



UNIVERSITÀ
DEGLI STUDI DI TRIESTE



Dipartimento di scienze economiche,
aziendali, matematiche e statistiche
“Bruno de Finetti”

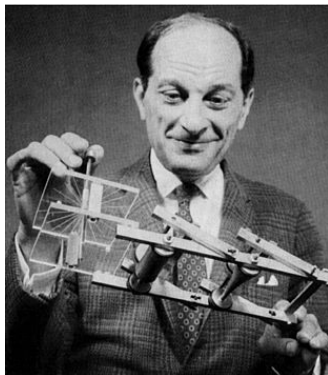
Bayesian Statistics

Introduction to Stan

Leonardo Egidi

A.A. 2018/19

Origins



Stanislaw Ulam (1909-1984): Manhattan project, H-Bomb experiments in Los Alamos, MCMC father jointly with John von Neumann.

Indice

- 1 What is Stan?
- 2 Why Stan?
- 3 Writing a Stan program
- 4 Linked package: bayesplot

What is Stan?

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

Why Stan?

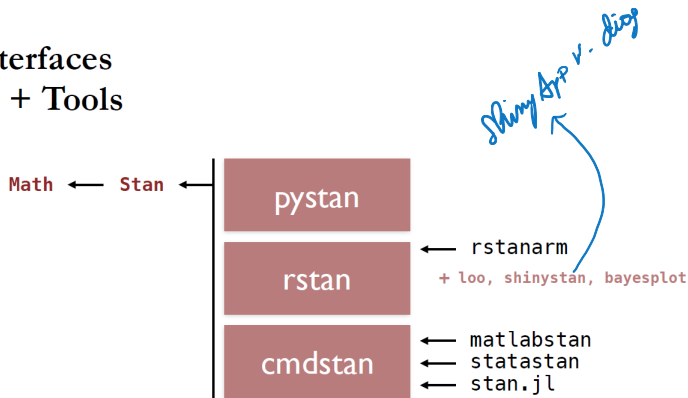
- Fit rich Bayesian statistical models. Close to the *big data* philosophy.
- Efficiency
 - Hamiltonian Monte Carlo + NUTS *NO-U-Turn Sampling*
 - Compiled to C++
- Flexible domain specific language
- “Freedom-respecting, open-source”
 - doc & written materials
 - interacting community
 - continuous development
- Interaction with some other R packages designed to explore the Stan output.

Who is using Stan?

- **Biological & physical sciences:** clinical trials, epidemiology, genomics, population ecology, entomology, ophthalmology, neurology, agriculture, fisheries, cancer biology, astrophysics & cosmology, molecular biology, oceanography, climatology.
- **Social sciences:** population dynamics, psycholinguistics, social networks, political science, human development, economics.
- **Many more:** sports analytics, public health, publishing, finance, pharma, actuarial, recommender systems, educational testing, materials engineering.

Interfaces

Interfaces + Tools



Indice

- 1 What is Stan?
- 2 Why Stan?**
- 3 Writing a Stan program
- 4 Linked package: bayesplot

Improving MCMC performance

With Stan, we aim to provide an MCMC implementation that works robustly for as many target distributions as possible

- Gibbs, RW Metropolis can be very inefficient, hard to diagnose.
- To explore complicated high-dimensional spaces we need to leverage what we know about the geometry of the typical set. *set of interest of our param. region*
- For such a reason, Stan **employs Hamiltonian Monte Carlo**.

We will have now only a brief sketch about how HMC works (further readings are mentioned later on). The Stan users, however, may use, analyze and interpret HMC outputs as they were *standard* MCMC outputs!

Moving to Hamiltonian Monte Carlo

Once we have built a model, Bayesian computation reduces to evaluating expectations, or integrals.

$$E_{\pi}(\theta|y) = \int \theta \pi(\theta|y) d\theta \quad (1)$$

How do we compute posterior expectations in practice?

