

Proof that Gibbs Sampler Will Converge:

Part B

Let $p(\underline{x})$ denote the probability density function of a vector of random variables x_1, \dots, x_p . We will assume that the Markov Chain that results from running the Gibbs Sampler is ergodic. We want to show that the stationary distribution of the p -component Gibbs sampler is indeed this target distribution $p(\underline{x})$.

To avoid having to write a bunch of stuff in L^AT_EX, let me introduce a couple pieces of notational shorthand. Let x^t denote the t -th iteration of the Gibbs Sampler (**not** x to the power t !). Additionally, $x_{[i:j]} = (x_i, x_{i+1}, \dots, x_{j-1}, x_j)$, just like the slice operation in Python or the colon function in R.

By the assumption of ergodicity, we know that the Markov Chain formed by the Gibbs Sampler, $\underline{x}^1, \underline{x}^2, \dots$, converges to a stationary distribution, denoted by $\pi(\underline{x})$ in the notes. In the continuous state space case, it is a property of a stationary Markov Chain that

$$\pi(\underline{y}) = \int_{\mathbb{R}^p} \pi(\underline{x}) p(\underline{x}, \underline{y}) d\underline{x}.$$

That is, we must have that the probability of arriving at state \underline{y} is the integral of the probability of being at state \underline{x} times the probability of moving from state \underline{x} to state \underline{y} .

This will require only two steps: deriving a closed form for the transition density $p(\underline{x}, \underline{y})$ and showing that the (multivariate) integral of the stationary distribution times the transition density gives us the stationary distribution. We will assert that one such stationary distribution that satisfies this property is the target distribution, then use the uniqueness of the stationary distribution to claim the target distribution is the unique stationary distribution.

Transition Density: What is the density of \underline{x}^{t+1} , given \underline{x}^t ? For this, by conditional expectation,

$$\begin{aligned} p(\underline{x}^{t+1} | \underline{x}^t) &= p(x_1^{t+1}, x_2^{t+1}, \dots, x_p^{t+1} | \underline{x}^t) \\ &= p(x_1^{t+1} | \underline{x}^t) \cdot p(x_2^{t+1}, \dots, x_p^{t+1} | x_1^{t+1}, \underline{x}_{[2:p]}^t) \\ &= p(x_1^{t+1} | \underline{x}^t) \cdot p(x_2^{t+1} | x_1^{t+1}, \underline{x}_{[2:p]}^t) \cdot p(x_3^{t+1}, \dots, x_p^{t+1} | x_{[1:2]}^{t+1}, \underline{x}_{[3:p]}^t) \\ &= \dots \text{(Repeated conditioning, unrolling each term conditioned on all the rest,)} \\ &= p(x_1^{t+1} | \underline{x}^t) \cdot p(x_2^{t+1} | x_1^{t+1}, \underline{x}_{[2:p]}^t) \dots p(x_p^{t+1} | \underline{x}_{[1:p-1]}^{t+1}, x_p^t) \end{aligned}$$

Intuitively, you can think of this as similar to what the Gibbs Sampler does. The distribution of the joint is the product of the conditionals of each random variable at the current iteration,

given the variables already drawn in the current iteration and the previous iteration for the variables not yet drawn in the current iteration.

Stationarity: The target distribution $p(\underline{x})$ is stationary iff

$$\begin{aligned} p(\underline{x}^{t+1}) &= \int_{\mathbb{R}^p} p(\underline{x}^t) p(\underline{x}^{t+1} | \underline{x}^t) d\underline{x}^t, \forall \underline{x}^t \in \mathbb{R}^p. \\ &= \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} p(x_1^t, x_2^t, \dots, x_p^t) p(x_1^{t+1}, x_2^{t+1}, \dots, x_p^{t+1} | x_1^t, x_2^t, \dots, x_p^t) dx_1^t dx_2^t \cdots dx_p^t, \end{aligned}$$

for all $x_i^t \in \mathbb{R}$, $i = 1, \dots, p$. Recall that $f(y) = \int_{\mathbb{R}} f(x, y) dx = \int_{\mathbb{R}} f(x) \cdot f(y|x) dx$ for arbitrary real-valued random variables y and x . The above integral can be seen to be the product of a marginal of \underline{x}^t and the conditional of \underline{x}^{t+1} given \underline{x}^t , integrated over the entire parameter space of \underline{x}^t . By similar logic to my univariate example earlier in the paragraph, we see that the result of this will be the marginal distribution of \underline{x}^{t+1} .

$$\Rightarrow \int_{\mathbb{R}^p} p(\underline{x}^t) p(\underline{x}^{t+1} | \underline{x}^t) d\underline{x}^t = p(\underline{x}^{t+1}).$$

Combining the uniqueness of the stationary distribution with the stationarity of the target distribution, we have that the Gibbs Sampler Markov Chain converges to the target distribution. ■

Part A

First, a story: An engineer and a mathematician shared an apartment. Their kitchen was equipped with an electric stove, and every morning someone had placed a pot of water on the back-right burner so they could make coffee. They both knew what knob turned on this burner. One morning the engineer came into the kitchen and found the pot was on the front-left burner. He got out the stove's schematics and followed the wiring diagram and finally figured out which knob turned on this burner and he then used that knob and made the coffee. The next morning the mathematician came in and also found the pot on the front-left burner. He moved the pot to the back-right burner, thereby reducing the problem to one which he had already solved.

Proof: Let $p = 2$. Refer to Part B. ■

Metropolis-Hastings with Simulated 2-dimensional Parameter

Being unable to directly sample the posterior distribution acquired by a multivariate normal prior and a binomial likelihood, we are reduced to more computational methods. Assume our likelihood and prior are

$$\begin{aligned} \beta &\sim N(\mu_0, \Sigma_0) \\ y_i | \beta &\sim \text{Bin}(m_i, \text{expit}(x_i^T \beta)), \quad i = 1, \dots, n, \end{aligned}$$

where $\text{expit}(u) = \exp(u)/(1 + \exp(u))$. For tractability, assume the responses are conditionally independent given β .

