Homework #02 Statistical Methods in Data Science II & Lab 2021/2022 Michelangelo Saveriano 06/15/2022 Michelangelo Saveriano Excercise 1 Illustrate the characteristics of the statistical model for dealing with the Dugong's data The data we are considering describe the lengths  $(Y_i)$  and ages  $(x_i)$  of 27 dugongs captured off the coast of Queensland. Dugong 0 0 2.6 2.4 0 0 2.0 <u>κ</u> 25 30 15 20 Age As we can see from the plot above age and length are **not linearly related**, that's why the following regression model is considered:  $Y_i \sim N(\mu_i, au^2) \ \mu_i = f(x_i) = lpha - eta \gamma^{x_i}$ Model parameters are  $lpha\in(1,\infty)$  ,  $eta\in(1,\infty)$  ,  $\gamma\in(0,1)$  ,  $au^2\in(0,\infty)$  . Let us consider the following prior distributions:  $lpha \sim N(0,\sigma_lpha^2)$  $eta \sim N(0,\sigma_{eta}^2)$  $\gamma \sim Unif(0,1)$  $au^2 \sim IG(a,b)(InverseGamma)$ Derive the corresponding likelihood function The Likelihood function, up to a proportionality constant, is:  $L_y(lpha,eta,\gamma, au^2) \ = \prod_{i=1}^n f(y_i|lpha,eta,\gamma, au^2)$  $\propto \prod_{i=1}^n rac{1}{ au^2} \, \exp\!\left(-rac{(y_i-\mu_i)^2}{2 au^2}
ight)$  $\propto au^{-2n} \; \expigg(-rac{1}{2 au^2} \sum_{i=1}^n (y_i - lpha + eta \gamma^{x_i})^2igg)$ Write down the expression of the joint prior distribution of the parameters at stake and illustrate your suitable choice for the hyperparameters. We can assume the parameter being independent each other and therefore write the joint prior distribution as the product of the marginal  $\pi(\alpha, \beta, \gamma, \tau^2) = \pi(\alpha)\pi(\beta)\pi(\gamma)\pi(\tau^2)$  $\sim rac{ au^{2(-a-1)}}{\sigma_{lpha}\sigma_{eta}} \mathrm{exp}igg(-rac{lpha^2}{2\sigma_{lpha}} - rac{eta^2}{2\sigma_{eta}} - rac{b}{ au^2}igg) 1_{[1,\infty]}(lpha) \ 1_{[1,\infty]}(eta) \ 1_{[0,1]}(\gamma) \ 1_{[0,\infty]}( au)$ The values we choose for the hyperparameters are the following: • Since we have no prior information on the parameter lpha we impose an high value for its variance:  $\sigma_lpha=1$ • Since we have no prior information on the parameter eta we impose an high value for its variance:  $\sigma_eta=1$ • Since  $au^2$  describes the length variance, assuming  $\mu_i=c\ orall\ i$  , we impose  $\mathbb{E}[ au^2]=\mathbb{V}ar[Y_i]$  and  $\mathbb{V}ar[ au^2]=1$  for the same reason above. Given these constraints we get a=2.00569, b=0.07586Derive the functional form (up to proportionality constants) of all full-conditionals. Which distribution can you recognize within standard parametric families so that direct simulation from full conditional can be easily implemented? •  $\pi(lpha|eta,\gamma, au^2,y^{obs}),\ lpha\in[1,+\infty]$  $\pi(lpha|eta,\gamma, au^2,y^{obs}) = \pi(lpha)L_{y^{obs}}(lpha,eta,\gamma, au^2)$  $\propto \expigg(-rac{lpha^2}{2\sigma_lpha}igg) \expigg(-rac{1}{2 au^2}\sum_{i=1}^n(y_i-lpha+eta\gamma^{x_i})^2igg)$  $=\expigg(-rac{lpha^2}{2\sigma_lpha}igg)\expigg(-rac{1}{2 au^2}\sum_{i=1}^n \left(y_i^2-2y_i(lpha-eta\gamma^{x_i})+(lpha-eta\gamma^{x_i})^2
ight)igg)$  $\propto \exp\!\left(-rac{lpha^2}{2\sigma_lpha}
ight) \exp\!\left(-rac{1}{2 au^2}\sum_{i=1}^n\left(-2y_ilpha+lpha^2-2lphaeta\gamma^{x_i}
ight)
ight)$  $=\expigg(-rac{lpha^2}{2\sigma_lpha}igg)\expigg(-rac{1}{ au^2}\Bigg[\Big(-\sum_{i=1}^n(y_i+eta\gamma^{x_i})\Big)lpha+rac{n}{2}lpha^2\Bigg]igg)$  $=\exp\!\left(lpha\left(rac{1}{ au^2}\sum_{i=1}^n(y_i+eta\gamma^{x_i})
ight)-rac{lpha^2}{2}igg(rac{n}{ au^2}+rac{1}{\sigma_lpha^2}igg)
ight)$ Please notice that if the heta's pdf can be expressed as  $p( heta)\sim \expig(-rac12a heta^2-b hetaig)$  then  $heta\sim N(\mu=rac{b}{a},\sigma^2=rac{1}{a})$  . Using the relation above we can derive:  $\pi(lpha|eta,\gamma, au^2,y^{obs}) \sim N\left(\mu=rac{b}{a},\sigma^2=rac{1}{a}
