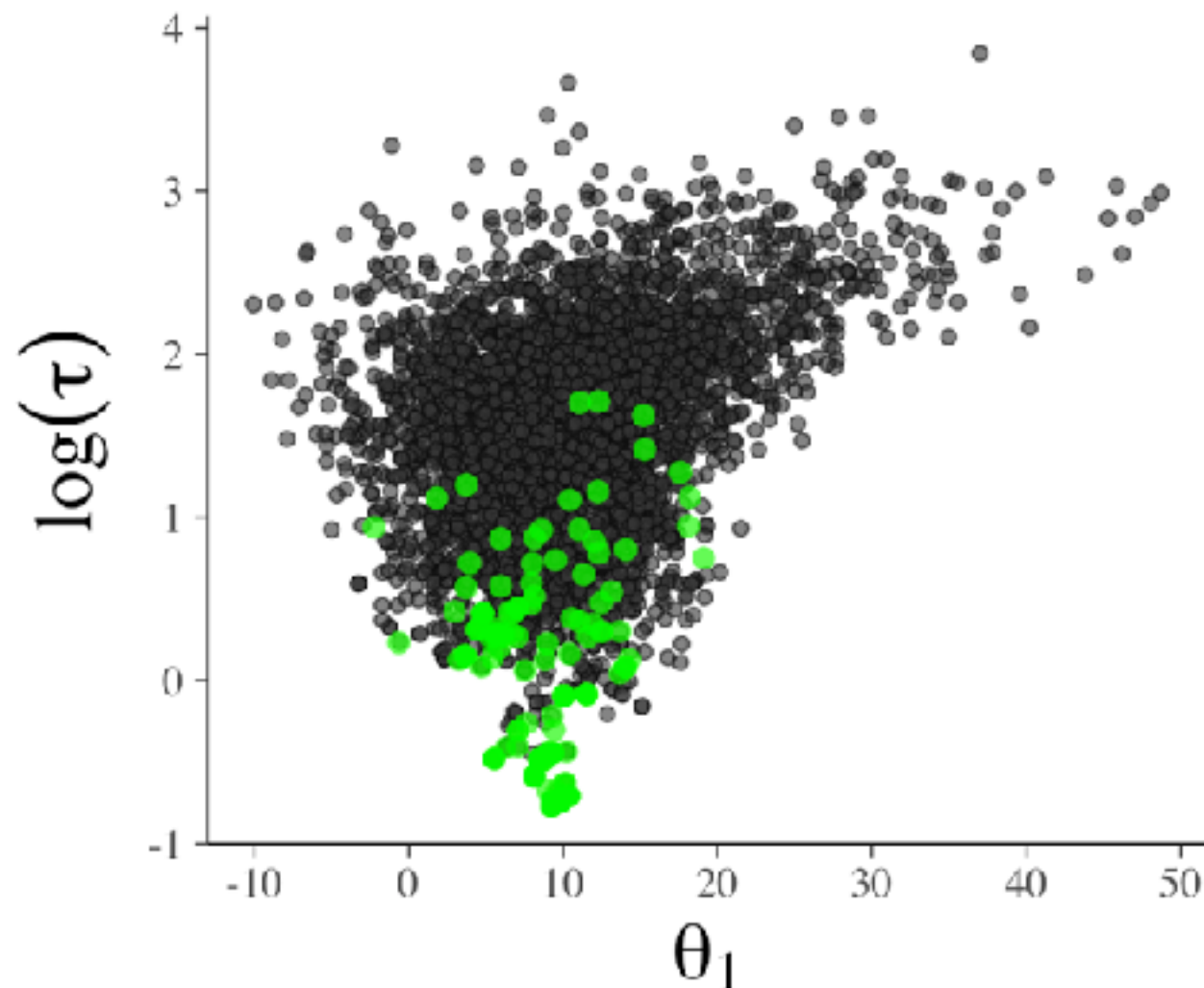
Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., and Gelman, A. (2018).  
**Visualization in Bayesian workflow.**  
*Journal of the Royal Statistical Society Series A*, accepted for publication.  
[arxiv.org/abs/1709.01449](https://arxiv.org/abs/1709.01449) | [github.com/jgabry/bayes-vis-paper](https://github.com/jgabry/bayes-vis-paper)

Betancourt, M. (2017).  
**A conceptual introduction to Hamiltonian Monte Carlo.**  
arXiv preprint:  
[arxiv.org/abs/1701.02434](https://arxiv.org/abs/1701.02434)



# MCMC diagnostics

## beyond trace plots

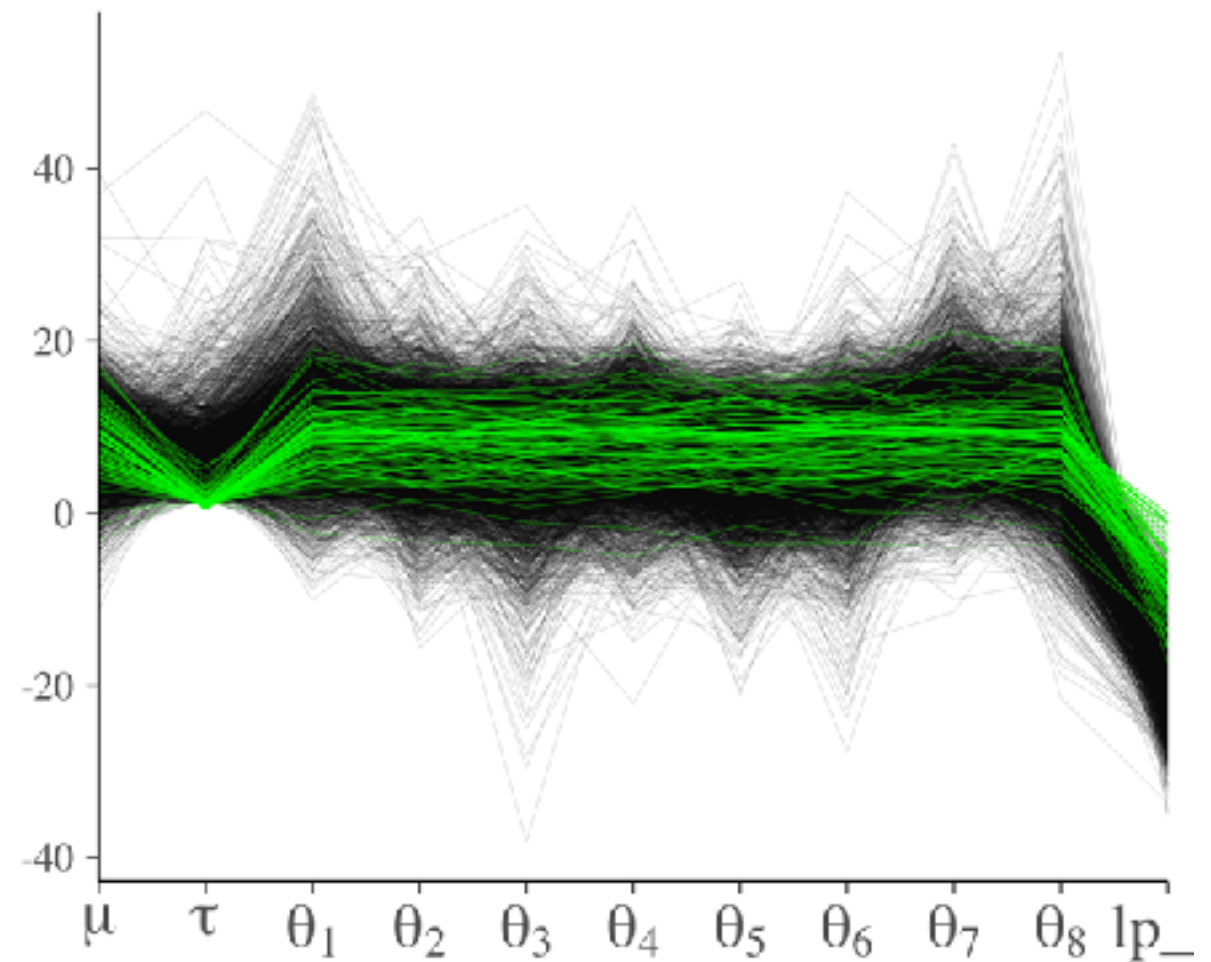
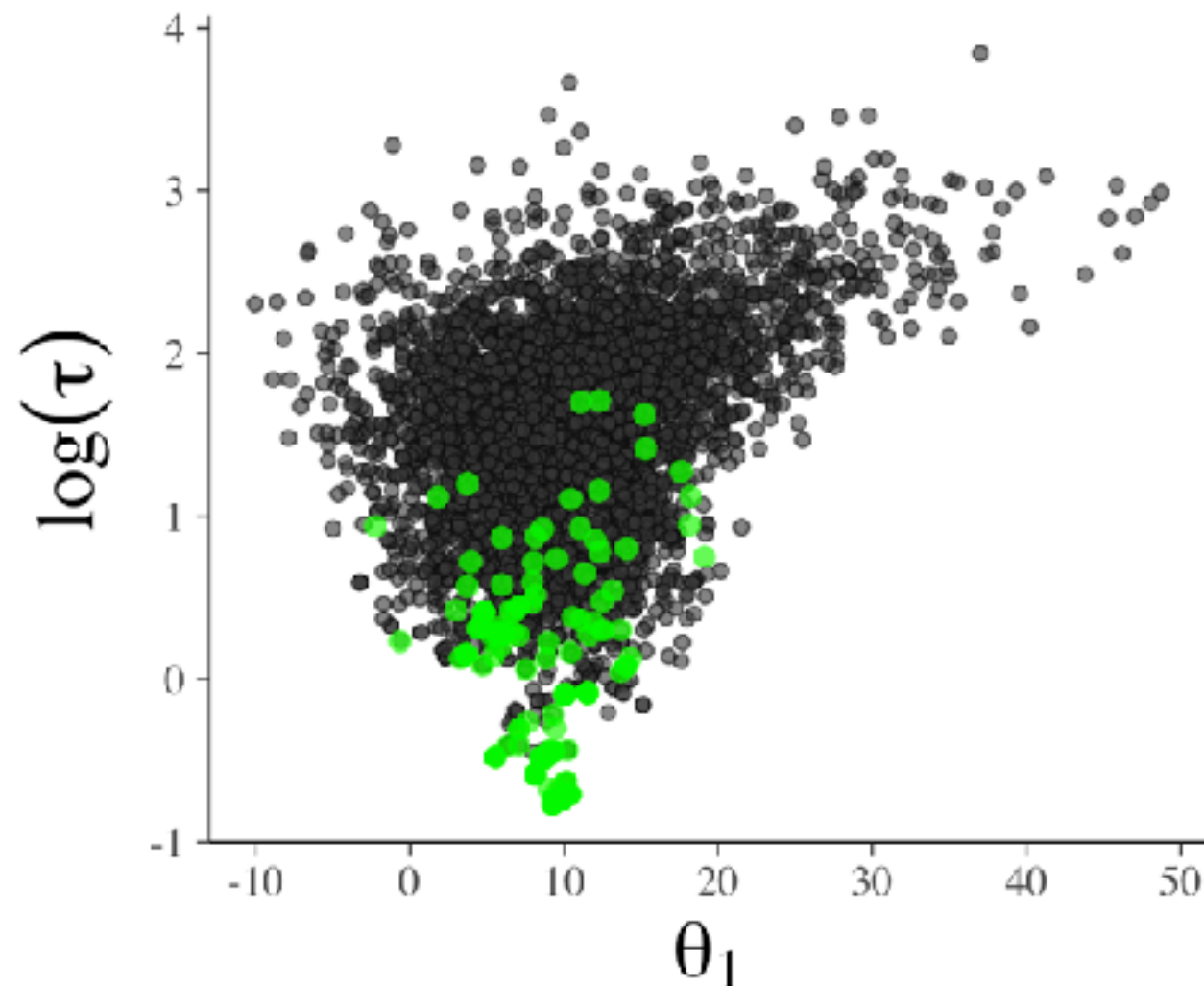


Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., and Gelman, A. (2018).  
**Visualization in Bayesian workflow.**  
*Journal of the Royal Statistical Society Series A*, accepted for publication.  
[arxiv.org/abs/1709.01449](https://arxiv.org/abs/1709.01449) | [github.com/jgabry/bayes-vis-paper](https://github.com/jgabry/bayes-vis-paper)

Betancourt, M. (2017).  
**A conceptual introduction to Hamiltonian Monte Carlo.**  
arXiv preprint:  
[arxiv.org/abs/1701.02434](https://arxiv.org/abs/1701.02434)

# MCMC diagnostics

## beyond trace plots



Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., and Gelman, A. (2018).  
**Visualization in Bayesian workflow.**  
*Journal of the Royal Statistical Society Series A*, accepted for publication.  
[arxiv.org/abs/1709.01449](https://arxiv.org/abs/1709.01449) | [github.com/jgabry/bayes-vis-paper](https://github.com/jgabry/bayes-vis-paper)

Betancourt, M. (2017).  
**A conceptual introduction to Hamiltonian Monte Carlo.**  
arXiv preprint:  
[arxiv.org/abs/1701.02434](https://arxiv.org/abs/1701.02434)

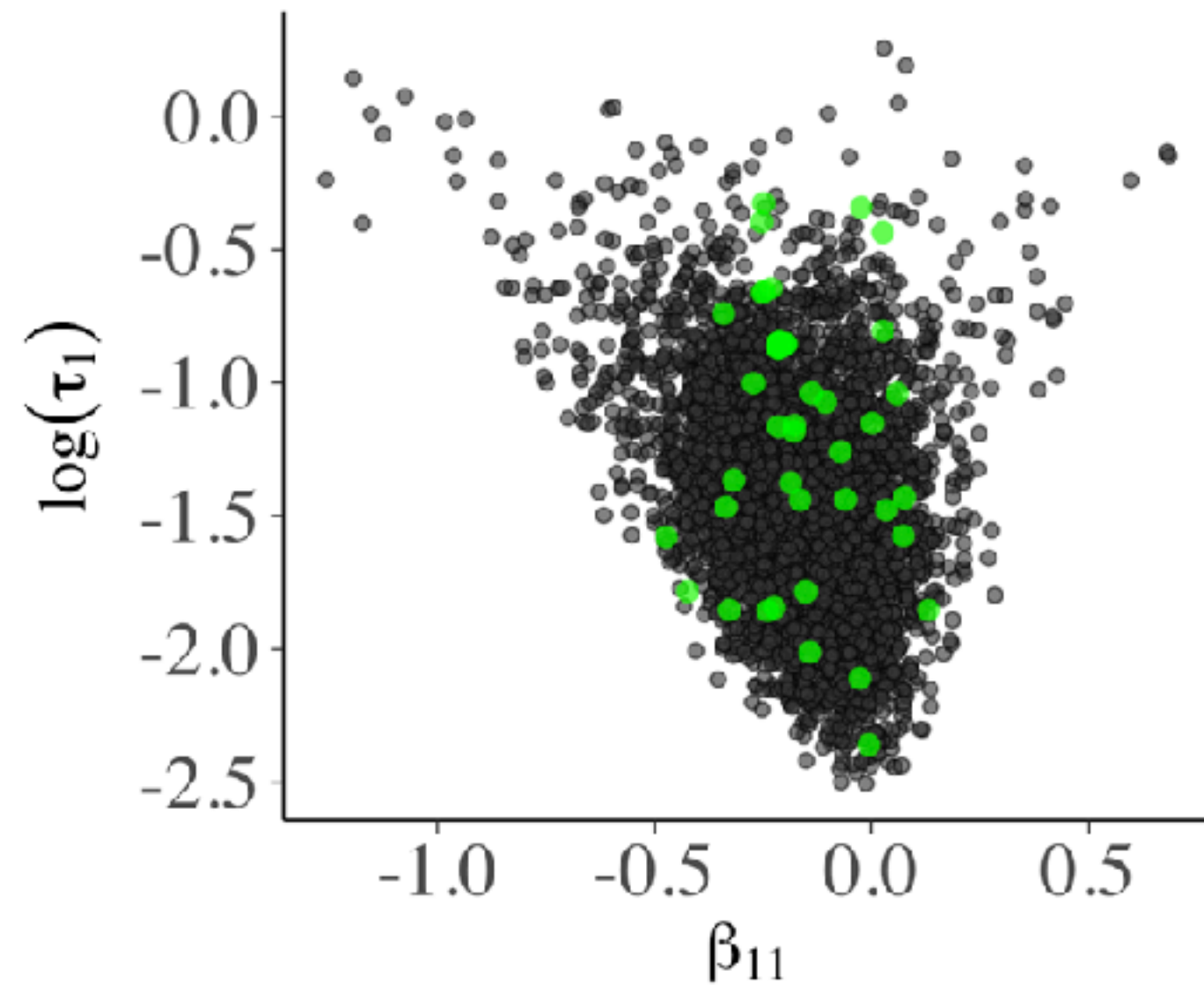


# **MCMC diagnostics**

## beyond trace plots

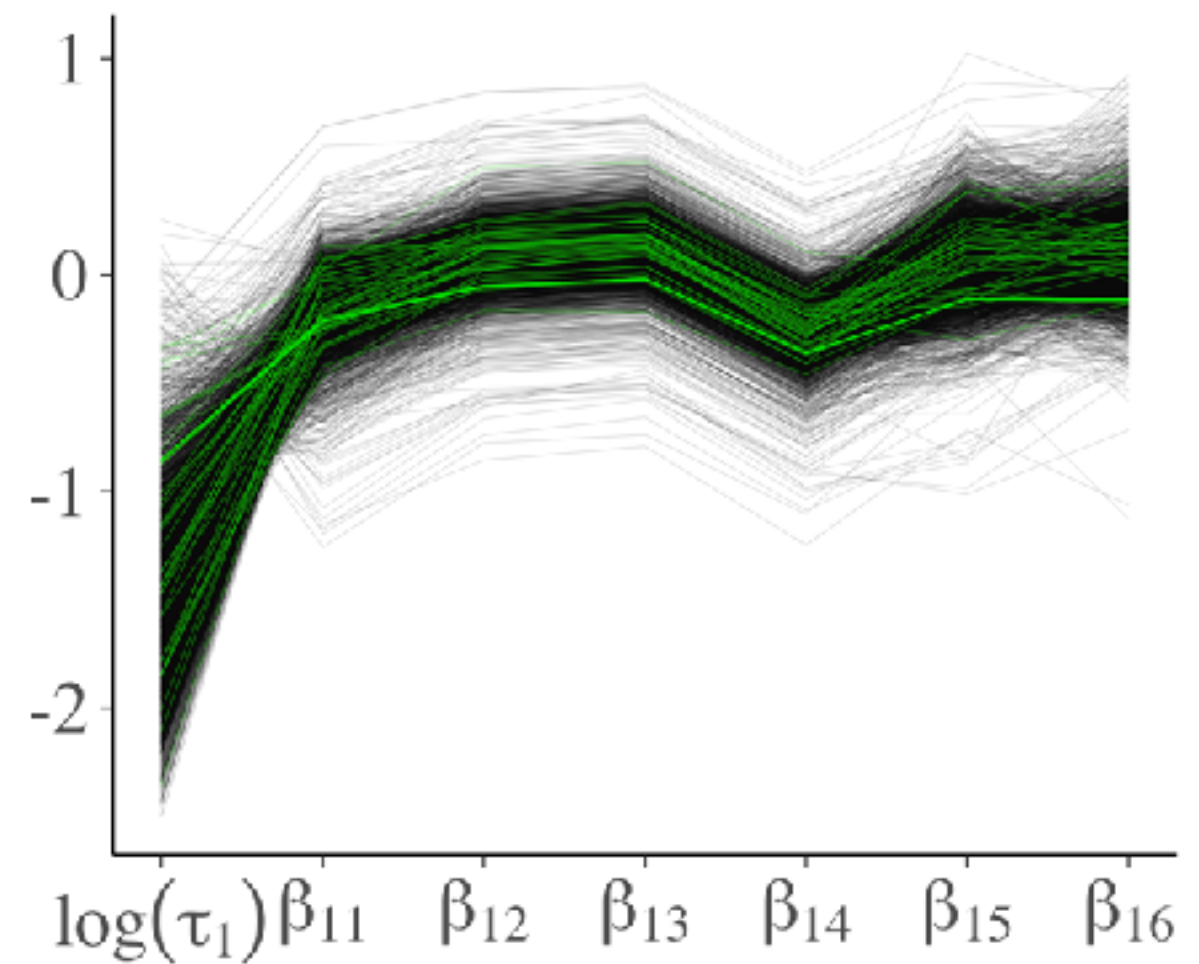
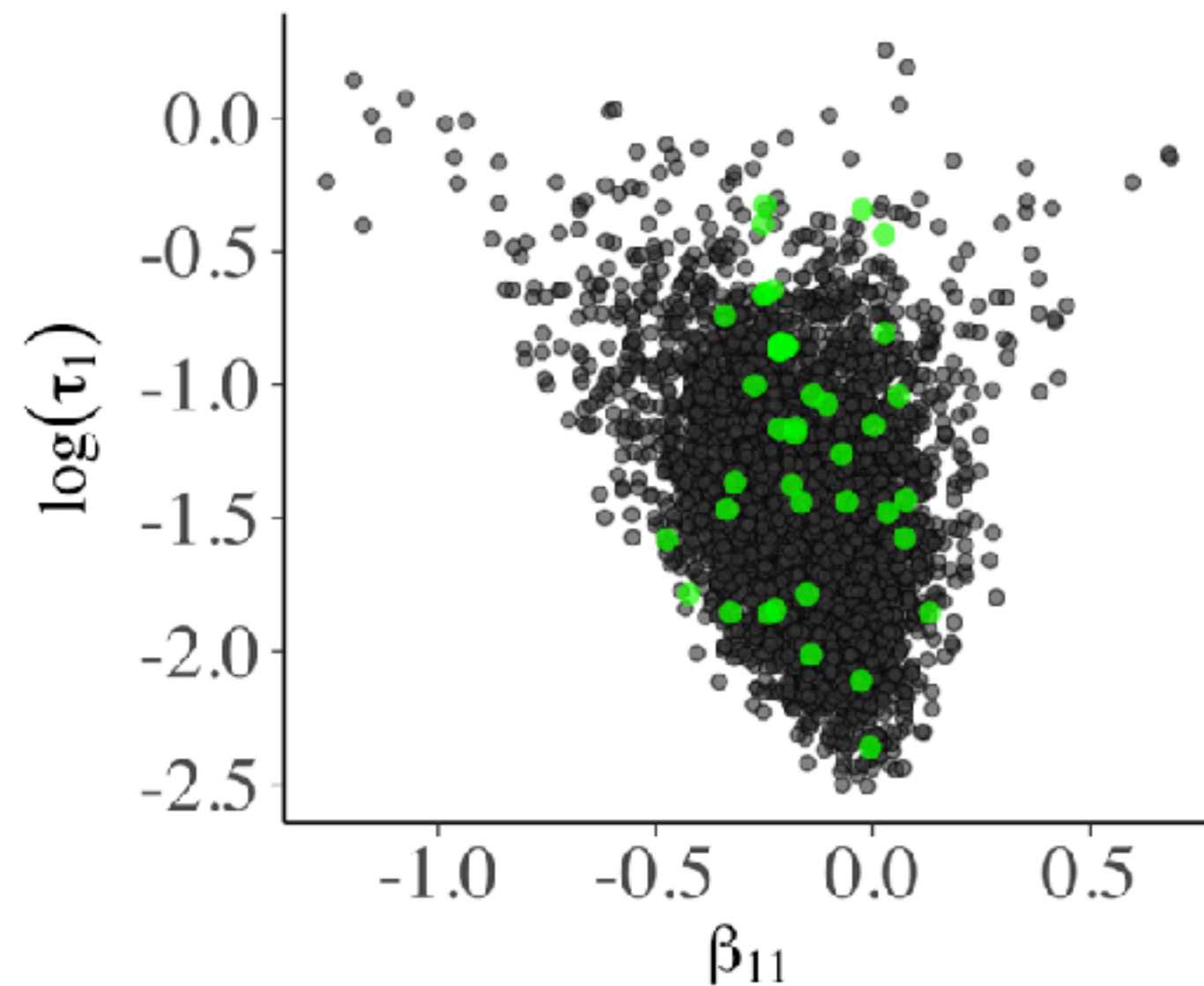
# MCMC diagnostics