The posterior is thus

$$\begin{aligned}
f(\underline{\beta}|\underline{y}) &\propto f(\underline{\beta}) \cdot \prod_{i=1}^n f(y_i|\beta) \quad (\text{by the cond. indep. of } y_i \text{ and Bayes' Theorem}) \\
&= N(\mu_0, \Sigma_0) \cdot \prod_{i=1}^n \text{Bin}(m_i, \text{expit}(x_i^T \beta)) \\
&\propto \exp\left(-\frac{1}{2}(\beta - \mu_0)^T \Sigma_0^{-1}(\beta - \mu_0)\right) \prod_{i=1}^n \binom{m_i}{y_i} \text{expit}(x_i^T \beta)^{y_i} (1 - \text{expit}(x_i^T \beta))^{m_i - y_i} \\
&\propto \exp\left(-\frac{1}{2}(\beta - \mu_0)^T \Sigma_0^{-1}(\beta - \mu_0)\right) \prod_{i=1}^n \frac{\exp(x_i^T \beta)^{y_i}}{1 + \exp(x_i^T \beta)^{y_i}} \frac{1}{1 + \exp(x_i^T \beta)^{m_i - y_i}} \\
&= \exp\left(-\frac{1}{2}(\beta - \mu_0)^T \Sigma_0^{-1}(\beta - \mu_0)\right) \prod_{i=1}^n \frac{\exp(x_i^T \beta)^{y_i}}{1 + \exp(x_i^T \beta)^{y_i}} \frac{1 + \exp(x_i^T \beta)^{y_i}}{1 + \exp(x_i^T \beta)^{m_i}} \\
&= \exp\left(-\frac{1}{2}(\beta - \mu_0)^T \Sigma_0^{-1}(\beta - \mu_0)\right) \prod_{i=1}^n \frac{\exp(x_i^T \beta)^{y_i}}{1 + \exp(x_i^T \beta)^{m_i}}
\end{aligned}$$

This can be simplified further, but there's not much point. This distribution has no closed form in β . We cannot directly sample from it, so we must apply Metropolis-Hastings.

1. Specify a starting point for the Markov Chain, β^0 .
2. At each increment i in the Markov Chain, propose a new value, β^* , centered on the previous one, β^i .
3. Compute the probability of the posterior at the candidate value $p(\beta^*|\underline{y})$ (derived up to proportionality above), as well as the posterior at the current value $p(\beta^i|\underline{y})$. In R, we work up to proportionality, so this is as simple as multiplying the prior by the likelihood, `dmvnorm` by `dbinom`. For stability reasons, we work on the log scale, so we'll add the log-prior and the log-likelihood to get the log-posterior.
4. Compute $\alpha^* = \frac{p(\beta^*|\underline{y})}{p(\beta^i|\underline{y})}$.
5. Draw $U \sim U(0, 1)$. If $U \leq \alpha^*$, set $\beta^{i+1} = \beta^*$. Else, $\beta^{i+1} = \beta^i$. Again, because of the log-scale, we can compare a log-uniform against the difference of the log-posterior at the candidate to the log-posterior at the current value of the Markov Chain, and the results will be equivalent with a greatly-reduced concern for numerical inaccuracy.
6. The vector $\beta^i, i = 1, \dots, n$ represents dependent draws from the target distribution $p(\beta|\underline{y})$.

This algorithm has a ton of knobs and switches to play with. The following summarizes some of the parameters I use to make the algorithm more efficient.

Starting Values: $\beta^0 = \hat{\beta}$, the MLE of the data, was chosen. We can think of this as being close to the posterior mode if we had an uninformative prior distribution. It can be shown that as the sample size increases, the effect of the prior diminishes for any non-degenerate prior, so for our sample size, the MLE will be quite close to the posterior mode. Using the MLE is a sensible choice if the sample size is large enough that the likelihood carries more weight in the posterior than the prior does.

Proposal Distribution The proposal is generated from the proposal distribution, which I chose to be a bivariate normal, centered on the previous value, with a covariance matrix $v^2\Sigma$, where v^2 is a tuning parameter that changes to give a suitable acceptance rate for samples. Σ is the estimated covariance matrix using a Frequentist logistic regression, which, for large enough sample size, closely models the structure of the posterior covariance. The reason for choosing a bivariate normal as my proposal distribution was out of simplicity. Using the bivariate normal, the probability of going from the candidate value to the current value is the same as going from current to candidate, a property called symmetry. This simplifies the calculation of the Metropolis-Hastings acceptance statistic. The choice of that strange form for the covariance matrix is an attempt at mimicking the variance of the posterior distribution. I acquired this matrix by fitting a logistic regression in R with `glm`, then using its estimate for the covariance as my proxy for the posterior variance.

