

Bayesian Data Analysis - Assignment 2 (S1889112)

The primary software used in this work was R, interfacing to JAGS via the libraries `rjags` and `runjags`. JAGS uses Gibbs sampling (a form of Markov Chain Monte Carlo (MCMC)).

1a: Exploratory Data Analysis

The results of the EDA are in Figure 1. Audouin abundance has fallen over time, while Yellowlegged abundance has increased over time. As such, we may expect a negative relationship between the two types of gulls. However looking at the correlation coefficient for all the data suggests no relationship between Audouin and Yellowlegged gulls. This is possibly due to the spikes in Audouin gulls in the mid to late 90s (before a sudden drop-off), where there is also a spike in Yellowlegged gulls. The small sample (26) means these values may have quite a large effect on the correlation. Sequentially removing the 2 outliers indicated by the boxplot leads to correlations of -0.07, and -0.30 for removing Audouin abundances of 625 and 525 respectively. The red and green lines are the cases with 625 and 525 removed respectively.

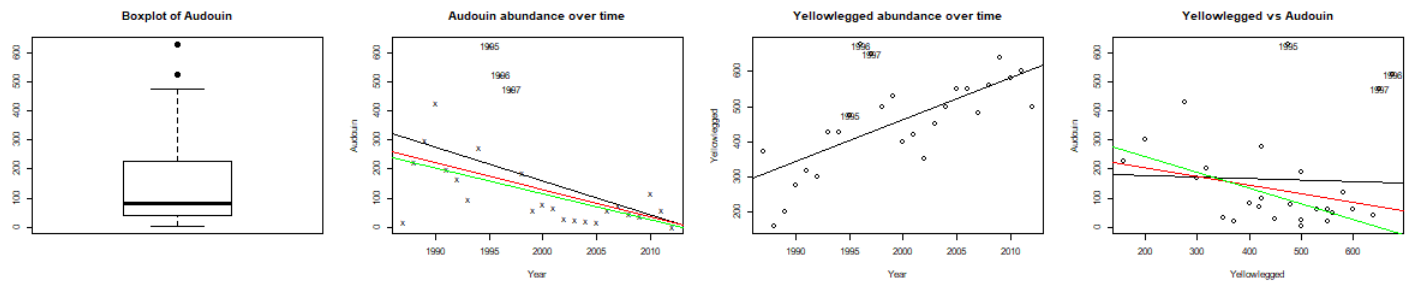


Figure 1: EDA Plots

1b: Initial Model

The only covariate is the year, which has been demeaned to assist convergence. The priors on the β_j are set to be $N(\mu_{\beta_j} = 0, \sigma_{\beta_j}^2 = 100)$ to reflect a lack of information about the parameters, while avoiding the issue of sparse priors causing numerical issues. Other priors were tested, all giving similar results. This was done by building the model for precision parameters (0.001, 0.01, 0.1, 1, 10). The code for this is in the “Consistency Check” section inside the R code for this question. The model is as follows:

$$\begin{aligned} \text{Likelihood:} \quad & y_i | \mu_i, \text{Year}_i \sim \text{Poisson}(\mu_i), \quad i = 1, \dots, n \\ \text{Link:} \quad & \log(\mu_i) = \beta_0 + \beta_1(\text{Year}_i - \overline{\text{Year}}) \\ \text{Priors:} \quad & \beta_j \sim N(0, 100), \quad j = 0, 1 \end{aligned}$$

The model was run with a burn-in of 20000, 20000 iterations, and a thinning interval of 4. Effective sample sizes are approximately 13500 for both β_0 and β_1 using the combined chains, with initial values determined by JAGS, due to the “Node inconsistent with parents” error when initialised manually or through a function call.

The Gelman-Rubin convergence diagnostic has point estimates and upper confidence intervals of 1. This is replicated by running the appended code. Trace and density plots are shown in Figure 2, showing good mixing and smoothness.

The results of the parameter estimates can be seen in Figure 3 to 3 decimal places. The 95% credible intervals for both parameters are very narrow. The value for β_0 suggests that in an average year, the log mean of Audouin abundance is 4.943, with a 95% credible interval of (4.909, 4.978). β_1 suggests that each extra year is associated with a fall in the log mean of -0.076, with a 95% credible interval of (-0.08, -0.071). In level terms this can be written $(\mu | \text{Year} = x + 1) = e^{-0.076}(\mu | \text{Year} = x)$.

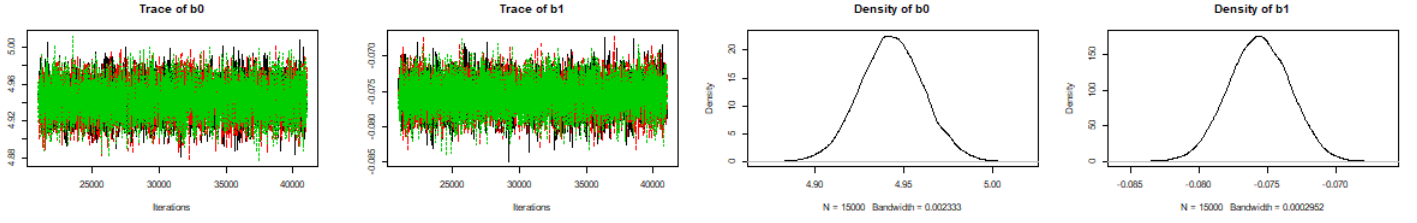


Figure 2: Trace and Density Plots

	2.5%	25%	50%	75%	97.5%	Mean	SD	Naive SE	Time-series SE	95% CI	95% CI Width
$b0$	4.909	4.931	4.943	4.955	4.978	4.943	0.018	0	0	(4.909, 4.978)	0.069
$b1$	-0.08	-0.077	-0.076	-0.074	-0.071	-0.076	0.002	0	0	(-0.08, -0.071)	0.009

Figure 3: Results of Poisson Regression $\log(\mu_i) = \beta_0 + \beta_1(\text{Year}_i - \overline{\text{Year}})$

1c: Extra-Variance

The objective of this section is to add variance to the model manually, so it is not solely determined by the mean, as is the case for a Poisson random variable. We augment the model from 1b to include a unit-specific error term in the link. That is:

$$\begin{aligned}
 \textbf{Likelihood:} \quad & y_i | \mu_i, \text{Year}_i \sim \text{Poisson}(\mu_i), \quad i = 1, \dots, n \\
 \textbf{Link:} \quad & \log(\mu_i) = \beta_0 + \beta_1(\text{Year}_i - \overline{\text{Year}}) + \epsilon_i \\
 \textbf{Priors:} \quad & \epsilon_i \sim N(0, \sigma_\epsilon^2), \quad \beta_j \sim N(0, 100), \quad \forall i, j \\
 \textbf{Hyperprior:} \quad & \sigma_\epsilon \sim U(0, 10)
 \end{aligned}$$

There is now much higher autocorrelation in the parameters. Hence we use a burn-in of 100000, 1 million iterations, and a thinning interval of 400. 3 chains are used, with initial values determined by JAGS. The thinning interval was determined using a 2-step process. Firstly, running a chain with no thinning, and seeing which parameter has the lowest effective sample size. Secondly, examining the effective sample size for this parameter for various thinning intervals by manually thinning the chain and computing the effective size on this thinned chain. This is shown in the `thincheck` function. We did not check multiple alternative priors, due to computational infeasibility. However changing the standard deviation of the β_j to 1 yielded no significant differences. The effective sample sizes are at least 4000. The Gelman-Rubin diagnostic is 1 for point estimates and upper confidence intervals, and the trace plots (Figure 5) suggest good mixing. The density plots are again smooth for all parameters.

The results can be seen in Figure 4. The credible intervals are now much wider, and the values for both β_j have shifted downwards. The extra-variance term in the linear component has approximately unit standard deviation, with a reasonably narrow 95% credible interval of (0.79, 1.44). Note that exponentiating this error, along with the multiplicative impact, suggests a reasonably large error with regards to the mean of raw series. β_0 has a mean of 4.5, with a 95% credible interval of (4.07, 4.92), while β_1 has a mean of -0.09 with a 95% credible interval of (-0.15, -0.04). Interpretation is the same as previously.

	2.5%	25%	50%	75%	97.5%	Mean	SD	Naive SE	Time-series SE	95% CI	95% CI Width
$b0$	4.071	4.365	4.502	4.645	4.921	4.503	0.215	0.002	0.003	(4.071, 4.921)	0.85
$b1$	-0.149	-0.112	-0.093	-0.074	-0.037	-0.093	0.029	0	0	(-0.149, -0.037)	0.112
sigma.epsilon	0.781	0.936	1.036	1.154	1.459	1.057	0.171	0.002	0.002	(0.781, 1.459)	0.678

Figure 4: Results of Poisson Regression $\log(\mu_i) = \beta_0 + \beta_1(\text{Year}_i - \overline{\text{Year}}) + \epsilon_i$

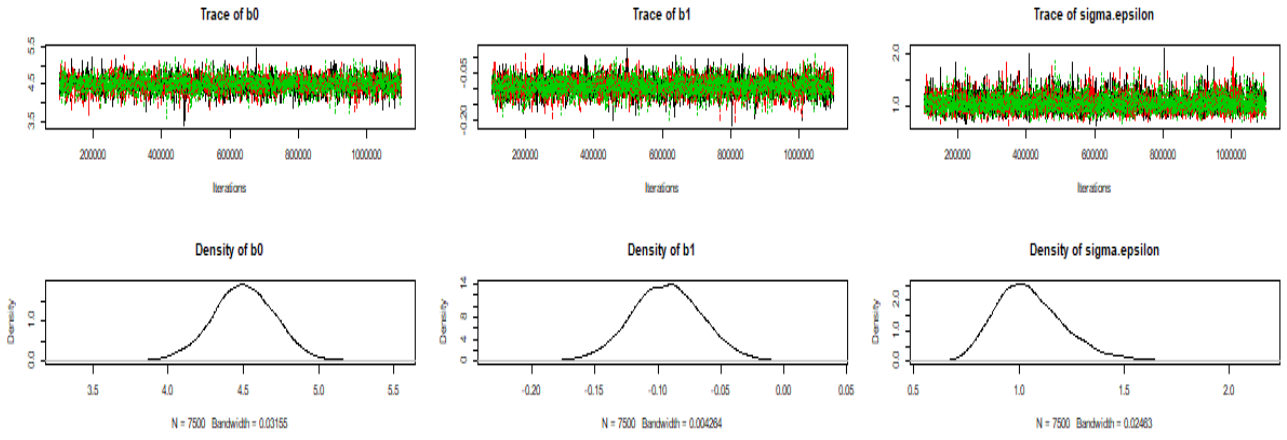


Figure 5: Trace and Density Plots

1d: Yellowlegged as a Covariate

Adding the abundance of Yellowlegged gulls is just a change in the link to include the covariate. Hence we have the same model as in 1c (with the same prior on β_2 as the other β_j) but with the following link:

Link:
$$\log(\mu_i) = \beta_0 + \beta_1(\text{Year}_i - \overline{\text{Year}}) + \beta_2(\text{Yellowlegged}_i - \overline{\text{Yellowlegged}}) + \epsilon_i$$

Using the same thinning evaluation as the previous section, we run 1.4 million iterations, thinning of 500, a burn-in of 200000, and 3 chains initialised by JAGS. Trace plots are in Figure 6, with full results in Figure 7. The trace plots show good mixing and the density plots are smooth. Effective sample sizes are at least 2500, with Gelman-Rubin diagnostics of 1 for all parameters.

The posterior parameter on β_2 has a mean of 0.004, and a 95% credible interval of (0, 0.008). This appears to be economically insignificant, in that the value of the coefficient is essentially 0. Based on the exploratory analysis this might be expected, given that the complete-case analysis (shown by the black line) suggests no relationship. However once we remove the 2 large outliers, there is a visible negative relationship. This would fit with intuition, since an increase in Yellowlegged gulls means food supplies are reduced once Audouin gulls arrived (since they have a later reproduction cycle). Given that all cases were used, we expected an insignificant relationship between the number of Yellowlegged gulls and the number of Audouin gulls in the same year, but should be wary that this lack of a relationship could be due to some environmental shock in 1996 and 1997 causing the abundance of both gulls to spike.