beyond trace plots

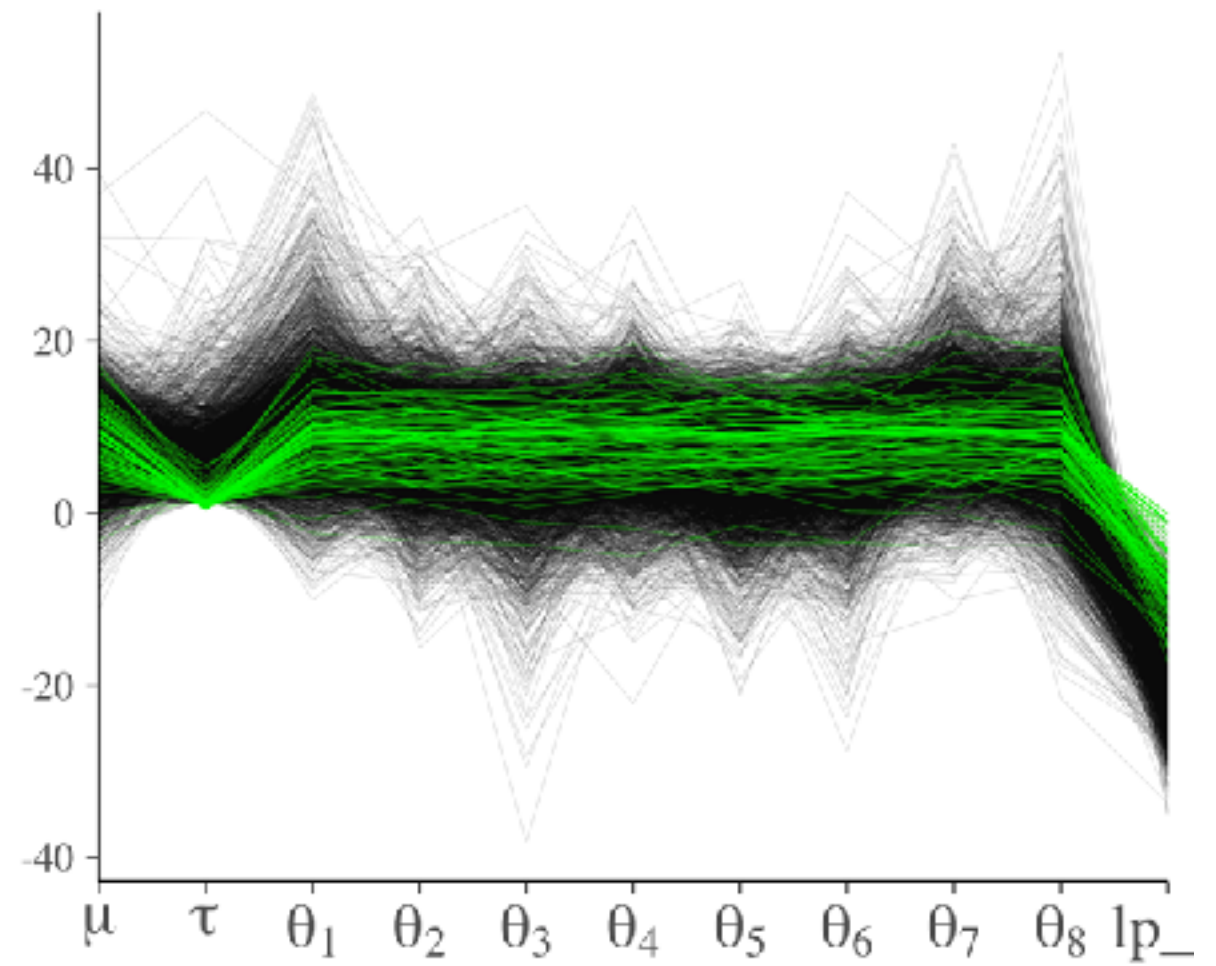
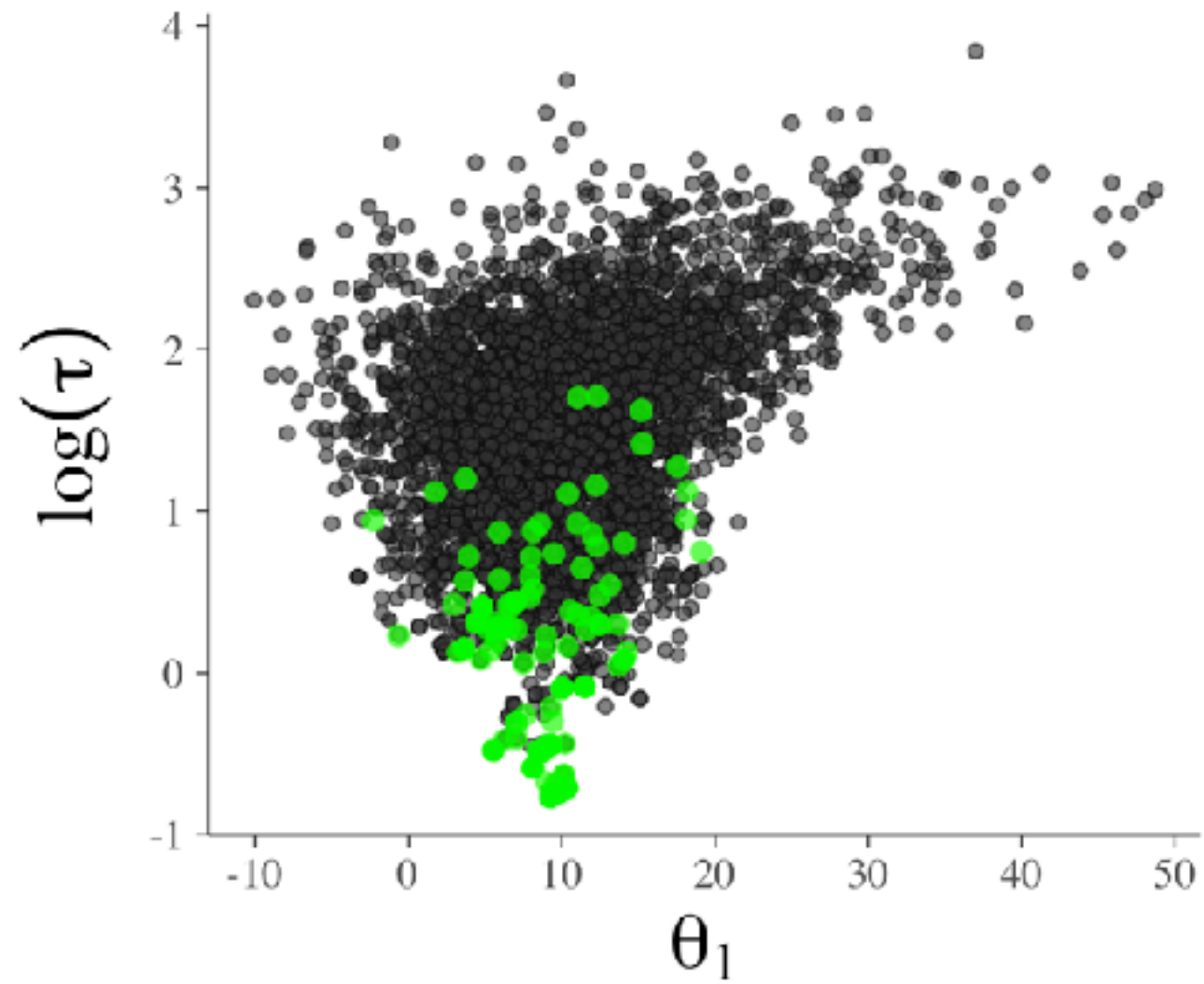