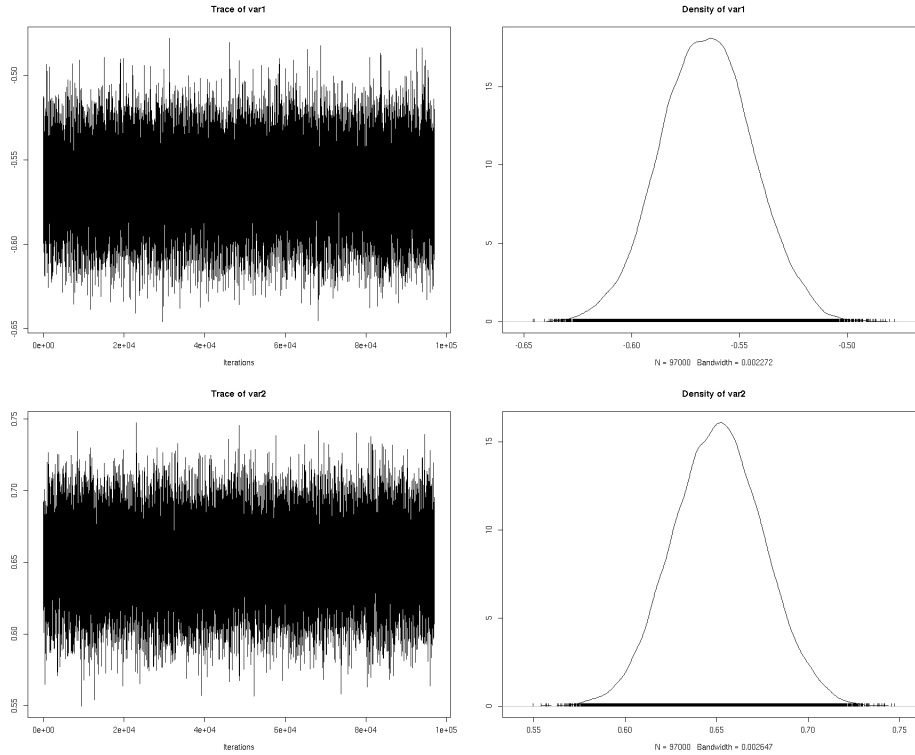
Tuning Process Every so often, we will go through a retuning process. It is possible that the proposal distribution does not have an appropriate variance. Too high a variance and we will propose bad values, rejecting too often and ending up with lots of redundant draws. If the variance is too small, we will stay in relatively the same spot, not exploring the space effectively, accepting too often. In both cases, we end up with a very small effective sample size, and gaining very little additional information on each draw. The solution to this problem is to keep track of our acceptance rates and make sure it stays within a good acceptance rate window. I chose 20% to 50% for this 2-dimensional example. This comes based on heuristic moreso than any analytical reason. This ended up being pretty reasonable, so I stuck with it. If the acceptance rate fell outside this band, I adjusted the proposal variance by a scalar factor of $v_i^2 = \exp(3 \cdot (\text{acceptance} - \text{ideal}))$, where v_i^2 is the scalar proposed during this tuning period (the previous tuning period times the scaling factor). Acceptance is the acceptance ratio, and ideal is the ideal accept rate. For the two-dimensional case, I chose 30%, which is a reasonable balance between accepting and rejecting. If the acceptance rate is higher than 30%, this scaling will increase the proposal variance (to make more adventurous jumps). If the acceptance is less, the scaling will decrease the variance. The coefficient of 3 was chosen as a factor of how much to scale by each time and was chosen by hand-tuning heuristically.

Number of Iterations Computing power is inexpensive in 2013. As long as I have other things to do, I would rather have the computer spend two hours running my code than have to spend 30 minutes of my time optimizing the program to make it run quicker.

Consistent with that philosophy, I spend very little time choosing a decent number of iterations. A couple million samples takes less than 10 minutes on my Linux machine, with a retuning every thousand observations up to the burn-in.

Burn-In Period Like the number of iterations, the burn-in is somewhat of a non-issue. I cleverly chose to start at the MLE of the likelihood, which means my Markov Chain was pointed at the correct position from the very beginning. I could've chosen a very small burn-in if I wanted, but I also needed to adjust the proposal variance, so I set the burn-in to be around 10,000. If I could've adjusted the variance on-the-fly *and* kept those draws instead of burning them, I could've made the burn-in period far less. I hit stationarity after less than a thousand draws most times.

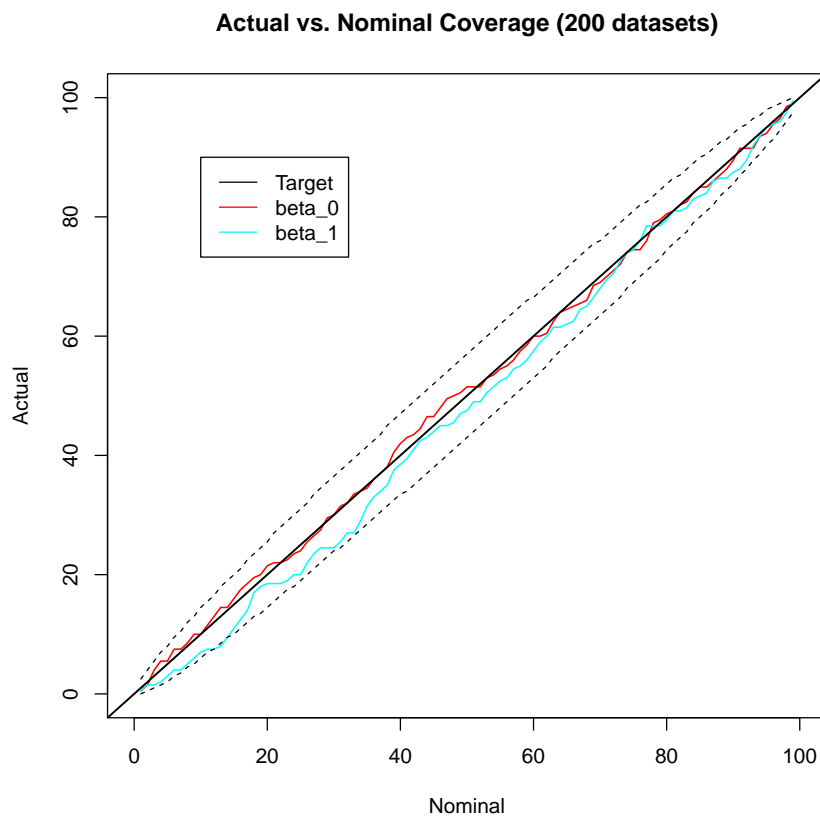
For this simple 2-dimensional problem, Metropolis-Hastings makes short work of sampling the posterior. While the effective sample size is only about 17% of the total sample size, this is good enough if we run enough iterations of the algorithm. We can thin the data (taking only every k^{th} iteration) if we are worried about autocorrelation between samples, but I didn't bother for only two dimensions. Looking at the traceplots, we see that stationarity has been attained.