ight)$  $where \left\{egin{array}{l} a = rac{n}{ au^2} + rac{1}{\sigma_lpha^2} \ b = rac{1}{ au^2} \sum_{i=1}^n (y_i + eta \gamma^{x_i}) \end{array}
ight.$ •  $\pi(eta|lpha,\gamma, au^2,y^{obs}),\ eta\in[1,+\infty]$  $\pi(eta|lpha,\gamma, au^2,y^{obs}) = \pi(eta)L_{y^{obs}}(lpha,eta,\gamma, au^2)$  $\propto \expigg(-rac{eta^2}{2\sigma_eta}igg) \expigg(-rac{1}{2 au^2}\sum_{i=1}^n(y_i-lpha+eta\gamma^{x_i})^2igg)$  $=\expigg(-rac{eta^2}{2\sigma_eta}igg)\expigg(-rac{1}{2 au^2}\sum_{i=1}^n \Big(y_i^2-2y_i(lpha-eta\gamma^{x_i})+(lpha-eta\gamma^{x_i})^2\Big)igg)$  $\propto \exp\!\left(-rac{eta^2}{2\sigma_eta}
ight) \exp\!\left(-rac{1}{2 au^2}\sum_{i=1}^n\left(2y_ieta\gamma^{x_i}-2lphaeta\gamma^{x_i}+eta^2\gamma^{2x_i}
ight)
ight)$  $=\expigg(-rac{eta^2}{2\sigma_eta}igg)\expigg(-rac{1}{ au^2}igg[\Big(-\sum_{i=1}^n(y_i\gamma^{x_i}-lpha\gamma^{x_i})\Big)eta+rac{\gamma^{2x_i}}{2}eta^2igg]igg)$  $=\expigg(etaigg(rac{1}{ au^2}\sum_{i=1}^n(lpha\gamma^{x_i}-y_i\gamma^{x_i})igg)-rac{eta^2}{2}igg(rac{1}{ au^2}\sum_{i=1}^n\gamma^{2x_i}+rac{1}{\sigma_eta^2}igg)igg)$ Using the relation above we can derive:  $\pi(eta|lpha,\gamma, au^2,y^{obs}) \sim N\left(\mu=rac{b}{a},\sigma^2=rac{1}{a}
ight)$  $where \left\{egin{array}{l} a = rac{1}{ au^2} \sum_{i=1}^n \gamma^{2x_i} + rac{1}{\sigma_eta^2} \ b = rac{1}{ au^2} \sum_{i=1}^n (lpha \gamma^{x_i} - y_i \gamma^{x_i}) \end{array}
ight.$ Please notice that, due to their support constraints, both  $\alpha, \beta$  follow a **truncated normal distribution**. •  $\pi(\gamma|\alpha,\beta,\tau^2,y^{obs})$  $\pi(\gamma|lpha,eta, au^2,y^{obs}) \ \propto \expigg(-rac{1}{2 au^2}\sum_{i=1}^n(y_i-\mu_i)^2igg) 1_{[0,1]}(\gamma)$  $=\expigg(-rac{1}{2 au^2}\sum_{i=1}^n(y_i-(lpha-eta\gamma^{x_i}))^2igg)1_{[0,1]}(\gamma)$  $\propto \expigg(-rac{1}{2 au^2}\sum_{i=1}^n(2y_ieta\gamma^{x_i}-2lphaeta\gamma^{x_i}+eta^2\gamma^{2x_i})igg)1_{[0,1]}(\gamma)$ Differently from the other parameters the  $\gamma$  full-conditional doesn't follow a standard distribution, so we have to use an algorithm like Metropolis Hasting to draw values from it. •  $\pi(\tau^2|\alpha,\beta,\gamma,y^{obs})$  $\pi( au^2|lpha,eta,\gamma,y^{obs}) \ \propto ( au^2)^{-a-1} \expigg(-rac{b}{ au^2}igg) au^{-2rac{n}{2}} \expigg(-rac{1}{2 au^2}igg(\sum_{i=1}^n (y_i-lpha+eta\gamma^{x_i})^2igg)igg)$  $=( au^2)^{-a-rac{n}{2}-1}\expigg(-rac{1}{2 au^2}igg(2b+\sum_{i=1}^n(y_i-lpha+eta\gamma^{x_i})^2igg)igg).$ From this we can clearly see that the  $au^2$  full-conditional has an **Inverse Gamma** shape:  $\pi( au^2|lpha,eta,\gamma,y^{obs}) \sim IG\left(a+rac{n}{2},b+rac{1}{2}\sum_{i=1}^n(y_i-lpha+eta\gamma^{x_i})^2
ight).$ Using a suitable Metropolis-within-Gibbs algorithm simulate a Markov chain (T=10000) to approximate the posterior distribution for the above model Gibbs Sampling The Gibbs sampling is an algorithm which allows us to simulate a sequence of dependent random variables starting from all the full conditionals. The algorithm: 1. fix the starting values of the parameter components at time t=0:  $oldsymbol{ heta}^0 = ( heta_1, heta, \dots, heta_k)$ 2. for  $t=1,\dots,T$  iterate the following cycle:  $heta_j^{t+1} \sim \pi( heta_j, heta_{(j)}) = \pi( heta_j | heta_1^{t+1}, \ldots, heta_{j-1}^{t+1}, heta_{j+1}^t, \ldots, heta_k^t) \;\; j = 1, \ldots, k$ Metropolis-Hastings Metropolis-Hastings is an algorithm which allows us to simulate a sequence of dependent random variables drawn from a proposal distribution  $p_x(y)$ , where x is the current state of the chain. Unlike the Gibbs sampler which relies on conditional distribution, the Metropolis-Hastings algorithm uses the joint distribution to generate a candidate draws. The algorithm works as follow: 1. Draw a candidate  $Y_{t+1}=y\sim p_x(y)$  2 Decide whether or not the candidate is accepted as the next state of the chain at time t+1 $X_{t+1} = egin{cases} y & ext{with probability } lpha(x,y) \ x & ext{with probability } 1-lpha(x,y) \end{cases}$  $lpha(x,y) = \min\left\{rac{\pi(y)p_y(x)}{\pi(x)p_x(y)},1