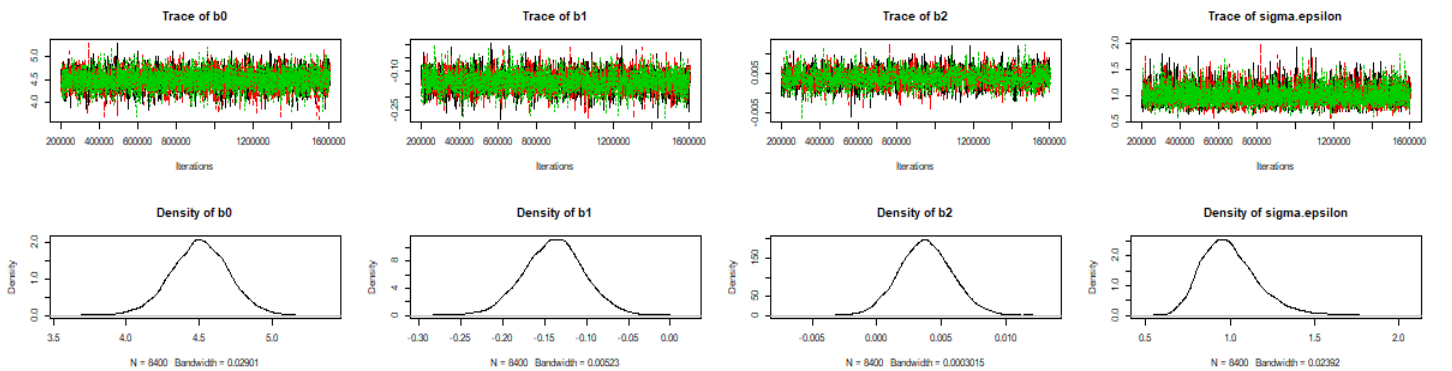


Figure 6: Trace and Density Plots

	2.5%	25%	50%	75%	97.5%	Mean	SD	Naive SE	Time-series SE	95% CI	95% CI Width
b_0	4.104	4.373	4.505	4.636	4.892	4.503	0.199	0.002	0.003	(4.104,4.892)	0.788
b_1	-0.211	-0.163	-0.139	-0.115	-0.067	-0.139	0.037	0	0.001	(-0.211,-0.067)	0.144
b_2	0	0.002	0.004	0.005	0.008	0.004	0.002	0	0	(0,0.008)	0.008
σ_{ϵ}	0.724	0.88	0.98	1.097	1.393	1	0.17	0.002	0.002	(0.724,1.393)	0.669

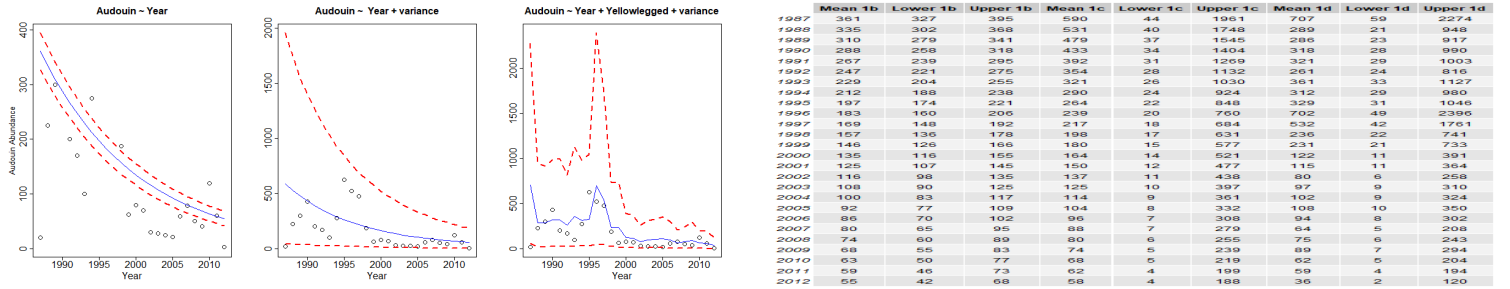
Figure 7: Results of Poisson Regression $\log(\mu_i) = \beta_0 + \beta_1(\text{Year}_i - \overline{\text{Year}}) + \beta_2(\text{Yellowlegged}_i - \overline{\text{Yellowlegged}}) + \epsilon_i$

1e: Posterior Predictive

Suppose we have S MCMC samples. For each observation i in the data, we compute $\mu_i = \exp(\beta_0 + \beta_1(\text{Year}_i - \overline{\text{Year}}) + \beta_2(\text{Yellowlegged}_i - \overline{\text{Yellowlegged}}) + \epsilon_i)$ repeatedly for each set of posterior samples of the β_j and ϵ . Therefore each observation has S predictions based on the different samples. To sample ϵ , we form a vector based on the samples of σ_{ϵ} . For each sample of σ_{ϵ} , we randomly draw from a normal distribution with mean 0 and standard deviation equal to the sample from σ_{ϵ} , and store the result in the vector. Then we can denote the s^{th} sample of the i^{th} observation as:

$$\mu_{is} = \exp(\beta_{0s} + \beta_{1s}(\text{Year}_i - \overline{\text{Year}}) + \beta_{2s}(\text{Yellowlegged}_i - \overline{\text{Yellowlegged}}) + \epsilon_s)$$

Where 's' is the index of the MCMC sample. Define the vector of posterior predictive samples for observation i as $\boldsymbol{\mu}_i = (\mu_{i1} \ \mu_{i2} \ \dots \ \mu_{iS})$. The mean can be calculated for each observation by taking the mean of this vector, and the 90% credible interval by passing the vector through the `quantile` function with `probs = c(0.05, 0.95)`. For each μ_{is} we sample from the Poisson distribution with parameter μ_{is} , creating the vector $\mathbf{y}_i = (y_{i1} \ \dots \ y_{iS})$ for each observation. That is $y_{is} \sim \text{Pois}(\mu_{is})$. The results of this analysis are shown in Figure 8. The dashed red lines represent the limits of the credible interval, and the blue line is the mean.



a: Posterior Predictive Mean and 90% Credible Interval

b: Full Results for Posterior Predictive

Figure 8: Posterior Predictive Graph and Full Results for all models

1f: Model Selection

Given the results in Figure 9, we can be reasonably confident that the model from 1b is not appropriate to model the data, caused by the lack of an extra-variation term leading to the 90% credible intervals capturing the observed data incredibly poorly since the Poisson distribution has variance equal to the mean. Just including this extra-variation allows the credible interval to capture the data well in terms of the credible interval. However the mean is quite poorly captured, with large residuals until about the year 2000. In terms of the mean, the data is captured better in model 1d, however this could be due to overfitting, since we have only 9 (26/3) observations per β_j .

It appears as though most of the relationship can be captured by a time trend, coinciding with the exploratory analysis done in 1a, and the small value of the coefficient on the demeaned Yellowlegged covariate in 1d.

Both of the models with the extra-variance term seem adequate to model the data, given the near-identical DIC, and the plots capturing the data. If we are looking to establish some semblance of causality, I would select model 1d, as the presence of Yellowlegged gulls dramatically improves how well the mean captures the data, despite the potential overfitting, and is intuitively an important variable to control for. Additionally, despite

mentioning the possibility of overfitting, the DIC is in fact the lowest among the models suggesting the improvement in predictive performance may be worth the loss in parsimony.

	Model 1B	Model 1C	Model 1D
DIC	2986.13	216.6	216.54

Figure 9: DIC Comparison

2a: EDA

The results of basic EDA are shown in Figure 10. The top left panel shows the correlation between the environment variables of the farm. In calculating this, we only took 1 row from each farm to avoid inflated correlations caused by an unequal number of observations for each farm. For instance suppose we have 2 farms with 100 observations each, and a 3rd farm with 1 observation. Correlation between any 2 variables will be almost perfect, as the first 2 farms dominate the sample. We see that temperature, height, and rain are all intricately related, with high correlations between each of the 3 covariates. Conversely, the slope and permeability have almost no linear relationship with the other covariates.

The other panels in this row calculate the parasite proportion within each farm, and for each cow age. These suggest that farm 1 is ‘safest’, given its large sample size. The farms with higher parasite proportions have smaller sample sizes, but farms 6, 7, and 13 appear to be the most parasitic farms. There appears to be a higher parasite probability in older cows.

The remaining panels show the distribution of each covariate given a specific parasite level. In addition, the Mann-Whitney p-value is reported. This is a nonparametric test for differences in the median of 2 groups, and does not lose much power from a t-test. Given this statistic, it appears as if there is clear discrepancy in the medians between groups for temperature, rainfall, and height. We conclude from this simple test that there is not much evidence for median difference between the groups for slope, age, and permeability. However we note that permeability only takes on 2 values, so a test for median differences should not be relied upon to draw conclusions. For instance, (0, 0, 0, 0, 0, 1, 1, 1, 1, 1) and (0, 0, 0, 0, 0, 0, 0, 0, 0, 1) will have the same median despite the groups clearly being different.

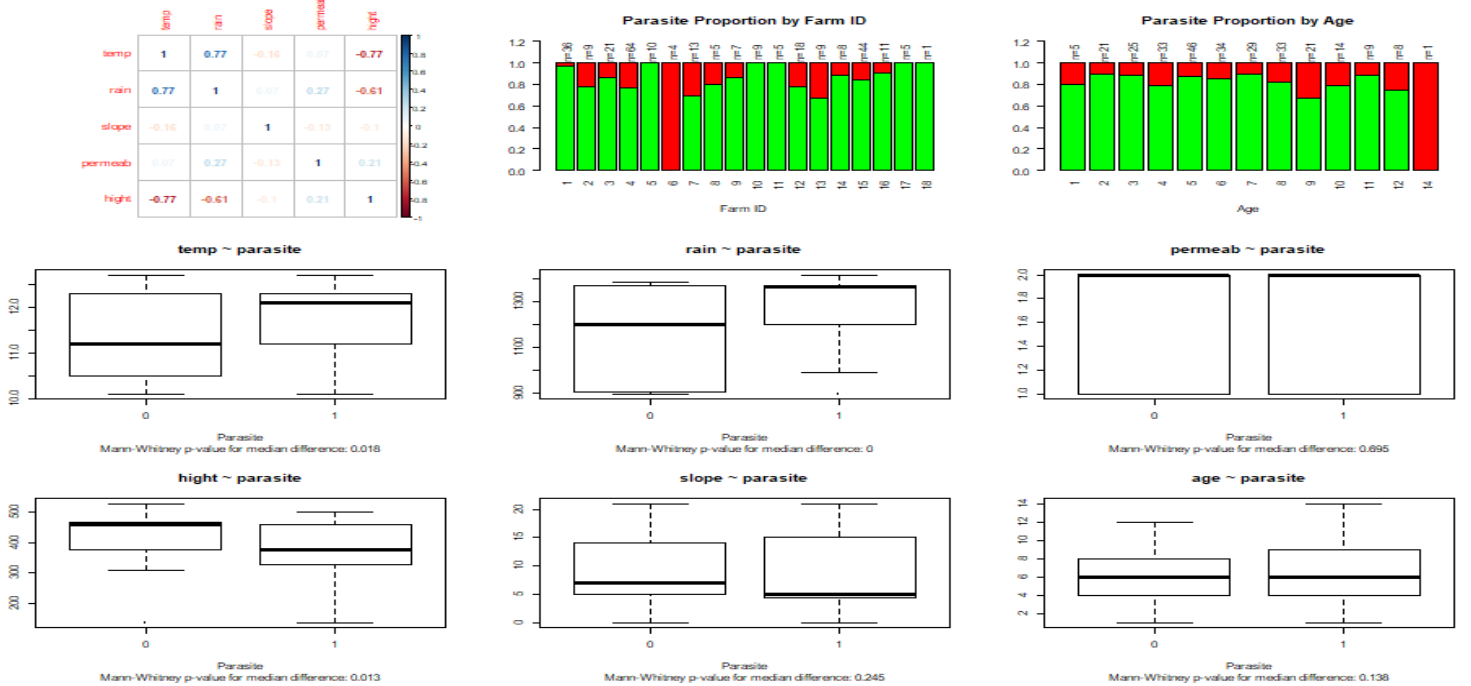


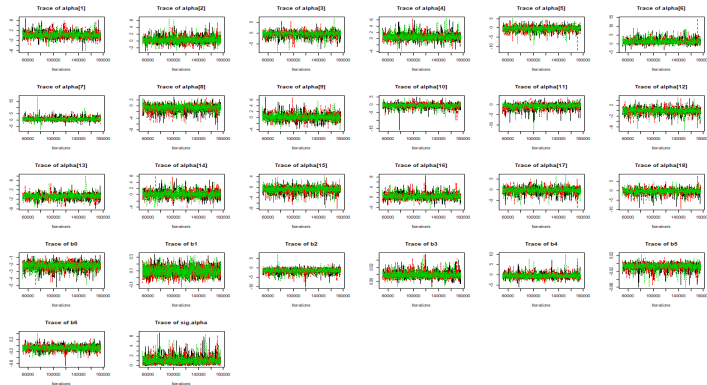
Figure 10: EDA Plots

2b: Model, All covariates

All covariates have been demeaned to aid convergence. The demeaned version of covariate x_k will be represented as \hat{x}_k . The vector of demeaned variables is $(\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4, \hat{x}_5, \hat{x}_6)$, corresponding respectively to age, temperature, rain, permeability, height, and slope. We represent the random effect as the sum of the common mean β_0 and a Gaussian white noise term α_j . The sum of these can be seen as the total random effect. In general, let us represent the total random intercept as $\gamma_j \sim N(\mu_\gamma, \sigma_\gamma^2)$, this is equivalent to $\mu_\gamma + \epsilon_j$, where $\epsilon_j \sim N(0, \sigma_\gamma^2)$, by properties of the normal distribution. In both cases, we place priors on μ_γ and σ_γ^2 , representing the common distribution of the random intercept. Referring back to our model, we have $\mu_\gamma = \beta_0$ and $\epsilon_j = \alpha_j$. β_0 has the same prior as the other β_k , where $\beta_k \sim N(0, 100)$, however this will be referred to as a hyperprior to make the relation to the random effect more clear. Instead of placing a prior on the variance, we place it on the standard deviation, such that $\sigma_\alpha \sim U(0, 10)$.

