

Homework 8 Solutions

Due: Wednesday 5/1/19 by 12:00pm (noon)

Note - problem 2. will require use of the `stochvol` package for R.

Grading Scheme:

- Maximum of 2 points for 1., determined as follows:
 - 0 points for no solutions or R output only;
 - 1 point if there are substantial mistakes and serious revisions needed for final project draft;
 - 2 points if mostly or entirely correct and only minor revisions needed for final project draft.
- Maximum of 3 points for problem 2., determined as follows:
 - 0 points for no solutions whatsoever or R output only;
 - 1 point for an honest effort but very few correct answers;
 - 2 points for mostly correct answers but some major mistakes;
 - 3 points for one or fewer mistakes.

1. Exploratory and State-Space Analysis of Project Data

For this problem, I'll ask you to select a data set from the following possibilities:

- Anomaly
 - Electricity
 - Stocks
 - Yields
 - Air
 - Beijing
- (a) Give the name of the dataset you've chosen. You'll have to stick with this dataset for the state-space part of the final project. All but one of the datasets are multivariate. For this problem, I will ask you to analyze the first time series in the dataset you've selected. For instance, the first time series in the `Anomaly` data is given by `Anomaly[, 1]`.
- (b) Plot the raw data.
- (c) All of the data sets have some type of "seasonal" aspect, i.e. they are measured quarterly, monthly, or daily and may have quarter-of-the-year, month-of-the-year, or day-of-the-week effects, respectively. What kind of seasonality might be present in the data you chose?
- (d) Define an $n \times s$ matrix \mathbf{Z} to capture seasonality, where s is the number of units of time per season minus one. Note that you'll have to transpose it when you pass it to the `MARSS` function(s). Fit four linear state-space models to the raw data minus the last 20 observations using the `MARSS` package:
- i. $y_t = ax_t + v_t, x_t = \phi x_{t-1} + w_t$
 - ii. $y_t = ax_t + \gamma' z_t + v_t, x_t = \phi x_{t-1} + w_t$
 - iii. $y_t = ax_t + v_t, x_t = \phi x_{t-1} + \mathbf{v}' z_t + w_t$
 - iv. $y_t = ax_t + \gamma' z_t + v_t, x_t = \phi x_{t-1} + \mathbf{v}' z_t + w_t$

Compute AIC for each, and indicate which model you would choose based on AIC alone.

- (e) Plot the last 40 observations from the raw data, the forecasts of the last 20 observations under each model, and 95% confidence intervals for each.

- (f) Compute the average squared forecast error for the last 20 observations under each of the four models. Indicate which would you choose based on squared forecast error alone.
- (g) In at most one sentence, indicate whether or not you would choose a model based on AIC or squared forecast error and explain why.