# MCMC diagnostics

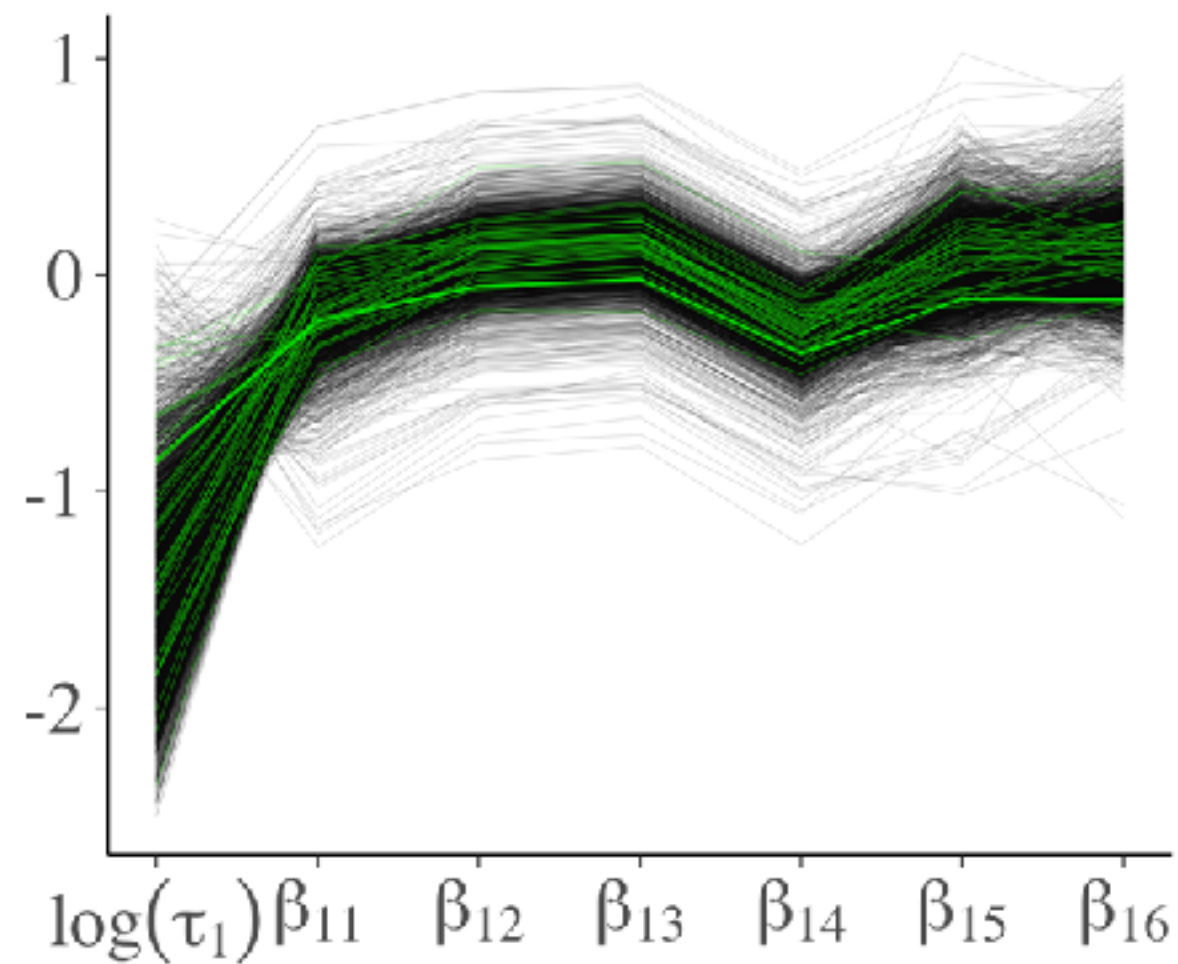
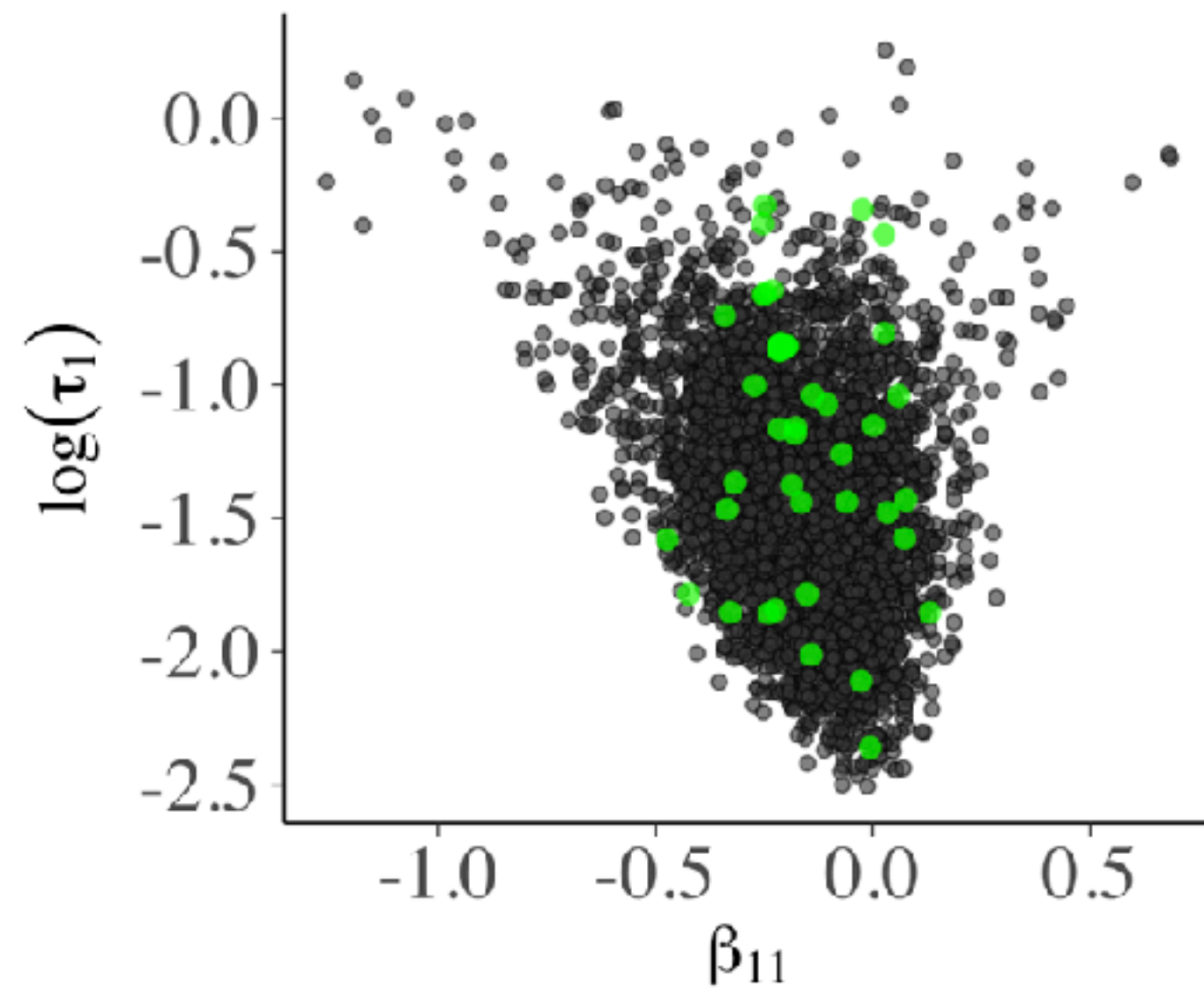
beyond trace plots



# Pathological geometry



# “False positives”



# Posterior predictive checks

*Visual model evaluation*

# Posterior predictive checking

visual model evaluation

The *posterior predictive distribution* is the average data generation process over the entire model

# Posterior predictive checking

visual model evaluation

The *posterior predictive distribution* is the average data generation process over the entire model

$$p(\tilde{y}|y) = \int p(\tilde{y}|\theta) p(\theta|y) d\theta$$



# Posterior predictive checking

## visual model evaluation

# Posterior predictive checking

visual model evaluation

- Misfitting and overfitting both manifest as tension between measurements and predictive distributions

# Posterior predictive checking

visual model evaluation

- Misfitting and overfitting both manifest as tension between measurements and predictive distributions
- Graphical posterior predictive checks visually compare the observed data to the predictive distribution

# Posterior predictive checking

## visual model evaluation

- Misfitting and overfitting both manifest as tension between measurements and predictive distributions
- Graphical posterior predictive checks visually compare the observed data to the predictive distribution

$$\theta^* \sim p(\theta|y)$$

# Posterior predictive checking

## visual model evaluation

- Misfitting and overfitting both manifest as tension between measurements and predictive distributions
- Graphical posterior predictive checks visually compare the observed data to the predictive distribution

$$\theta^* \sim p(\theta|y)$$



$$\tilde{y} \sim p(y|\theta^*)$$

# Posterior predictive checking

## visual model evaluation

- Misfitting and overfitting both manifest as tension between measurements and predictive distributions
- Graphical posterior predictive checks visually compare the observed data to the predictive distribution

