

## Chapter 6

- 6.1 The place of model checking in applied Bayesian statistics
- 6.2 Do the inferences from the model make sense?
- 6.3 Posterior predictive checking
- 6.4 Graphical posterior predictive checks
  - this can be skimmed, see instead the paper  
Gabry et al. (2019). *Visualization in Bayesian workflow*
- 6.5 Model checking for the educational testing example

## Model checking

- demo6\_1: Posterior predictive checking - light speed
- demo6\_2: Posterior predictive checking - sequential dependence
- demo6\_3: Posterior predictive checking - poor test statistic
- demo6\_4: Posterior predictive checking - marginal predictive p-value

## Model checking – overview

- Sensibility with respect to additional information not used in modeling
  - e.g., if posterior would claim that hazardous chemical decreases probability of death

## Model checking – overview

- Sensibility with respect to additional information not used in modeling
  - e.g., if posterior would claim that hazardous chemical decreases probability of death
- External validation
  - compare predictions to completely new observations
  - cf. relativity theory predictions

## Model checking – overview

- Sensibility with respect to additional information not used in modeling
  - e.g., if posterior would claim that hazardous chemical decreases probability of death
- External validation
  - compare predictions to completely new observations
  - cf. relativity theory predictions
- Internal validation
  - posterior predictive checking
  - cross-validation predictive checking

## Posterior predictive checking – example

- Newcomb's speed of light measurements
  - model  $y \sim \text{normal}(\mu, \sigma)$  with prior  $(\mu, \log \sigma) \propto 1$

## Posterior predictive checking – example

- Newcomb's speed of light measurements
  - model  $y \sim \text{normal}(\mu, \sigma)$  with prior  $(\mu, \log \sigma) \propto 1$
- Posterior predictive replicate  $y^{\text{rep}}$

## Posterior predictive checking – example

- Newcomb's speed of light measurements
  - model  $y \sim \text{normal}(\mu, \sigma)$  with prior  $(\mu, \log \sigma) \propto 1$
- Posterior predictive replicate  $y^{\text{rep}}$ 
  - draw  $\mu^{(s)}, \sigma^{(s)}$  from the posterior  $p(\mu, \sigma|y)$

## Posterior predictive checking – example

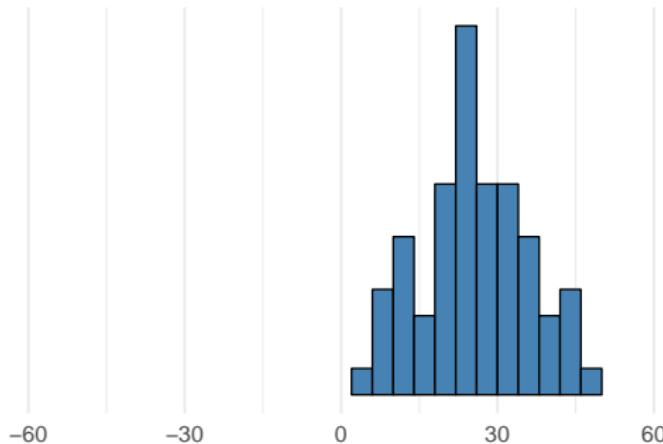
- Newcomb's speed of light measurements
  - model  $y \sim \text{normal}(\mu, \sigma)$  with prior  $(\mu, \log \sigma) \propto 1$
- Posterior predictive replicate  $y^{\text{rep}}$ 
  - draw  $\mu^{(s)}, \sigma^{(s)}$  from the posterior  $p(\mu, \sigma|y)$
  - draw  $y^{\text{rep}(s)}$  from  $\text{normal}(\mu^{(s)}, \sigma^{(s)})$

## Posterior predictive checking – example

- Newcomb's speed of light measurements
  - model  $y \sim \text{normal}(\mu, \sigma)$  with prior  $(\mu, \log \sigma) \propto 1$
- Posterior predictive replicate  $y^{\text{rep}}$ 
  - draw  $\mu^{(s)}, \sigma^{(s)}$  from the posterior  $p(\mu, \sigma|y)$
  - draw  $y^{\text{rep}(s)}$  from  $\text{normal}(\mu^{(s)}, \sigma^{(s)})$
  - repeat  $n$  times to get  $y^{\text{rep}}$  with  $n$  replicates

# Posterior predictive checking – example

- Newcomb's speed of light measurements
  - model  $y \sim \text{normal}(\mu, \sigma)$  with prior  $(\mu, \log \sigma) \propto 1$
- Posterior predictive replicate  $y^{\text{rep}}$ 
  - draw  $\mu^{(s)}, \sigma^{(s)}$  from the posterior  $p(\mu, \sigma|y)$
  - draw  $y^{\text{rep}(s)}$  from  $\text{normal}(\mu^{(s)}, \sigma^{(s)})$
  - repeat  $n$  times to get  $y^{\text{rep}}$  with  $n$  replicates



## Replicates vs. future observation

- Predictive  $\tilde{y}$  is the next not yet observed possible observation.  
 $y^{\text{rep}}$  refers to replicating the whole experiment (potentially with same values of  $x$ ) and obtaining as many replicated observations as in the original data.

## Posterior predictive checking – example

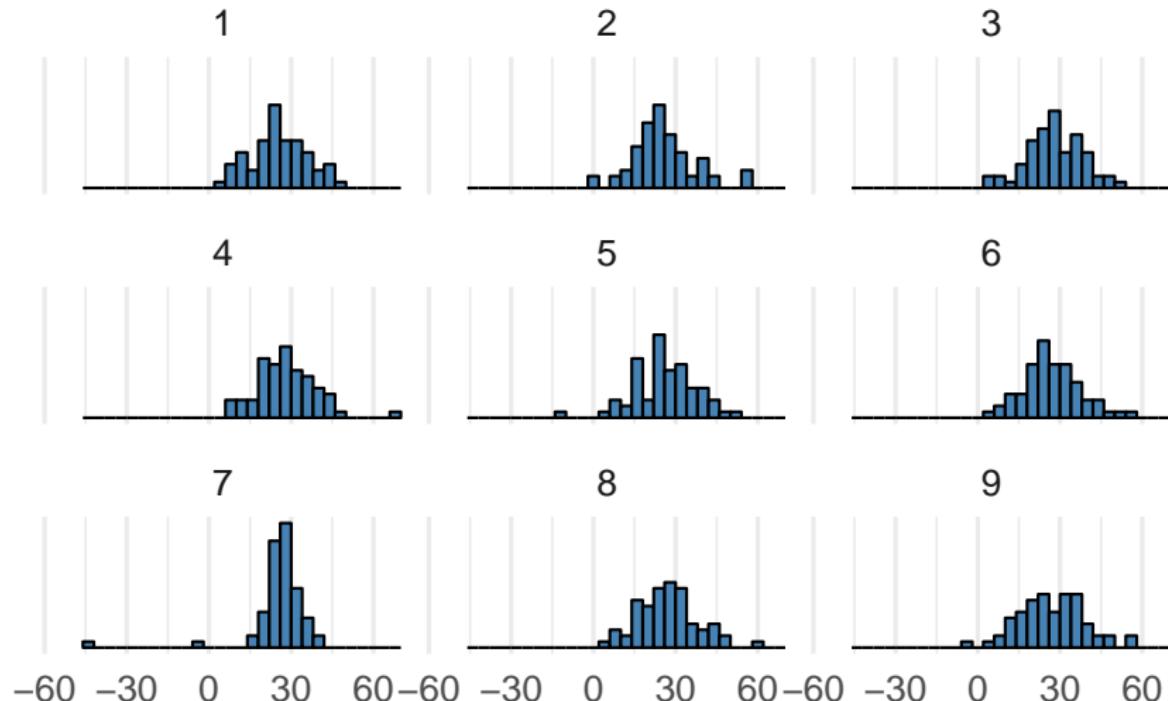
- Generate several replicated datasets  $y^{\text{rep}}$

## Posterior predictive checking – example

- Generate several replicated datasets  $y^{\text{rep}}$
- Compare to the original dataset

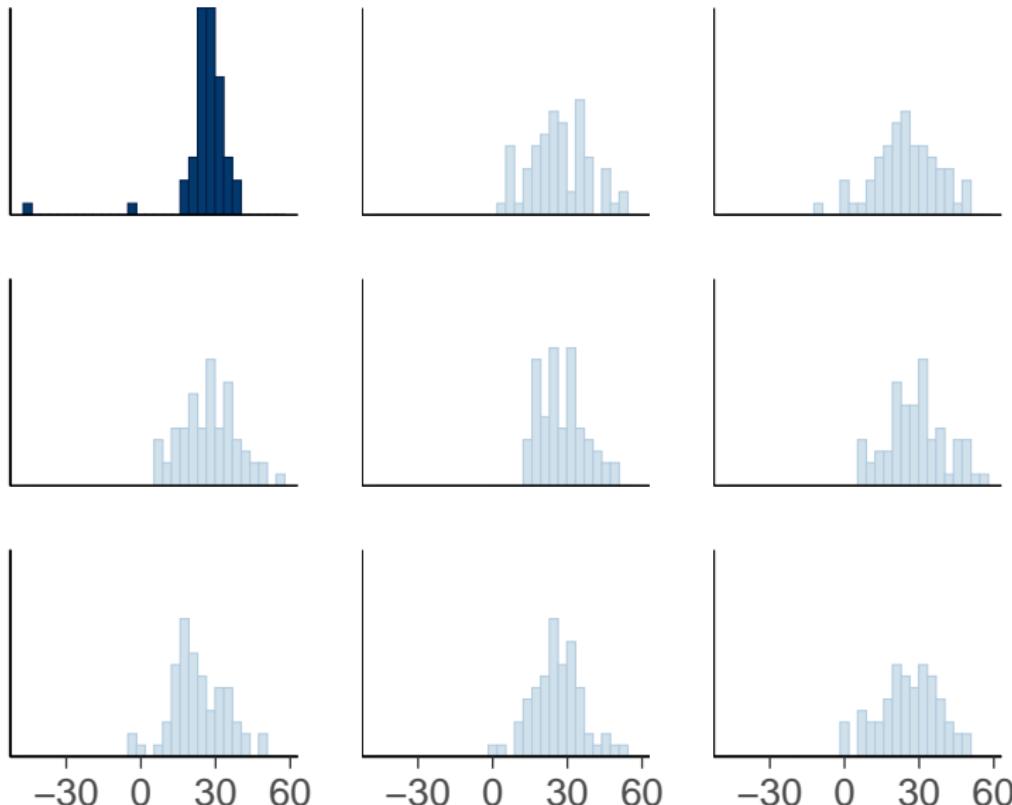
## Posterior predictive checking – example

- Generate several replicated datasets  $y^{\text{rep}}$
- Compare to the original dataset



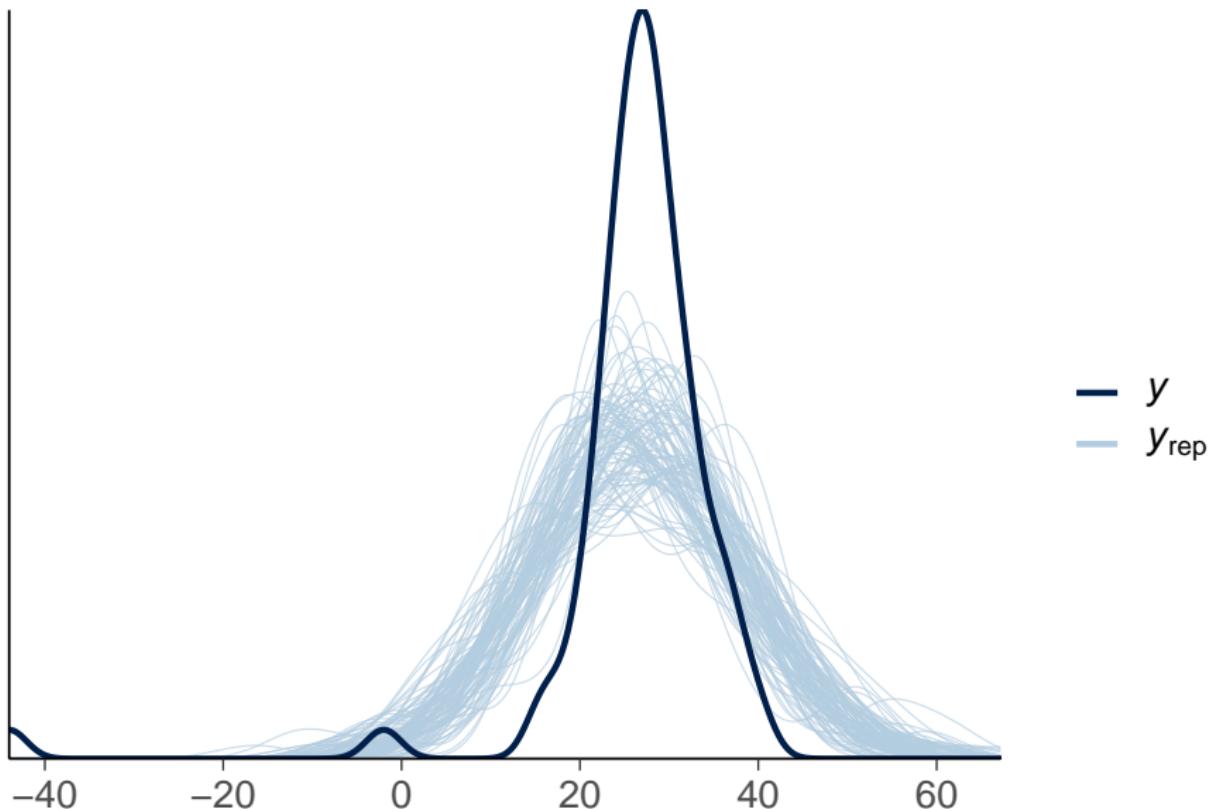
## Posterior predictive checking – bayesplot

```
ppc_hist(y, yrep[1:8,])
```