Likelihood:	$y_{ij} p_{ij}, x_i \sim \text{Bernoulli}(p_{ij}) \quad \forall i, j$
Link:	$\text{Logit}(p_{ij}) = \beta_0 + \alpha_j + \beta_1 \hat{x}_{1,ij} + \beta_2 \hat{x}_{2,ij} + \beta_3 \hat{x}_{3,ij} + \beta_4 \hat{x}_{4,ij} + \beta_5 \hat{x}_{5,ij} + \beta_6 \hat{x}_{6,ij}$
Priors:	$\alpha_j \sim N(0, \sigma_\alpha^2), \quad \beta_k \sim N(0, 100), \quad i \in (1, \dots, n), j \in (1, \dots, 18), k \in (1, \dots, 6)$
Hyperpriors:	$\sigma_\alpha \sim U(0, 10) \quad \beta_0 \sim N(0, 100)$

The thinning interval was determined in the same manner as specified in 1c. A burn-in of 50000 was used, with 125000 iterations and a thinning interval of 100. Effective sample sizes were at least 1800 for all parameters. Trace plots and results are in Figure 11. Trace plots show good mixing, and the highest Gelman-Rubin upper confidence interval was 1.05. Initial values were determined by JAGS, due to the same error as previously. For all future models, initial values will also be determined by JAGS. β_0 has a mean value of -2.144, with a 95% credible interval of (-3.2, -1.4), suggesting there is quite a lot of uncertainty around the common mean of the random effect. In general, there is a high degree of uncertainty around almost all of the parameters, with the width of the credible intervals being around the range (3,5) for all parameters except the overall mean (β_0), and $\beta_1, \beta_3, \beta_5$ and β_6 . These are the coefficients on age, rain, height, and slope respectively. Of these, β_1 (age) seems to be the most economically significant, with a mean of 0.099, but its credible interval (-0.03, 0.23) contains 0 casting doubt on the significance of age. However there is a 93% probability this parameter is greater than 0, suggesting it is likely to be positive (see R code). Of the parameters with narrow credible intervals, only β_5 has a credible interval not containing 0, suggesting it is the most statistically significant. As such, of the covariates, it seems age (economic significance) and height (statistical significant) have the strongest relationship with the parasite probability.



a: Trace Plots

	2.5%	25%	50%	75%	97.5%	Mean	SD	Naive SE	Time-series SE	95% CI	95% CI Width
alpha[1]	-1.768	-0.342	0.064	0.621	2.566	0.176	1.051	0.017	0.017	(-1.768, 2.566)	4.334
alpha[2]	-1.146	-0.152	0.209	0.79	2.33	0.352	0.872	0.014	0.016	(-1.146, 2.33)	3.476
alpha[3]	-3.245	-0.849	-0.215	0.201	1.563	-0.4	1.173	0.019	0.023	(-3.245, 1.563)	4.808
alpha[4]	-1.208	-0.047	0.354	0.974	2.938	0.514	1.003	0.016	0.023	(-1.208, 2.938)	4.146
alpha[5]	-3.114	-0.765	-0.14	0.278	1.978	-0.282	1.206	0.02	0.02	(-3.114, 1.978)	5.092
alpha[6]	-0.651	0.214	0.828	1.693	4.425	1.104	1.31	0.021	0.023	(-0.651, 4.425)	5.076
alpha[7]	-1.011	0.022	0.505	1.294	3.626	0.76	1.224	0.02	0.022	(-1.011, 3.626)	4.637
alpha[8]	-3.152	-1.018	-0.34	0.065	1.078	-0.551	1.045	0.017	0.02	(-3.152, 1.078)	4.23
alpha[9]	-1.742	-0.321	0.084	0.629	2.37	0.173	0.993	0.016	0.017	(-1.742, 2.37)	4.112
alpha[10]	-3.277	-0.816	-0.167	0.227	1.666	-0.361	1.225	0.02	0.022	(-3.277, 1.666)	4.943
alpha[11]	-3.384	-1.023	-0.34	0.074	1.204	-0.568	1.163	0.019	0.019	(-3.384, 1.204)	4.588
alpha[12]	-2.072	-0.355	0.04	0.508	1.963	0.041	0.967	0.016	0.017	(-2.072, 1.963)	4.035
alpha[13]	-2.898	-1.058	-0.413	0.003	1.168	-0.56	1.005	0.016	0.018	(-2.898, 1.168)	4.056
alpha[14]	-1.692	-0.316	0.079	0.646	2.234	0.167	0.946	0.015	0.016	(-1.692, 2.234)	3.926
alpha[15]	-3.371	-1.203	-0.625	-0.027	1.02	-0.703	1.103	0.018	0.02	(-3.371, 1.102)	4.391
alpha[16]	-1.335	-0.169	0.224	0.86	2.882	0.391	1.03	0.017	0.019	(-1.335, 2.882)	4.217
alpha[17]	-2.994	-0.664	-0.047	0.429	2.157	-0.159	1.225	0.02	0.02	(-2.994, 2.157)	5.151
alpha[18]	-3.047	-0.706	-0.082	0.369	2.048	-0.214	1.217	0.02	0.019	(-3.047, 2.048)	5.095
beta[0]	-3.243	-2.367	-2.091	-1.847	-1.336	-2.142	0.474	0.008	0.009	(-3.243, -1.336)	1.907
beta[1]	-0.028	0.054	0.098	0.141	0.229	0.098	0.065	0.001	0.001	(-0.028, 0.229)	0.257
beta[2]	-4.139	-1.995	-1.257	-0.66	0.519	-1.392	1.183	0.019	0.026	(-4.139, 0.519)	4.658
beta[3]	0.003	0.007	0.009	0.012	0.02	0.01	0.004	0	0	(0.003, 0.02)	0.017
beta[4]	-2.347	-1.143	-0.597	0.081	2.002	-0.478	1.119	0.018	0.022	(-2.347, 2.002)	4.349
beta[5]	-0.028	-0.011	-0.006	-0.002	0.004	-0.007	0.008	0	0	(-0.028, 0.004)	0.032
beta[6]	-0.244	-0.115	-0.071	-0.029	0.078	-0.074	0.079	0.001	0.001	(-0.244, 0.078)	0.322
sig.alpha	0.088	0.568	0.948	1.411	2.997	1.085	0.764	0.012	0.017	(0.088, 2.997)	2.909

b: Full Results

Figure 11: Results and Convergence for model 2b

2c: Model simplification

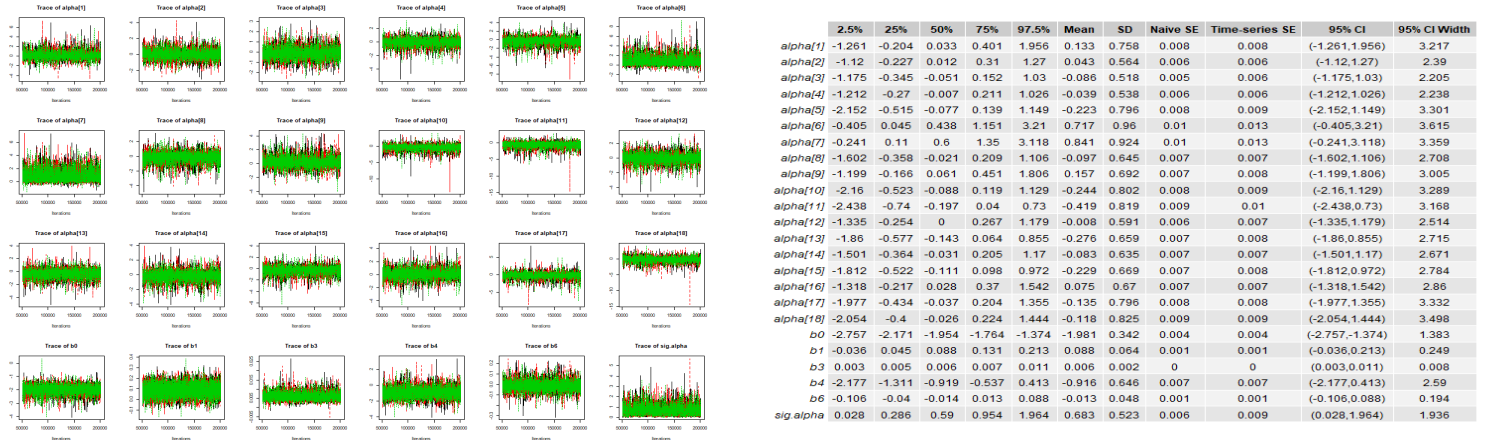
Now we look at the correlations between the environment variables to determine some model simplifications, and therefore refer back to Figure 10, looking at the top left correlation plot. We can see that temperature is highly correlated with both rain and height, and therefore remove it, as we might suggest that we can explain most of the influence of temperature by these 2 variables. Hence we try a model without temperature. We then see that rain and height are highly correlated with each other, and therefore remove both variables individually, as one may capture the effect of the other. As such we try 3 simplifications. Firstly, a model without temperature. Secondly, a model without temperature or height. Thirdly, a model without temperature or rain. The coefficient indices will be maintained for consistency. E.g. the coefficient on permeability remains β_4 throughout. We rely on DIC to select the model. For computational speed, each model was run with only 2 chains for the DIC calculation. A burn-in of 50000 was used, with 100000 iterations for each DIC sample. The results of can be seen in Figure 12. We have assumed convergence after a burn-in of 10000, and that there is a reasonable effective sample size. This is perhaps not unreasonable, given the effective sizes for previous and future models, but is by no means a guarantee. This was not done due to time constraints. We select the model which has removed temperature and height as covariates, as it has the lowest DIC.

	Full Model	Temp Removed	Temp, Rain Removed	Temp, Height Removed
DIC	242.65	243.45	244.04	242.24

Figure 12: DIC of considered models

We recompile the model with 3 chains, a burn-in of 50000, 150000 iterations, and a thinning interval of 50. The effective sizes are at least 4000. The Gelman-Rubin diagnostic upper intervals are below 1.02, and the trace plots in Figure 13 show good mixing.

In this same Figure, we see the width of credible intervals for the parameters has generally fallen. The mean for β_1 (age) is mostly unchanged, but the mean for β_4 (permeability) has increased in magnitude (decreased in value) from -0.467 to -0.917, suggesting some of the predictive power lost by dropping β_2 and β_5 has been captured. The coefficient on rain β_3 is largely unchanged, falling from 0.01 to 0.006. This could be due to the omitted variables pulling the coefficient in different directions as a result of the opposite sign correlations (rain is positively correlated with temperature, negatively with height). The coefficient on the slope, β_6 , shrinks in importance, falling in magnitude from -0.075 to -0.014. The parameters on many of the random effects have significantly changed. For example farms 2 to 4 now have farm-specific effects much closer to 0.



a: Trace Plots

b: Full Results

Figure 13: Results and Convergence for no temperature, no height model

2d: Posterior for farms 1 and 6

The calculation method for the posterior is very similar to that used in 1e, but we are reporting on the link component, and not sampling from the distribution of the outcome variable. That is we are reporting p (the probability of a cow from farm j having a parasite), not y (the classification of a cow). We know the environmental characteristics of a specific farm are fixed, and assume the age to be equal to the mean for the following calculations. We calculate the posterior probability for farms 1 and 6 for each of the S MCMC samples we have, in the following manner:

$$\text{logit}(p_{js}) = \beta_{0s} + \alpha_{js} + \beta_{1s}\hat{x}_{1j} + \beta_{3s}\hat{x}_{3j} + \beta_{4s}\hat{x}_{4j} + \beta_{6s}\hat{x}_{6j}$$

Where the \hat{x} are defined as in 2b, s denotes the MCMC sample index, and the subscript i is dropped since we assume age equal to the mean, hence the 3rd term on the right hand side is zero by definition of \hat{x} , which removes all cow individuality in the model. After we calculate the linear combination, we apply `inverse.logit` from the `boot` package to calculate the probabilities. We then have a vector \mathbf{p}_j with S probabilities corresponding to each MCMC sample. The mean is calculated by taking the average of this vector. The 95% symmetric credible interval is calculated by using the `quantile` function with argument `probs = c(0.025, 0.975)`, and the probability of an epidemic is equivalent to the proportion of elements of the relevant \mathbf{p}_j greater than 0.2.

A histogram for each farm is illustrated, as well as the exact figures. These are in Figures 15 and 14. The mean is marked by a blue vertical line, and the 95% credible interval by red lines. Note that figures are rounded to 3 decimal places, and we used a limited number of samples. Hence an epidemic is not ‘impossible’ for farm 1, but is very unlikely. This makes sense given the fairly large sample for farm 1 ($n = 36$ as per 2a), with only 1 parasitic cow. The expected proportion of cows on farm 1 (with an average age) with a parasite is 2.5%, and there is a 95% probability this proportion is between 0.2% and 7.6%.