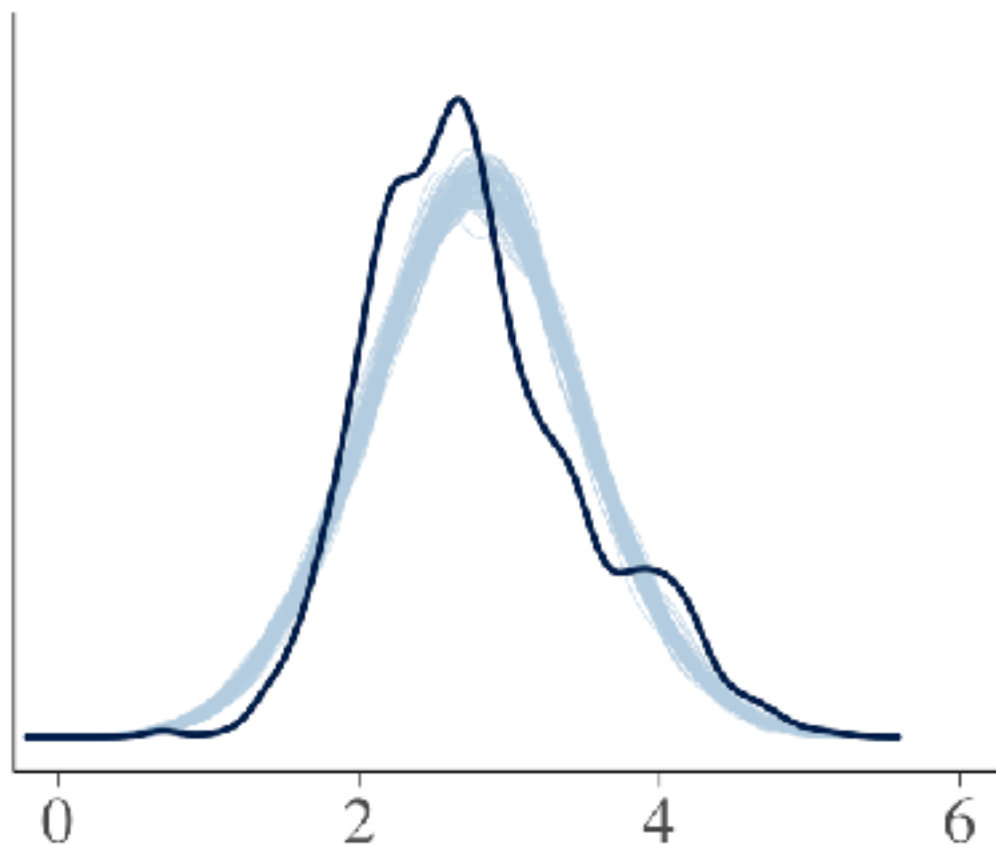
$$\begin{array}{ccc} \theta^* \sim p(\theta|y) & & \\ \downarrow & \longleftrightarrow & \tilde{y} \sim p(\tilde{y}|y) \\ \tilde{y} \sim p(y|\theta^*) & & \end{array}$$