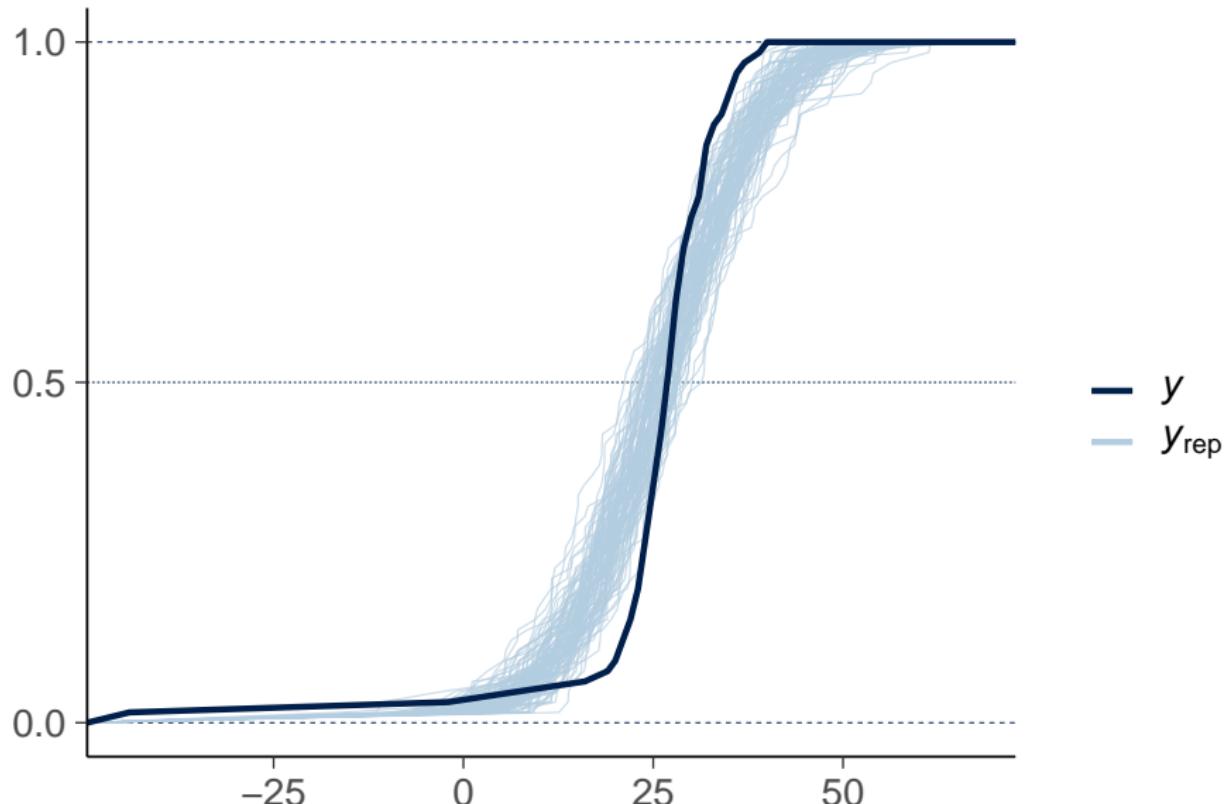
## Posterior predictive checking – bayesplot

```
ppc_dens_overlay(y, yrep[1:100, ])
```



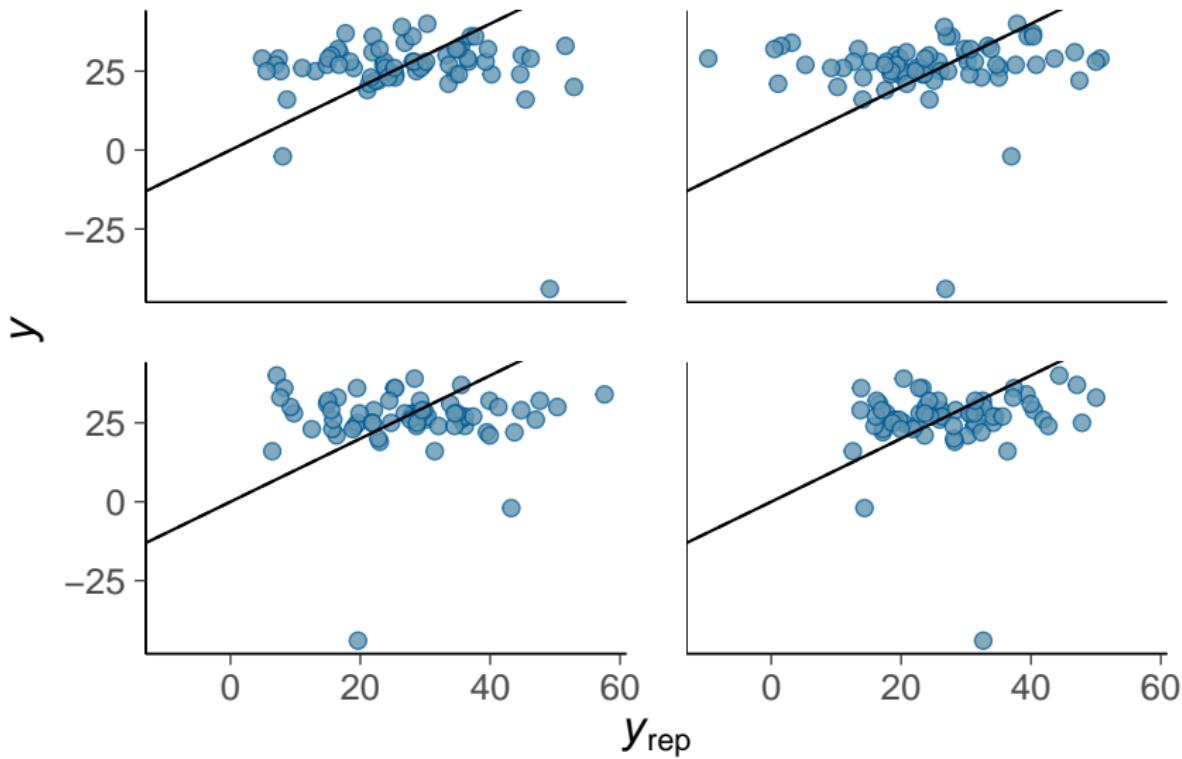
## Posterior predictive checking – bayesplot

```
ppc_ecdf_overlay(y, yrep[1:100, ])
```



## Posterior predictive checking – bayesplot

```
ppc_scatter(y, yrep[1:4,]) + geom_abline()
```



## Posterior predictive checking with test statistic

- Replicated data sets  $y^{\text{rep}}$
- Test quantity (or discrepancy measure)  $T(y, \theta)$ 
  - summary quantity for the observed data  $T(y, \theta)$
  - summary quantity for a replicated data  $T(y^{\text{rep}}, \theta)$
  - can be easier to compare summary quantities than data sets

## Posterior predictive checking – example

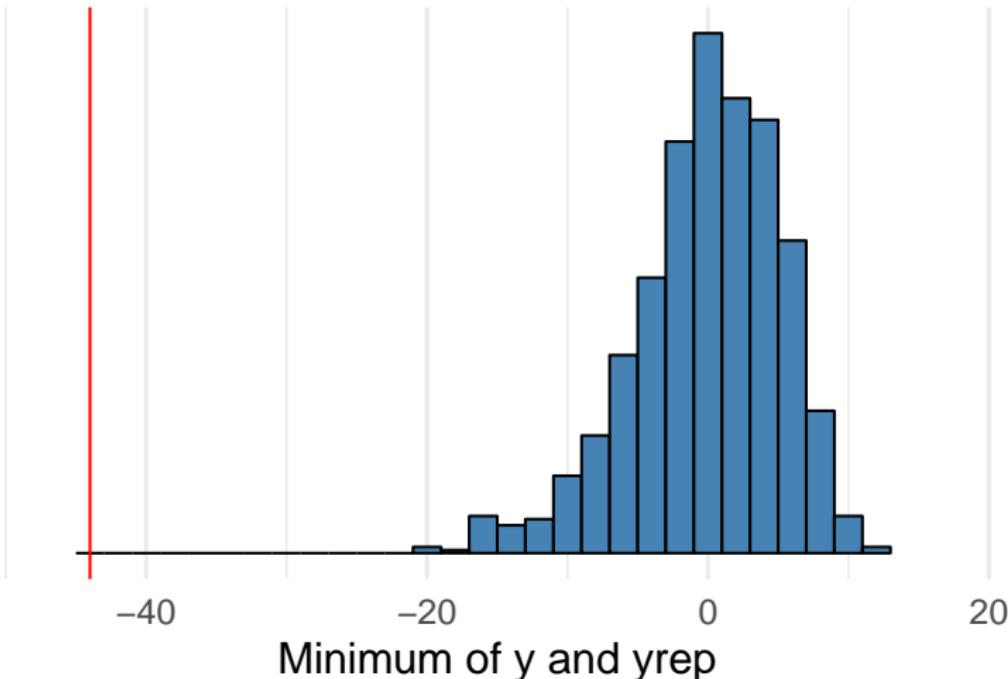
- Compute test statistic for data  $T(y, \theta) = \min(y)$

## Posterior predictive checking – example

- Compute test statistic for data  $T(y, \theta) = \min(y)$
- Compute test statistic  $\min(y^{\text{rep}})$  for many replicated datasets

## Posterior predictive checking – example

- Compute test statistic for data  $T(y, \theta) = \min(y)$
- Compute test statistic  $\min(y^{\text{rep}})$  for many replicated datasets



## Posterior predictive checking – example

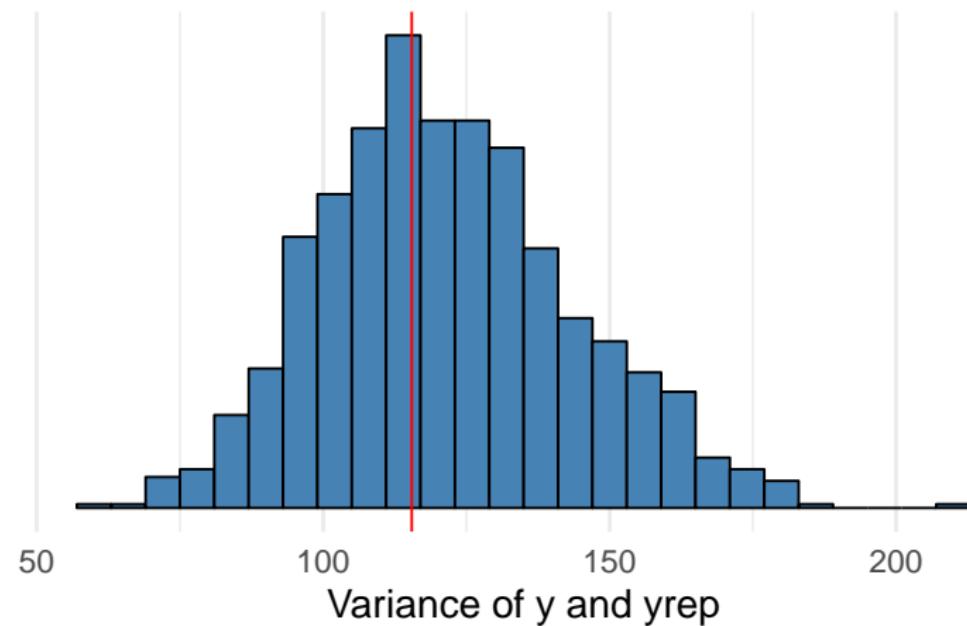
- Good test statistic is ancillary (or almost)
  - ancillary if it depends only on observed data and if its distribution is independent of the parameters of the model

## Posterior predictive checking – example

- Good test statistic is ancillary (or almost)
  - ancillary if it depends only on observed data and if its distribution is independent of the parameters of the model
- Bad test statistic is highly dependent of the parameters
  - e.g. variance for normal model

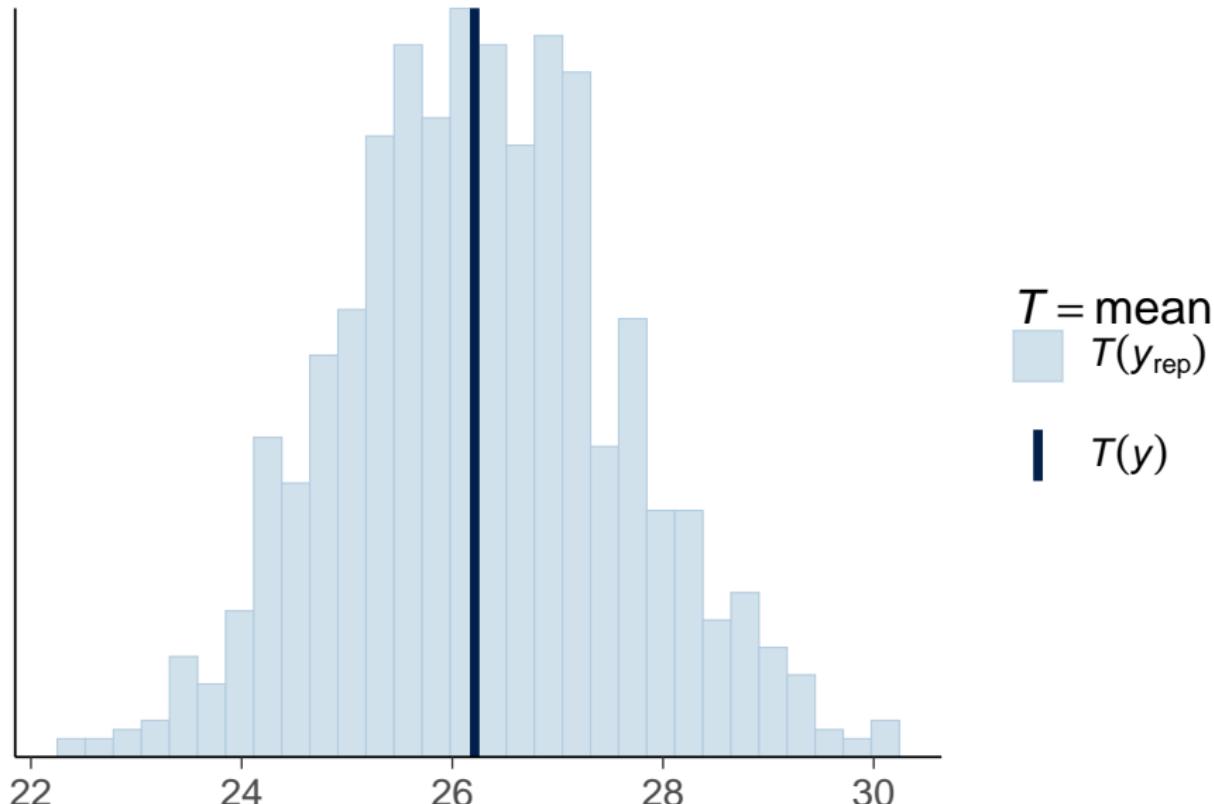
## Posterior predictive checking – example