Similar logic applies to farm 6. All 4 cows on the farm have a parasite, which would be very unlikely with a probability of below 20%. With a probability of exactly 20%, all 4 cows having a parasite has a probability of 0.16% assuming independence (`dbinom(4, 4, 0.2)` in R). The low sample size leads to a wide credible interval of (0.223, 0.963), suggesting there is a 95% probability that the proportion of cows on farm 6 with the parasite (with an average age) is between 22.3% and 96.1%.

	Mean	95% CI	P(Epidemic)
Farm 1 (Rand)	0.025	(0.002, 0.074)	0
Farm 6 (Rand)	0.593	(0.226, 0.961)	0.989

Figure 14: Posterior Parasite Probabilities (Random Effects)

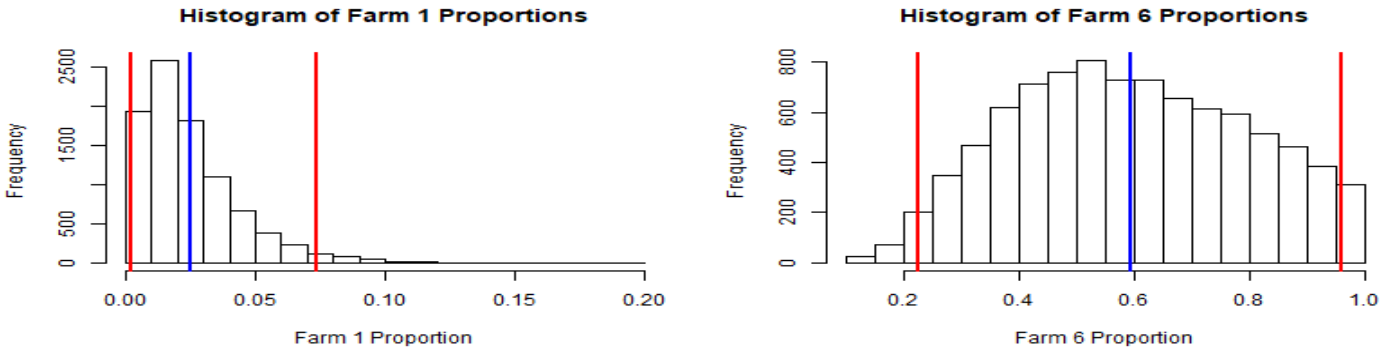


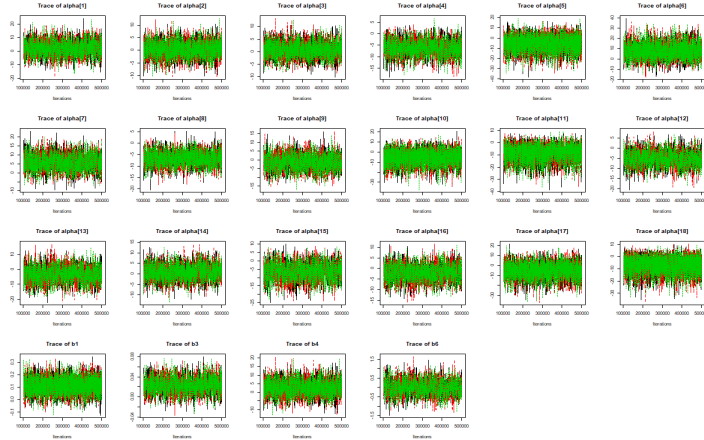
Figure 15: Histogram of Posterior Parasite Probabilities (Random Effects)

Q2e: Fixed Effects

In estimating a fixed effects model, each farm-specific intercept has its own hyperparameters, meaning the model has the same likelihood, but the following link and priors:

$$\begin{aligned} \text{Likelihood:} \quad & y_{ij}|p_{ij}, x_i \sim \text{Bernoulli}(p_{ij}) \quad \forall i, j \\ \text{Link:} \quad & \text{Logit}(p_{ij}) = \alpha_j + \beta_1 \hat{x}_{1,ij} + \beta_3 \hat{x}_{3,ij} + \beta_4 \hat{x}_{4,ij} + \beta_6 \hat{x}_{6,ij} \\ \text{Priors:} \quad & \alpha_j \sim N(0, 100), \quad \beta_k \sim N(0, 100), \quad \forall j, k \end{aligned}$$

The model was run with burn-in 100000, 400000 iterations, thinning of 200, and 3 chains. Effective sizes were at least 1500, and Gelman-Rubin diagnostics were 1 for all upper intervals. The relevant trace plots are in Figure 16, along with the full results. The trace plots show good mixing. The posterior summaries are in Figure 17. The credible intervals for the parameters are now much wider, since there is no information sharing via the common distribution, with 95% credible interval widths of between 10 and 30 for the farm-specific intercepts. The expected value of the proportion of parasitic cows has increased for both farms. From 2.5% to 2.9% for farm 1, and drastically from 59% to 98% for farm 6. The reason for farm 6 shooting up is likely because of the lack of information sharing between the groups, as they no longer come from a common distribution. To 3 decimal places, there is almost total certainty about the probability of an epidemic with fixed effects (of these farms). The credible interval is now slightly wider for farm 1, and much narrower for farm 6.



	2.5%	25%	50%	75%	97.5%	Mean	SD	Naive SE	Time-series SE	95% CI	95% CI Width
alpha[1]	-8.5	-1.502	2.139	5.704	12.782	2.151	5.392	0.07	0.113	(-8.5, 12.782)	21.282
alpha[2]	-5.687	-1.756	0.363	2.455	6.443	0.362	3.106	0.04	0.058	(-5.687, 6.443)	12.13
alpha[3]	-5.079	-1.259	0.697	2.789	6.848	0.76	3.048	0.039	0.061	(-5.079, 6.848)	11.927
alpha[4]	-13.04	-8.441	-6.035	-3.781	0.506	-6.131	3.428	0.044	0.081	(-13.04, 0.506)	13.546
alpha[5]	-21.932	-9.947	-4.41	0.522	8.526	-5.012	7.79	0.101	0.106	(-21.932, 8.526)	30.458
alpha[6]	-6.782	2.33	7.253	12.418	23.463	7.494	7.627	0.098	0.127	(-6.782, 23.463)	30.245
alpha[7]	-2.972	2.748	5.941	9.092	15.011	5.957	4.608	0.059	0.104	(-2.972, 15.011)	17.983
alpha[8]	-12.826	-8.872	-6.719	-4.554	-0.617	-6.723	3.141	0.041	0.063	(-12.826, -0.617)	12.209
alpha[9]	-9.915	-4.216	-1.324	1.825	7.506	-1.223	4.421	0.057	0.106	(-9.915, 7.506)	17.421
alpha[10]	-21.462	-9.933	-4.31	0.504	8.414	-4.969	7.666	0.099	0.103	(-21.462, 8.414)	29.876
alpha[11]	-22.608	-11.983	-7.262	-3.249	3.066	-7.953	6.611	0.085	0.094	(-22.608, 3.066)	25.674
alpha[12]	-12.763	-7.786	-5.087	-2.576	2.2	-5.166	3.834	0.049	0.095	(-12.763, 2.2)	14.963
alpha[13]	-13.952	-7.059	-3.374	0.338	7.104	-3.386	5.427	0.07	0.127	(-13.952, 7.104)	21.066
alpha[14]	-6.283	-1.653	0.662	3.21	7.855	0.767	3.617	0.047	0.071	(-6.283, 7.855)	14.138
alpha[15]	-17.164	-10.197	-6.53	-2.945	3.425	-6.577	5.284	0.068	0.131	(-17.164, 3.425)	20.589
alpha[16]	-9.611	-4.599	-1.945	0.851	5.442	-2.006	3.842	0.05	0.096	(-9.611, 5.442)	15.053
alpha[17]	-21.096	-9.728	-4.307	0.71	9.324	-4.712	7.828	0.101	0.108	(-21.096, 9.324)	30.42
alpha[18]	-22.131	-10.778	-5.736	-1.67	4.488	-6.59	6.817	0.088	0.093	(-22.131, 4.488)	26.619
b1	-0.019	0.066	0.11	0.153	0.238	0.109	0.065	0.001	0.001	(-0.019, 0.238)	0.257
b3	-0.007	0.013	0.023	0.034	0.055	0.024	0.016	0	0	(-0.007, 0.055)	0.062
b4	-6.044	-0.382	2.68	5.677	11.561	2.695	4.498	0.058	0.087	(-6.044, 11.561)	17.605
b6	-0.798	-0.272	-0.004	0.282	0.825	0.005	0.413	0.005	0.011	(-0.798, 0.825)	1.623

a: Trace Plots

b: Full Results

Figure 16: Results and Convergence for no temperature, no height model (Fixed Effects)

	Mean	95% CI	P(Epidemic)
Farm 1 (Fixed)	0.02892	(0.00072, 0.1025)	0.00017
Farm 6 (Fixed)	0.97996	(0.74967, 1)	1

Figure 17: Posterior Parasite Probabilities (Fixed Effects)

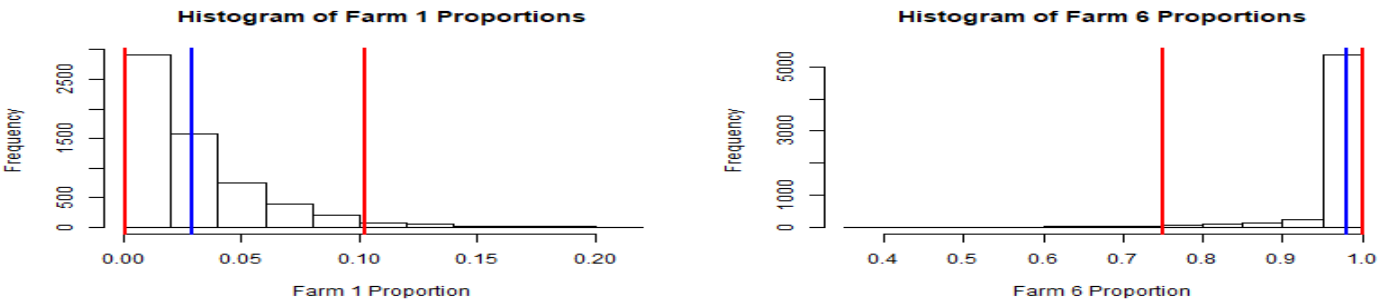


Figure 18: Posterior Parasite Probabilities (Fixed Effects)

Q3: A New Model

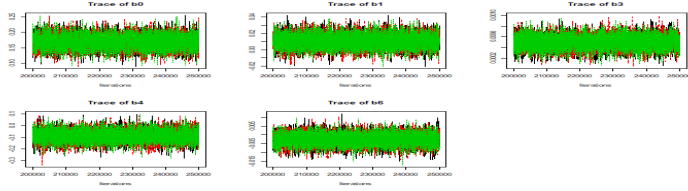
This section explores the performance of a Bayesian linear probability model (BLPM) applied to the cows data. We will use the reduced form of the model from 2c, and use a fixed intercept. That is, we are using a very simple model to examine a baseline for the model. Results will be compared using DIC. Additionally, a normal likelihood will be used. We place a conjugate, pseudo-uninformative gamma prior on $\tau = \frac{1}{\sigma^2}$.

$$y_i | \mu_i, x_i, \sigma^2 \sim N(\mu_i, \sigma^2) \quad \forall i$$

$$\mu_i = \beta_0 + \beta_1 \hat{x}_{1i} + \beta_3 \hat{x}_{3i} + \beta_4 \hat{x}_{4i} + \beta_6 \hat{x}_{6i}$$

$$\tau \sim \Gamma(0.01, 0.01) \quad \beta_j \sim N(0, 0.01) \quad j \in \{0, 1, 3, 4, 6\}$$

Here, μ_i represents the probability observation i has a parasite. To calculate the DIC, 3 chains were used with a burn-in of 50000, and 150000 iterations. The results were then sampled using 50000 further iterations, with thinning of 10 for computational ease. The DIC was 236.6, which is lower than all of the models under the hierarchical GLM framework (Figure 12). There is no reason to suspect a lack of convergence based on the trace plots in Figure 19, while the Gelman-Rubin statistics were 1 across the board. Effective sample sizes were all around 15000.



a: Trace Plots

	2.5%	25%	50%	75%	97.5%	Mean	SD	Naive SE	Time-series SE	95% CI	95% CI Width
b0	0.1262	0.1539	0.1683	0.183	0.2117	0.1685	0.0219	0.0002	0.0002	(0.1262, 0.2117)	0.0855
b1	-0.0054	0.0049	0.0104	0.0159	0.0269	0.0105	0.0082	0.0001	0.0001	(-0.0054, 0.0269)	0.0323
b3	0.0003	0.0004	0.0005	0.0006	0.0007	0.0005	0.0001	0	0	(0.0003, 0.0007)	0.0004
b4	-0.1859	-0.1232	-0.0902	-0.0563	0.0101	-0.0896	0.0502	0.0004	0.0004	(-0.1859, 0.0101)	0.196
b6	-0.0095	-0.0051	-0.0028	-0.0005	0.0038	-0.0028	0.0034	0	0	(-0.0095, 0.0038)	0.0133

b: Full Results

Figure 19: Results and Convergence for no temperature, no height model (BLPM)

With the BLPM, β_0 represents the proportion of cows with a parasite. The mean is 16.8% (the proportion in the dataset), with a 95% credible interval of (12.5%, 21.2%). β_3 , the coefficient on rainfall is 0.0005, suggesting that ceteris paribus, an increase in rainfall of 100ml is associated with a 0.05 (5%) increase in the probability a cow has a parasite. For this parameter, all samples were positive suggesting a very high probability that, when controlling for the variables in the model, there is a positive association between rainfall and the probability a cow has a parasite. This positive association is consistent with the results from Figure 13. The probability each parameter has the same sign as its mean is (1, 0.898, 1, 0.9654, 0.8) for $(\beta_0, \beta_1, \beta_3, \beta_4, \beta_5)$. These same probabilities from the model in 2c are (0.923, 0.999, 0.927, 0.65) for β_1 to β_6 (β_0 not included due to the random intercept distorting these probabilities). There appears to be consistency here, with β_6 having the most uncertain sign consistency across both models, while the positive effect of β_3 remains near certain.