- Construct a Markov chain that explores the parameter space.
- Anything you would want to do if you could write it analytically, you can do to any accuracy with the draws (history) of the chain

$$\lim_{S \rightarrow \infty} \frac{1}{S} \sum_{s=1}^S \theta^{(s)} \rightarrow E_{\pi}(\theta|y)$$

Moving to Hamiltonian Monte Carlo

To be efficient we need to focus computation on the **relevant neighborhoods** of parameter space. Relevant neighborhoods, however, are defined not by probability density but rather by probability mass.



But exactly which neighborhoods end up contributing most to arbitrary expectations?



The neighborhoods around the maxima of probability distributions feature a lot of probability density, but, especially in a large number of dimensions, or in long tailed distributions, they do not feature much volume. In other words, the *sliver* size $d\theta$ tends to be small there.

The Geometry of High-Dimensional Spaces

Consider a rectangular partitioning centered around a distinguished point, such as the mode (example from Betancourt, 2017):

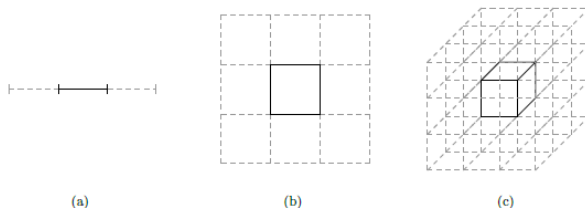
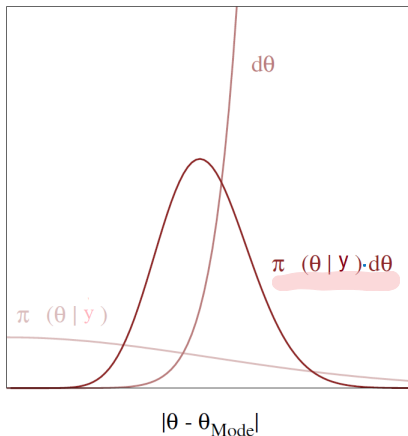


FIG 1. To understand how the distribution of volume behaves with increasing dimension we can consider a rectangular partitioning centered around a distinguished point, such as the mode. (a) In one dimension the relative weight of the center partition is $1/3$, (b) in two dimensions it is $1/9$, (c) and in three dimensions it is only $1/27$. Very quickly the volume in the center partition becomes negligible compared to the neighboring volume.

One of the characteristic properties of high-dimensional spaces is that there is much more volume outside any given neighborhood than inside of it!

Typical set

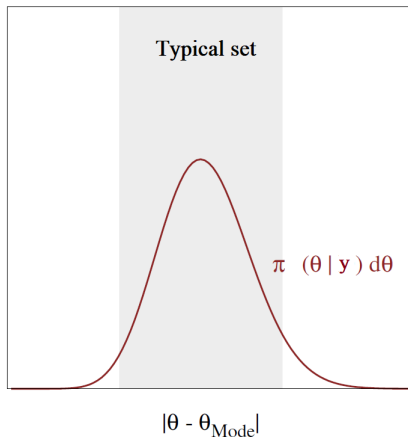
Thus, relevant neighborhoods are defined not by probability density but rather by probability mass.



← Extreme Case

Typical set

Probability mass concentrates on a hypersurface called the **typical set** that surrounds the mode.



Moving to Hamiltonian Monte Carlo

To accurately estimate expectations we need a method for numerically finding and then exploring the typical set.



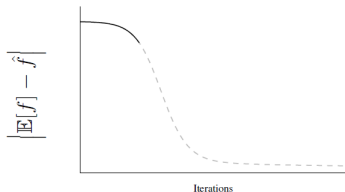
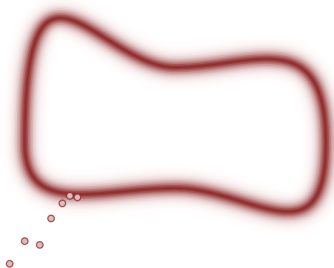
A Markov transition that targets our desired distribution naturally concentrates towards probability mass. An inherent inefficiency in the Gibbs sampler and in the random walk Metropolis Hastings is their random walk behaviour: the simulations can take a long time *zigging* and *zagging* while moving through the target distribution.