- Good test statistic is ancillary (or almost)
  - ancillary if it depends only on observed data and if its distribution is independent of the parameters of the model
- Bad test statistic is highly dependent of the parameters
  - e.g. variance for normal model



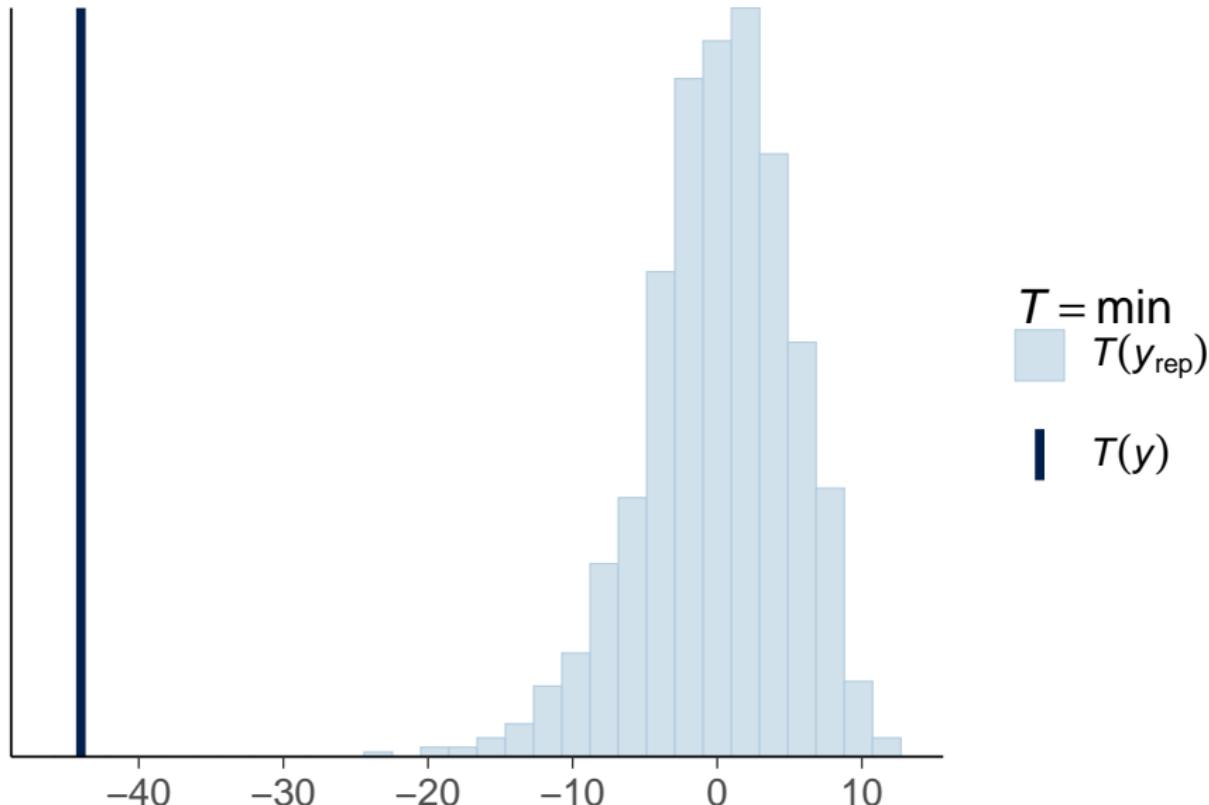
## Posterior predictive checking – bayesplot

`ppc_stat(y, yrep)`, the default statistic "mean" is usually bad



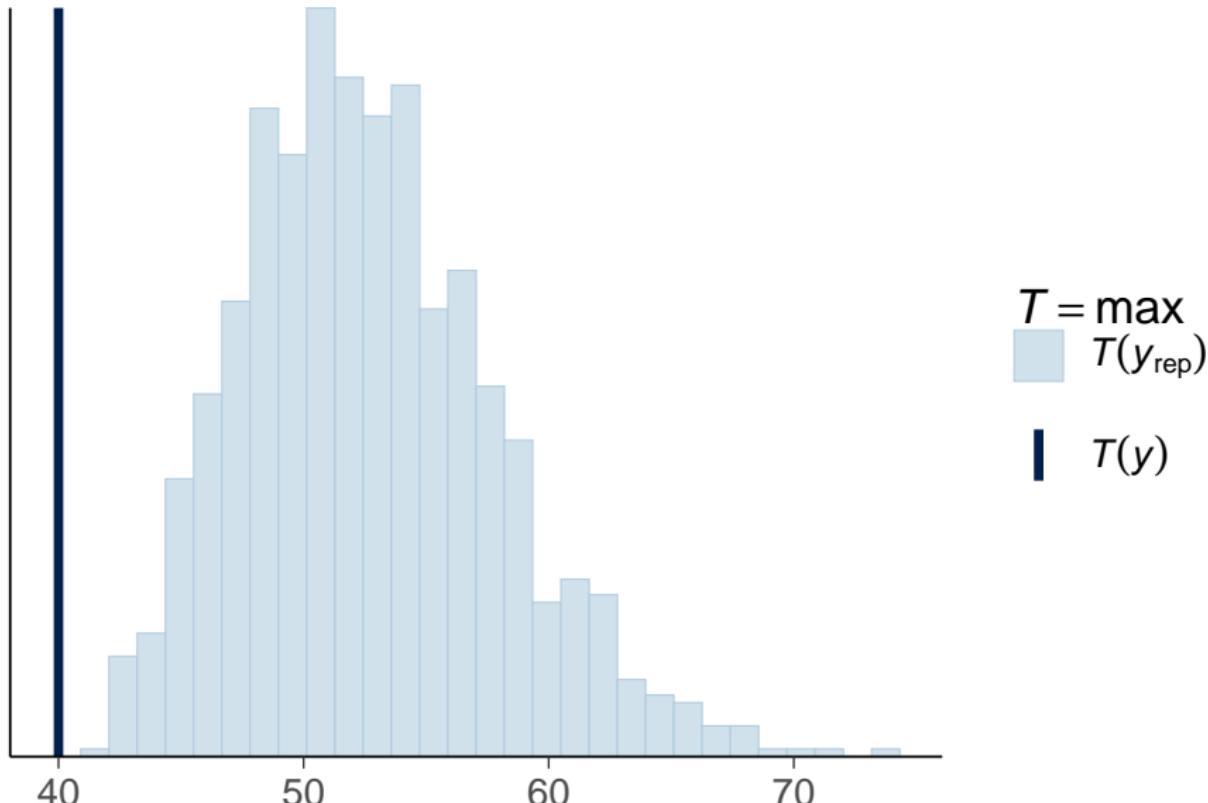
## Posterior predictive checking – bayesplot

```
ppc_stat(y, yrep, stat="min")
```



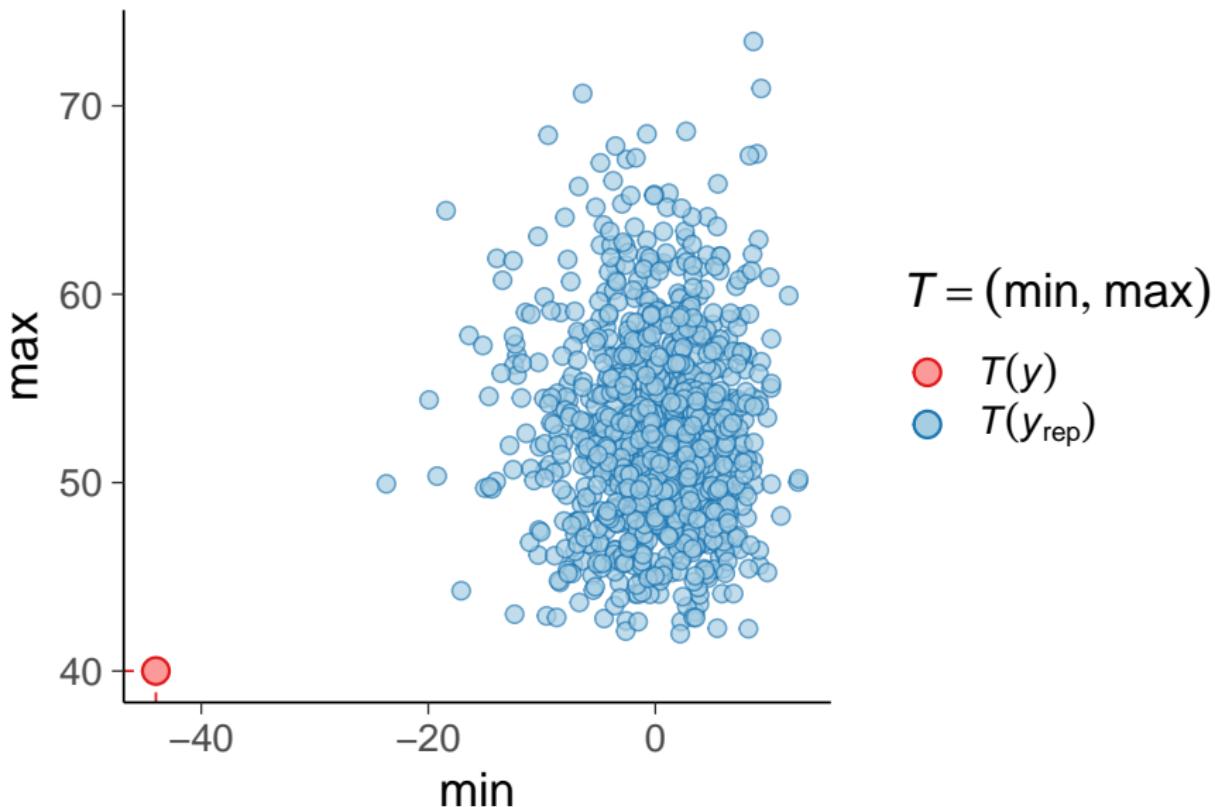
## Posterior predictive checking – bayesplot

```
ppc_stat(y, yrep, stat="max")
```



## Posterior predictive checking – bayesplot

```
ppc_stat2d(y, yrep, stat=c("min", "max"))
```



## Posterior predictive checking

- *Posterior predictive p-value*

$$\begin{aligned} p &= \Pr(T(y^{\text{rep}}, \theta) \geq T(y, \theta) | y) \\ &= \int \int I_{T(y^{\text{rep}}, \theta) \geq T(y, \theta)} p(y^{\text{rep}} | \theta) p(\theta | y) dy^{\text{rep}} d\theta \end{aligned}$$

where  $I$  is an indicator function

# Posterior predictive checking

- *Posterior predictive p-value*

$$\begin{aligned} p &= \Pr(T(y^{\text{rep}}, \theta) \geq T(y, \theta) | y) \\ &= \int \int I_{T(y^{\text{rep}}, \theta) \geq T(y, \theta)} p(y^{\text{rep}} | \theta) p(\theta | y) dy^{\text{rep}} d\theta \end{aligned}$$

where  $I$  is an indicator function

- having  $(y^{\text{rep}(s)}, \theta^{(s)})$  from the posterior predictive distribution, easy to compute

$$T(y^{\text{rep}(s)}, \theta^{(s)}) \geq T(y, \theta^{(s)}), \quad s = 1, \dots, S$$

# Posterior predictive checking

- *Posterior predictive p-value*

$$\begin{aligned} p &= \Pr(T(y^{\text{rep}}, \theta) \geq T(y, \theta) | y) \\ &= \int \int I_{T(y^{\text{rep}}, \theta) \geq T(y, \theta)} p(y^{\text{rep}} | \theta) p(\theta | y) dy^{\text{rep}} d\theta \end{aligned}$$

where  $I$  is an indicator function

- having  $(y^{\text{rep}(s)}, \theta^{(s)})$  from the posterior predictive distribution, easy to compute

$$T(y^{\text{rep}(s)}, \theta^{(s)}) \geq T(y, \theta^{(s)}), \quad s = 1, \dots, S$$

- Posterior predictive p-value (ppp-value) estimates whether difference between the model and data could arise by chance

# Posterior predictive checking

- Posterior predictive *p*-value

$$\begin{aligned} p &= \Pr(T(y^{\text{rep}}, \theta) \geq T(y, \theta) | y) \\ &= \int \int I_{T(y^{\text{rep}}, \theta) \geq T(y, \theta)} p(y^{\text{rep}} | \theta) p(\theta | y) dy^{\text{rep}} d\theta \end{aligned}$$

where  $I$  is an indicator function

- having  $(y^{\text{rep}(s)}, \theta^{(s)})$  from the posterior predictive distribution, easy to compute

$$T(y^{\text{rep}(s)}, \theta^{(s)}) \geq T(y, \theta^{(s)}), \quad s = 1, \dots, S$$