Intuitively it seems problematic that this model has a lower DIC, since the output ‘probability’ can leave the (0,1) bound, and there is no distinction between the farms. The reason for this is due to the unbalanced data, in that the number of cows without parasites far outweighs those with parasites. As such, by setting all parameters to 0 (except the intercept), it is still possible to attain an 83% classification accuracy, given a standard cutoff of 0.5 as everything will be classified as 0 (no parasite). This is what the BLPM is doing (classifying everything as 0). See the PREDICTIONS section in the R code for details. As such, it appears that the hierarchical random effects model is hardly outperforming this basic model in terms of classification, with a higher penalty due to complexity, leading to higher DIC. As such the random effects model seems reasonably poor given this basic metric, and we need more data to form a reasonable model for the parasite probabilities, since at the moment (in terms of classification) even just a linear intercept model would hardly lose any predictive power vs the random effects model.

A Appendix: Functions

```
1 ##### Demean #####
2
3 demean <- function(x) {x - mean(x)}
4
5
6
7
8 ##### Density Plotter #####
9
10
11 mcmc.dens <- function(combined.results) {
12   #
13   # combined results is of class "mcmc".
14   # Output you would get from mcmc.combine(coda.samples(...))
15   #
16
17   # Dataframe conversion
18   mcmc.df <- as.data.frame(combined.results)
19
20   # Parameter names
21   param.names <- colnames(mcmc.df)
22
23   for (name in param.names){
24     # Parameter vector
25     param.samples <- mcmc.df[, name]
26     main = paste("Density of", name)
27     plot(density(param.samples), main = main)
28   }
29 }
30
31
32 ##### Thinning Checker #####
33
34 thincheck <- function(results.obj, var.idx, thinhigh, thinby, chainsize){
35   # var.idx: Specify variable to observe ESS for
36   # chainsize: Size of each individual chain in results.obj
37   #
38   #
39   #
40   #
41   #
42
43   # Thinning intervals to consider, slice to remove 0
44   thins <- seq(from = 0, to = thinhigh, by = thinby)
45   thins <- thins[2:length(thins)]
46
47   # Vector to store effective sample sizes corresponding to
48   # various thinning intervals
49   ess <- numeric(length(thins))
50
51   # i: Index of vector add to
```

```

52 i      <- 1
53
54 for (thin in thins) {
55   # Create indices to take from posterior sample
56   thin.vec <- seq(from=1, to=chainsize, by = thin)
57   # Get thinned sample by slicing posterior sample using
58   # above indices
59   tempsamp <- results$obj[[1]][, var.idx][thin.vec]
60   # Multiply by 3 to get ESS of combined chain
61   ess[i]    <- 3*effectiveSize(tempsamp)
62   i = i + 1
63 }
64 plot(thins, ess)
65 }
66
67
68 ##### RESULTS TABLE #####
69
70 results.table <- function(combres, dig = 3){
71   sum.res <- summary(combres)
72   quants.res <- sum.res$quantiles
73   stats.res <- sum.res$statistics
74   rep.res <- cbind(quants.res, stats.res)
75   rep.res <- round(rep.res, dig)
76
77   CI.95 <- paste('(', rep.res[, '2.5%'], ',', rep.res[, '97.5%'], ')', sep='')
78   CI.95.width <- round(rep.res[, '97.5%'] - rep.res[, '2.5%'], dig)
79
80   rep.res <- cbind(rep.res, CI.95)
81   colnames(rep.res) <- replace(colnames(rep.res), length(colnames(rep.res)), '95% CI')
82
83   rep.res <- cbind(rep.res, CI.95.width)
84   colnames(rep.res) <- replace(colnames(rep.res), length(colnames(rep.res)), '95% CI Width')
85   return(rep.res)
86 }
87
88
89 ##### QUESTION 2 EXCLUSIVE #####
90
91 ##### BARPLOTS FOR EDA #####
92
93 barplot.2 <- function(height.name, xlab, col = c('Green', 'Red'), srt = 90){
94   ### Create barplot table
95   t.height <- table(cows$parasite, cows[, height.name])
96   t.height.pc <- prop.table(t.height, margin = 2)
97   t.height.pc <- round(t.height.pc, 2)
98   # Get sample size to put at top of barplot
99   sample.size <- apply(t.height, MARGIN = 2, sum)
100
101   ### Create barplot ###
102
103   main <- paste('Parasite Proportion by', Hmisc::capitalize(xlab))
104

```

```

105 mybar <- barplot(t.height.pc, main = main,
106                 xlab = xlab, col = col,
107                 ylim = c(0, 1.2), las = 2)
108
109 ### Label with sample size
110 text(mybar, y = 1.1, label = paste('n=', sample.size, sep = ''), srt = srt)
111
112 }
113
114 ##### Mann-Whitney U-statistic with boxplot #####
115
116 boxplot.mw <- function(df, y.name, var.names){
117   # y must be categorical, with 0 and 1
118
119   ## Split dataset into 2 for readability
120   df.0 <- df[df[, y.name] == 0, ]
121   df.1 <- df[df[, y.name] == 1, ]
122
123   for (i in 1:length(var.names)) {
124
125     var.name <- var.names[i]
126     ## Different vectors
127     zero.vec <- df.0[, var.name]
128     one.vec <- df.1[, var.name]
129     ## Mann-Whitney
130     MW.p <- wilcox.test(zero.vec, one.vec)$p.value
131     MW.p <- round(MW.p, 3)
132
133     ## Boxplot
134     main <- paste(var.name, '~', y.name)
135     sub <- paste('Mann-Whitney p-value for median difference:', MW.p)
136     xlab <- Hmisc::capitalize(y.name)
137     formula <- as.formula(main)
138     boxplot(formula = formula, data = df, xlab = xlab, main = main, sub = sub)
139   }
140 }
141
142
143
144
145
146
147 ##### FUNCTIONS FOR QUESTION 2D AND E #####
148
149 ### EXTRACT POSTERIOR SAMPLES ###
150
151 farmprobs <- function(xt, Bt, alpha){
152   # xt: row of data of class matrix: Dimensions 1 x k (k parameters)
153   # Bt: Matrix of coefficient samples: Dimensions k x S (S MCMC samples)
154   # alpha: row of samples of farm-specific intercept, class matrix: Dimensions 1 x S
155
156   # Add 1 on for random intercept influence to Xt, and add alpha to Bt.
157   xt.a <- cbind(xt, 1)

```

```

158 Bt.a <- rbind(Bt, alpha)
159
160 ## Linear combination
161 lc.farm <- xt.a%*%Bt.a
162
163 ## Probability, plogis is the inverse logit function
164 pr.farm <- plogis(lc.farm)
165
166 ## Return probability vector
167 return(pr.farm)
168 }
169
170
171 #### SUMMARISE POSTERIOR SAMPLES ####
172
173 results.post <- function(farmID, probvec, digits=2) {
174   #### Get mean, CI, probability of epidemic
175   mean.farm <- mean(probvec)
176   ci.farm <- as.vector(quantile(probvec, probs = c(0.02,0.975)))
177   pr.epi.farm <- mean(probvec > 0.2)
178
179   # Get confidence interval as a string
180   ci.farm.string <- paste('(', round(ci.farm[1], digits=digits), ', ',
181                           round(ci.farm[2], digits=digits), ')', sep='')
182
183   # Concatenate all results
184   res <- round(cbind(ci.farm[1], ci.farm[2], pr.epi.farm, mean.farm), digits=digits)
185   res <- cbind(res, ci.farm.string)
186   colnames(res) <- c('2.5%', '97.5%', 'P(Epidemic)', 'Mean', '95% CI')
187   rownames(res) <- paste('Farm', farmID)
188   return(res)
189 }

```


B Appendix: Question 1

B.1 1a

```
1 ##### Load packages and data
2
3 require(rjags)
4 require(runjags)
5 require(gridExtra)
6
7
8 gulls <- read.csv('Q1/Data/gulls_data.csv')
9
10
11 year.raw <- gulls$year
12
13 png('Q1/Q1a.png', width = 1200, height = 250)
14 par(mfrow = c(1, 4))
15
16 ## Univariate plot of Audouin
17 boxplot(gulls$audouin, main = 'Boxplot of Audouin', pch = 16, cex = 1.5)
18
19
20 ## Plot Audouin gulls per year
21
22 plot(gulls$year, gulls$audouin, xlab = 'Year', ylab = 'Audouin',
23       main = 'Audouin abundance over time', pch = 'x')
24
25 abline(lm(audouin~year, data = gulls))
26 abline(lm(audouin~year, data = gulls[gulls$audouin < 525, ]), col = 'red')
27 abline(lm(audouin~year, data = gulls[gulls$audouin < 475, ]), col = 'green')
28
29
30 # Label years 1995 to 1997 for both time series plots
31 # Done by slicing the data to only include these years for the text addition
32
33 years <- c(1995, 1996, 1997)
34
35 text(x      = gulls$year[gulls$year%in%years],
36      y      = gulls$audouin[gulls$year%in%years],
37      labels = year.raw[gulls$year%in%years])
38
39
40
41 ## Plot yellowlegged gulls against year
42
43 plot(gulls$year, gulls$yellowlegged, xlab = 'Year', ylab = 'Yellowlegged',
44       main = 'Yellowlegged abundance over time')
45
46 abline(lm(yellowlegged~year, data = gulls))
47
48
49
```

```

50 text(x      = gulls$year[gulls$year%in%years],
51      y      = gulls$yellowlegged[gulls$year%in%years],
52      labels = year.raw[gulls$year%in%years])
53
54
55 ## Plot yellowlegged against Audouin
56
57 plot(gulls$yellowlegged, gulls$audouin, xlab = 'Yellowlegged', ylab = 'Audouin',
58      main = 'Yellowlegged vs Audouin')
59 abline(lm(audouin~yellowlegged, data = gulls))
60 abline(lm(audouin~yellowlegged, data = gulls[gulls$audouin < 525, ]), col='red')
61 abline(lm(audouin~yellowlegged, data = gulls[gulls$audouin < 475, ]), col='green')
62
63 text(x      = gulls$yellowlegged[gulls$year%in%years],
64      y      = gulls$audouin[gulls$year%in%years],
65      labels = year.raw[gulls$year%in%years])
66
67
68
69 par(mfrow = c(1,1))
70 dev.off()
71
72 ##### Get correlations
73
74 ## All data
75
76 cor(gulls$audouin, gulls$year)
77 cor(gulls$audouin, gulls$yellowlegged)
78
79 ##### Remove outliers for audouin
80
81 sort(gulls$audouin)
82
83 # [1]   3  20  21  25  28  30  41  50  59  60  62  70  79
84 # [14]  80 100 120 170 187 201 225 275 300 430 476 525 625
85
86 slice <- function(boundary) {gulls$audouin < boundary}
87
88
89 upp <- 525
90 cor(gulls$audouin[slice(upp)], gulls$yellowlegged[slice(upp)])
91 cor(gulls$audouin[slice(upp)], gulls$year[slice(upp)])
92
93 upp <- 625
94 cor(gulls$audouin[slice(upp)], gulls$yellowlegged[slice(upp)])
95 cor(gulls$audouin[slice(upp)], gulls$year[slice(upp)])