2. Stochastic Volatility

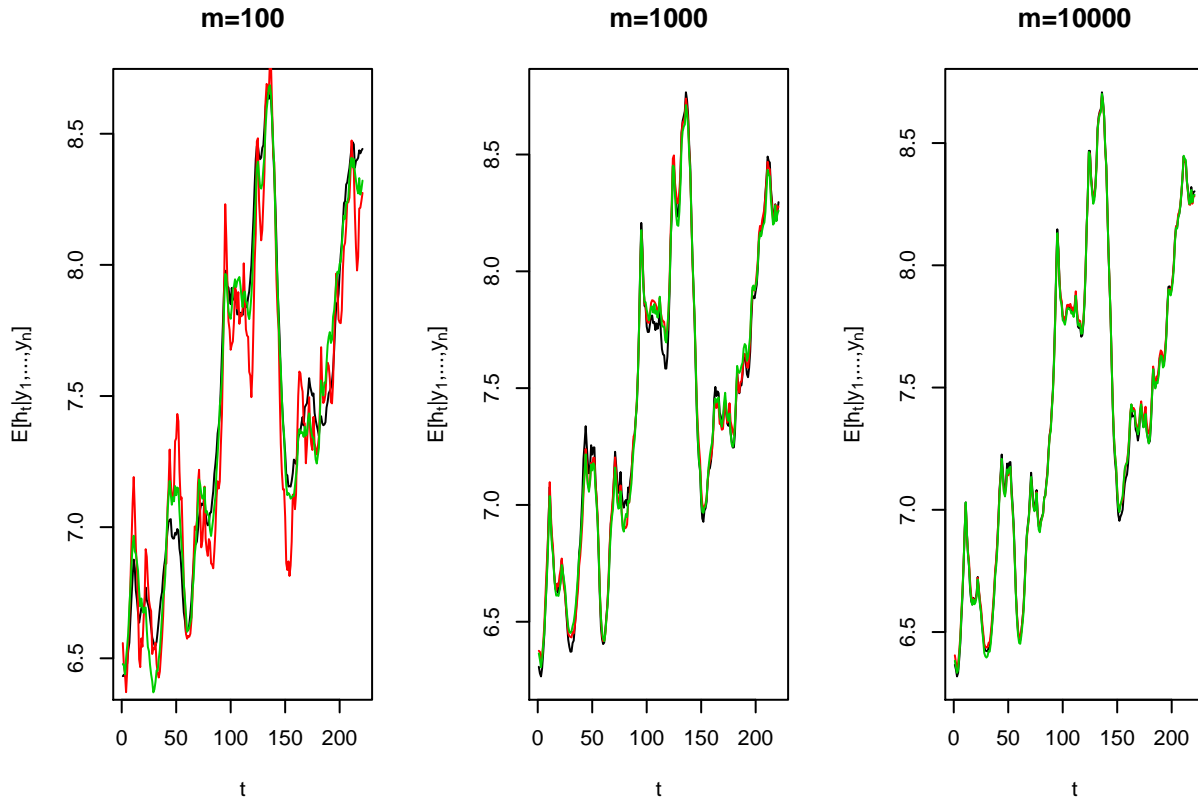
On the second exam, we applied a **GARCH**($m, 0$) model to the second differences of the demeaned **gnp** data. We're going to apply a stochastic volatility model to the same data. You'll want to start with the following code to load the packages we need and the data:

```
library(astsa)
library(stochvol)
data(gnp)
y <- (diff(gnp, d = 2) - mean(diff(gnp, d = 2)))
n <- length(y)
```

- (a) Fit stochastic volatility model with the default prior specifications to **y** using the **svsample** function **three times for each** of the following values of m , the number of simulated values of the states and parameters drawn from the posterior distribution:
 - i. $m = 100$;
 - ii. $m = 1000$;
 - iii. $m = 10000$.

Make a plot with three panels, one for each value of m . In each panel, plot the estimates of the posterior means for the latent states \mathbf{h} for each run of **svsample**. You will have three lines per panel.

```
ms <- c(100, 1000, 10000)
times <- 3
par(mfrow = c(1, 3))
for (m in ms) {
  for (t in 1:times) {
    res <- svsample(y, draws = m)
    if (t == 1) {
      plot(colMeans(res$latent), type = "l", col = t,
           ylab = expression(paste("E[", h[t], "|", y[1], ", ..., ", y[n], "]", sep = "")),
           xlab = "t",
           main = paste("m=", m, sep = ""))
    } else {
      lines(colMeans(res$latent), type = "l", col = t)
    }
  }
}
```



(b) For this data, which value of m seems reasonable to use in practice? Answer in at most one sentence and base your answer on your plots from (a).

I would use $m = 10000$, because that is **certainly** enough simulated values to get a good estimate of the posterior means of the states. If you said that you would use $m = 1000$ and justified it based on computation time, that is also ok.

- (c) Using the value of m you argued for in (b), fit the stochastic volatility models to the data with the last 20 observations held out using the following priors:
- Default priors for μ_h , ϕ and σ_w^2 ;
 - Default priors for ϕ and σ_w^2 , normal prior for μ_h with mean 0 and variance 1;
 - Default priors ϕ and σ_w^2 , normal prior for μ_h with mean 0 and variance 1000000;
 - Default priors for μ_h and σ_w^2 , beta prior for $(\phi + 1)/2$ with $a_0 = 1$ and $b_0 = 1$;
 - Default priors for μ_h and σ_w^2 , beta prior for $(\phi + 1)/2$ with $a_0 = 10$ and $b_0 = 10$;
 - Default priors for μ_h and ϕ , gamma prior for σ_w^2 with shape 1/2 and rate 1/20.
 - Default priors for μ_h and ϕ , gamma prior for σ_w^2 with shape 1/2 and rate 1/0.02.