ight\}$ Metropolis-within-Gibbs Gibbs sampling requires the knowledge of the full conditionals to work, unfortunately, since  $\gamma$  does not follow a well-known parametric distribution, we're lacking this information. To overcome this issue we can use the Metropolis-within-Gibbs algorithm which allows us to replace the original i-th kernel  $K_i$  with one related to the Metropolis algorithm  ${ ilde K}_i^{MET_i}$  . To perform such simulations we'll use the Jags library. parameters <- c("alpha", "beta", "gamma", "tau2",</pre> "Ypred\_20", "cond\_exp\_20", "Ypred\_30", "cond\_exp\_30") dugongjags <- jags(data=mydata,</pre> # let JAGS choose the initial values parameters.to.save=parameters, model.file="dugong\_jags\_model.txt", DIC = F,n.chains=1, n.burnin=1, n.iter=10000) ## module glm loaded ## module dic loaded ## Compiling model graph Resolving undeclared variables Allocating nodes ## Graph information: Observed stochastic nodes: 27 Unobserved stochastic nodes: 6 Total graph size: 134 ## Initializing model print(dugongjags) ## Inference for Bugs model at "dugong\_jags\_model.txt", fit using jags, ## 1 chains, each with 10000 iterations (first 1 discarded), n.thin = 9 ## n.sims = 1111 iterations saved mu.vect sd.vect 2.5% 25% 50% 75% 97.5% ## Ypred\_20 2.592 0.135 2.301 2.507 2.594 2.678 2.853 ## Ypred\_30 2.620 0.142 2.345 2.527 2.622 2.714 2.893 ## alpha 2.642 0.087 2.479 2.584 2.637 2.697 2.822 ## beta ## tau2 0.004 0.009 0.012 0.014 0.017 0.025 0.015 Show the 4 univariate trace-plots of the simulations of each parameter As we can see from the trace-plots below all the parameters converge to some stationary distribution. OL. 3.0 -2.8 -2.6 -2.4 -2693 8093 5393 2.00 -1.75 -1.50 -1.25 -1.00 -Chain alpha 2693 5393 8093 beta gamma 0.8 tau2 0.6 -0.4 -2693 5393 8093 tau2 0.03 -0.02 -0.01 -2693 5393 8093 Iteration Evaluate graphically the behaviour of the empirical averages  $\hat{I}_t$  with growing  $t=1,\dots,T$ The parameter convergence can also be observed through the empirical cumulative average. OL. 2.65 --2.60 -2.55 -2.50 -2.45 -2.40 -2693 8093 5393 1.4 -1.3 -Chain 1.2 alpha 2693 5393 8093 beta gamma 0.85 0.75 tau2 0.65 -0.55 -0.45 -2693 8093 5393 tau2 0.030 -0.025 -0.020 -0.015 -2693 5393 8093 Iteration Provide estimates for each parameter together with the approximation error and explain how you have evaluated such error The estimates for each parameter can be easily computed using the empirical mean. To evaluate the approximation error we have to take into account the temporal dependency between the samples: contrary to standard MC, where the samples are iid, in a MCMC setting the current state depends on the previous one. Therefore we have to compute the MCMC variance as the MC variance divided by a penalisation term that depends on the correlation between the samples:  $\mathbb{V}ar_{MCMC}[\hat{I}\,] = rac{\mathbb{V}ar[\hat{I}\,]}{S_{eff}}$ where  $S_{eff}$  measures the effective sample size. This quantity can be interpreted as the number of independent Monte Carlo samples necessary to carry the same amount of information. **Parameters Estimates** Approximation\_Error Variance **ESS** alpha 2.6424961 0.0047367 0.0076045 338.9453 1.0645952 0.0027528 0.0044558 587.9958 beta 0.0030258 443.7499 0.8387293 0.0026113 gamma