# Posterior predictive checking

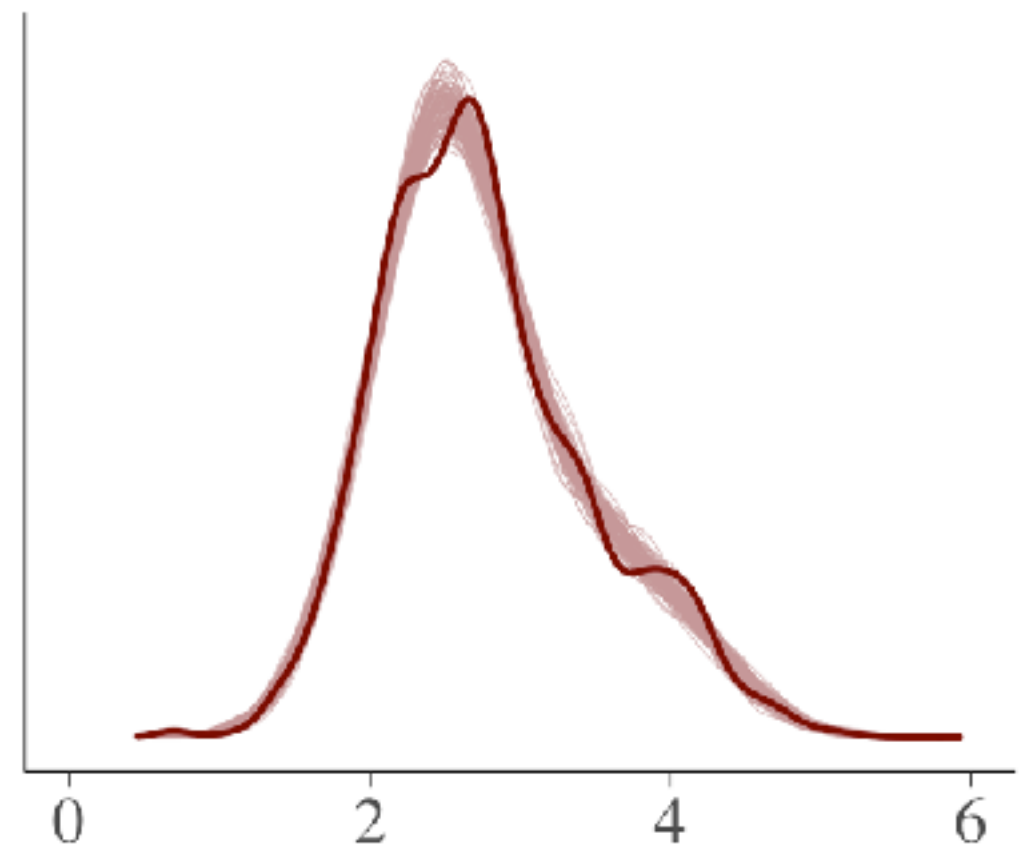
visual model evaluation

Observed data vs posterior predictive simulations

**Model 1 (single level)**



**Model 3 (multilevel)**

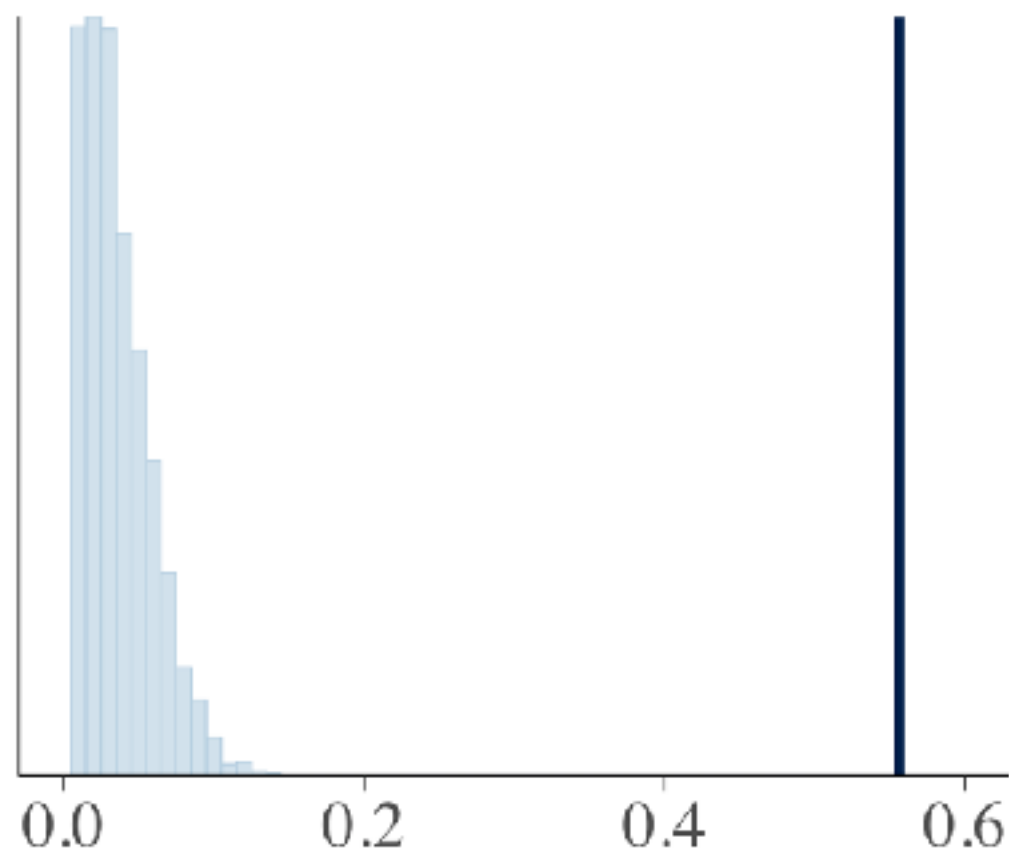


# Posterior predictive checking

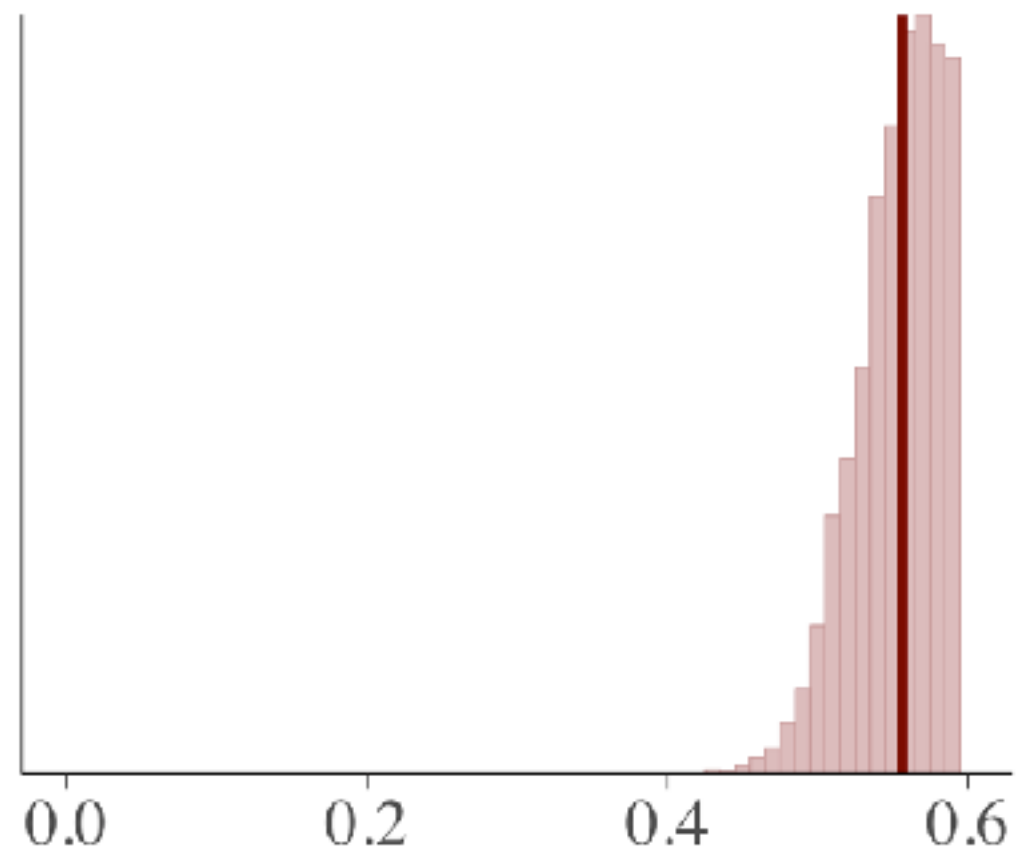