Trace and density plots for (β_1, β_2) , $n = 10^5$, burn = 3000, retune = 100.

A similar picture can be seen on the other 200 datasets generated on the compute cluster. On the topic of coverage probabilities, you may take a look at the coverage summaries.

Things turned out okay. Keep in mind that there is randomness associated with the Markov Chain. Running the simulation another time changes the coverage probabilities.



Coverage Line Plot for (β_1, β_2) , $n = 5 \cdot 10^5$, burn = 10^4 , retune = 10^3 .

	β_0	β_1
p_01	0.01	0.01
p_05	0.06	0.03
p_10	0.10	0.07
p_25	0.24	0.20
p_50	0.52	0.47
p_75	0.74	0.74
p_90	0.90	0.88
p_95	0.94	0.95
p_99	0.99	0.99

Coverage Probabilities for (β_0, β_1) .

Metropolis-Hastings with 11-dimensional Parameter on Real Data

The previous approach did not work very well on this particular problem (wait for the punchline after the table). I tried do it in the same fashion, with 100,000 iterations, a burn-in of 20,000, an acceptance range between 10% and 30% to accommodate the more frequent number of rejections, the algorithm finished quickly, but the effective sample size for the Texture variable was not much larger than 5. In the interest of not re-writing my Metropolis-Hastings code from Problem 2, I ramped up the iterations to 5,000,000 (*five million!*). The effective sample size of the Texture variable (X11) was still only 10.

The dimension of the parameter space is partially to blame. It's difficult to propose plausible values in 11-dimensional space and have them be more likely than another previously-selected value. But the "Curse of Dimensionality" cannot fully explain the difficulty of sampling the β s in this problem. Looking at the covariance matrix of the 11 parameters, we can see where we might have trouble.

β_1	β_2	β_3	β_4	β_5	β_6	β_7	β_8	β_9	β_{10}	β_{11}
165.19	0.0003	413.82	813.91	65.94	7319.94	0.26	13.81	1021.12	113.01	0.004

Variances for MLEs $\hat{\beta}_i$, $i = 1, \dots, 11$.

The scale is simply ridiculous! Variables X2 (area), X7 (perimeter), and X11 (texture) are on the order of tenths to thousandths, while X6 (fracdim) and X9 (smoothness) are in the thousands. Between the smallest-varied and largest-varied variables, we're looking at a six order of magnitude difference! This makes it very hard to tune our proposal distribution. However, a clever choice of proposal distribution solves a number of problems. I was initially unable to get the MH algorithm to converge correctly with this complicated space. Using the estimated covariance matrix created from a logistic regression solved a large part of the problem, but was still not enough until about 10pm sitting in my office, I realized I had been using the wrong prior ($\frac{1}{1000} \cdot I_{11} \neq 1000 \cdot I_{11}$)! Changing the code drastically increased my effective sample size to the point that Metropolis-Hastings became a viable option for sampling an 11-dimensional parameter space. The convergence was not as nice as the 2-dimensional case, but people who don't want to spend extra time re-factoring optimized code can't be choosers. Running the chain a long time handily solves this problem anyway. Let's talk about the autocorrelation. Using built-in R function `acf()`, computing the lag-1 autocorrelation is trivial.

β_1	β_2	β_3	β_4	β_5	β_6	β_7	β_8	β_9	β_{10}	β_{11}
0.95	0.96	0.94	0.93	0.96	0.83	0.96	0.96	0.93	0.96	0.96

The Lag-1 autocorrelation for each β_i .

Terrible. Each draw is very correlated to the previous draw. One simple and slightly naive

way to handle this problem is to thin the samples, removing every k^{th} observation. Thinning will reduce the effective sample size, but since I ran the chain for half a million iterations, I can afford to lose 98% of my Markov Chain iterations to cut down the autocorrelation.

β_1	β_2	β_3	β_4	β_5	β_6	β_7	β_8	β_9	β_{10}	β_{11}
0.16	0.19	0.12	0.05	0.14	0.05	0.20	0.18	0.01	0.20	0.14

