In [46]: `using Distributions, Plots, KernelDensity`

# Bayesian Methods

A **data generating process** is a model characterized by a parameter vector, and the data is generated from the model at a **fixed parameter** vector, $\theta_0$.

Bayesians view **data as given**, and **update beliefs about a parameter** using the information about the parameter contained in the data.

There's nothing obviously contradictory in these views. Nevertheless, it's not hard to find discussions where there are disagreements.

To compare frequentist approach vs. Bayesian approach, see this video.

Here, I'm trying to address a model with a fixed non-random parameter about which we would like to learn. As long as the object of interest (the dgp and it's parameter) is agreed upon, then we can contemplate using **any** convenient tool.

- A pure frequentist ''full information'' approach would require writing the joint likelihood of all samples, which would almost certainly constitute an impossible task.
- The Bayesian approach concentrates all of the information coming from previous work in the form of a prior. A fairly simple, easy to use prior may not exactly capture all previous information, but it could offer a handy and reasonably accurate summary, and it's almost certainly better than simply pretending that all of that previous information simply doesn't exist. So, the idea of a prior as a summary of what we have learned may simply be viewed as a practical solution to the problem of using all the available information. Given that it's a summary, one may as well use a convenient form, as long as it's plausible and the results don't depend too exaggeratedly particular form used.

As long as one takes the view that there is a fixed unknown parameter value $\theta_0$ which generates all samples, then frequentist and Bayesian methods are trying to inform us about the same object, and the choice between tools may become one of convenience.

## Definitions

The Bayesian approach summarizes beliefs about parameters using a density function:

- There is a true unknown parameter vector, $\theta_0$, and the density, $\pi(\theta)$, which is known as the **prior**, reflects current beliefs about the parameter, before observing the sample. It is assumed that the econometrician can provide this density.
- We also have sample information, $y=\{y_1, y_2, \ldots y_n\}$. We're already familiar with the likelihood function, $f(y|\theta)$, which is the density of the sample given a parameter value.

Given these two pieces, we can write the joint density of the sample and the beliefs:

$$f(y, \theta) = f(y|\theta)\pi(\theta)$$

We can get the marginal likelihood by integrating out the parameter, integrating over its support $\Theta$:

$$f(y) = \int_\Theta f(y, \theta) d\theta$$

The last step is to get the **posterior** of the parameter. This is simply the density of the parameter conditional on the sample, and we get it in the normal way we get a conditional density, using **Bayes' theorem**:

$$f(\theta|y) = \frac{f(y, \theta)}{f(y)} = \frac{f(y|\theta)\pi(\theta)}{f(y)}$$

- The movement from the prior to the posterior reflects the **learning** that occurs about the parameter when one receives the sample information.
- The sources of information used to make the posterior are the prior and the likelihood function.
- Once we have the posterior, one can provide a complete probabilistic description about our updated beliefs about the parameter, using quantiles or moments of the posterior.

So far, this is pretty straightforward. The complications are mostly computational. To illustrate, the posterior mean is

$$E(\theta|y) = \int_\Theta \theta f(\theta|y) d\theta = \frac{\int_\Theta \theta f(y|\theta)\pi(\theta) d\theta}{\int_\Theta f(y, \theta) d\theta}$$

- One can see that a means of integrating will be needed.
- Only in very special cases will the integrals have analytic solutions.
- Otherwise, computational methods will be needed. Advances in computational methods are what have lead to the increased use of Bayesian methods.

## Extensions

Chernozhukov and Hong (2003) "An MCMC Approach to Classical Estimation" is a very interesting article that shows how Bayesian methods may be used with criterion functions that are associated with classical estimation techniques.

For example, when the criterion function $L_n(\theta)$ in their paper is set to the log-likelihoodfunction and the pseudo-prior $\pi(\theta)$ is a real Bayesian prior, the class of estimators discussed by CH reduces to the ordinary Bayesian posterior mean. this estimator is consistent and asymptotically normally distributed. In particular, the Bayesian posterior mean has the same asymptotic distribution as does the ordinary maximum likelihood estimator.