- Posterior predictive *p*-value (ppp-value) estimates whether difference between the model and data could arise by chance
- Not commonly used, as
  - the distribution of test statistic has more information
  - not calibrated in case of non-ancillary statistic

## Posterior predictive checking – example

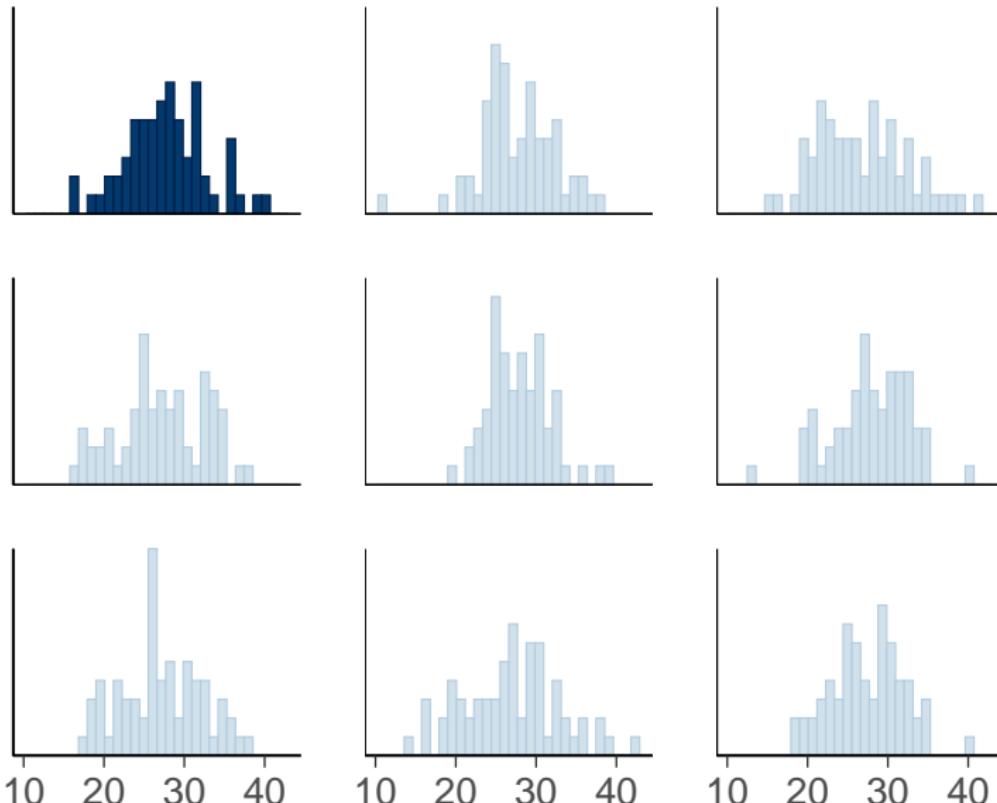
- In general it is often better to improve the model than remove data
- There can be cases where after the first model fit, it is clear that some data are clearly incorrectly measured, but there has to be really strong justifications for dropping out observations

## Posterior predictive checking – example

- In general it is often better to improve the model than remove data
- There can be cases where after the first model fit, it is clear that some data are clearly incorrectly measured, but there has to be really strong justifications for dropping out observations
- Let's assume that in Newcomb experiment, the two observations with negative values are clearly incorrect and let's refit by removing them

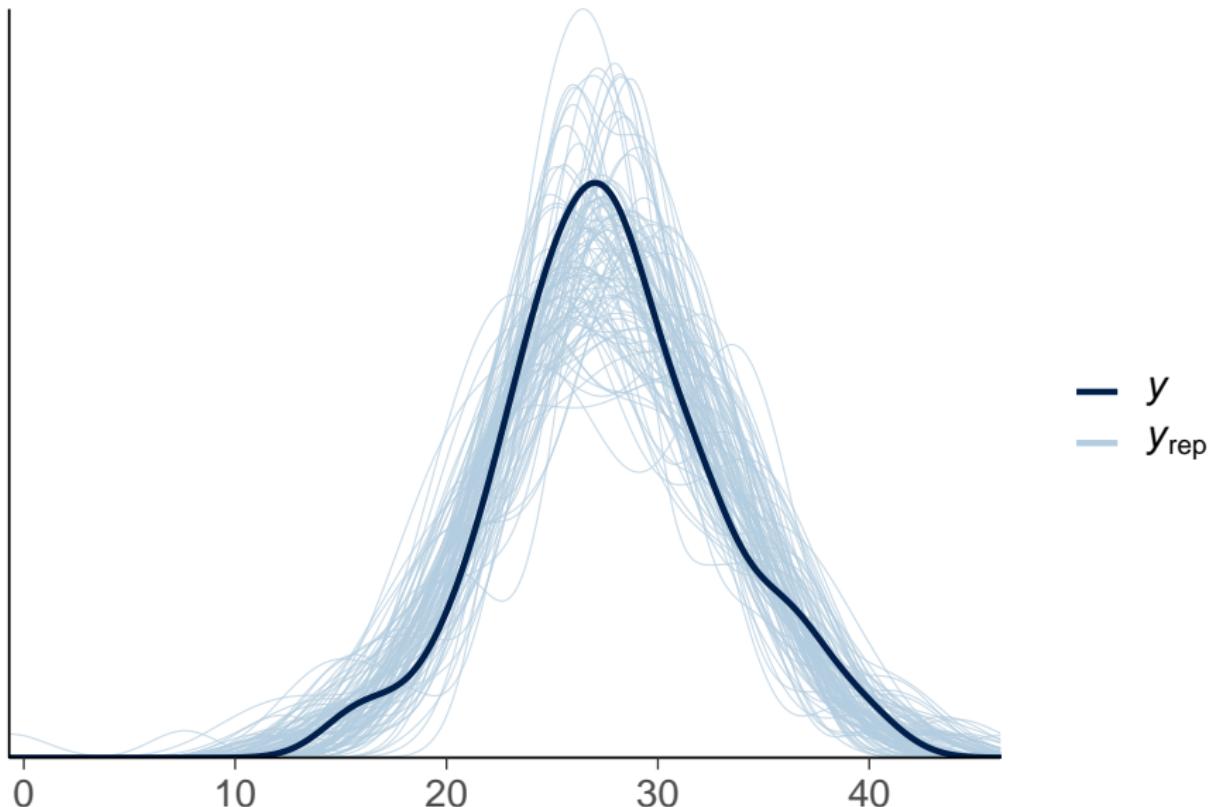
## Posterior predictive checking (good fit)

```
ppc_hist(y2, y2rep[1:8,])
```



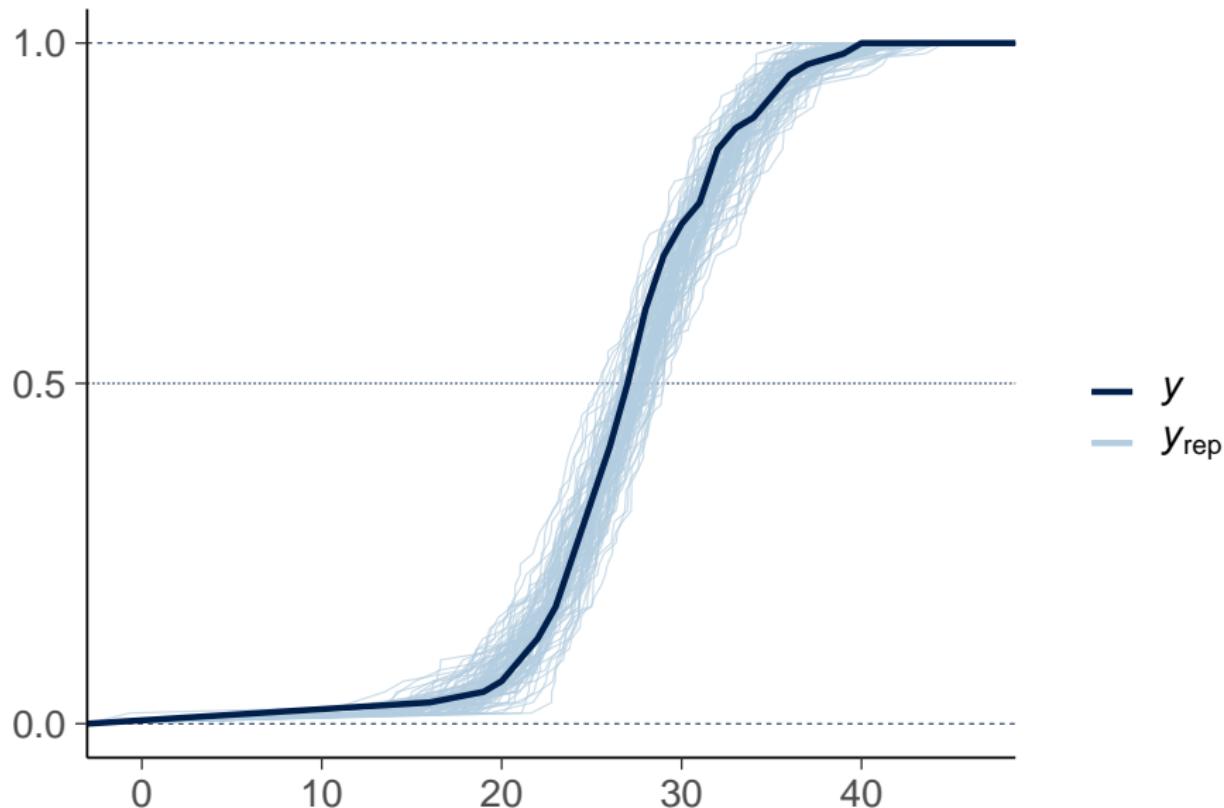
## Posterior predictive checking (good fit)

```
ppc_dens_overlay(y2, y2rep[1:100,])
```



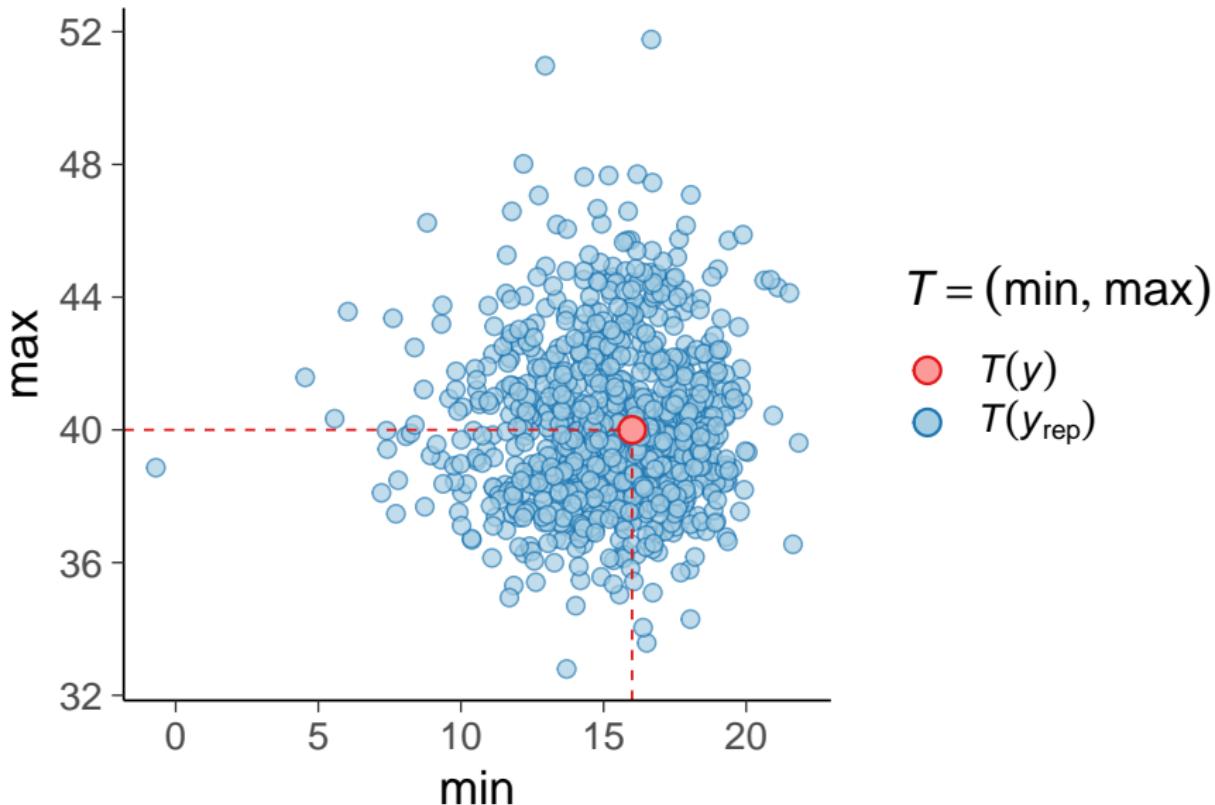
## Posterior predictive checking (good fit)

```
ppc_ecdf_overlay(y2, y2rep[1:100,])
```



## Posterior predictive checking (good fit)

```
ppc_stat2d(y2, y2rep, stat=c("min", "max"))
```



## Marginal and CV predictive checking

- Consider marginal predictive distributions  $p(\tilde{y}_i|y)$  and each observation separately
  - marginal posterior p-values

$$p_i = \Pr(T(y_i^{\text{rep}}) \leq T(y_i)|y)$$

if  $T(y_i) = y_i$

$$p_i = \Pr(y_i^{\text{rep}} \leq y_i|y)$$

## Marginal and CV predictive checking

- Consider marginal predictive distributions  $p(\tilde{y}_i|y)$  and each observation separately
  - marginal posterior p-values

$$p_i = \Pr(T(y_i^{\text{rep}}) \leq T(y_i)|y)$$

if  $T(y_i) = y_i$

$$p_i = \Pr(y_i^{\text{rep}} \leq y_i|y)$$

- if  $\Pr(\tilde{y}_i|y)$  well calibrated, distribution of  $p_i$  would be uniform between 0 and 1
  - holds better for cross-validation predictive tests (cross-validation BDA3 Ch 7)