HMC borrows strengths from physics to suppress the random walk behaviour in the Metropolis algorithm, thus allowing it to move much more rapidly through the target distribution. The method enjoys the gradient of the log-posterior distribution, $\frac{d \log(\pi(\theta|y))}{d\theta}$ for a sort of adjustment of the algorithm towards the typical set area.

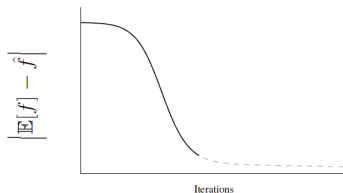
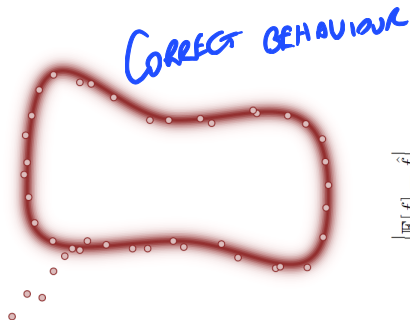
When something goes wrong

Under ideal conditions, MCMC estimators converge to the true expectations in a very practical progression.



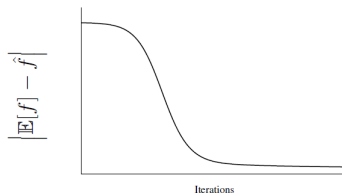
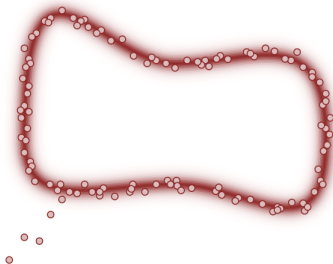
When something goes wrong

Under ideal conditions, MCMC estimators converge to the true expectations in a very practical progression.



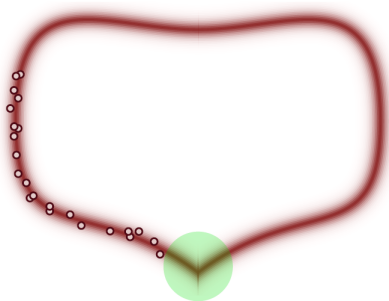
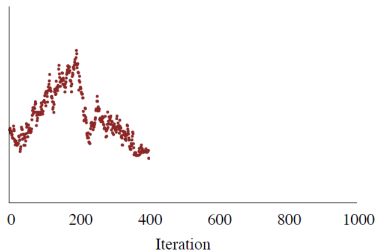
When something goes wrong

Under ideal conditions, MCMC estimators converge to the true expectations in a very practical progression.



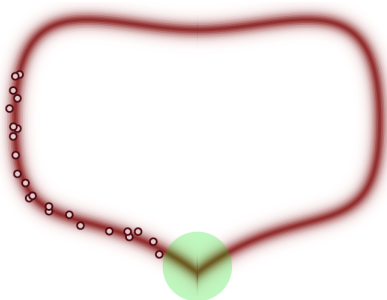
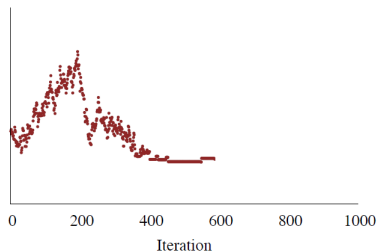
When something goes wrong

There are many pathological posterior geometries, however, that spoil these ideal conditions.



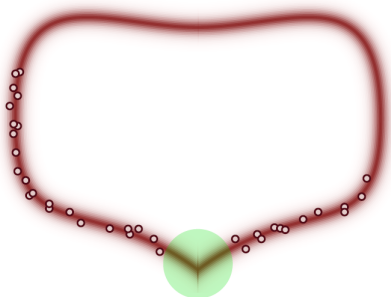
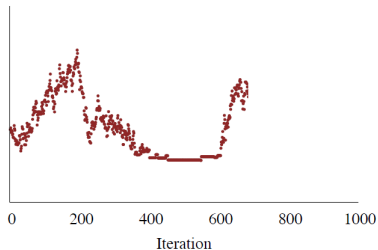
When something goes wrong

There are many pathological posterior geometries, however, that spoil these ideal conditions.