```

B.2 1b

```
1 ## Demean covariates
2 gulls$yellowlegged <- demean(gulls$yellowlegged)
3 gulls$year          <- demean(gulls$year)
4
5 ##### DATA #####
6 ## Dataset
7 n      <- nrow(gulls)
8 aud    <- gulls$audouin
9 year   <- gulls$year
10
11 ## Prior
12 b0.mu  <- 0
13 b1.mu  <- 0
14 b0.tau <- 0.01
15 b1.tau <- 0.01
16
17 ## data list
18
19 data <- list(n = n, aud = aud, year = year,
20             b0.mu = b0.mu, b0.tau = b0.tau,
21             b1.mu = b1.mu, b1.tau = b1.tau)
22
23
24 ##### MODEL #####
25
26 modstr.1b <- "model{
27
28   b0 ~ dnorm(b0.mu, b0.tau)
29   b1 ~ dnorm(b1.mu, b1.tau)
30
31   # likelihood
32
33   for (i in 1:n){
34     aud[i] ~ dpois(mu[i])
35     log(mu[i]) = b0 + b1*(year[i])
36   }
37
38 }"
39
40
41
42 m1.b <- jags.model(textConnection(modstr.1b), data = data, n.chains = 3)
43
44 update(m1.b, 20000)
45
46
47 ##### RESULTS #####
48 res.1b <- coda.samples(m1.b, c('b0', 'b1'), n.iter = 20000, thin = 4) ## thin to make following
   steps quicker
49
50 ## COMBINE CHAINS
51 combres.1b <- combine.mcmc(res.1b)
```

```

52
53
54 ##### CONVERGENCE, ESS #####
55
56 png( 'Q1/Q1bConvergence.png', width = 1200, height = 200)
57
58 par(mfrow = c(1, 4))
59
60 traceplot(res.lb)
61 mcmc.dens(combres.lb)
62 par(mfrow=c(1,1))
63
64 dev.off()
65
66 autocorr.plot(combres.lb)
67 effectiveSize(combres.lb)
68 gelman.diag(res.lb)
69 gelman.plot(res.lb)
70
71 ##### CONSISTENCY CHECK #####
72
73 ## Change priors of b0 and b1, and then compare quantiles
74 ## Reduced iterations for speed
75
76 modstr.lb <- modstr.lb
77
78 taus <- c(0.001, 0.01, 0.1, 1, 10)
79
80 sumlist <- list()
81
82 # Loop over index
83 for (tau.idx in 1:length(taus)) {
84
85   # Extract the prior being tested
86   tau.prior <- taus[tau.idx]
87   # Create label for the results to add to sumlist. of the form 'tau=0.01'
88   listlabel <- paste('tau', '=', tau.prior, sep='')
89
90   # Alter data to have new priors, and create model
91   data$b0.tau <- tau.prior
92   data$b1.tau <- tau.prior
93   ml.b.check <- jags.model(textConnection(modstr.lb),
94                             data = data,
95                             n.chains = 1)
96   # Update and extract results, create summary object from results
97   update(ml.b.check, 20000)
98
99   res.lb.check <- coda.samples(ml.b.check, c('b0', 'b1'), n.iter = 30000, thin = 4)
100   sum.check <- summary(res.lb.check)
101
102   # Store results in sumlist
103   sumlist[[listlabel]] <- sum.check$quantiles
104 }

```

```

105
106 # View results
107 sumlist
108
109 ##### REPORT RESULTS #####
110
111 ## EXTRACT RESULTS AS SUMMARY ##
112
113 restab.1b <- results.table(combres.1b)
114
115 png('Q1/Q1bResults.png', width = 720, height = 70)
116 grid.table(restab.1b)
117 dev.off()

```

B.3 1c

```
1 ##### DATA RESPECIFICATION #####
2
3 data$b0.tau <- 0.01
4 data$b1.tau <- 0.01
5
6 ##### MODEL #####
7
8 modstr.1c <- "model{
9
10
11 ### Priors on beta
12 b0 ~ dnorm(b0.mu, b0.tau)
13 b1 ~ dnorm(b1.mu, b1.tau)
14
15 ### Hyperprior information
16 sigma.epsilon ~ dunif(0, 10)
17 tau.epsilon = pow(sigma.epsilon, -2)
18
19
20
21 # likelihood
22 for (i in 1:n){
23     aud[i] ~ dpois(mu[i])
24     log(mu[i]) = b0 + b1*(year[i]) + epsilon[i]
25     epsilon[i] ~ dnorm(0, tau.epsilon)
26 }
27
28 }"
29
30 ml.c <- jags.model(textConnection(modstr.1c), data = data, n.chains = 3)
31
32 update(ml.c, 100000)
33
34
35 ## 5:20 mins for 1000000 iterations, 3 chains
36 # Thin for ease of following steps
37 start_time <- Sys.time()
38 res.1c <- coda.samples(ml.c,
39                        c('b0', 'b1', 'sigma.epsilon'),
40                        n.iter = 1000000,
41                        thin = 400)
42
43 end_time <- Sys.time()
44 end_time - start_time
45
46
47 # Combine results to 1 chain
48 combres.1c <- combine.mcmc(res.1c)
49
50
51 ##### CHECK HOW MUCH TO THIN #####
52 # only run if this is 1
```



```

53
54 runthincheck <- 0
55
56 if (runthincheck == 1){
57
58 thincheck(res.lc, 1, 1000, 50, dim(res.lc[[1]])[1])
59 abline(v = c(100, 200, 300, 400, 500), col = c('red', 'blue', 'green', 'purple', 'pink'))
60
61 }
62
63
64
65
66 ##### CONVERGENCE, ESS #####
67
68 ### TRACE, DENSITY PLOTS ###
69
70 png('Q1/Q1cConvergence.png', width = 1200, height = 270)
71 # 2 x 3 plot window
72 # 2: Trace and density
73 # 3: Parameters
74 par(mfrow = c(2, 3))
75
76 # All the trace plots
77 traceplot(res.lc)
78
79 # Density plot of each parameter
80 # Check functions for documentation
81 mcmc.dens(combres.lc)
82 par(mfrow=c(1,1))
83
84 dev.off()
85
86 ### AUTOCORR, GELMAN, ESS ###
87 autocorr.plot(combres.lc)
88 effectiveSize(combres.lc)
89 gelman.diag(res.lc)
90 gelman.plot(res.lc)
91
92
93 ##### REPORT RESULTS #####
94
95 ## EXTRACT RESULTS AS SUMMARY ##
96 sum.lc <- summary(combres.lc)
97
98
99 ## REPORT TABLE ##
100
101 restab.lc <- results.table(combres.lc)
102
103 png('Q1/Q1cResults.png', width = 750, height=85)
104 grid.table(restab.lc)
105 dev.off()

```

B.4 1d

```
1 ##### DATA RESPECIFICATION #####
2
3 # Already been demeaned
4 yellowlegged <- gulls$yellowlegged
5
6 # Add to data list
7 # Add after year, so next to covariates
8 #data <- append(data, list(yel = yellowlegged), 3)
9
10 data$yel <- yellowlegged
11 # Add priors to list
12 data$b2.mu <- 0
13 data$b2.tau <- 0.01
14
15
16 ##### MODEL #####
17
18 modstr.1d <- "model{
19
20 # Parameter priors
21
22 b0 ~ dnorm(b0.mu, b0.tau)
23 b1 ~ dnorm(b1.mu, b1.tau)
24 b2 ~ dnorm(b2.mu, b2.tau)
25
26 # Extra variation hyperprior details
27 # sigma uniformly distributed
28 # tau is inverse sigma squared (inverse variance)
29
30 sigma.epsilon ~ dunif(0, 10)
31 tau.epsilon <- pow(sigma.epsilon, -2)
32
33
34 # likelihood
35 for (i in 1:n){
36   aud[i] ~ dpois(mu[i])
37   log(mu[i]) = b0 + b1*(year[i]) + b2*(yel[i]) + epsilon[i]
38   epsilon[i] ~ dnorm(0, tau.epsilon)
39 }
40
41 }"
42
43
44 m1.d <- jags.model(textConnection(modstr.1d), data = data, n.chains = 3)
45
46 update(m1.d, 200000)
47
48
49 ## RUNTIME: 2000000: 13m
50 ## RUNTIME: 1500000: 10m
51 ## RUNTIME: 1000000: 7m
52 start_time <- Sys.time()
```

```

53 res.ld <- coda.samples(m1.d, c('b0', 'b1', 'b2', 'sigma.epsilon'),
54                          n.iter = 1400000, thin = 500)
55
56 end_time <- Sys.time()
57 end_time - start_time
58
59 effectiveSize(res.ld)
60
61 ### Combined results
62 combres.ld <- combine.mcmc(res.ld)
63
64 ##### CHECK HOW MUCH TO THIN #####
65 runthincheck <- 0
66
67 if (runthincheck == 1){
68
69   thincheck(res.ld, 3, 1000, 50, dim(res.ld[[1]])[1])
70   abline(v = c(100, 200, 300, 400, 500), col = c('red', 'blue', 'green', 'purple', 'pink'))
71
72 }
73
74
75
76
77
78 ##### CONVERGENCE #####
79 ## Save trace and densityplots
80 png('Q1/Q1dConvergence.png', width = 1200, height = 300)
81 par(mfrow = c(2, 4))
82 traceplot(res.ld)
83 mcmc.dens(combres.ld)
84 par(mfrow=c(1,1))
85 dev.off()
86
87 ### AUTOCORR, GELMAN, ESS ###
88 autocorr.plot(combres.ld)
89 effectiveSize(combres.ld)
90 gelman.diag(res.ld)
91 gelman.plot(res.ld)
92
93
94
95 ## REPORT TABLE ##
96
97 restab.ld <- results.table(combres.ld)
98
99 png('Q1/Q1dResults.png', width = 750, height=110)
100 grid.table(restab.ld)
101 dev.off()

```

B.5 1e

```
1 sampler.1e.2 <- function(X, B, eps){
2   # X: nxp matrix
3   # B: pxS matrix
4   # eps: 1xS matrix, each row same
5
6   S <- ncol(B)
7   X <- cbind(X, 1)
8   Be <- rbind(B, eps)
9
10  ## Linear combination and sampling
11  lc <- X%*%Be
12  mu <- exp(lc)
13  samp <- t(apply(mu, 1, FUN=rpois, n = S))
14
15
16  ## Summary statistics
17  means <- apply(samp, 1, mean)
18  lower <- apply(samp, 1, quantile, probs = 0.05)
19  upper <- apply(samp, 1, quantile, probs = 0.95)
20
21  ## Results reporting
22  ret.mat <- cbind(means, lower, upper)
23  colnames(ret.mat) <- c('mean', 'lower', 'upper')
24  ret.df <- as.data.frame(ret.mat)
25
26
27  return(ret.df)
28 }
29
30 ##### PLOTTER.1E #####
31 plotter.1e <- function(ret.df, cols=c('blue', 'red', 'red'), main = '', cex = 1){
32
33   upperlim <- max(ret.df$upper)
34   ## BASE PLOT
35   plot(x = year.raw, y = aud, ylim = c(0, upperlim),
36        main = main, xlab = 'Year', ylab = '',
37        cex = cex, cex.main = cex, cex.axis = cex, cex.lab = cex)
38
39   ## MEAN PLOT
40   lines(x = year.raw, y = ret.df$mean, col = cols[1], pch = 'x', cex = cex)
41   lines(x = year.raw, y = ret.df$lower, col = cols[2], lwd = 2, lty = 2)
42   lines(x = year.raw, y = ret.df$upper, col = cols[3], lwd = 2, lty = 2)
43 }
44
45
46
47
48
49
50 ##### MODEL 1B #####
51
52 Xb <- cbind(1, gulls[, 'year'])
```

```

53 Bb <- t(as.matrix(combres.1b))
54 epsb <- 0
55
56 ret.dfb <- sampler.1e.2(X = Xb, B=Bb, eps = 0)
57
58
59 ##### MODEL 1C #####
60 Sc <- nrow(combres.1c)
61
62
63 Xc <- cbind(1, gulls[, 'year'])
64 Bc <- t(as.matrix(combres.1c[, c('b0', 'b1')]))
65 eps <- rbind(rnorm(Sc, 0, sd = combres.1c[, 'sigma.epsilon']), NULL)
66
67 ret.dfc <- sampler.1e.2(Xc, Bc, eps)
68
69
70 ##### MODEL 1D #####
71
72 Sd <- nrow(combres.1d)
73
74 Xd <- as.matrix(cbind(1, gulls[, c('year', 'yellowlegged')]))
75 Bd <- t(as.matrix(combres.1d[, c('b0', 'b1', 'b2')]))
76 epsd <- rbind(rnorm(Sd, 0, sd = combres.1d[, 'sigma.epsilon']), NULL)
77
78 ret.dfd <- sampler.1e.2(Xd, Bd, epsd)
79
80 ##### PLOTS #####
81
82 cex <- 2
83
84 png('Q1/Q1eFinalPlots.png', width = 1200, height = 500)
85 par(mfrow = c(1,3),
86     mar = c(5, 4,4,2)+ 1)
87 plotter.1e(ret.df = ret.dfb, main = 'Audouin ~ Year', cex = cex)
88 title(ylab = 'Audouin Abundance', cex.lab = 1.7)
89 plotter.1e(ret.df = ret.dfc, main = 'Audouin ~ Year + variance', cex = cex)
90 plotter.1e(ret.df = ret.dfd, main = 'Audouin ~ Year + Yellowlegged + variance', cex = cex)
91 par(mfrow=c(1,1))
92 dev.off()
93
94
95
96
97 ##### REPORT FIGURES #####
98
99 # Capitalise first letter of a string
100 capitalise <- function(x) {paste(toupper(substr(x,1,1)), substr(x,2,nchar(x)),sep='')}
101
102 df1 <- cbind.data.frame(ret.dfb, ret.dfc, ret.dfd)
103
104 model.ind <- paste(c(rep('1b', 3), rep('1c', 3), rep('1d', 3)))
105