Plot kernel density estimates of $p(\mu_h | \mathbf{y})$ from i.-iii. in the first panel, kernel density estimates of $p(\phi | \mathbf{y})$ from i., iv.-v. in the second panel, and $p(\sigma_w^2 | \mathbf{y})$ from i., vi.-vii. in the last panel. Kernel density estimates can be computed using the `density` function in R applied to simulated values of the corresponding parameter returned by `svsample`.

```
par(mfrow = c(1, 3))
res <- svsample(y[1:(n-20)], draws = 10000, priormu = c(0, 100))
plot(density(res$para[, 1]),
     main = "", xlab = expression(mu[h]),
     ylab = expression(paste("p(", mu[h], "|", y[1],
                             ", ..., ", y[n], ")")),
     xlim = c(-10, 20))
res <- svsample(y[1:(n-20)], draws = 10000, priormu = c(0, 1))
```

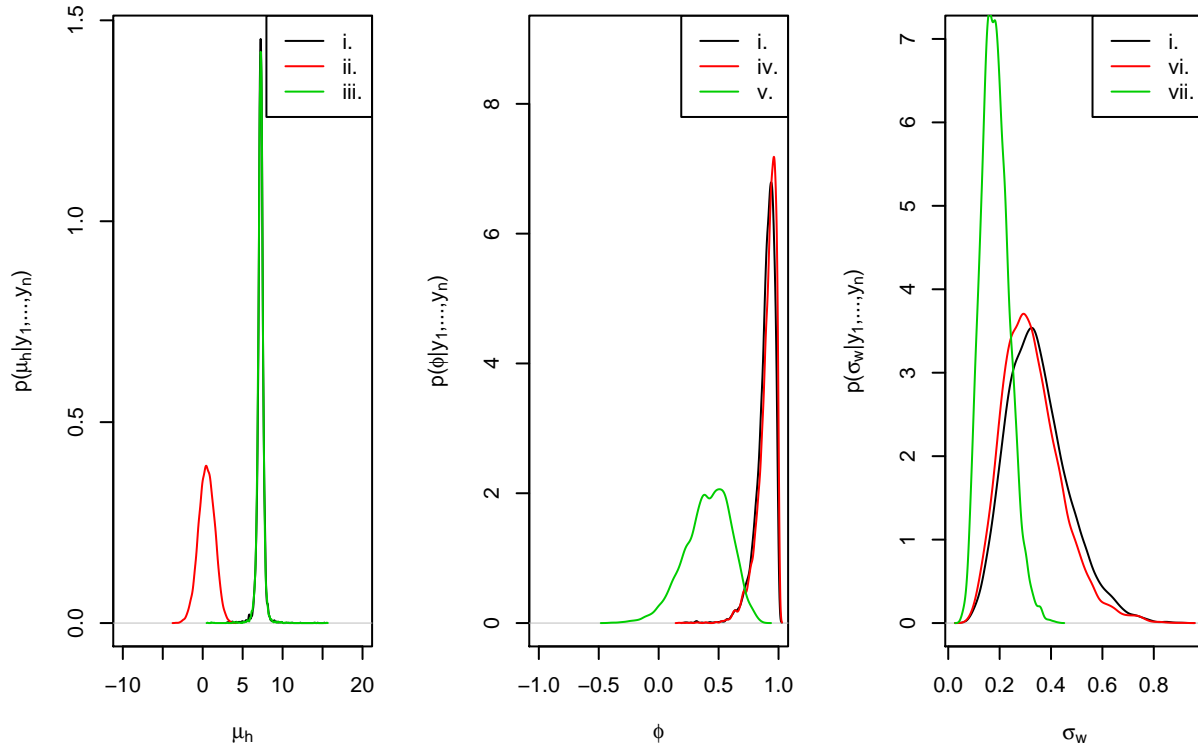
```

lines(density(res$para[, 1]), col = 2)
res <- svsample(y[1:(n-20)], draws = 10000, priormu = c(0, 1000))
lines(density(res$para[, 1]), col = 3)
legend("topright", lty = c(1, 1, 1),
      col = c(1, 2, 3), legend = c("i.", "ii.", "iii."))

res <- svsample(y[1:(n-20)], draws = 10000, priorphi = c(5, 1.5))
plot(density(res$para[, 2]), main = "", xlab = expression(phi),
     ylab = expression(paste("p(", phi, "|", y[1],
                             "...", y[n], ")"), sep =)),
     xlim = c(-1, 1),
     ylim = c(0, 9))
res <- svsample(y[1:(n-20)], draws = 10000, priorphi = c(1, 1))
lines(density(res$para[, 2]), col = 2)
res <- svsample(y[1:(n-20)], draws = 10000, priorphi = c(10, 10))
lines(density(res$para[, 2]), col = 3)
legend("topright", lty = c(1, 1, 1),
      col = c(1, 2, 3), legend = c("i.", "iv.", "v."))

res <- svsample(y[1:(n-20)], draws = 10000, priorsigma = 1)
plot(density(res$para[, 3]),
     main = "", xlab = expression(sigma[w]),
     ylab = expression(paste("p(", sigma[w], "|",
                             y[1], "...", y[n], ")"),
                             sep =)),
     ylim = c(0, 7))
res <- svsample(y[1:(n-20)], draws = 10000, priorsigma = 10)
lines(density(res$para[, 3]), col = 2)
res <- svsample(y[1:(n-20)], draws = 10000, priorsigma = 0.01)
lines(density(res$para[, 3]), col = 3)
legend("topright", lty = c(1, 1, 1),
      col = c(1, 2, 3), legend = c("i.", "vi.", "vii."))