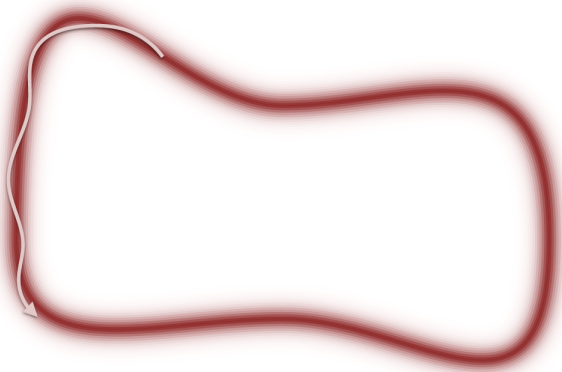
When something goes wrong

There are many pathological posterior geometries, however, that spoil these ideal conditions.



Hamiltonian Monte Carlo

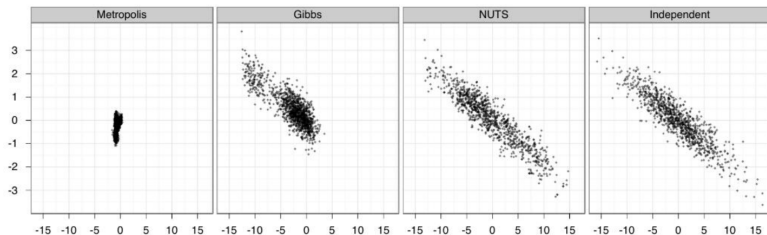
- Hamiltonian Monte Carlo yields fast, and **robust, exploration of the distributions** common in practice.



Hamiltonian Monte Carlo: bivariate Gaussian

Comparison of algorithms on **highly correlated** 250-dimensional Gaussian distribution

- Do **1,000,000** draws with both Random Walk Metropolis and Gibbs, thinning by 1000
- Do **1,000** draws using Stan's NUTS algorithm (no thinning)
- Do 1,000 independent draws (we can do this for multivariate normal)



Indice

- 1 What is Stan?
- 2 Why Stan?
- 3 Writing a Stan program**
- 4 Linked package: bayesplot

Before starting

What is a Bayesian model?

- Building a Bayesian model forces us to build a model for how the data is generated
- We often think of this as specifying a prior and a likelihood, as if these are two separate things
- They are not!

Generative models

The philosophy behind Stan is to think **generatively**.



The model is expressed as a joint probability distribution of observed and unobserved variables, which may be decomposed as follows:

$$p(y, \theta) = p(y|\theta)\pi(\theta) \quad (2)$$

The posterior of interest is then proportional to the joint distribution (2):

$$p(\theta|y) \propto p(y|\theta)\pi(\theta) \quad (3)$$

Generative models

A Bayesian modeller commits to an a priori joint distribution:

$$p(y, \theta) = \underbrace{p(y|\theta)\pi(\theta)}_{\text{Likelihood} \times \text{Prior}} = \underbrace{\pi(\theta|y)p(y)}_{\text{Posterior} \times \text{Marginal Likelihood}} \quad (4)$$

Generative models and vague priors

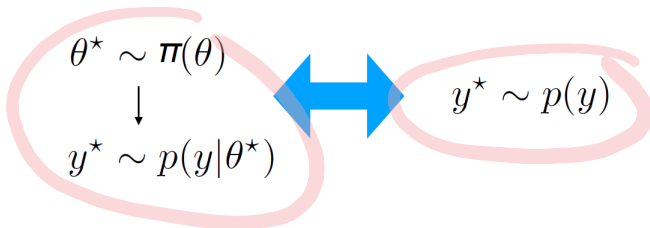
What is the problem with *vague/diffuse* priors?

- If we use an improper prior, then we do not specify a joint model for our data and parameters.
- More importantly, we do not specify a data generating mechanism $p(y)$.
- By construction, these priors do not regularize inferences, which is quite often a bad idea
- Proper but diffuse is better than improper but is still often problematic.