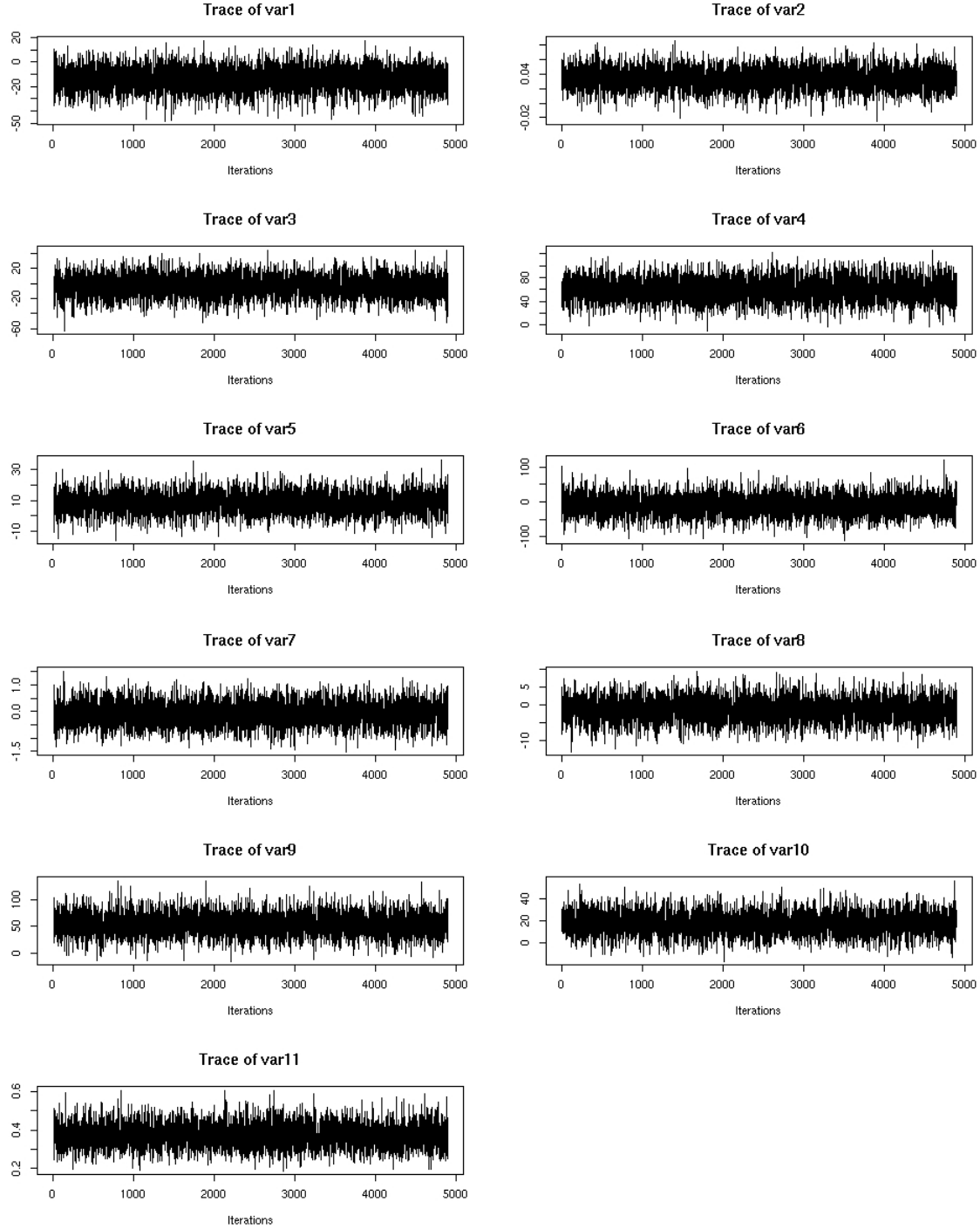
The Lag-1 autocorrelation for each β_i using 50-thinning.

β_1	β_2	β_3	β_4	β_5	β_6	β_7	β_8	β_9	β_{10}	β_{11}
0.08	0.05	-0.04	-0.02	0.08	-0.05	0.02	0.02	-0.01	0.08	-0.01

The Lag-1 autocorrelation for each β_i using 100-thinning.

The 100-thinned version is basically uncorrelated at this point. The effective sample size demonstrated this, as most variables had effective sample sizes that were very close to actual number of thinned samples. For my large iteration choice, this yielded 5,000 uncorrelated samples from the posterior, which was good enough for me. In hindsight, I could've been fine running 50,000 iterations and still had 500 uncorrelated samples.

Convince yourself from my traceplots that the chain has good convergence.



Trace and density plots for $(\beta_1, \dots, \beta_{11})$, $n = 5 \cdot 10^5$, burn = 10^4 , retune = 100.

Next, we consider the problem of variable selection. Bayesian variable selection works in a different way from Frequentist variable selection, since we don't have the same kind of

Neyman-Pearson formulation that we did before—the concept of controlling a Type 1 error for a model isn’t as meaningful. One approach is simply to look at the posterior credible intervals and see which ones contain zero, because while the magnitude of each posterior parameter does depend on the scaling of each variable, whether the posterior intervals contain zero is only very weakly dependent on the scaling. From the table, it seems the most important variables (the ones that don’t contain zero) are β_4 , β_9 , and β_{11} (concave points, smoothness, and texture, respectively).

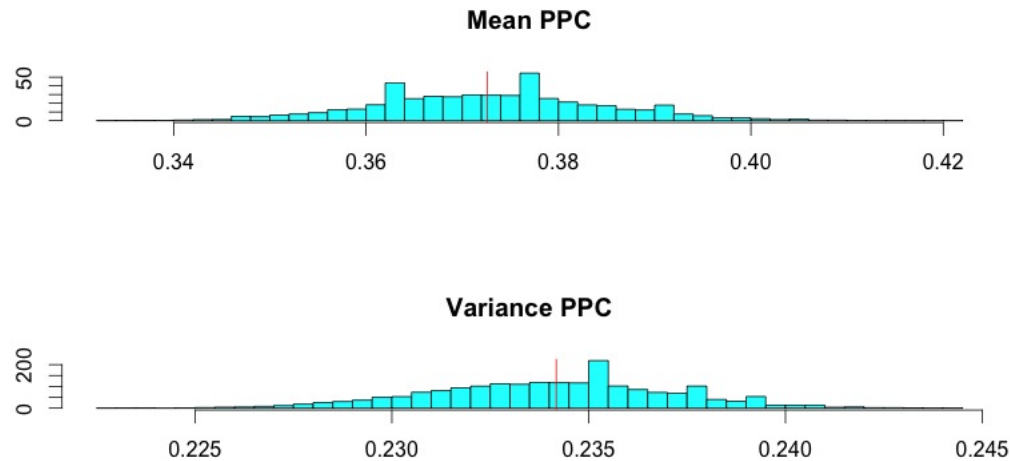
	β_1	β_2	β_3	β_4	β_5	β_6	β_7	β_8	β_9	β_{10}	β_{11}
2.5%	-36.31	0.00	-30.85	19.69	-6.33	-67.32	-0.85	-8.15	15.52	-0.66	0.26
97.5%	4.10	0.06	22.50	98.34	23.22	51.30	0.84	4.49	87.96	36.64	0.49