- the intuition is clear: as the amount of information coming from the sample increases, the likelihood function brings an increasing amount of information, relative to the prior. Eventually, the prior is no longer important for determining the shape of the posterior.
- when the sample is large, the shape of the posterior depends on the likelihood function. The likelihood function collapses around $\theta_0$ when the sample is generated at $\theta_0$. The same is true of the posterior, it narrows around $\theta_0$. This causes the posterior mean to converge to the true

parameter value. In fact, all quantiles of the posterior converge to $\theta_0$. Chernozhukov and Hong discuss estimators defined using quantiles.

- For an econometrician coming from the frequentist perspective, this is attractive. The Bayesian estimator has the same asymptotic behavior as the MLE. There may be computational advantages to using the Bayesian approach, because there is no need for optimization. If the objective function that defines the classical estimator is irregular (multiple local optima, nondifferentiabilities, noncontinuities...), then optimization may be very difficult. However, Bayesian methods that use integration may be more tractable. This is the main motivation of CH's paper. Additional advantages include the benefits if an informative prior is available. When this is the case, the Bayesian estimator can have better small sample performance than the maximum likelihood estimator.

## Example (Exponential Distribution)

Suppose data is generated by i.i.d. sampling from an exponential distribution with mean $\theta$. An exponential random variable takes values on the positive real numbers. Waiting times are often modeled using the exponential distribution.

- The density of a typical sample element is $f(y|\theta) = \frac{1}{\theta}e^{-y/\theta}$. The likelihood is simply the product of the sample contributions.
- Suppose the prior for $\theta$ is $\theta \sim$ lognormal(1,1). This means that the logarithm of $\theta$ is standard normal. We use a lognormal prior because it enforces the requirement that the parameter of the exponential density be positive.

```
In [47]:   n = 10    # sample size
           theta = 3 # true theta
           y = rand(Exponential(theta), n) # sample from exponential(theta)
```

```
Out[47]:  10-element Array{Float64,1}:
           15.694341912542694
            0.7458285508366813
            8.741647661375469
            0.9163465573331584
            0.9154960432774923
            5.544738338157552
            2.348936698893298
            3.430901323005183
            0.6597096037830681
            1.8874567713487362
```

```
In [48]:   likelihood(y,theta)
```

```
Out[48]:  0-dimensional Array{Float64,0}:
           2.041827654972068e-11
```

```
In [49]:   # the prior is lognormal(1,1)
           function prior(theta)
               d = LogNormal(1.0,1.0)
               p = pdf.(Ref(d), theta)
               pmean = exp(1.5) # mean of lognormal is exp(mu+sig/2)
               return p, pmean
           end
```

```julia
# the likelihood function
function likelihood(y, theta)
    dens = zeros(size(theta))
    for i = 1:size(theta,1)
        d = Exponential(theta[i])
        dens[i] = prod(pdf.(Ref(d), y))
    end
    return dens
end


# joint is prior X likelihood
function joint(y, theta)
    l = likelihood(y, theta)
    p, junk = prior(theta)
    dens = l.*p
end

# compute marginal likelihood of Y by integrating out theta (crude, only illustrative)
function marginal(y)
    dens = 0.0
    theta = 0.0
    delta = 0.01
    # evaluate joint over grid
    for r = 1:1000
        theta += delta
        dens += joint(y, theta)
    end
    # marginalize by integrating the joint (sum up height X width)
    dens = dens*delta
end

# the posterior, by Bayes' Law
function posterior(y, theta)
    m = marginal(y)
    j = joint(y, theta)
    dens = j ./ m
    thetas = range(0.01,stop=10,length=1000)
    pmean = sum(dens.*thetas.*0.01)
    return dens, pmean
end
```