Try using informative priors

Generative models

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



Stan computations

Stan works in logarithmic terms: all the computations are actually done on log-scale. So, for the posterior we have.

$$\log(\pi(\theta|y)) = \log(\pi(\theta)) + \log(p(y|\theta)) + \text{constant} \quad (5)$$

Products become sums of logs:

$$p(y|\theta) = \prod_{i=1}^n p(y_i|\theta) \rightarrow \log(p(y|\theta)) = \sum_{i=1}^n \log(p(y_i|\theta)).$$

Starting point

We are now going to write a Stan program together:

- Open a new empty file in RStudio
- Save it as `linear_regression.stan`

Blocks strategy

Stan programs are organized into **blocks**:

- **data block**: declare data types, sizes, and constraints. Read from data source and constraints validated. Evaluated: once.
- **parameters block**: declare parameter types, sizes, and constraints. Evaluated: every log prob evaluation.
- **transformed parameters block**: declare those parameters transformed from the original ones declared in the parameters block. Evaluated: every log prob evaluation.
- **model block**: statements defining the posterior density in log scale. Evaluated: every log prob evaluation.
- **generated quantities**: declare and define derived variables. (P)RNGs, predictions, event probabilities, decision making. Constraints validated. Evaluated: once per draw.

Data block

```
data {
```

```
  // Dimensions
```

```
  int<lower=1> N;
```

```
  int<lower=1> K;
```

```
  // Variables
```

```
  matrix[N, K] X;
```

```
  vector[N] y;
```

```
}
```

integer

use capitals for sample sizes (convention)

C++ style

```
// single line comment
```

```
/* multiple lines of  
comments */
```

Parameters' block

```
parameters {  
  real alpha;  $\rightarrow$  INTERCEPT  
  vector[K] beta;  $\rightarrow$  SLOPES  
  real<lower=0> sigma;  
}
```

$$y \sim N(\alpha + X\beta, \sigma^2)$$

constraints *required* in
parameters block

Model block

```

model {
  // priors (flat, uniform, if omitted)
  sigma ~ exponential(1);
  alpha ~ normal(0, 10);
  for (k in 1:K) beta[k] ~ normal(0, 5);
  for (n in 1:N) {
    y[n] ~ normal(X[n, ] * beta + alpha, sigma);
  }
}

```

Handwritten notes:
 - An arrow points from `normal(0, 10)` to the text "This is SD $\Rightarrow \sigma^2 = 100$ ".
 - The line `for (k in 1:K)` is underlined in blue.
 - Next to `normal(0, 5)` is the text "OR beta ~ normal(0, 5)".
 - Below the underlined `for (k in 1:K)` is the text "No for cycle".

Why is the default automatically uniform?

- $\pi(\theta) \propto 1$ (0 on log scale)
- Nothing added to log prob

Generated quantities block

Exploit generability

```
generated quantities {  
  vector[N] y_rep;  
  for (n in 1:N) {  
    real y_hat = X[n,] * beta + alpha; // local/temp  
    y_rep[n] = normal_rng(y_hat, sigma);  
  }  
}
```

Complete Stan model

data {	
int<lower=1> N;	
int<lower=1> K;	
matrix[N, K] X;	
vector[N] y;	Observed variables
}	
parameters {	
real alpha;	
vector[K] beta;	
real<lower=0> sigma;	Unobserved variables
}	
model {	
sigma ~ exponential(1);	$\log \pi(\theta)$
alpha ~ normal(0, 10);	
for (k in 1:K) beta[k] ~ normal(0, 5);	+
for (n in 1:N)	$\log p(y \theta)$
y[n] ~ normal(alpha + X[n,] * beta, sigma);	
}	
generated quantities {	
vector[N] y_rep;	
for (n in 1:N)	Simulate from
y_rep[n] = normal_rng(alpha + X[n,] * beta, sigma);	generative model
}	

Launching the Stan model from R

Now we may launch the Stan program directly in R:

```
library(rstan)

# passing the data (already stored)
data <- list(N=N, K=K, X=X, y=y)

# fitting the model
fit1 <- stan(
  file = 'linear_regression.stan',
  data = data,
  iter = 2000,
  chains = 4)

# extracting the estimates
sims <- extract(fit1)
```

Handwritten blue notes:

- An arrow points from `iter = 2000` to the word "default".
- An arrow points from `chains = 4` to the phrase "(parallel) default".