95% Central Credible Interval

Finally, I will talk about checking our model. Model checking in Bayesian inference is done through posterior predictive checking. This procedure involves drawing parameters $\underline{\beta}$ from the posterior distribution, then generating simulated data from the assumed likelihood, given those posterior values as the parameters. If the model is an accurate representation of the data, the posterior predictive distributions should not be all that different from the true data, which we can corroborate using summary statistics (mean, median, variance, and the like). Our algorithm is as follows:

1. Draw $\underline{\beta}^i$ from the posterior, $p(\underline{\beta}|y)$, for $i = 1, \dots, M$.
2. Simulate M datasets \tilde{y}^i , where $\tilde{y}^i \sim \text{Bin}(m_i, \text{expit}(X^T \beta^i))$.
3. Compute a summary statistic on each dataset, $T(\tilde{y}^i)$, then create a histogram with the M summary statistics.

If the model is correct, the true summary statistic should sit somewhere near the center of the histogram (or at least not in the far tails of the histogram). Since the likelihood is Bernoulli data, minimum and maximum statistics are not useful. I went with the mean and the variance, and we can see from the plots that the binomial likelihood model is quite sensible in this case.



Posterior Predictive Checks for the Mean and Variance

The red line indicates the true mean and variance. The mean was chosen as a measure of center, and the variance was the choice for measuring dispersion. There is an inflexible relationship between the mean and the variance with the binomial, so if another model was more appropriate, checking both variance and mean could have given some insight. As it turns out, the lines sit very closely to the middle of both histograms, leading us to believe that the model is quite sound as a way of modeling both the mean and the variance. It is possible that the model is not a good fit through some other summary statistic. A more thorough method would be to look at a number of quantiles as the summary statistics. This provides a very accurate picture of both the tails and the center of the model. Similar to the QQ Plot, the posterior predictive quantiles ought to match up against the true corresponding quantiles if the model is valid (think of this as as many histograms as there are quantiles). In the interest of parsimony, I chose only a measure of center and a measure of dispersion.

R Source Code

Helper_functions is a script containing the majority of the algorithms I used to do Bayesian Logistic regression. It contains basic functions, as well as my Metropolis Hastings implementation, computation of quantiles, and posterior predictive checks.

```

1 expit = function(X, Beta){
2   exp(X %*% Beta) / {1+exp(X %*% Beta)}
3 }
4 expit = cmpfun(expit)
5
6 #Compute posterior of MVN prior and Binom likelihood. Return log posterior.
7 posterior.prob = function(Beta, y, beta.0, sigma.0, m, X){
8   #Computes sum of Pr(Beta)^N(beta.0, sigma)
9   logprior = dmnorm(Beta, mean=beta.0, sigma=sigma.0, log=TRUE)

```

```

10   loglikelihood = sum(dbinom(y, size=m, prob=expit(X,Beta), log=TRUE))
11   logprior+loglikelihood
12 }
13 posterior.prob = cmpfun(posterior.prob)
14
15 MH = function(current, candidate, uniform, m, y, sigma.0, beta.0, X){
16   #Computes the Metropolis-Hastings criterion. If the criterion exceeds a uniform (everything logged, btw)
17   ,
18   #return the new value. Otherwise, return the old one.
19   r.top = posterior.prob(candidate, y, beta.0, sigma.0, m, X)
20   r.bottom = posterior.prob(current, y, beta.0, sigma.0, m, X)
21   r = r.top - r.bottom
22   if (uniform < r) return(list(sample=candidate, accept=TRUE))
23   else return( list(sample=current, accept=FALSE) )
24 }
25 MH = cmpfun(MH)
26
27 "bayes.logreg.samples" <- function(m, y, X, beta.0=NA, sigma.0=NA, niter=10000, burnin=1000, print.every
28   =1000, retune=100, verbose=TRUE)
29 {
30   p = ncol(X)
31   n = nrow(X)
32   if (all(is.na(beta.0))) beta.0 = rep(0, p)
33   if (all(is.na(sigma.0))) sigma.0 = diag(p)
34   Sigma.0.inv = solve(sigma.0) #Compute the inverse of precision--covariance. Speeds up inner loop
35
36   beta.samples = matrix(NA, nrow=p, ncol=niter)
37   acceptances = rep(NA, niter)
38   acceptances[1] = TRUE
39
40   uniform = log(runif(niter)) #Generate all uniforms at once, to utilize R's vectorization.
41   logist_reg = glm(cbind(y, m-y) ~ X-1, family="binomial") #Run a logistic regression on the data.
42   beta.samples[,1] = logist_reg$coefficients #Start at the Logistic Regression MLE
43   covar = vcov(logist_reg) #Covariance of Frequentist logistic regression. To be scaled by a constant
44   later.
45
46   #Used to adjust the variance in the burn-in period. Jump less if space is large.
47   var.scale = ifelse(p > 3, .01, 1)
48
49   #Generate a lot of proposal jumps at once. If  $Z \sim N(0, S)$ , then  $Z + b \sim N(b, S)$  for const  $b$ !
50   #Vectorization makes things fast, yo.
51   perturbations = matrix(nrow=p, ncol=niter)
52   perturbations[,1:retune] = t(rmvnorm(retune, mean=rep(0,p), sigma=var.scale * covar))
53
54   #If parameter space is big (p>3), I would expect to be rejected more.
55   low = ifelse(p>3, .10, .20)
56   high = ifelse(p>3, .30, .50)
57   ideal = ifelse(p>3, .20, .30)
58
59   for(l in 1:(niter-1)) {
60     #Retune loop:
61     if (l %% retune == 0 & l <= burnin) {
62       acceptance = mean(acceptances[(l-retune):(l-1)])
63
64       if (acceptance <= low | acceptance >= high) {
65         scaling = exp(3 * (acceptance-ideal))
66         var.scale = var.scale * scaling #If acceptance is greater than ideal, scale it down.
67         if (verbose) {
68           print(paste("Acceptance rate in iterations", l-retune+1, "through", l, "was", round(acceptance
69             ,3), sep=" "))
70           print(paste("Scaling variance by:", round(scaling, 3), sep=" "))
71         }
72       }
73     }
74     perturbations[, (l+1):(burnin+1)] = t(rmvnorm( burnin , mean=rep(0,p), sigma=var.scale * covar))
75   }
76
77   #First legit obs loop:
78   if ( l == (burnin+1) ) {
79     #Generate all proposal perturbations by the final variance, then proceed as usual.
80     if (verbose) print(paste("Starting the first non-burned observation at l=", l, sep=" "))
81     perturbations[, 1:niter] = t(rmvnorm( niter-l+1 , mean=rep(0,p), sigma=var.scale * covar))
82   }
83
84   current = beta.samples[,1]