```



- (d) Give the average squared forecast error for the last 20 observations for all of the models fit in (c) in a table. For your forecasts, use the average simulated value of each future y_{n+k} , which can be obtained using the the `predict` function.

```

mses <- c()
draws <- 10000
res <- svsample(y[1:(n-20)], draws = draws)
pred <- predict(res, 20)
mses <- c(mses, mean((y[(n - 20 + 1):n] - colMeans(pred$y))^2))
res <- svsample(y, draws = draws, priormu = c(0, 1))
pred <- predict(res, 20)
mses <- c(mses, mean((y[(n - 20 + 1):n] - colMeans(pred$y))^2))
res <- svsample(y, draws = draws, priormu = c(0, 10000))
pred <- predict(res, 20)
mses <- c(mses, mean((y[(n - 20 + 1):n] - colMeans(pred$y))^2))

res <- svsample(y, draws = draws, priorphi = c(1, 1))
pred <- predict(res, 20)
mses <- c(mses, mean((y[(n - 20 + 1):n] - colMeans(pred$y))^2))
res <- svsample(y, draws = draws, priorphi = c(10, 10))
pred <- predict(res, 20)
mses <- c(mses, mean((y[(n - 20 + 1):n] - colMeans(pred$y))^2))

res <- svsample(y, draws = draws, priorsigma = 10)
pred <- predict(res, 20)
mses <- c(mses, mean((y[(n - 20 + 1):n] - colMeans(pred$y))^2))
res <- svsample(y, draws = draws, priorsigma = 0.01)
pred <- predict(res, 20)
mses <- c(mses, mean((y[(n - 20 + 1):n] - colMeans(pred$y))^2))

```

Model	Average MSE
i.	4651.266
ii.	4689.445
iii.	4669.992
iv.	4629.529
v.	4661.115
vi.	4636.418
vii.	4622.868

Full credit was given to students who got similar but not exactly the same numbers here - because these are approximations computed from simulations, they may differ a bit depending on the seed used (which means we should have used even more simulated values).

- (e) Based on what you observe (c) and (d), explain in one sentence which prior specification(s) you prefer. You don't have to choose a single one, but you should comment on whether or not any seem like especially good or bad choices.

I would prefer the prior specifications used in models i., iii., iv., and vi. because they are uninformative (have high variance), and give posterior distributions for the parameters that reflect the data well, and I would not use the forecast mean squared error to compare models because they are all very large and indicate that the stochastic volatility model does not help us forecast values of y_t , regardless of which priors are chosen.