First example: 8 schools

This example studied coaching effects from eight schools.

We denote with y_{ij} the result of the i -th test in the j -th school. We assume the following model:

$$y_{ij} \sim \mathcal{N}(\theta_j, \sigma_y^2)$$

$$\theta_j \sim \mathcal{N}(\mu, \tau^2)$$

Do some schools perform better/worse according to these coaching effects?

Here is the data, already aggregated by schools:

```
schools_dat <- list(J = 8,
  y = c(28, 8, -3, 7, -1, 1, 18, 12),
  sigma = c(15, 10, 16, 11, 9, 11, 10, 18))
```

↳ VARIES \forall obs.

First Stan model: 8 schools

```
// saved as 8schools.stan
data {
  int<lower=0> J;           // number of schools
  real y[J];               // estimated treatment effects
  real<lower=0> sigma[J];  // standard error of effect estimates
}
parameters {
  real mu;                // population treatment effect
  real<lower=0> tau;       // standard deviation in treatment effects
  vector[J] eta;          // unscaled deviation from mu by school
}
transformed parameters {
  vector[J] theta = mu + tau * eta;    // school treatment effects
}
model {
  eta ~ normal(0,1);           // prior
  y ~ normal(theta, sigma);    //likelihood
}
```


First example: 8 schools

To fit the model and visualize the estimates, it is sufficient to type in R the following commands (with 2000 iterations and 4 chains as a default):

```
fit_8schools <- stan(file = '8schools.stan', data = schools_dat)
print(fit_8schools, pars=c("mu", "tau", "theta"))
```

	mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
mu	7.89	5.04	-2.31	4.74	7.92	11.05	18.05	2352	1
tau	6.70	5.71	0.24	2.61	5.43	9.19	21.16	1480	1
theta[1]	11.36	8.23	-2.25	6.18	10.29	15.46	31.15	3161	1
theta[2]	7.89	6.21	-4.43	3.96	7.83	11.78	20.47	4923	1
theta[3]	6.05	7.59	-10.81	1.92	6.56	10.81	20.25	4057	1
theta[4]	7.60	6.44	-5.36	3.74	7.63	11.73	20.57	5055	1
theta[5]	5.13	6.23	-8.45	1.35	5.60	9.26	16.37	4346	1
theta[6]	5.95	6.68	-8.21	1.99	6.30	10.21	18.21	4313	1
theta[7]	10.62	6.93	-1.58	6.12	10.14	14.53	25.63	3381	1
theta[8]	8.40	7.77	-7.18	3.84	8.26	12.63	25.78	3854	1

Indice

- 1 What is Stan?
- 2 Why Stan?
- 3 Writing a Stan program
- 4 Linked package: bayesplot**

Posterior graphical analysis with bayesplot

Once we fit a model, it is to vital check it via graphical inspection. The bayesplot package (for any help, see the [vignette](#)) is designed to this task.



The package allows to display:

- Posterior uncertainty intervals
- Univariate marginal posterior distributions
- Bivariate plots
- Trace plots
- Posterior predictive plots