## Marginal predictive checking (good fit)

- Marginal tail area or Probability integral transform (PIT)

$$p_i = p(y_i^{\text{rep}} \leq y_i | y)$$

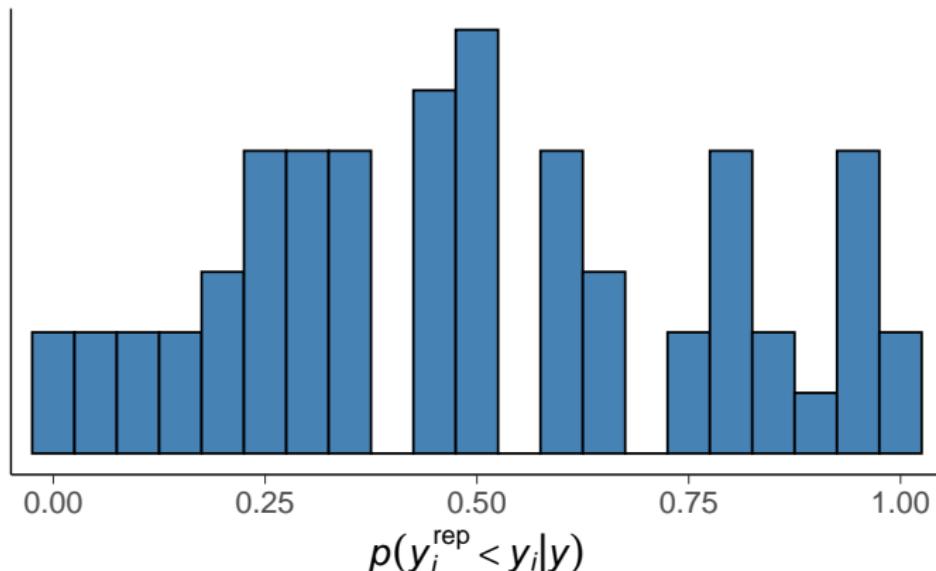
- if  $p(\tilde{y}_i | y)$  is well calibrated, distribution of  $p_i$ 's would be uniform between 0 and 1

## Marginal predictive checking (good fit)

- Marginal tail area or Probability integral transform (PIT)

$$p_i = p(y_i^{\text{rep}} \leq y_i | y)$$

- if  $p(\tilde{y}_i | y)$  is well calibrated, distribution of  $p_i$ 's would be uniform between 0 and 1

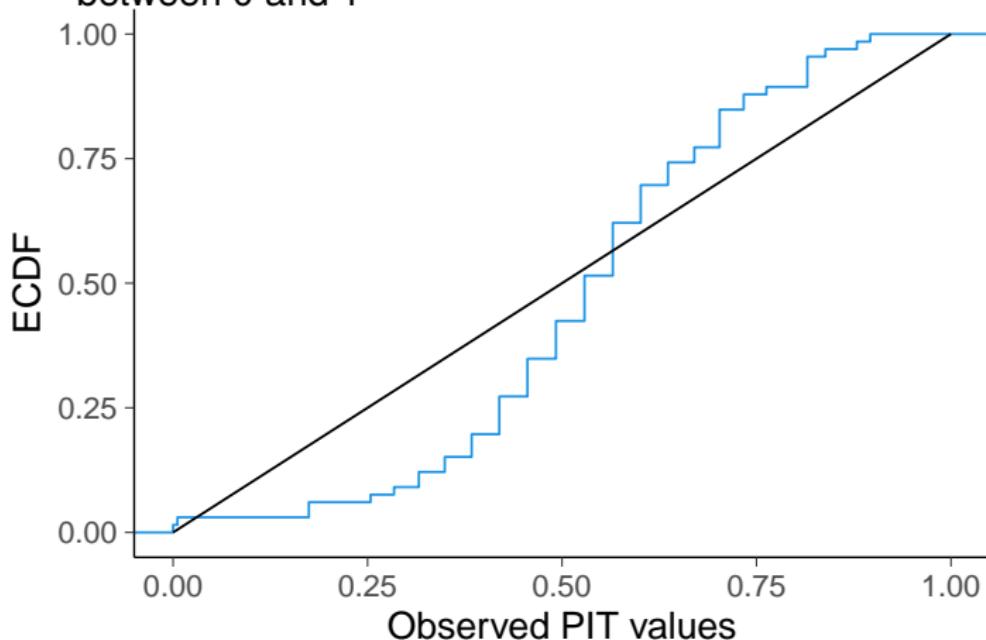


## Marginal predictive checking (good fit)

- Marginal tail area or Probability integral transform (PIT)

$$p_i = p(y_i^{\text{rep}} \leq y_i | y)$$

- if  $p(\tilde{y}_i | y)$  is well calibrated, distribution of  $p_i$ 's would be uniform between 0 and 1

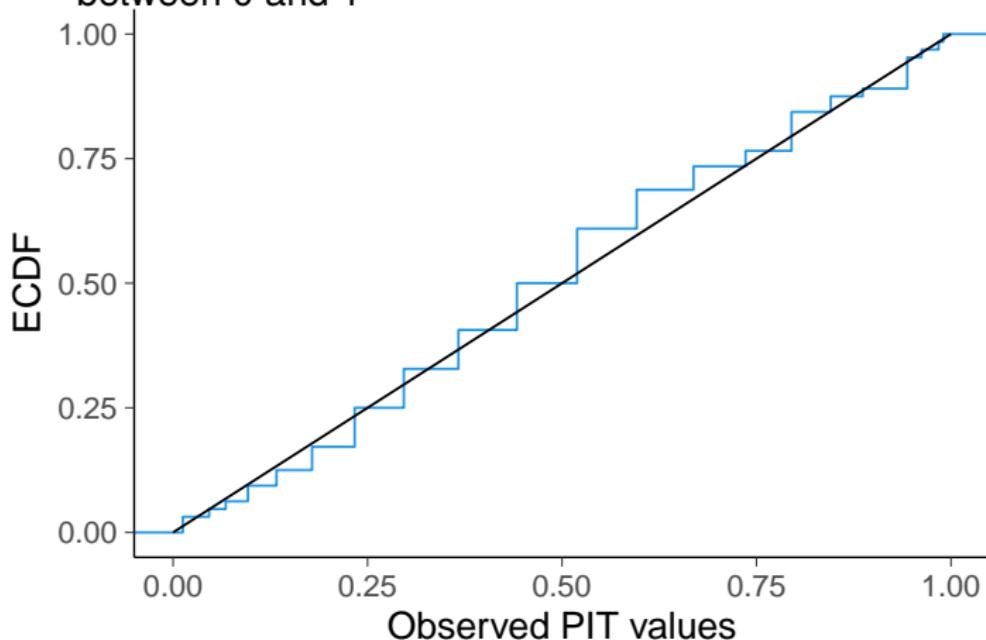


## Marginal predictive checking (good fit)

- Marginal tail area or Probability integral transform (PIT)

$$p_i = p(y_i^{\text{rep}} \leq y_i | y)$$

- if  $p(\tilde{y}_i | y)$  is well calibrated, distribution of  $p_i$ 's would be uniform between 0 and 1

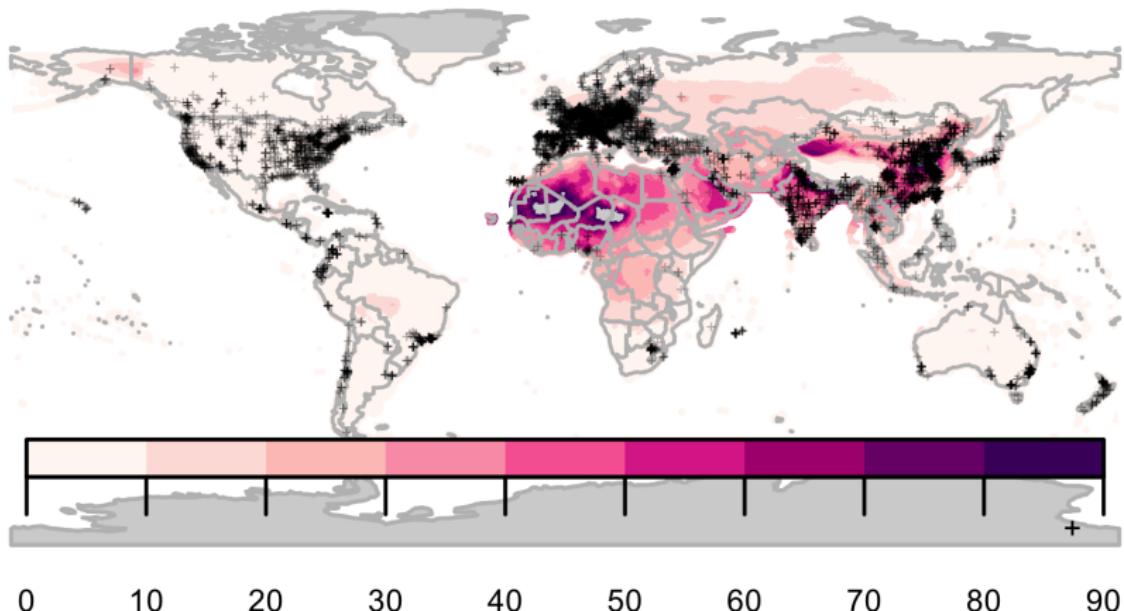


## Example: Exposure to air pollution

- Example from Jonah Gabry, Daniel Simpson, Aki Vehtari, Michael Betancourt, and Andrew Gelman (2019). Visualization in Bayesian workflow. <https://doi.org/10.1111/rssc.12378>
- Estimation of human exposure to air pollution from particulate matter measuring less than 2.5 microns in diameter (PM<sub>2.5</sub>)
  - Exposure to PM<sub>2.5</sub> is linked to a number of poor health outcomes and a recent report estimated that PM<sub>2.5</sub> is responsible for three million deaths worldwide each year (Shaddick et al., 2017)
  - In order to estimate the public health effect of ambient PM<sub>2.5</sub>, we need a good estimate of the PM<sub>2.5</sub> concentration at the same spatial resolution as our population estimates.

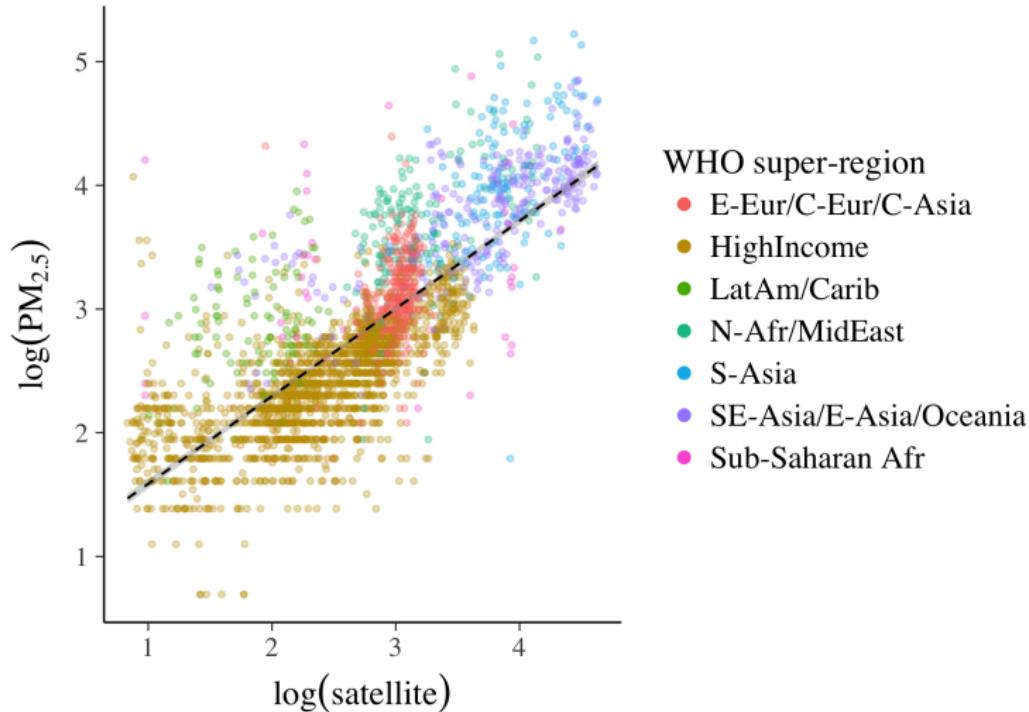
## Example: Exposure to air pollution

- Direct measurements of PM 2.5 from ground monitors at 2980 locations
- High-resolution satellite data of aerosol optical depth



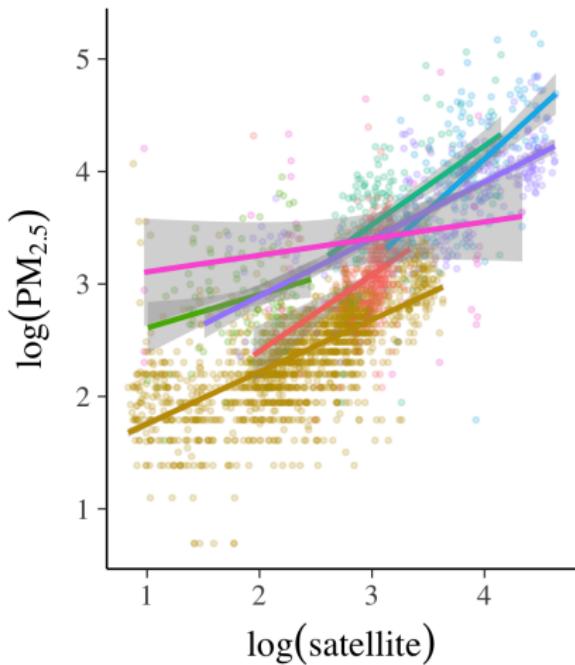
## Example: Exposure to air pollution

- Direct measurements of PM 2.5 from ground monitors at 2980 locations
- High-resolution satellite data of aerosol optical depth