visual model evaluation

## Observed statistics vs posterior predictive statistics

**Model 1 (single level)**



**Model 3 (multilevel)**

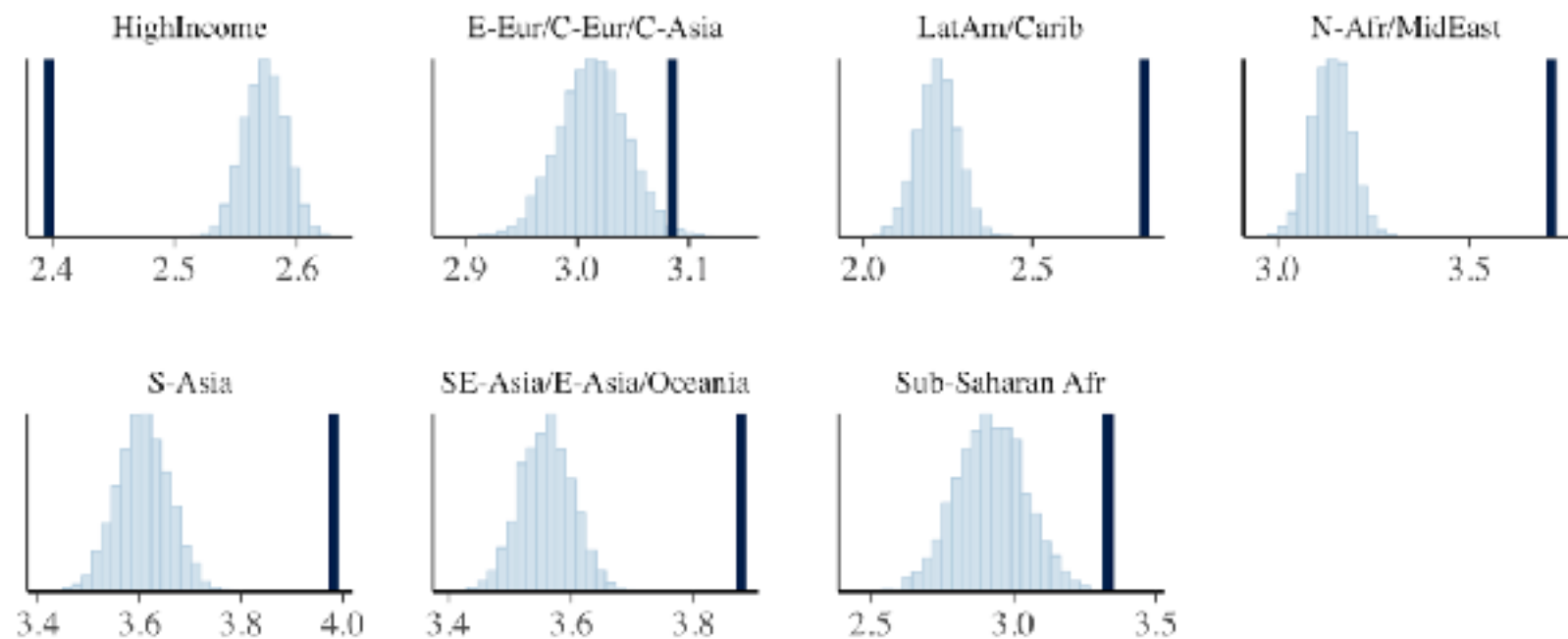


$$T(y) = \text{skew}(y)$$



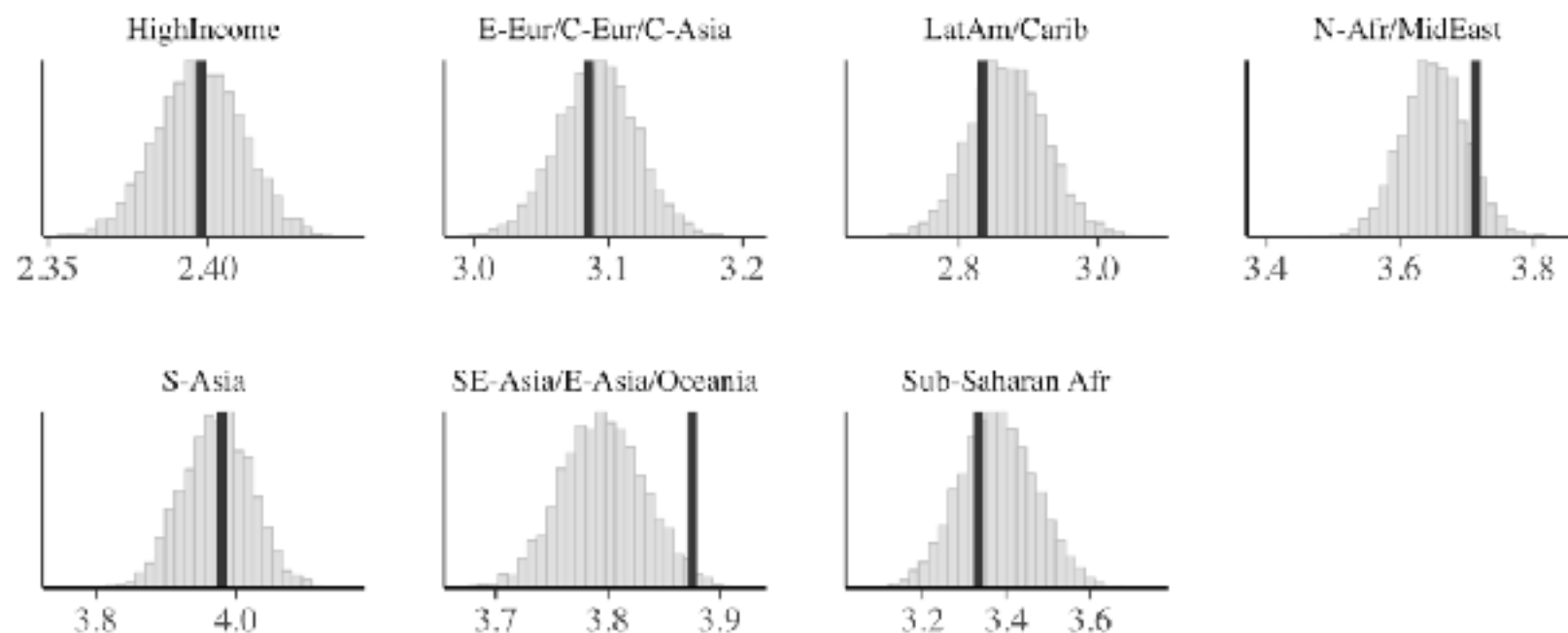
# Posterior predictive checking: visual model evaluation