Posterior graphical analysis with bayesplot

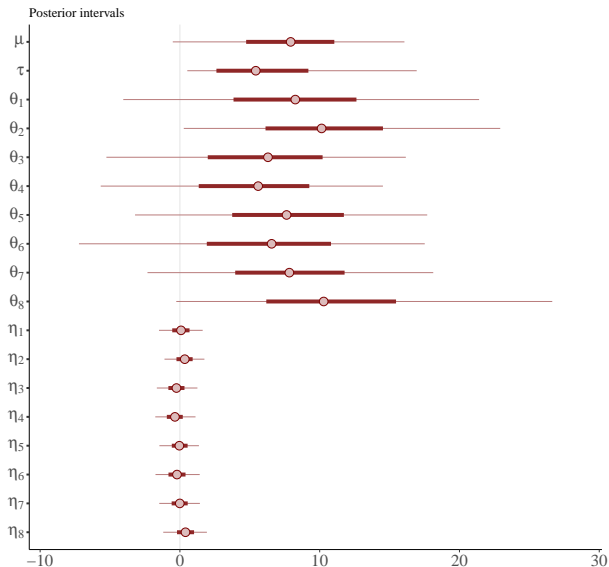
The first step is to save the posterior. Then you have many choices:

```
library(bayesplot)
posterior <- as.array(fit_8schools)

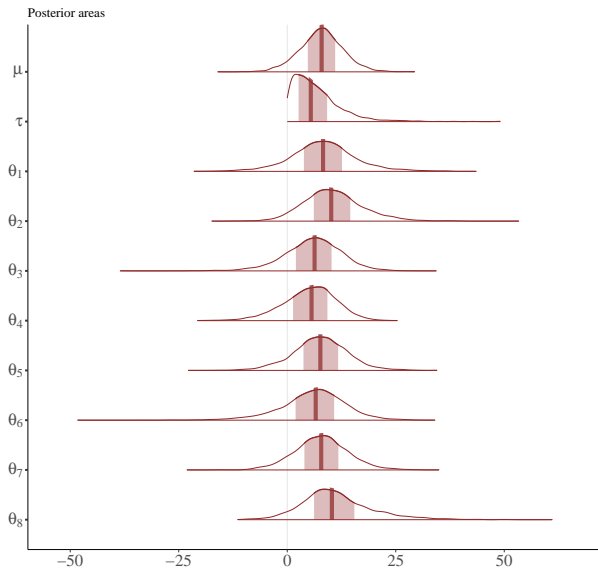
mcmc_intervals(posterior)      # posterior intervals
mcmc_areas(posterior)         # posterior areas
mcmc_dens(posterior)          # marginal posteriors
mcmc_pairs(posterior)         # bivariate plots
mcmc_trace(posterior)         # trace plots
```

With the arguments `pars` or `regex_pars` you may select the desired parameters.