## Example: Exposure to air pollution

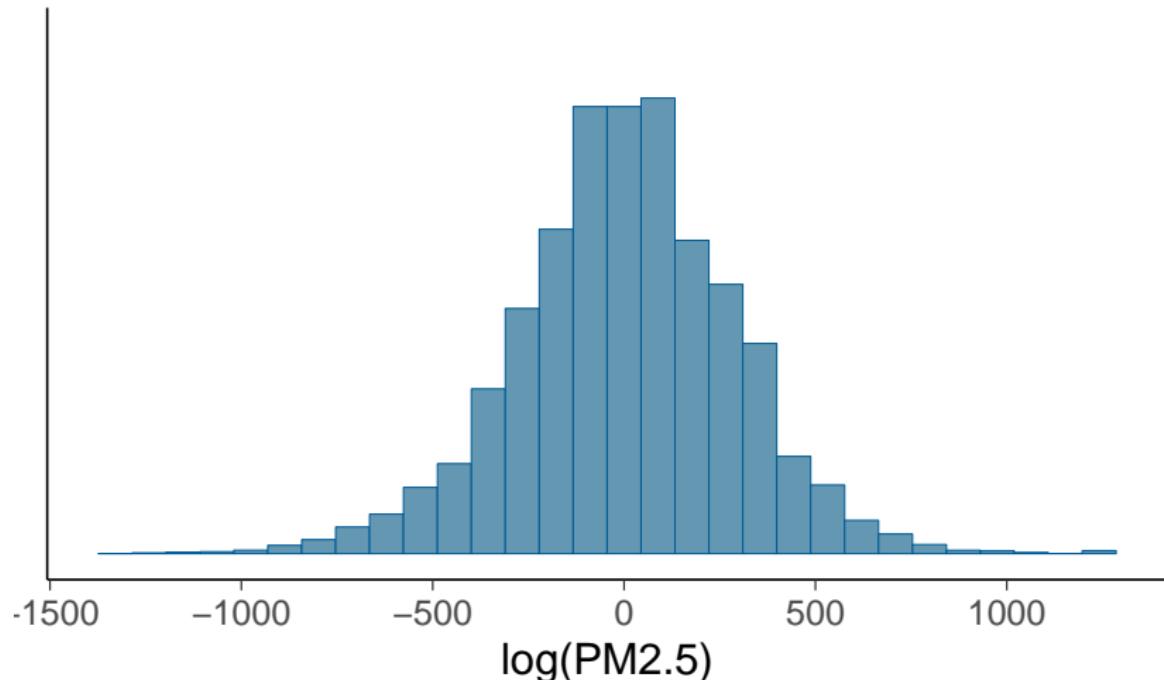
- Direct measurements of PM 2.5 from ground monitors at 2980 locations
- High-resolution satellite data of aerosol optical depth



## Example: Exposure to air pollution

Prior predictive checking

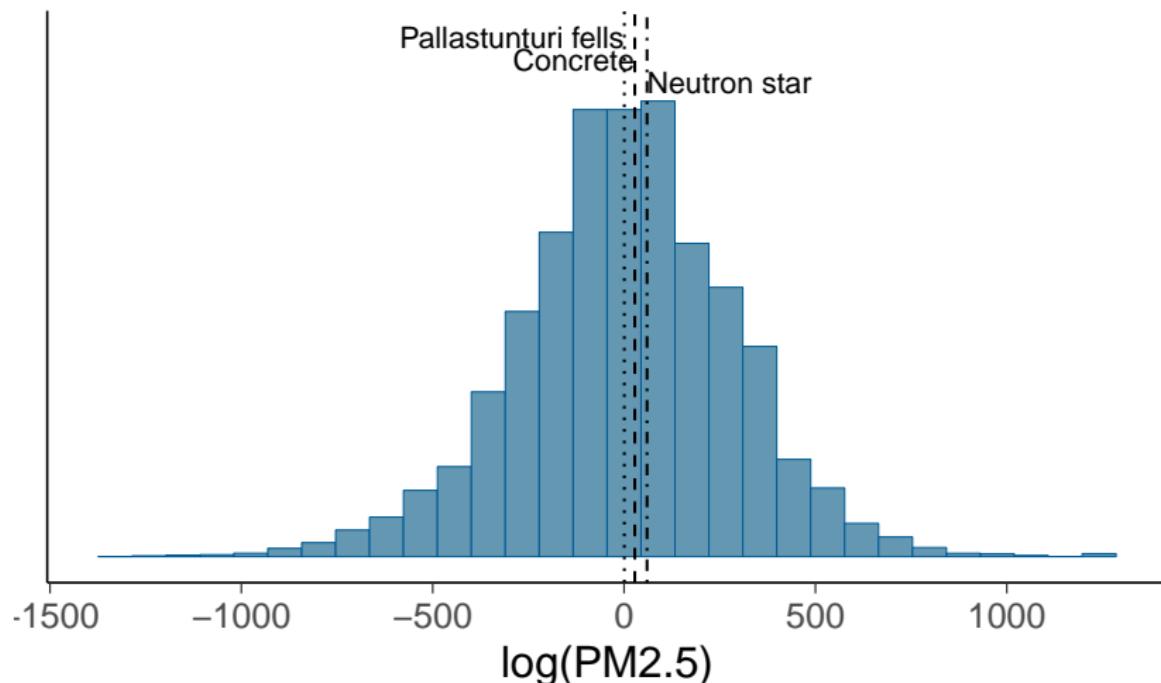
Prior predictive distribution with vague prior



## Example: Exposure to air pollution

Prior predictive checking

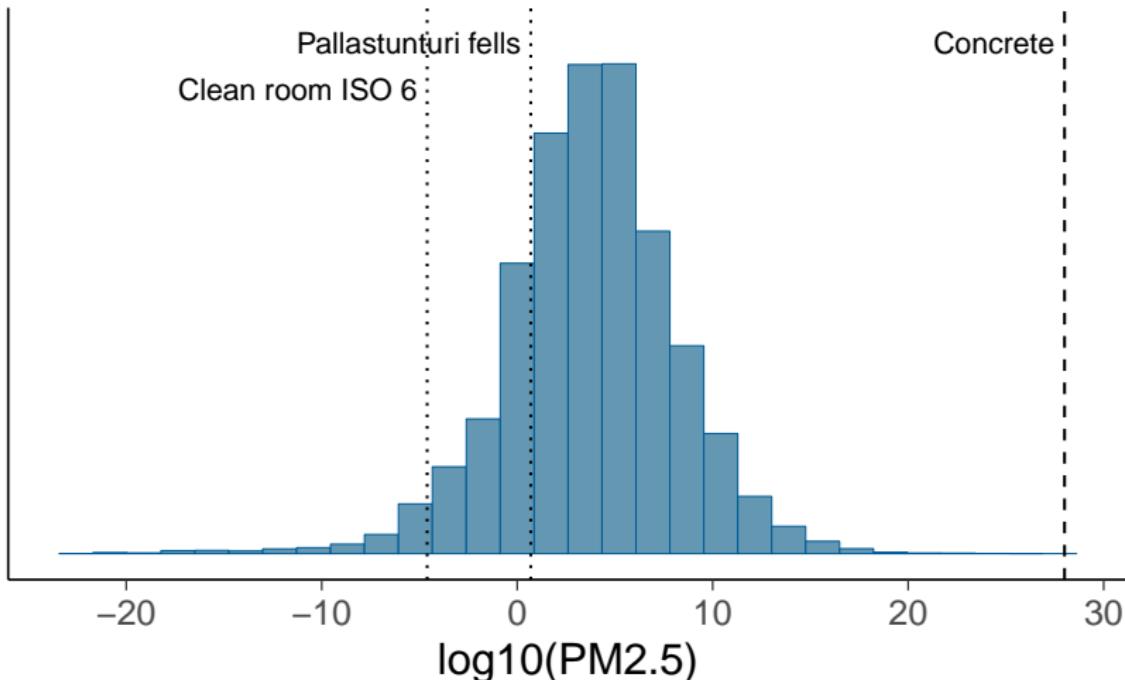
Prior predictive distribution with vague prior



## Example: Exposure to air pollution

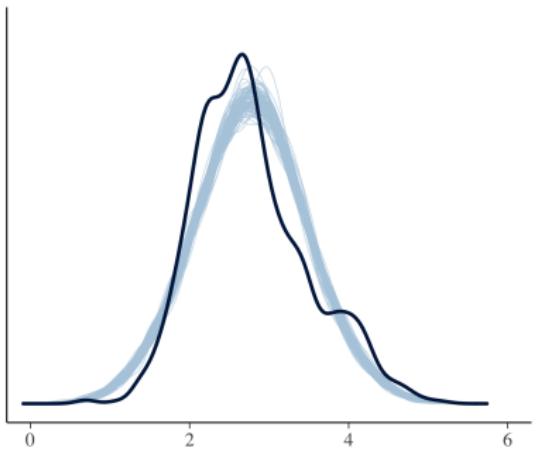
Prior predictive checking

Prior predictive distribution with weakly informative

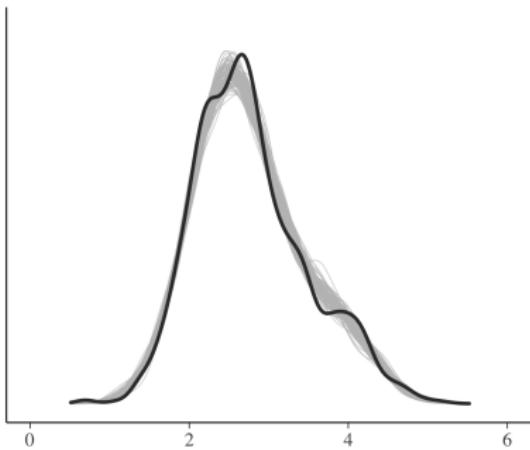


# Example: Exposure to air pollution

Posterior predictive checking – marginal predictive distributions



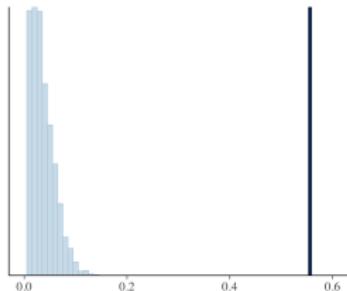
(a) Model 1



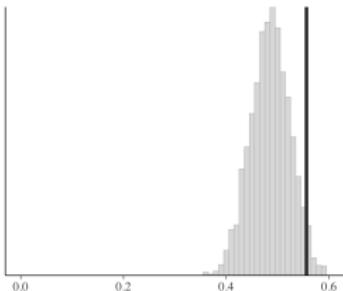
(b) Model 2

# Example: Exposure to air pollution

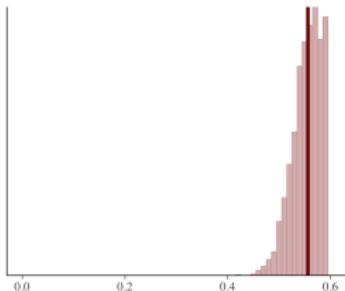
Posterior predictive checking – test statistic (skewness)



(a) Model 1



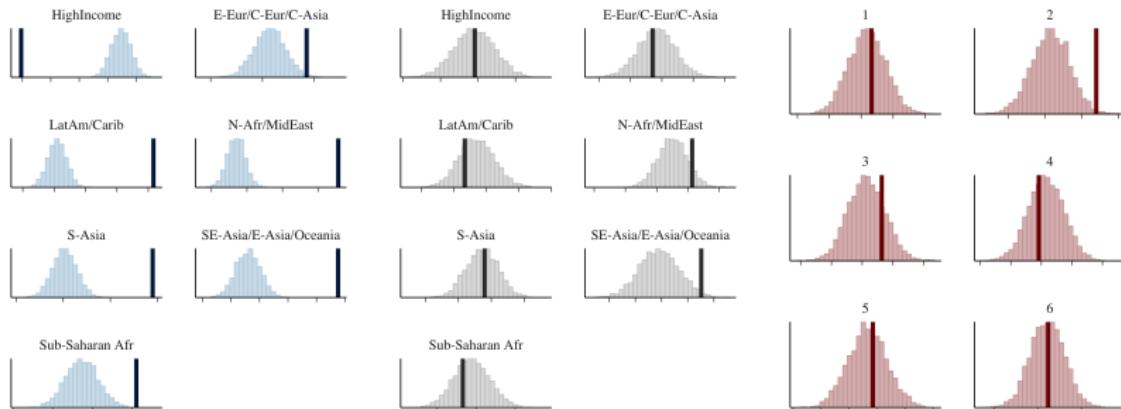
(b) Model 2



(c) Model 3

# Example: Exposure to air pollution

Posterior predictive checking – test statistic (median for groups)