```

```
106 df1.names <- colnames(df1)
107 df1.names <- paste(df1.names, model.ind)
108 colnames(df1) <- capitalise(df1.names)
109 df1 <- as.data.frame(round(as.matrix(df1), 0))
110 rownames(df1) <- year.raw
111
112 png('Q1/Q1eResults.png', width = 650, height = 600)
113 grid.table(df1)
114 dev.off()
```


B.6 1f

```
1 ##### DIC #####
2
3 # No need to update, already burned in.
4
5 ## 16m
6
7
8 start <- Sys.time()
9 dic.1b <- dic.samples(model = m1.b, n.iter = 20000, type='pD')
10 dic.1c <- dic.samples(model = m1.c, n.iter = 1000000, type='pD')
11 dic.1d <- dic.samples(model = m1.d, n.iter = 1400000, type='pD')
12 end <- Sys.time()
13 runtime <- end-start
14 runtime
15
16
17 c1 <- sum(dic.1b$deviance) + sum(dic.1b$penalty)
18 c2 <- sum(dic.1c$deviance) + sum(dic.1c$penalty)
19 c3 <- sum(dic.1d$deviance) + sum(dic.1d$penalty)
20
21 dic.mat <- cbind(c1, c2, c3)
22 colnames(dic.mat) <- c('Model 1B', 'Model 1C', 'Model 1D')
23 rownames(dic.mat) <- 'DIC'
24 dic.mat <- round(dic.mat, 2)
25 dic.df <- as.data.frame(dic.mat)
26
27 png('Q1/Q1fDIC.png', width = 250, height = 50)
28 grid.table(dic.df)
29 dev.off()
```

C Appendix: Question 2

C.1 2a

```
1 #### Funcs from Q1 ####
2
3 demean <- function(x) {x - mean(x)}
4
5
6 ##### Load #####
7
8 require(rjags)
9 require(runjags)
10 require(gridExtra)
11 require(corrplot)
12 require(dplyr)
13
14
15 cows <- read.csv('Q2/Data/cows.csv')
16
17 # --- Correlation between variables of farm environment
18
19 #### To avoid inflation of correlation, get unique
20 env.names <- c('temp', 'rain', 'permeab', 'hight', 'slope')
21 env.df <- cows[, env.names]
22 env.df.unique <- distinct(env.df)
23
24 png('Q2/Q2aEDA.png', width = 930, height = 600)
25 par(mfrow = c(3, 3))
26 ## Correlation Plot ##
27 corrplot(corr = cor(env.df.unique), type = 'lower', method = 'number', order = 'hclust')
28
29 # --- Farm barplot data
30
31 barplot.2('farmID', xlab = 'Farm ID', srt = 90)
32
33 # --- Age barplot data
34
35 barplot.2('age', xlab = 'Age', srt = 90)
36
37 # --- Boxplots
38 # Remove cowID, farmID, and parasite for boxplot variables
39 var.names <- colnames(cows)[!(colnames(cows) %in% c('cowID', 'farmID', 'parasite'))]
40 boxplot.mw(cows, 'parasite', var.names = var.names)
41 par(mfrow = c(1,1))
42 dev.off()
```

C.2 2b

```

1 ## ===== ##
2 ##### 2b: Mod #####
3 ## ===== ##
4
5
6 ##### Extract as vectors #####
7
8 para <- cows$parasite
9
10 ### Fixed ###
11 # Specific #
12 age <- demean(cows$age)
13
14 # Environment #
15 temp <- demean(cows$temp)
16 rain <- demean(cows$rain)
17 perm <- demean(cows$permeab)
18 height <- demean(cows$hight)
19 slop <- demean(cows$slope)
20
21 # Random
22 farm <- cows$farmID
23
24
25 ##### Data #####
26 ## Dataset ##
27 n <- nrow(cows)
28 J <- max(farm)
29 para <- para
30 age <- age
31 temp <- temp
32 rain <- rain
33 perm <- perm
34 height <- height
35 slop <- slop
36 farm <- farm
37
38 ## Prior ##
39
40 ##### Same priors for each beta
41 beta.mu <- 0
42 beta.tau <- 0.01
43
44 ## Hyperpriors ##
45 # Priors on the random parameter alpha
46 sig.alpha.ub <- 20
47
48 ## DATA LIST ##
49
50 data <- list(n = n, J = J, # Loop idx
51             para = para, age = age, temp = temp, rain = rain, # Covariates
52             perm = perm, height = height, slop = slop,

```

```

53     farm = farm ,
54     beta.mu = beta.mu, beta.tau = beta.tau ,           # Priors
55     sig.alpha.ub = sig.alpha.ub) # Hyperpriors
56
57
58 ##### MODEL #####
59 modstr.2b <- "model{
60
61   # Likelihood
62   for (i in 1:n) {
63     para[i] ~ dbern(p[i])
64     # alpha is the random farm-specific intercept.
65     logit(p[i]) = b0 + alpha[farm[i]] + b1*age[i] + b2*temp[i] + b3*rain[i] +
66                 b4*perm[i] + b5*height[i] + b6*slop[i]
67   }
68
69   # Priors
70   b0 ~ dnorm(beta.mu, beta.tau)
71   b1 ~ dnorm(beta.mu, beta.tau)
72   b2 ~ dnorm(beta.mu, beta.tau)
73   b3 ~ dnorm(beta.mu, beta.tau)
74   b4 ~ dnorm(beta.mu, beta.tau)
75   b5 ~ dnorm(beta.mu, beta.tau)
76   b6 ~ dnorm(beta.mu, beta.tau)
77
78   for (j in 1:J){
79     alpha[j] ~ dnorm(0, tau.alpha)
80   }
81
82   # Hyperpriors #
83   sig.alpha ~ dunif(0, sig.alpha.ub)
84   tau.alpha = pow(sig.alpha, -2)
85
86
87 }"
88
89 m.2b <- jags.model(textConnection(modstr.2b), data = data, n.chains = 3)
90
91
92 var.names <- c('b0', 'b1', 'b2', 'b3', 'b4', 'b5', 'b6', 'alpha', 'mu.alpha', 'sig.alpha')
93
94
95
96 update(m.2b, 50000)
97
98
99 # 12m: 150000 iterations
100 # 8.7m: 125000?
101 start_time <- Sys.time()
102
103 res.2b <- coda.samples(m.2b, variable.names = var.names, n.iter = 125000, thin = 100)
104
105 end_time <- Sys.time()

```

```

106 end_time - start_time
107
108 ## Combine ##
109 # From runjags
110 combres.2b <- combine.mcmc(res.2b)
111
112
113
114 ##### CHECK HOW MUCH TO THIN #####
115 # only run if this is 1
116
117 runthincheck <- 0
118
119 if (runthincheck == 1){
120
121     thincheck(res.2b, 'b2', 500, 50, dim(res.2b[[1]])[1])
122     abline(v = c(100, 200, 300, 400, 500), col = c('red', 'blue', 'green', 'purple', 'pink'))
123
124 }
125
126
127
128
129
130
131 ##### CONVERGENCE #####
132 gelman.diag(res.2b) # All 1, upper CI 1.02
133 effectiveSize(combres.2b) # Minimum of 1800
134
135 png('Q2/Q2bTrace.png', width = 1200, height = 800)
136 par(mfrow = c(5, 6))
137 traceplot(res.2b)
138 par(mfrow=c(1,1))
139 dev.off()
140
141
142
143 ##### RESULTS #####
144
145 restab.2b <- results.table(combres.2b)
146 png('Q2/Q2bResults.png', width = 800, height = 600)
147 grid.table(restab.2b)
148 dev.off()
149
150
151 ##### Further #####
152
153 mean(combres.2b[, 'b1'] > 0)

```

C.3 2c

```
1 ##### TEMP REMOVED MODEL STRING #####
2
3 modstr.2c.rtemp <- "model{
4
5   # Likelihood
6   for (i in 1:n) {
7 para[i] ~ dbern(p[i])
8 # alpha is the random farm-specific intercept.
9 logit(p[i]) = b0 + alpha[farm[i]] + b1*age[i] +
10               #b2*temp[i] +
11               b3*rain[i] + b4*perm[i] + b5*height[i] + b6*slop[i]
12 }
13
14 # Priors
15 b0 ~ dnorm(beta.mu, beta.tau)
16 b1 ~ dnorm(beta.mu, beta.tau)
17 #b2 ~ dnorm(beta.mu, beta.tau)
18 b3 ~ dnorm(beta.mu, beta.tau)
19 b4 ~ dnorm(beta.mu, beta.tau)
20 b5 ~ dnorm(beta.mu, beta.tau)
21 b6 ~ dnorm(beta.mu, beta.tau)
22
23 for (j in 1:J){
24 alpha[j] ~ dnorm(0, tau.alpha)
25 }
26
27 # Hyperpriors #
28 sig.alpha ~ dunif(0, sig.alpha.ub)
29 tau.alpha = pow(sig.alpha, -2)
30
31 }"
32
33 ##### TEMP, RAIN REMOVED MODEL STRING #####
34
35 modstr.2c.rtemprain <- "model{
36
37
38   # Likelihood
39   for (i in 1:n) {
40 para[i] ~ dbern(p[i])
41 # alpha is the random farm-specific intercept.
42 logit(p[i]) = b0 + alpha[farm[i]] + b1*age[i] +
43               #b2*temp[i] + b3*rain[i] +
44               b4*perm[i] + b5*height[i] + b6*slop[i]
45 }
46
47 # Priors
48 b0 ~ dnorm(beta.mu, beta.tau)
49 b1 ~ dnorm(beta.mu, beta.tau)
50 #b2 ~ dnorm(beta.mu, beta.tau)
51 #b3 ~ dnorm(beta.mu, beta.tau)
52 b4 ~ dnorm(beta.mu, beta.tau)
```



```

53 b5 ~ dnorm(beta.mu, beta.tau)
54 b6 ~ dnorm(beta.mu, beta.tau)
55
56 for (j in 1:J){
57   alpha[j] ~ dnorm(0, tau.alpha)
58 }
59
60 # Hyperpriors #
61 sig.alpha ~ dunif(0, sig.alpha.ub)
62 tau.alpha = pow(sig.alpha, -2)
63
64 }"
65
66 ##### TEMP, HEIGHT REMOVED MODEL STRING #####
67
68 modstr.2c.rtempheight <- "model{
69
70
71   # Likelihood
72   for (i in 1:n) {
73     para[i] ~ dbern(p[i])
74   # alpha is the random farm-specific intercept.
75   logit(p[i]) = b0 + alpha[farm[i]] + b1*age[i] +
76                 #b2*temp[i] +
77                 b3*rain[i] + b4*perm[i] +
78                 #b5*height[i] +
79                 b6*slop[i]
80   }
81
82   # Priors
83   b0 ~ dnorm(beta.mu, beta.tau)
84   b1 ~ dnorm(beta.mu, beta.tau)
85   #b2 ~ dnorm(beta.mu, beta.tau)
86   b3 ~ dnorm(beta.mu, beta.tau)
87   b4 ~ dnorm(beta.mu, beta.tau)
88   #b5 ~ dnorm(beta.mu, beta.tau)
89   b6 ~ dnorm(beta.mu, beta.tau)
90
91   for (j in 1:J){
92     alpha[j] ~ dnorm(0, tau.alpha)
93   }
94
95   # Hyperpriors #
96   sig.alpha ~ dunif(0, sig.alpha.ub)
97   tau.alpha = pow(sig.alpha, -2)
98
99   }"
100
101 ##### MODEL INITIALISATION #####
102
103 ## Note: Should be warnings about unused variables in data
104
105 m2c.rtemp <- jags.model(textConnection(modstr.2c.rtemp), data = data, n.chains = 2)

```

```

106
107 m2c.rtemprain <- jags.model(textConnection(modstr.2c.rtemprain), data = data, n.chains = 2)
108
109 m2c.rtempheight <- jags.model(textConnection(modstr.2c.rtempheight), data = data, n.chains = 2)
110
111 ##### BURN IN #####
112 update(m2c.rtemp, 20000)
113 update(m2c.rtemprain, 20000)
114 update(m2c.rtempheight, 20000)
115
116
117 ##### DIC SAMPLES #####
118
119 n.iter <- 100000 ## 18 minutes
120
121 start <- Sys.time()
122 dic.full <- dic.samples(m.2b, n.iter = n.iter)
123 dic.rtemp <- dic.samples(m2c.rtemp, n.iter = n.iter)
124 dic.rtemprain <- dic.samples(m2c.rtemprain, n.iter = n.iter)
125 dic.rtempheight <- dic.samples(m2c.rtempheight, n.iter = n.iter)
126 end <- Sys.time()
127 end-start
128
129
130
131 dic.full.val <- sum(dic.full$deviance) + sum(dic.full$penalty)
132 dic.rtemp.val <- sum(dic.rtemp$deviance) + sum(dic.rtemp$penalty)
133 dic.rtemprain.val <- sum(dic.rtemprain$deviance) + sum(dic.rtemprain$penalty)
134 dic.rtempheight.val <- sum(dic.rtempheight$deviance) + sum(dic.rtempheight$penalty)
135
136 ## Results matrix
137 dic.all <- cbind(dic.full.val, dic.rtemp.val, dic.rtemprain.val, dic.rtempheight.val)
138 colnames(dic.all) <- c('Full Model', 'Temp Removed', 'Temp, Rain Removed', 'Temp, Height Removed')
139 rownames(dic.all) <- c('DIC')
140 dic.all <- round(dic.all, 2)
141 dic.all <- as.data.frame(dic.all)
142
143
144
145 png('Q2/Q2cDIC.png', width = 500, height = 50)
146 grid.table(dic.all)
147 dev.off()
148
149 ##### PART 2, PARAMETERS OF CHOSEN MODEL #####
150
151 ## Recompile model with 3 chains
152 m2c.rtempheight <- jags.model(textConnection(modstr.2c.rtempheight), n.chains = 3, data = data)
153
154 ## Burn-in
155 update(m2c.rtempheight, 50000)
156
157 ## Determine variables to keep (remove b2 and b5)