Out[49]: posterior (generic function with 1 method)

```julia
n = 10   # sample size
theta = 3 # true theta
y = rand(Exponential(theta), n) # sample from exponential(theta)
# make plots
thetas = range(0.01,stop=10,length=1000)
p, priormean = prior(thetas)
post, postmean = posterior(y, thetas);
```
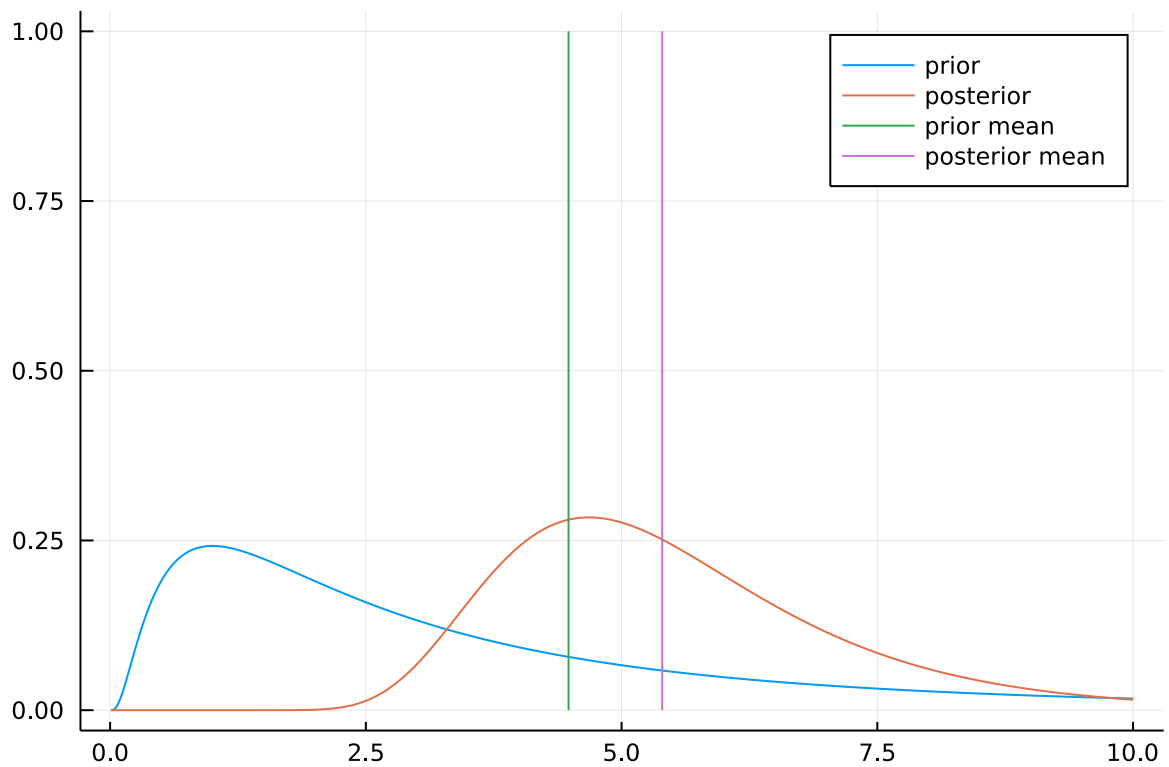
In [51]:
```julia
plot(thetas, [p post], label = ["prior" "posterior"])
plot!([priormean, priormean], [0.0, 1.0], label = "prior mean")
plot!([postmean, postmean], [0.0, 1.0], label = "posterior mean")
```
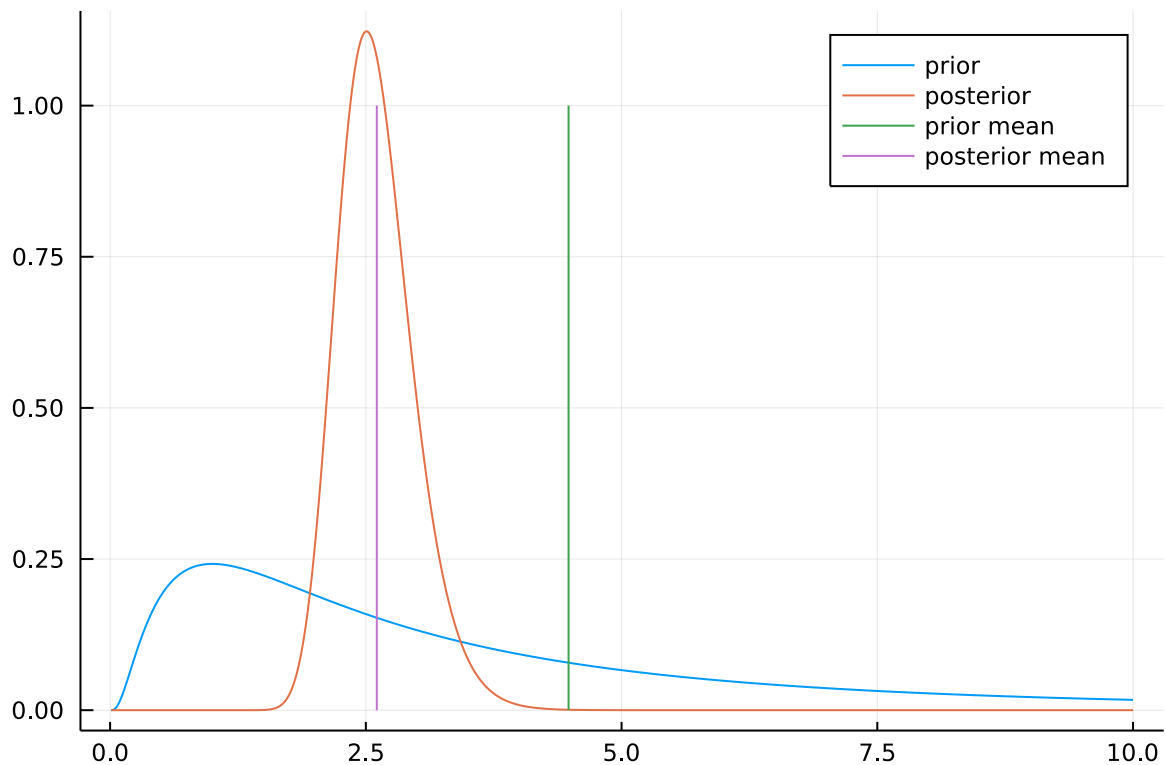
Out[51]:

```
In [52]:   n = 50   # sample size
           theta = 3 # true theta
           y = rand(Exponential(theta), n) # sample from exponential(theta)
           # make plots
           thetas = range(0.01,stop=10,length=1000)
           p, priormean = prior(thetas)
           post, postmean = posterior(y, thetas)
           plot(thetas, [p post], label = ["prior" "posterior"])
           plot!([priormean, priormean], [0.0, 1.0], label = "prior mean")
           plot!([postmean, postmean], [0.0, 1.0], label = "posterior mean")
```