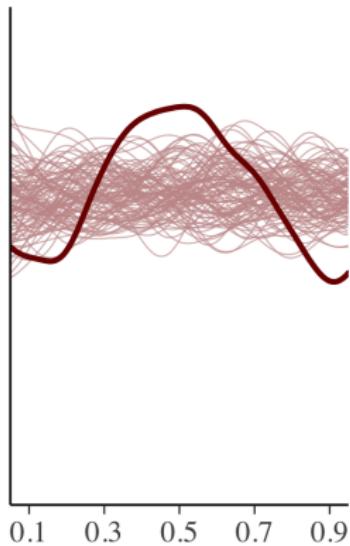
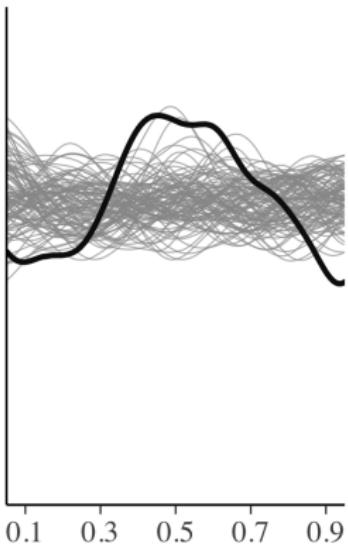
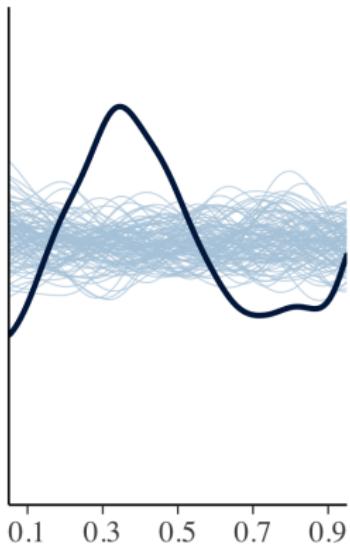
(a) Model 1

(b) Model 2

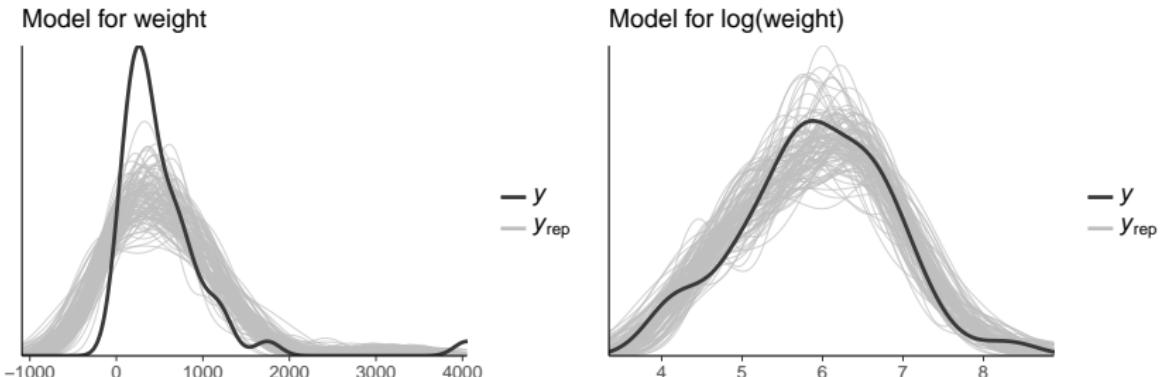
(c) Model 3

# Example: Exposure to air pollution

LOO predictive checking – LOO-PIT



# Positive target

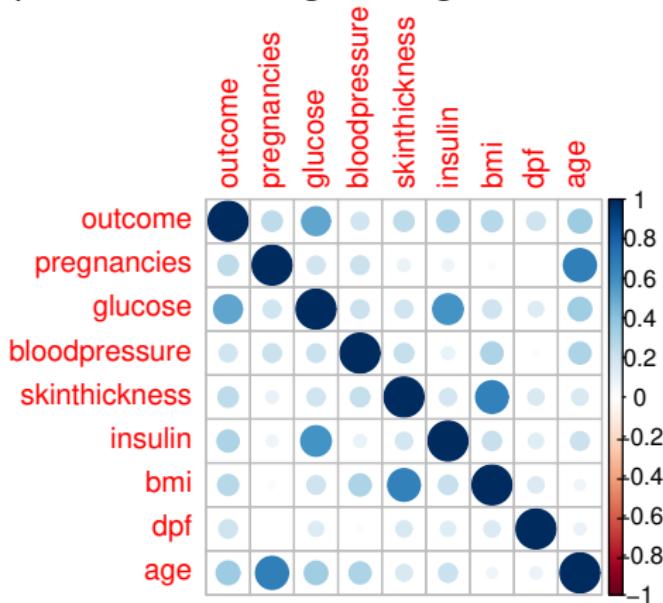


Predicting the yields of mesquite bushes.

Gelman, Hill & Vehtari (2020): Regression and Other Stories, Chapter 11.

## PPC for binary target

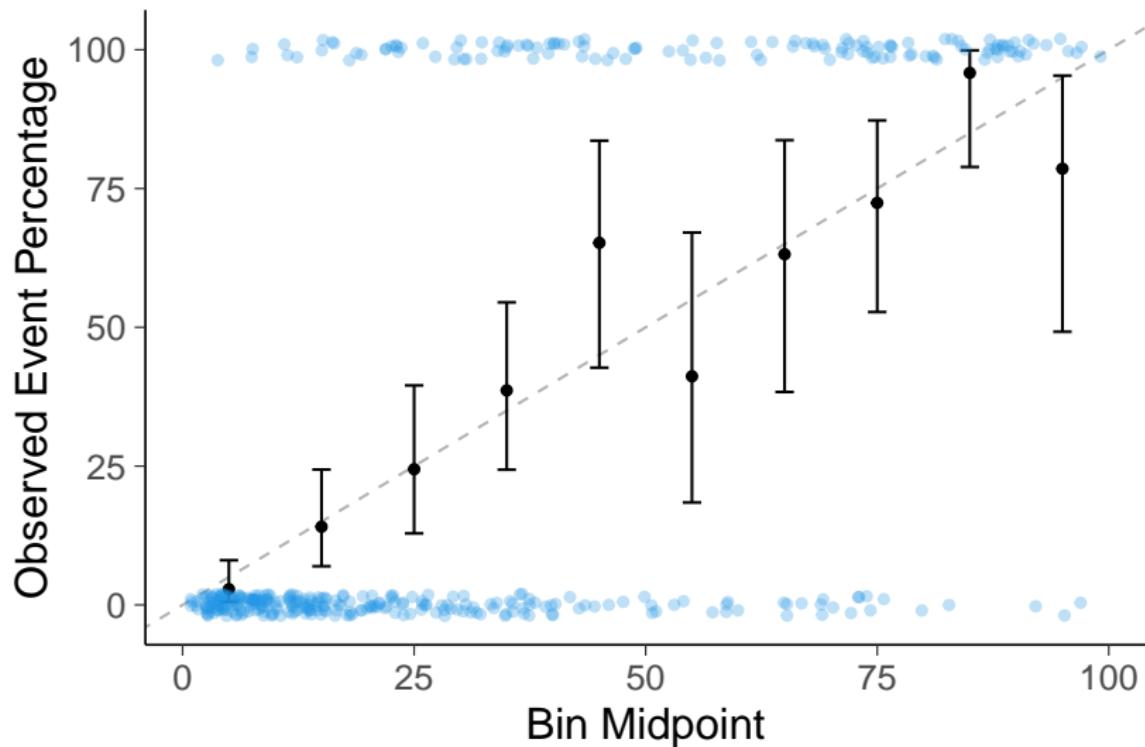
Diabetes prediction with logistic regression - diabetes demo



# PPC for binary target

Diabetes prediction with logistic regression - diabetes demo

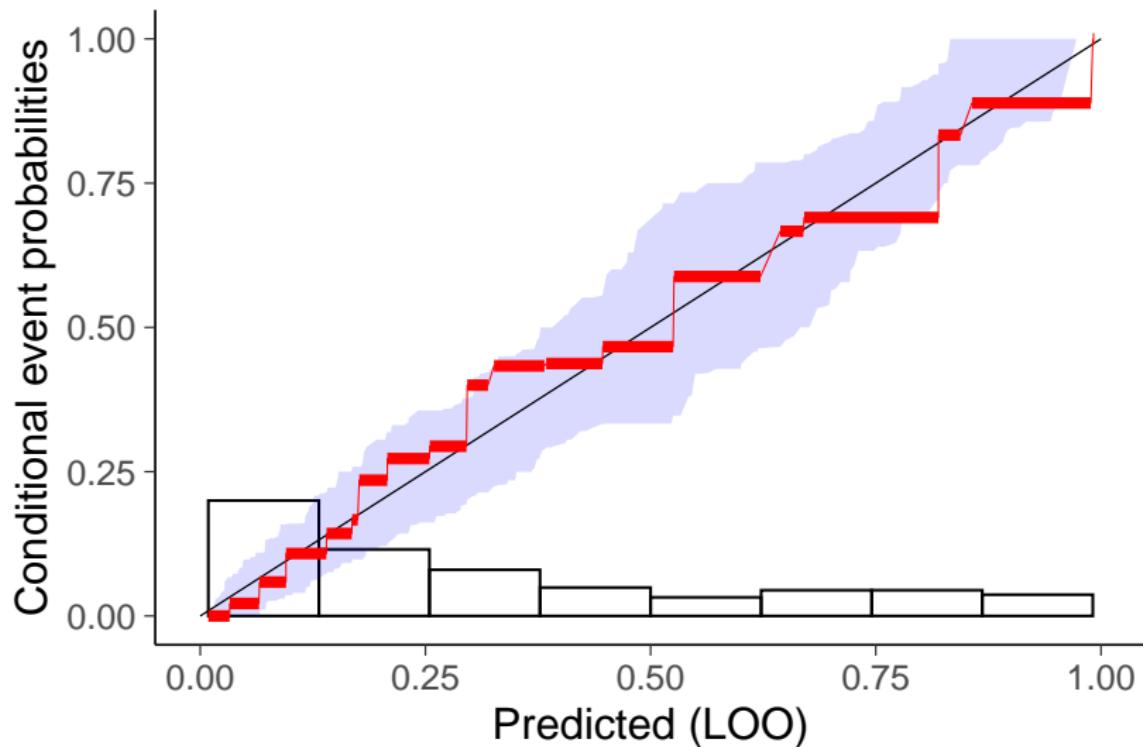
PPC with binning for binary data



## PPC for binary target

Diabetes prediction with logistic regression - diabetes demo

PPC with monotonic regression for binary data



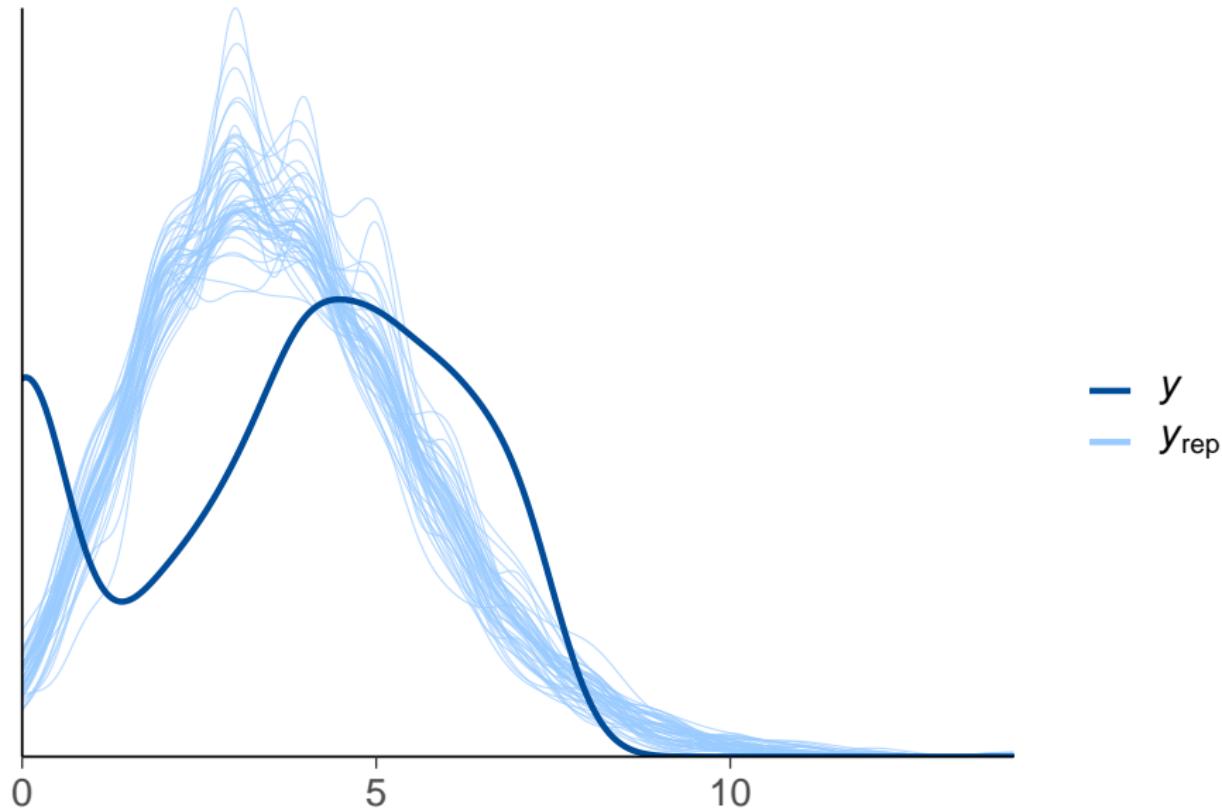
# Posterior predictive checking

- demo demos\_rstan/ppc/poisson-ppc.Rmd

```
data {
  int <lower=1> N;
  int <lower=0> y[N];
}
parameters {
  real <lower=0> lambda;
}
model {
  lambda ~ exponential(0.2);
  y ~ poisson(lambda);
}
generated quantities {
  real log_lik[N];
  int y_rep[N];
  for (n in 1:N) {
    y_rep[n] = poisson_rng(lambda);
    log_lik[n] = poisson_lpmf(y[n] | lambda);
  }
}
```