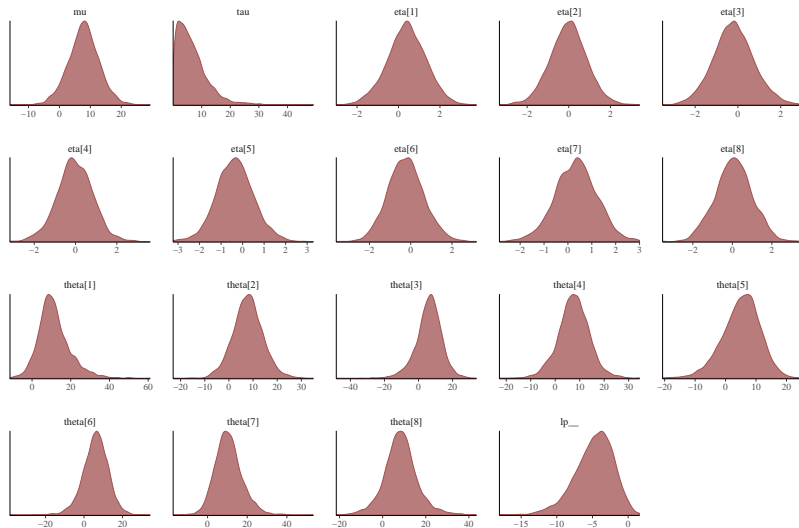
Posterior uncertainty intervals



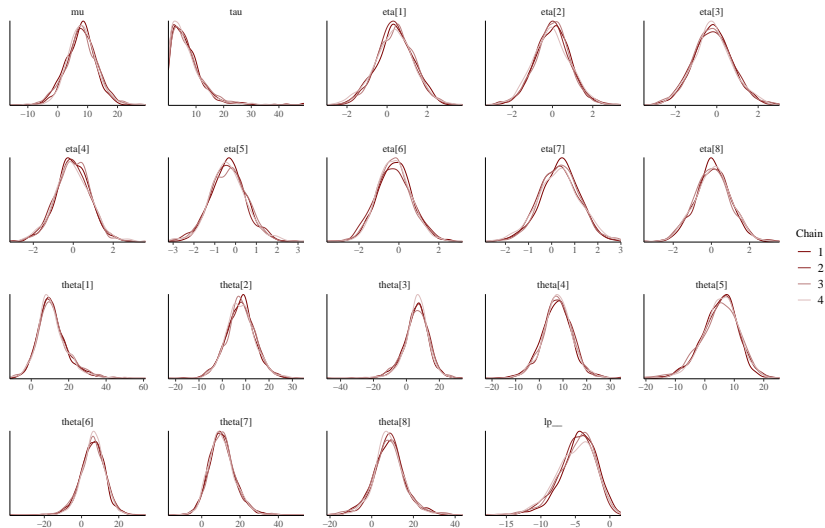
Posterior uncertainty areas



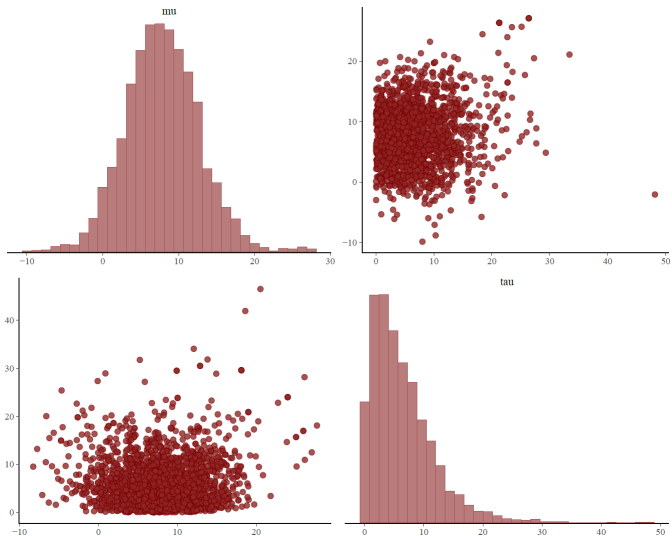
Marginal posteriors



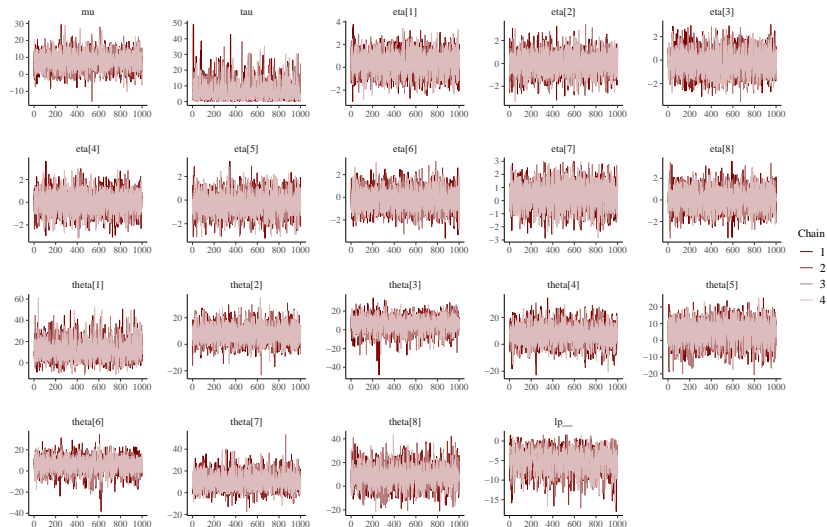
Marginal posteriors separated for each chain



Bivariate posterior plots



Trace plots for the Markov chains



Our challenge with Stan

The Stan shuttle is ready to start! We will learn to:

- **write** simple and more complex model in Stan: lm, glm, hierarchical models.
- **analyze** the posterior summaries.
- **criticize** the model and, eventually, change/reparametrize it.

Further reading

Further reading:

- Carpenter, B, and Gelman, A, Hoffman, M.D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P., Riddell, A. (2017). Stan: A Probabilistic Programming Language, *Journal of statistical software* 76(1). Here the [pdf](#)

Further optional reading about Hamiltonian Monte Carlo:

- Betancourt, M. (2017) A conceptual introduction to Hamiltonian Monte Carlo. Here the [pdf](#)