```

```

158 var.names <- c('b0', 'b1', 'b2', 'b3', 'b4', 'b5', 'b6', 'alpha', 'sig.alpha')
159 var.names.rtempheight <- var.names[!(var.names %in% c('b2', 'b5'))]
160
161
162 ## Generate samples
163
164 start <- Sys.time()
165 res.2c <- coda.samples(m2c.rtempheight, var.names.rtempheight, n.iter = 150000, thin = 50)
166 end <- Sys.time()
167 end - start
168
169 combres.2c <- combine.mcmc(res.2c)
170
171
172 ##### CONVERGENCE #####
173 ## Traceplots
174 png('Q2/Q2cTrace.png', width = 1200, height=800)
175 par(mfrow = c(4, 6))
176 traceplot(res.2c)
177 par(mfrow=c(1,1))
178 dev.off()
179
180
181 effectiveSize(res.2c)
182 gelman.diag(res.2c)
183 gelman.plot(res.2c)
184
185 ##### RESULTS #####
186
187 restab.2c <- results.table(combres.2c)
188 png('Q2/Q2cResults.png', width = 800, height = 600)
189 grid.table(restab.2c)
190 dev.off()
191
192 ##### RESULTS EVALUATION #####
193 b6.2c <- combres.2c[, 'b6']
194 mean(b6.2c < 0.01 & b6.2c > -0.01)
195
196 ##### RESULTS COMPARISON #####
197 restab.2b <- results.table(combres.2b)
198 restab.2b <- as.data.frame(restab.2b)
199 restab.2c <- as.data.frame(restab.2c)
200
201 ### JOIN ###
202 joined <- merge(restab.2c, restab.2b, by = 0)
203
204 ### Meandif
205
206 meandif <- cbind.data.frame(joined$Mean.x, joined$Mean.y)
207 meandif
208 View(meandif)

```

C.4 2d

```
1 #### Get common function inputs for both farms
2
3 # Data matrix and coefficients
4 # 1 (intercept), age, rain, permeability, slope
5 # age = 0 since assuming mean age (mean(age) - mean(age)) since demeaned variable
6 X <- cbind(1, age=0, rain, perm, slop)
7 B <- combres.2c[, c('b0', 'b1', 'b3', 'b4', 'b6')]
8 Bt <- t(as.matrix(B))
9
10 ##### Farm 1 #####
11 ID <- 1
12
13 # Get farm 1 matrix
14 X.2d.1 <- X[cows$farmID == ID, ]
15 # Extract 1st row, since all values same for a given farm
16 xt.2d.1 <- t(as.matrix(X.2d.1[1, ]))
17
18 # Extract farm-specific effect.
19 a.2d.1 <- combres.2c[, 'alpha[1]']
20
21 probs.2d.1 <- farmprobs(xt = xt.2d.1, Bt = Bt, alpha = a.2d.1)
22 post.2d.1 <- results.post(1, probs.2d.1, 3)
23
24 ##### Farm 6 #####
25
26 ID <- 6
27
28 # Get farm 1 matrix
29 X.2d.6 <- X[cows$farmID == ID, ]
30 # Extract 1st row, since all values same for a given farm
31 xt.2d.6 <- t(as.matrix(X.2d.6[1, ]))
32
33 # Extract farm-specific effect.
34 a.2d.6 <- combres.2c[, 'alpha[6]']
35
36 probs.2d.6 <- farmprobs(xt = xt.2d.6, Bt = Bt, alpha = a.2d.6)
37 post.2d.6 <- results.post(6, probs.2d.6, 3)
38
39 ##### SUMMARISE RESULTS #####
40
41 ## Combined posterior summary ##
42
43 Q2d.results <- rbind(post.2d.1, post.2d.6)
44 rownames(Q2d.results) <- paste(rownames(Q2d.results), '(Rand)')
45 Q2d.results <- Q2d.results[, c('Mean', '95% CI', 'P(Epidemic)')]
46
47 ##### REPORT RESULTS #####
48 png('Q2/Q2dResultsFARM.png', height = 80, width = 300)
49 grid.table(Q2d.results)
50 dev.off()
51
52 ## Histograms ##
```

```

53
54
55 png('Q2/Q2dHists.png', width = 800, height = 300)
56 par(mfrow = c(1, 2))
57 hist(probs.2d.1, main = 'Histogram of Farm 1 Proportions', xlab = 'Farm 1 Proportion')
58 abline(v = c(post.2d.1[, '2.5%'], post.2d.1[, '97.5%'], post.2d.1[, 'Mean']),
59        col = c('Red', 'Red', 'Blue'), lwd = 2)
60
61 hist(probs.2d.6, main = 'Histogram of Farm 6 Proportions', xlab = 'Farm 6 Proportion')
62 abline(v = c(post.2d.6[, '2.5%'], post.2d.6[, '97.5%'], post.2d.6[, 'Mean']),
63        col = c('Red', 'Red', 'Blue'), lwd = 2)
64 par(mfrow = c(1,1))
65 dev.off()

```

```

1 modstr.2e <- "model{
2
3
4 # Likelihood
5   for (i in 1:n) {
6     para[i] ~ dbern(p[i])
7     # alpha is the FIXED farm-specific intercept.
8     logit(p[i]) = alpha[farm[i]] + b1*age[i] +
9     #b2*temp[i] +
10    b3*rain[i] + b4*perm[i] +
11    #b5*height[i] +
12    b6*slop[i]
13   }
14
15 # Priors
16 b1 ~ dnorm(beta.mu, beta.tau)
17 #b2 ~ dnorm(beta.mu, beta.tau)
18 b3 ~ dnorm(beta.mu, beta.tau)
19 b4 ~ dnorm(beta.mu, beta.tau)
20 #b5 ~ dnorm(beta.mu, beta.tau)
21 b6 ~ dnorm(beta.mu, beta.tau)
22
23   for (j in 1:J){
24     alpha[j] ~ dnorm(beta.mu, beta.tau)
25   }
26 }"
27
28
29
30
31 mod.2e <- jags.model(textConnection(modstr.2e), data = data, n.chains = 3)
32
33 update(mod.2e, 100000)
34
35 var.names.2e <- c('b0', 'b1', 'b3', 'b4', 'b6', 'alpha')
36
37
38 start <- Sys.time()
39 res.2e <- coda.samples(mod.2e, variable.names = var.names.2e, n.iter = 400000,
40                       thin = 200)
41 Sys.time() - start
42
43 combres.2e <- combine.mcmc(res.2e)
44
45 ##### CONVERGENCE #####
46 effectiveSize(res.2e)
47
48 png('Q2/Q2eTrace.png', width = 1200, height=800)
49 par(mfrow = c(4, 6))
50 traceplot(res.2e)
51 par(mfrow=c(1,1))
52 dev.off()

```

```

53
54 gelman.diag(res.2e)
55
56
57 ##### RESULTS #####
58 restab.2e <- results.table(combres = combres.2e)
59
60
61 ##### REPORT RESULTS OF PARAMETERS #####
62 png('Q2/Q2eResults.png', width = 800, height = 600)
63 grid.table(restab.2e)
64 dev.off()
65
66
67
68 ##### FARM SPECIFIC #####
69
70 ## Recall objects of interest:
71 # X: Same rows are used, so just copy
72 # No b0, so remove the first element with '-1' as column index
73 xt.2e.1 <- t(as.matrix(xt.2d.1[-1]))
74 xt.2e.6 <- t(as.matrix(xt.2d.6[-1]))
75
76 # B, alpha: Need to extract new parameters
77 B <- combres.2e[, c('b1', 'b3', 'b4', 'b6')]
78 Bt <- t(as.matrix(B))
79
80 ##### FARM 1
81 ID <- 1
82
83 # Extract farm-specific effect.
84 a.2e.1 <- combres.2e[, 'alpha[1]']
85
86 probs.2e.1 <- farmprobs(xt = xt.2e.1, Bt = Bt, alpha = a.2e.1)
87 post.2e.1 <- results.post(1, probs.2e.1, 5)
88
89
90 ##### FARM 6
91
92 ID <- 6
93
94 # Extract farm-specific effect.
95 a.2e.6 <- combres.2e[, 'alpha[6]']
96
97 probs.2e.6 <- farmprobs(xt = xt.2e.6, Bt = Bt, alpha = a.2e.6)
98 post.2e.6 <- results.post(6, probs.2e.6, 5)
99
100 ##### SUMMARISE RESULTS #####
101
102 ## Combined posterior summary ##
103
104 Q2e.results <- rbind(post.2e.1, post.2e.6)
105 rownames(Q2e.results) <- paste(rownames(Q2e.results), '(Fixed)')

```

```

106 Q2e.results <- Q2e.results[, c('Mean', '95% CI', 'P(Epidemic)')]
107
108 ##### REPORT RESULTS #####
109 png('Q2/Q2eResultsFARM.png', height = 80, width = 340)
110 grid.table(Q2e.results)
111 dev.off()
112
113 ## Histograms ##
114
115
116 png('Q2/Q2eHists.png', width = 800, height = 300)
117 par(mfrow = c(1, 2))
118 hist(probs.2e.1, main = 'Histogram of Farm 1 Proportions', xlab = 'Farm 1 Proportion')
119 abline(v = c(post.2e.1[, '2.5%'], post.2e.1[, '97.5%'], post.2e.1[, 'Mean']),
120        col = c('Red', 'Red', 'Blue'), lwd = 2)
121
122 hist(probs.2e.6, main = 'Histogram of Farm 6 Proportions', xlab = 'Farm 6 Proportion')
123 abline(v = c(post.2e.6[, '2.5%'], post.2e.6[, '97.5%'], post.2e.6[, 'Mean']),
124        col = c('Red', 'Red', 'Blue'), lwd = 2)
125 par(mfrow = c(1,1))
126 dev.off()

```


D Appendix: Question 3

```
1 ##### MODEL STRING #####
2
3 modstr.3 <- "model{
4
5
6 # Likelihood
7 for (i in 1:n) {
8   para[i] ~ dnorm(p[i], para.tau)
9
10   p[i] = b0 + b1*age[i] +
11           b3*rain[i] +
12           b4*perm[i] +
13           b6*slop[i]
14
15 }
16
17 # Priors
18 b0 ~ dnorm(beta.mu, beta.tau)
19 b1 ~ dnorm(beta.mu, beta.tau)
20 b3 ~ dnorm(beta.mu, beta.tau)
21 b4 ~ dnorm(beta.mu, beta.tau)
22 b6 ~ dnorm(beta.mu, beta.tau)
23 para.tau ~ dgamma(0.01, 0.01)
24
25 }"
26
27 m.3 <- jags.model(textConnection(modstr.3), data = data, n.chains = 3)
28
29 ##### DIC #####
30 update(m.3, 50000)
31 dic.3 <- dic.samples(m.3, 150000)
32
33 dic.3
34
35
36 ##### MODEL #####
37 res.3 <- coda.samples(m.3, c('b0', 'b1', 'b3', 'b4', 'b6'), n.iter = 50000, thin = 10)
38
39 combres.3 <- combine.mcmc(res.3)
40
41
42 ##### CONVERGENCE #####
43
44 png('Q2/Q3Trace.png', width = 700, height = 400)
45 par(mfrow = c(2, 3))
46 traceplot(res.3)
47 par(mfrow = c(1, 1))
48 dev.off()
49
50 ##### RESULTS #####
51 restab.3 <- results.table(combres.3, dig = 4)
```

```

52
53 png('Q2/Q3Results.png', width = 750, height = 150)
54 grid.table(restab.3)
55 dev.off()
56
57 ### COMP ###
58
59 mean(combres.3[, 'b0'] > 0)
60 mean(combres.3[, 'b1'] > 0)
61 mean(combres.3[, 'b3'] > 0)
62 mean(combres.3[, 'b4'] < 0)
63 mean(combres.3[, 'b6'] < 0)
64
65 mean(combres.2c[, 'b0'] > 0)
66 mean(combres.2c[, 'b1'] > 0)
67 mean(combres.2c[, 'b3'] > 0)
68 mean(combres.2c[, 'b4'] < 0)
69 mean(combres.2c[, 'b6'] < 0)
70
71
72 ### PREDICTIONS ###
73 X <- cbind(1, age, rain, perm, slop)
74 Bt <- t(as.matrix(combres.3))
75
76 lc <- X%*%Bt
77 pr <- apply(lc, 1, mean)
78 table(pr>0.5, para)

```