Out[52]:

Note how the posterior is more concentrated around the true parameter value when sample has size 50. Also note how the posterior mean is closer to the prior mean when the sample is small. When the sample is small, the likelihood function has less weight, and more of the information comes from the prior. When the sample is larger, the likelihood function will have more weight, and its effect will dominate the prior's.

# Computational Methods: MCMC

To compute the posterior mean, we need to evaluate

$$E(\theta|y) = \int_{\Theta} \theta f(\theta|y)d\theta$$
$$= \frac{\int_{\Theta} \theta f(y|\theta)\pi(\theta)d\theta}{\int_{\Theta} f(y,\theta)d\theta}.$$

- Note that both of the integrals are multiple integrals, with the dimension given by that of the parameter, $\theta$.
- Under some special circumstances, the integrals may have analytic solutions: e.g., Gaussian likelihood with a Gaussian prior leads to a Gaussian posterior.
- When the dimension of the parameter is low, quadrature methods may be used. What we have done with marginal(y) and posterior(y, theta) is an unsophisticated example of this. More sophisticated methods use an intelligently chosen grid to reduce the number of function evaluations. Still, these methods only work for dimensions up to 3 or so.
- Otherwise, some form of simulation-based "Monte Carlo" integration must be used. The basic idea is that $E(\theta|y)$ can be approximated by $(1/S)\sum_{s=1}^{S}\theta^s$, where $\theta^s$ is a random draw from the posterior distribution $f(\theta|y)$.