```

```

81     candidate = current+perturbations[,1]
82     foo = MH(current, candidate, uniform[1], m, y, sigma.0, beta.0, X)
83     beta.samples[,1+1] = foo$sample
84     acceptances[1+1] = foo$accept
85
86     if (1 %% print.every == 0) print(paste("Current iteration:", 1, sep=" "))
87   }
88   beta.samples.burned = t(beta.samples[, (burnin+1):niter])
89 }
90 bayes.logreg.samples = cmpfun(bayes.logreg.samples)
91
92
93 "bayes.logreg" <- function(m, y, X, beta.0=NA, sigma.0.inv=NA, niter=10000, burnin=1000, print.every=1000,
94   retune=100, verbose=TRUE, thin=1) {
95   # m: Vector containing the number of trials for each observation (of length n)
96   # y: Vector containing the number of successes for each observation (of length n)
97   # X: Design matrix (of dimension n p)
98   # beta.0: Prior mean for beta (of length p)
99   # Sigma.0.inv: Prior precision (inverse covariance) matrix for beta (of dimension p p)
100  # niter: Number of iterations to run the MCMC after the burnin period
101  # burnin: Number of iterations for the burnin period (draws will not be saved)
102  # print.every: Print an update to the user after every period of this many iterations
103  # retune: Retune the proposal parameters every return iterations. No tuning should be done after the
104    burnin period is completed
105  # verbose: If TRUE then print lots of debugging output, else be silent
106  # thin: integer -- Take every "thin"th observation from the chain.
107  if (all(!is.na(sigma.0.inv))) {sigma.0 = solve(sigma.0.inv)}
108  samples = bayes.logreg.samples(m, y, X, beta.0, sigma.0, niter, burnin, print.every, retune, verbose)
109  Posterior_Draws = samples[seq.int(1, (niter-burnin), by=thin), ]
110
111  cbind(quantile(Posterior_Draws[,1], probs=seq(.01, .99, length=99)),
112    quantile(Posterior_Draws[,2], probs=seq(.01,.99, length=99)))
113 }
114
115 PPC.Graphs = function(Y, X, m=1E5, mi=rep(1,length(Y)), posterior=chain_thinned){
116   #Y: Response variable
117   #X: Design matrix (explanatory variables)
118   #m: Number of PPC datasets to generate.
119   #mi: Number of trials for each level of the design matrix.
120   #posterior: MCMC object or matrix (n x p) of posterior draws.
121
122   n = length(Y)
123   niter = nrow(posterior) #Number of MCMC draws from posterior
124
125   ppc.datasets <- matrix(nrow=m, ncol=n)
126   rownames(ppc.datasets) <- paste("PPC_", 1:m, sep="")
127
128   index = sample.int(niter, m, replace=TRUE) #Sample the posteriors we want.
129
130   #Calculate p_i using posterior draw.
131   #Rows: Probability at each X_r; Columns: PP dataset m_i.
132   p.i = sapply(index, function(i) expit(X, posterior[i,]))
133
134   PPC_Data = sapply(1:m, function(i) rbinom(n, mi, p.i[,i]))
135   PPC_mean = apply(PPC_Data, 2, mean) #Apply the mean to each PPC Dataset
136   PPC_var = apply(PPC_Data, 2, var) #And the variance
137
138   true_mean = mean(Y)
139   true_var = var(Y)
140
141   par(mfrow=c(2,2))
142   truehist(PPC_mean, main="Posterior Predictive Check on Mean", xlab="", ylab="")
143   abline(v=true_mean, col="red")
144
145   truehist(PPC_var, main="Posterior Predictive Check on Variance", xlab="", ylab="")
146   abline(v=true_var, col="red")
147
148   NULL
149 }
150 PPC.Graphs = cmpfun(PPC.Graphs)