**Model 1 (single level)**



$$T(y) = \text{med}(y|\text{region})$$

**Model 2 (multilevel)**



# Model comparison

*Pointwise predictive comparisons & LOO-CV*

# **Model comparison**

pointwise predictive comparisons & LOO-CV

# Model comparison

pointwise predictive comparisons & LOO-CV

- Visual PPCs can also identify unusual/influential (outliers, high leverage) data points

# Model comparison

## pointwise predictive comparisons & LOO-CV

- Visual PPCs can also identify unusual/influential (outliers, high leverage) data points
- We like using cross-validated leave-one-out predictive distributions

# Model comparison

## pointwise predictive comparisons & LOO-CV

- Visual PPCs can also identify unusual/influential (outliers, high leverage) data points
- We like using cross-validated leave-one-out predictive distributions

$$p(y_i | y_{-i})$$

# Model comparison

## pointwise predictive comparisons & LOO-CV

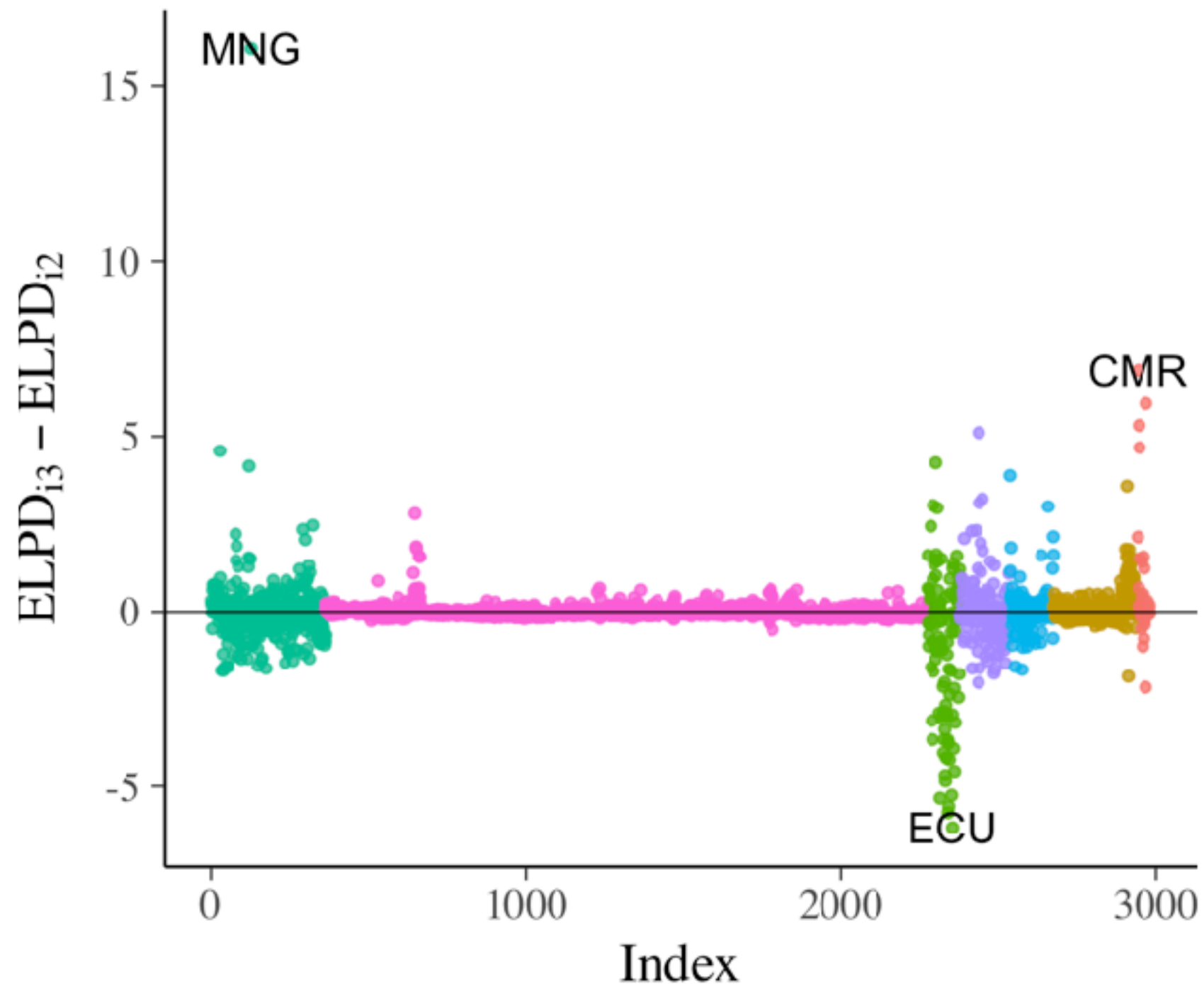
- Visual PPCs can also identify unusual/influential (outliers, high leverage) data points
- We like using cross-validated leave-one-out predictive distributions

$$p(y_i | y_{-i})$$

- Which model best predicts each of the data points that is left out?

# Model comparison

pointwise predictive comparisons & LOO-CV





# Model comparison

## Efficient approximate LOO-CV

---

Vehtari, A., Gelman, A., and Gabry, J. (2017).  
**Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC.**  
*Statistics and Computing*. 27(5), 1413–1432.  
doi: [10.1007/s11222-016-9696-4](https://doi.org/10.1007/s11222-016-9696-4)

Vehtari, A., Gelman, A., and Gabry, J. (2017).  
**Pareto smoothed importance sampling.**  
working paper  
arXiv: [arxiv.org/abs/1507.02646/](https://arxiv.org/abs/1507.02646/)

# Model comparison

## Efficient approximate LOO-CV

- How do we compute LOO-CV without fitting the model  $N$  times?

---

Vehtari, A., Gelman, A., and Gabry, J. (2017).  
**Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC.**  
*Statistics and Computing*. 27(5), 1413–1432.  
doi: [10.1007/s11222-016-9696-4](https://doi.org/10.1007/s11222-016-9696-4)

Vehtari, A., Gelman, A., and Gabry, J. (2017).  
**Pareto smoothed importance sampling.**  
working paper  
arXiv: [arxiv.org/abs/1507.02646/](https://arxiv.org/abs/1507.02646/)

# Model comparison

## Efficient approximate LOO-CV

- How do we compute LOO-CV without fitting the model  $N$  times?
- Fit once, then use Pareto smoothed importance sampling (PSIS-LOO)

# Model comparison

## Efficient approximate LOO-CV

- How do we compute LOO-CV without fitting the model  $N$  times?
- Fit once, then use Pareto smoothed importance sampling (PSIS-LOO)
- Has finite variance property of truncated IS

# Model comparison

## Efficient approximate LOO-CV

- How do we compute LOO-CV without fitting the model  $N$  times?
- Fit once, then use Pareto smoothed importance sampling (PSIS-LOO)
- Has finite variance property of truncated IS
- And less bias (replace largest weights with order stats of generalized Pareto)

# Model comparison

## Efficient approximate LOO-CV

- How do we compute LOO-CV without fitting the model  $N$  times?
- Fit once, then use Pareto smoothed importance sampling (PSIS-LOO)
- Has finite variance property of truncated IS
- And less bias (replace largest weights with order stats of generalized Pareto)
- Assumes posterior not highly sensitive to leaving out single observations

# Model comparison

## Efficient approximate LOO-CV

- How do we compute LOO-CV without fitting the model  $N$  times?
- Fit once, then use Pareto smoothed importance sampling (PSIS-LOO)
- Has finite variance property of truncated IS
- And less bias (replace largest weights with order stats of generalized Pareto)
- Assumes posterior not highly sensitive to leaving out single observations
- Asymptotically equivalent to WAIC

# Model comparison

## Efficient approximate LOO-CV

- How do we compute LOO-CV without fitting the model  $N$  times?
- Fit once, then use Pareto smoothed importance sampling (PSIS-LOO)
- Has finite variance property of truncated IS
- And less bias (replace largest weights with order stats of generalized Pareto)
- Assumes posterior not highly sensitive to leaving out single observations
- Asymptotically equivalent to WAIC
- Advantage: PSIS-LOO CV more robust + has diagnostics (check assumptions)



# Diagnostics

Pareto shape parameter & influential observations