**How to make draws from the posterior** when in general we can't compute the posterior?

- the law of large numbers tells us that this average will converge to the desired expectation as $S$ gets large
- convergence will be more rapid if the random draws are independent of one another, but insisting on independence may have computational drawbacks.

Monte Carlo methods include importance sampling, Markov chain Monte Carlo (MCMC) and sequential Monte Carlo (SMC, also known as particle filtering). The great expansion of these methods over the years has caused Bayesian econometrics to become much more widely used than it was in the not so distant (for some of us) past. There is much literature - here we will only look at a basic example that captures the main ideas.

Let's consider the basic **Metropolis-Hastings MCMC** algorithm. We will generate a long realization of a Markov chain process for $\theta$, as follows:

The prior density is $\pi(\theta)$, as above.

Let $g(\theta^*; \theta^s)$ be a proposal density, which describes the density of a trial value $\theta^*$ conditional on starting at $\theta^s$. It must be possible to sample from the proposal. This gives a new trial parameter value $\theta^*$, given the most recently accepted parameter value $\theta^s$. A proposal will be accepted if

$$\frac{f(\theta^*|y)}{f(\theta^s|y)} \frac{g(\theta^s; \theta^*)}{g(\theta^*; \theta^s)} > \alpha$$

where $\alpha$ is a $U(0, 1)$ random variate.

There are two parts to the numerator and denominator: the posterior, and the proposal density.

- The numerator, when the trial value of the proposal has a higher posterior, acceptance is favored.
- The density associated with returning to $\theta^s$ when starting at $\theta^*$, which has to do with the reversibility of the Markov chain. If this is too low, acceptance is not favored. We don't want to jump to a new region if we will never get back, as we need to sample from the entire support of the posterior.

The two together mean that we will jump to a new area only if we are able to eventually jump back with a reasonably high probability. The probability of jumping is higher when the new area has a higher posterior density, but lower if it's hard to get back.

The idea is to sample from all regions of the posterior, those with high and low density, sampling more heavily from regions of high density. We want to go occasionally to regions of low density, but it is important not to get stuck there.

Note that the ratio of posteriors is equal to the ratio of likelihoods times the ratio of priors:

$$\frac{f(\theta^*|y)}{f(\theta^s|y)} = \frac{f(y|\theta^*)}{f(y|\theta^s)} \frac{\pi(\theta^*)}{\pi(\theta^s)}$$

because the marginal likelihood $f(y)$ is the same in both cases. We don't need to compute that integral! We don't need to know the posterior, either. The acceptance criterion can be written as: accept if

$$\frac{f(y|\theta^*)}{f(y|\theta^s)} \frac{\pi(\theta^*)}{\pi(\theta^s)} \frac{g(\theta^s; \theta^*)}{g(\theta^*; \theta^s)} > \alpha$$

otherwise, reject

From this, we see that the information needed to determine if a proposal is accepted or rejected is the prior, the proposal density, and the likelihood function $f(y|\theta)$.

- in principle, the prior is non-negotiable. In practice, people often chose priors with convenience in mind
- the likelihood function is what it is
- the place where artistry comes to bear is the choice of the proposal density
- when the proposal density is \emph{symmetric, }so that $g(\theta^s; \theta^*) = g(\theta^*; \theta^s)$, the acceptance criterion simplifies to

$$\frac{f(y|\theta^*)}{f(y|\theta^s)} \frac{\pi(\theta^*)}{\pi(\theta^s)} > \alpha$$

  A random walk proposal, where the trial value is the current value plus a shock that doesn't depend on the current value, satisfies symmetry.