```

The BLR_Fit file contains the code to run the simulations when the true value was known. It outputs a matrix containing the quantiles for each of the two parameters.

```

1  library(mvtnorm)
2  library(coda)
3  library(compiler)
4  suppressWarnings(library(mvtnorm))
5  suppressWarnings(library(coda))
6  suppressWarnings(library(compiler))
7
8  setwd("/home/caden11/sta250/Explorations-into-Computational-Statistics/HW1/BayesLogit/")
9
10 #####
11 #####
12 ## Handle batch job arguments:
13 args <- commandArgs(TRUE) #l-indexed version is used now.
14
15 cat(paste0("Command-line arguments:\n"))
16 print(args)
17
18 #####
19 # sim_start ==> Lowest simulation number to be analyzed by this particular batch job
20 ###
21
22 #####
23 sim_start <- 1000
24 length.datasets <- 200
25 beta.0=c(0,0)
26 sigma.0.inv = diag(2)
27 #####
28
29 if (length(args)==0){
30   sinkit <- FALSE
31   sim_num <- sim_start + 1
32   set.seed(1330931)
33 } else {
34   # Sink output to file?
35   sinkit <- TRUE
36   # Decide on the job number, usually start at 1000:
37   sim_num <- sim_start + as.numeric(args[1])
38   # Set a different random seed for every job number!!!
39   set.seed(762*sim_num + 1330931)
40 }
41
42 # Simulation datasets numbered 1001-1200
43
44 #####
45 #####
46
47 # Read data corresponding to appropriate sim_num:
48 Data = read.csv(paste("./data/blr_data_", sim_num, ".csv", sep=""))
49 pars = read.csv(paste("./data/blr_pars_", sim_num, ".csv", sep=""))
50
51 # Extract X and y:
52 X = as.matrix(Data[,3:4], ncol=2)
53 y = Data$y
54 m = Data$n
55
56 #####
57 # Set up the specifications:
58 p = ncol(X)
59 beta.0 <- matrix(rep(0,p))
60 sigma.0 <- sigma.0.inv <- diag(p)
61
62 niter <- 5E5
63 burnin = 1E4
64 print.every = niter/10
65 retune = 1E3
66 verbose = TRUE
67 local = FALSE #If running a cluster job, don't the piece to get raw samples--just quantiles.
68 thin = 100 #Take only every 100 observations to generate the quantiles.
69 #####
70
71 # Fit the Bayesian model:
72 source("./helper_functions.R")
73 if (local) {

```

```

74 Posterior_Draws = bayes.logreg.samples(m, y, X, beta.0, sigma.0, niter, burnin, print.every, retune,
    verbose)
75 chain = mcmc(Posterior_Draws) #Markov Chain objects, to make plots and ESS easy.
76 effectiveSize(chain) #Effective Sample Size for each variable.
77 plot(chain) #Convergence plots
78 }
79
80 # Extract posterior quantiles...
81 posterior_quantiles = bayes.logreg(m, y, X, beta.0, sigma.0.inv, niter, burnin, print.every, retune,
    verbose, thin)
82
83 # Write results to a (99 x p) csv file...
84 write.table(posterior_quantiles, file=paste(getwd(), "/results/blr_res_", sim_num, ".csv", sep=""), row.
    names=FALSE, col.names=FALSE, sep=",")
85
86 # Go celebrate.
87 #stop("You can't tell me what to do--this is America!")

```

The Cancer file handles the import and cleaning of the cancer data, along with the creation of graphs and the MCMC for the data set. At the end, it produces effective sample size, convergence plots, and does a little bit of model checking using posterior predictive checks.

```

1  library(mvtnorm) #Generate MV Normals
2  library(coda) #MCMC datatype and Eff Samp Size
3  library(compiler) #Bytecode Compiler--speeds up for loops
4  library(MASS) #For Truehist
5  library(xtable) #To display Latex for variable selection
6
7  if (Sys.info()['sysname'] %in% c("Linux", "Darwin")) {
8    setwd("/Dropbox/sta250/Assignments/HW1/BayesLogit/")
9  } else {
10   stop("You're using a different architecture than I am. setwd() so that the code is in the pwd!")
11 }
12
13 # Read data corresponding to appropriate sim_num:
14 cancer = read.table("breast_cancer.txt", sep="", header=TRUE)
15
16
17 X = model.matrix(diagnosis ~ ., data=cancer) #Making the design matrix
18 p = ncol(X)
19 y = as.numeric(cancer$diagnosis)-1 #Benign=0, Malignant=1;
20 m = rep(1, length(y))
21
22 #####
23 # Set up the specifications:
24 beta.0 <- matrix(rep(0,p))
25 sigma.0 <- diag(p)*1000
26 niter <- 5E5
27 burnin = 1E4
28 print.every = niter/10
29 retune = 1E2
30 verbose = TRUE
31 #####
32
33 # Fit the Bayesian model:
34 source("../helper_functions.R")
35 Posterior_Draws = bayes.logreg.samples(m, y, X, beta.0, sigma.0, niter, burnin, print.every, retune,
    verbose)
36
37 chain = mcmc(Posterior_Draws)
38 ESS = effectiveSize(chain)
39 corrs = acf(chain, lag.max=1, plot=FALSE) #Compute lag-1 autocorrelations
40 diag_acf = numeric(p)
41 for (i in 1:p) { diag_acf[i] = corrs$acf[2,i,i] } #Take only the non cross-autocorrs.
42
43 #Let's thin the draws to reduce the autocorrelation.
44 chain_thinned = mcmc(Posterior_Draws[seq.int(1, (niter-burnin), by=100), ])
45 ESS_thinned = effectiveSize(chain_thinned)
46 corrs_thinned = acf(chain_thinned, lag.max=1, plot=FALSE) #Compute lag-1 autocorrelations
47 diag_acf_thinned = numeric(p)

```

```

48 for (i in 1:p) { diag_acf_thinned[i] = corrs_thinned$acf[2,i,i] } #Take only the non cross-autocorrs.
49
50 #Plotting the thinned, final chain.
51 plot(chain_thinned, density=FALSE, ask=FALSE)
52 Posterior_Quantiles = apply(chain_thinned, 2, quantile, probs=c(.025, .975))
53 xtable(Posterior_Quantiles, caption="95% Central Credible Interval", align="|cccccccccc|")
54
55 #Making the Posterior Predictive Check graph.
56 PPC.Graphs(Y, X)

```