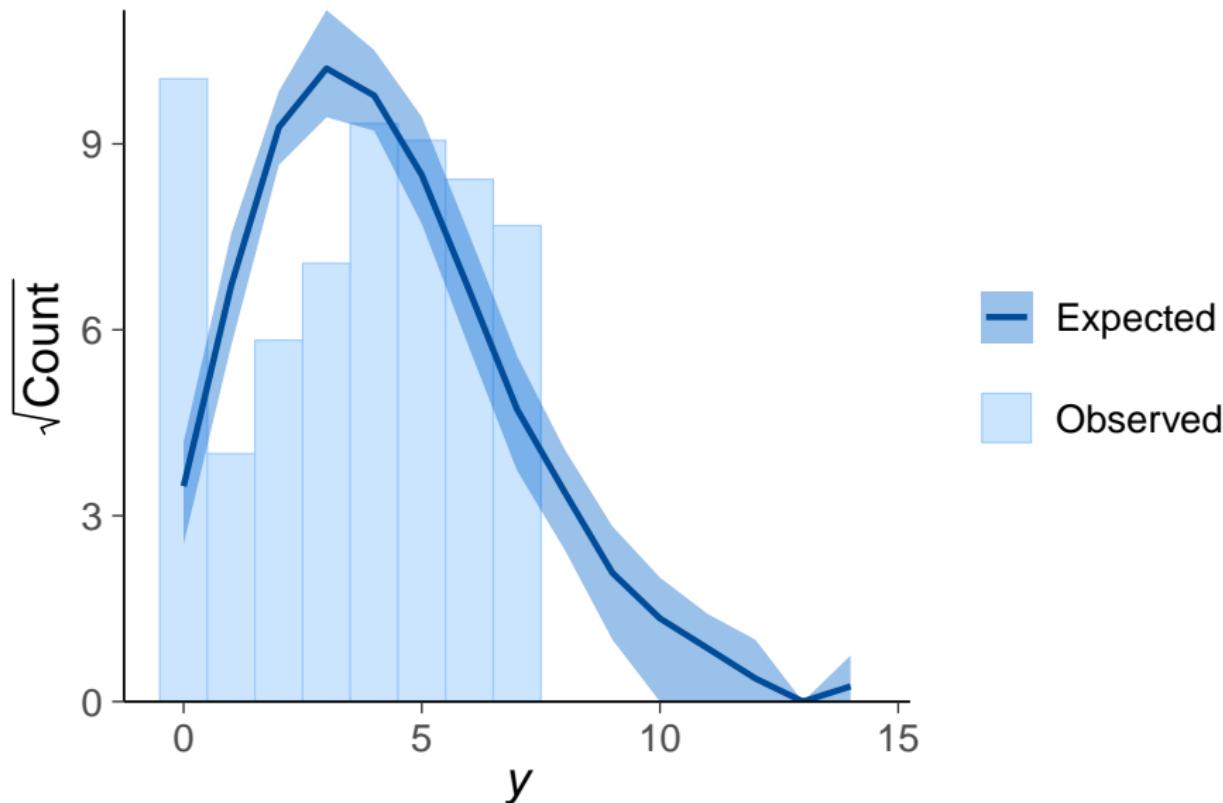
## PPC for count data – Poisson model

```
ppc_dens_overlay(y, yrep[1:50,])
```



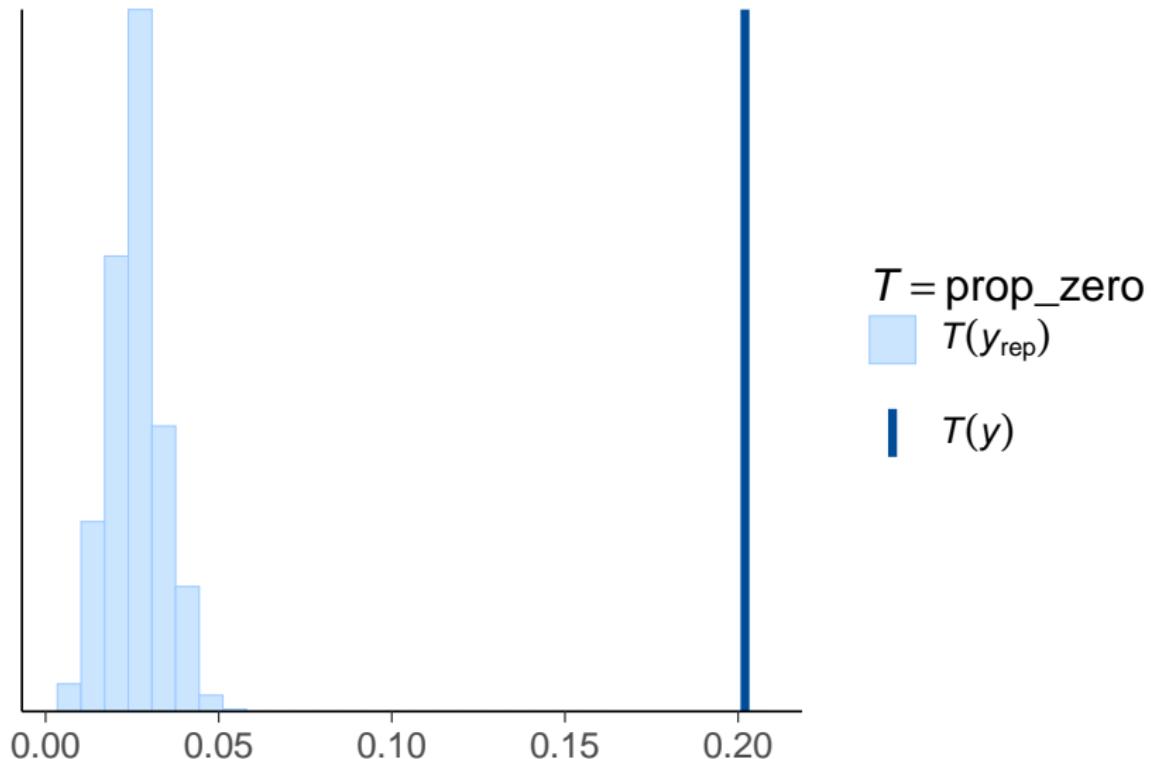
## PPC for count data – Poisson model

`ppc_rootogram(y, yrep)`



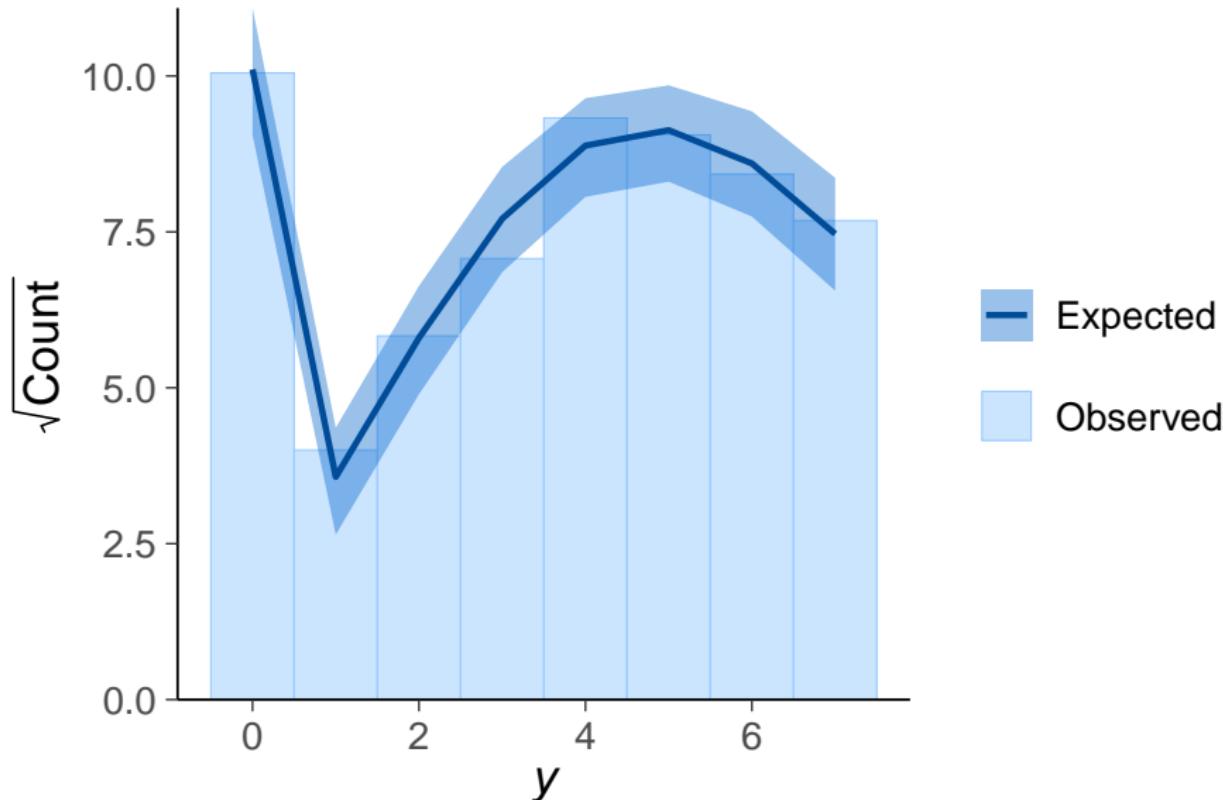
## PPC for count data – Poisson model

```
prop_zero <- function(x) mean(x == 0)  
ppc_stat(y, yrep, stat = "prop_zero")
```



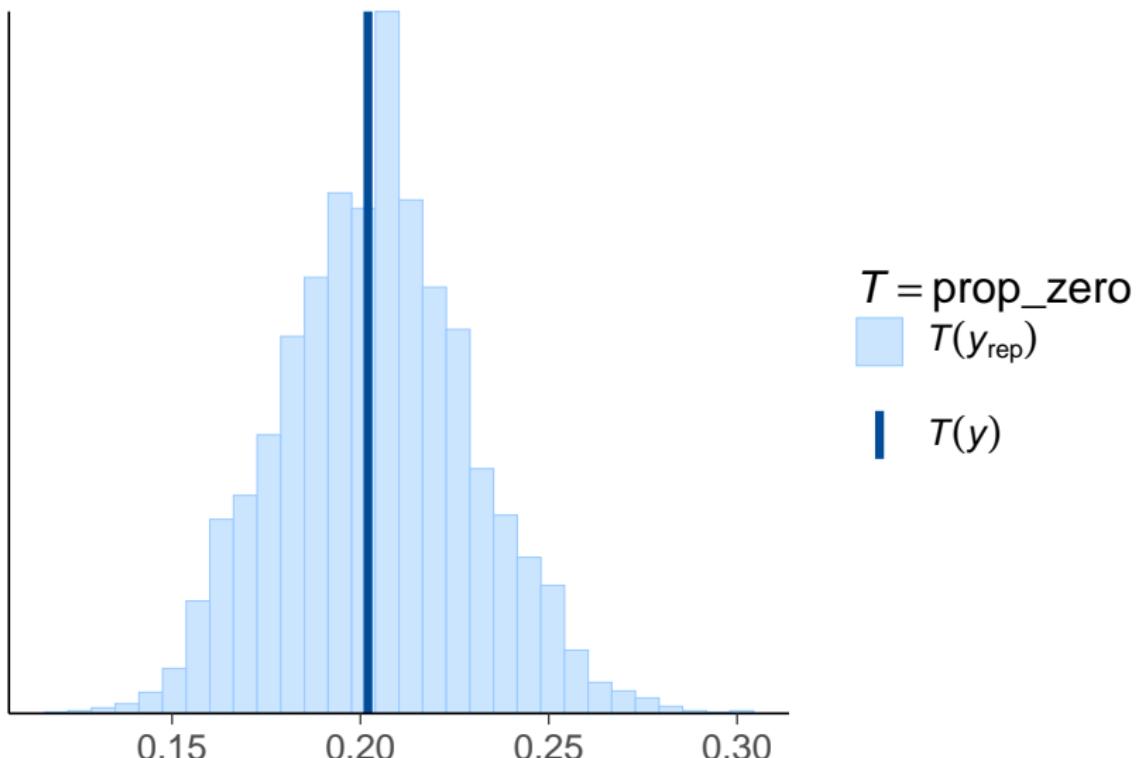
## PPC for count data – hurdle truncated Poisson model

ppc\_rootogram(y, yrep2)



## PPC for count data – hurdle truncated Poisson model

```
prop_zero <- function(x) mean(x == 0)  
ppc_stat(y, yrep2, stat = "prop_zero")
```



## Further reading and examples

- Gabry, Simpson, Vehtari, Betancourt, and Gelman (2019).  
Visualization in Bayesian workflow.  
<https://doi.org/10.1111/rssc.12378>.
- Graphical posterior predictive checks using the bayesplot package  
<http://mc-stan.org/bayesplot/articles/graphical-ppcs.html>
- Another demo `demos_rstan/ppc/poisson-ppc.Rmd`

## Sensitivity analysis

- How much different choices in model structure and priors affect the results

## Sensitivity analysis

- How much different choices in model structure and priors affect the results
  - test different models and priors

# Sensitivity analysis

- How much different choices in model structure and priors affect the results
  - test different models and priors
  - alternatively combine different models to one model
    - e.g. hierarchical model instead of separate and pooled
    - e.g.  $t$  distribution contains Gaussian as a special case
  - robust models are good for testing sensitivity to “outliers”
    - e.g.  $t$  instead of Gaussian

# Sensitivity analysis

- How much different choices in model structure and priors affect the results
  - test different models and priors
  - alternatively combine different models to one model
    - e.g. hierarchical model instead of separate and pooled
    - e.g.  $t$  distribution contains Gaussian as a special case
  - robust models are good for testing sensitivity to “outliers”
    - e.g.  $t$  instead of Gaussian
- Compare sensitivity of essential inference quantities
  - extreme quantiles are more sensitive than means and medians
  - extrapolation is more sensitive than interpolation