The steps are:

1. the algorithm is initialized at some $\theta^1$
2. for $s = 2, \ldots, S$,
    A. draw $\theta^*$ from $g(\theta^*; \theta^s)$
    B. according to the acceptance/rejection criterion, if the result is acceptance, set $\theta^{s+1} = \theta^*$, otherwise set $\theta^{s+1} = \theta^s$
    C. iterate

Once the chain is considered to have stabilized, say at iteration $r$, the values of $\theta^s$ for $s > r$ are taken to be draws from the posterior. The posterior mean is computed as the simple average of the value. Quantiles, etc., can be computed in the appropriate fashion.

The art of applying these methods consists of providing a good proposal density so that the acceptance rate is reasonably high, but not too high. There is a vast literature on this, and the vastness of the literature should serve as a warning that getting this to work in practice is not necessarily a simple matter. If it were, there would be fewer papers on the topic.

- too high acceptance rate: this is usually due to a proposal density that gives proposals very close to the current value, e.g, a random walk with very low variance. This means that the posterior is being explored inefficiently, we travel around through the support at a very low rate, which means the chain will have to run for a (very, very...) long time to do a thorough exploration.

- too low acceptance rate: this means that the steps are too large, and we attempt to move to low posterior density regions too frequently. The chain will become highly autocorrelated, as it stays in the same place due to rejections, so long periods convey little additional information relative to a subset of the values in the interval

## Example

The simple exponential example with log-normal prior can be implemented using MH MCMC

In [53]:
```julia
function mcmc(θ, reps::Int64, burnin::Int64, Prior::Function, lnL::Function, Proposal::
    reportevery = Int((reps+burnin)/10)
    lnLθ = lnL(θ)
    chain = zeros(reps, size(θ,1)+1)
    naccept = zeros(size(θ))
    for rep = 1:reps+burnin
        θᵗ = Proposal(θ) # new trial value
        if report
            changed = Int.(.!(θᵗ .== θ)) # find which changed
        end
        # MH accept/reject: only evaluate logL if proposal is in support of prior (avoi
        pt = Prior(θᵗ)
        accept = false
        if pt > 0.0
            lnLθᵗ = lnL(θᵗ)
            accept = rand() < exp(lnLθᵗ-lnLθ) * pt/Prior(θ)
            if accept
                θ = θᵗ
                lnLθ = lnLθᵗ
            end
        end
        if report
            naccept = naccept .+ changed .* Int.(accept)
        end
        if (mod(rep,reportevery)==0 && report)
            println("current parameters: ", round.(θ,digits=3))
            println("  acceptance rates: ", round.(naccept/reportevery,digits=3))
            naccept = naccept - naccept
        end
        if rep > burnin
            chain[rep-burnin,:] = [θ; accept]
        end
    end
    return chain
end
```

Out[53]: mcmc (generic function with 2 methods)

In [54]:
```julia
function npdensity(z)
    #pyplot()
    n = size(z,2)
    p = zeros(n)
    for i = 1:n
        x = z[:,i]
        y = kde(x)
        q05 = quantile(x,0.05)
        Plots.plot(range(minimum(x), stop=q05, length=100),z->pdf(y,z), color=:orange,
        q95 = quantile(x,0.95)
        Plots.plot!(range(q05, stop=q95, length=100),z->pdf(y,z), color=:green, fill=(0
```

```julia
        Plots.plot!(range(q95, stop=maximum(x), length=100),z->pdf(y,z), color = :orang
        m = mean(x)
        Plots.plot!([m,m],[0,pdf(y,m)],color=:blue, label="mean")
        m = median(x)
        if i == 1
            p = Plots.plot!([m,m],[0,pdf(y,m)],color=:black, label="median")
        else
            p = [p, Plots.plot!([m,m],[0,pdf(y,m)],color=:black, label="median")]
        end
    end
    return p
end
```

Out[54]: npdensity (generic function with 1 method)

In [55]:
```julia
# sample is from exponential, prior is lognormal, proposal is random walk lognormal
y = rand(Exponential(3.0),30)
# set prior, likelihood and proposal
Prior = θ -> pdf(LogNormal(1.0,1.0), θ)
lnL = θ -> sum(logpdf.(Ref(Exponential(θ)),y))
tuning = 0.5
Proposal = θ -> rand(LogNormal(log(θ),tuning))
# get the chain, plot posterior, and descriptive stats
chain = mcmc(1.0, 100000, 10000, Prior, lnL, Proposal) # start value, chain length, and
p = npdensity(chain[:,1]) # nonparametric plot of posterior density
plot!(p, title="posterior density: true value = 3.0") # add a title
display(p)
```
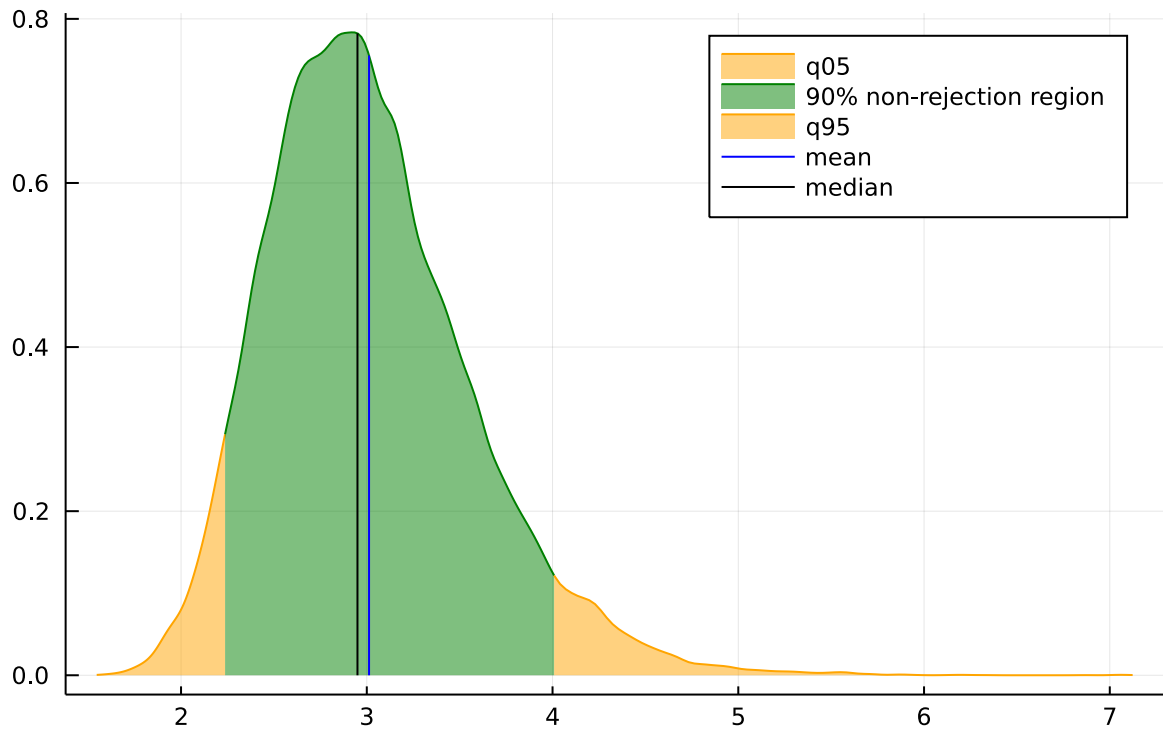
```
current parameters: 3.093
  acceptance rates: 0.396
current parameters: 2.361
  acceptance rates: 0.384
current parameters: 3.531
  acceptance rates: 0.395
current parameters: 2.096
  acceptance rates: 0.4
current parameters: 2.499
  acceptance rates: 0.387
current parameters: 3.673
  acceptance rates: 0.404
current parameters: 2.544
  acceptance rates: 0.392
current parameters: 3.94
  acceptance rates: 0.393
current parameters: 3.139
  acceptance rates: 0.391
```

posterior density: true value = 3.0

current parameters: 3.13
 acceptance rates: 0.394