1.a

Length

1.b

1.c

1.d/e

1.f

where

1.g

value

1.h

Running Mean

1.i

tau2

1.l

alpha

beta

tau2

alpha

beta

tau

gamma

1.n/o/p

Ypred\_20

7.5 -

2.5 -

Ypred\_20

Ypred\_30

density

2.25

2.50

Below we can see how the two predictions are distributed.

2.75

3.00

From the table above we see that the prediction Ypred\_20 is slightly less precise than the one for Ypred\_30.

value

In order to compare the precision of the two predictions we can compare the length of 95% confidence interval as we did before.

gamma

0.0147211

Which parameter has the largest posterior uncertainty? How did you measure it?

To measure the posterior uncertainty we can use the equal-tails 95% confidence interval.

lower\_limit

2.4791811

1.0020391

0.7087873

0.0087796

Which couple of parameters has the largest correlation (in absolute value)?

alpha

1.0000000

-0.1148434

0.8734139

-0.1863524

ullet cond\_exp\_20 , the conditional expectation regressed for x=20

ullet cond\_exp\_30 , the conditional expectation regressed for x=30

mean

2.591852

From the correlation matrix below we can see that the most correlated couple is  $(\alpha, \gamma)$ .

Provide the prediction of a different dugong with age 30. Which prediction is less precise?

• Ypred\_20, the predicted dugong length, normally distributed, centered in cond\_exp\_20

• Ypred\_30, the predicted dugong length, normally distributed, centered in cond\_exp\_30

sd

0.1349927

In order to perform such predictions we ask Jags to provide us samples drawn from the following distributions:

0.0001357

As we can see the parameter with the highest uncertainty (the parameter whose 95% confidence interval length is the largest) is  $\alpha$ .

beta

-0.1148434

1.0000000

-0.4312648

0.2778105

Use the Markov chain to approximate the posterior predictive distribution of the length of a dugong with age of 20 years.

upper\_limit

2.8224715

1.2295500

0.9135787

0.0249764

gamma

0.8734139

-0.4312648

1.0000000

-0.2769318

50%

2.593742

75%

2.677527

Parameter

cond\_exp\_20

cond\_exp\_30

Ypred\_20

Ypred\_30

length

0.5516561

0.5480885

Parameter

Ypred\_20

Ypred\_30

0.0000180

976.8022

length

0.3432904

0.2275109

0.2047914

0.0161968

tau

-0.1863524

0.2778105

-0.2769318

1.0000000

97.5%

2.852933

2.589237 0.0471591 2.477502 2.563777 2.669320 cond\_exp\_20 2.594298 2.622006 Ypred\_30 2.620153 0.1418062 2.345267 2.526597 2.621536 2.713773 2.893355 cond\_exp\_30 2.626529 0.0702877 2.479060 2.582269 2.628604 2.677605 2.756374 As expected Ypred\_20 and cond\_exp\_20 are centered in almost the same value. The same consideration holds for Ypred\_30 and cond\_exp\_30 as well. 20 30

2.5%

2.301277

25%

2.507194

2.50

2.75

3.00

2.25	2.50	2.75	3.00
	value		

Excercise 2 Let us consider a Markov chain  $(X_t)_{t\geq 0}$  defined on the state space  $\mathcal{S}=\{1,2,3\}$  with the following transition and the corresponding transition probability matrix: Starting at time t=0 in the state  $X_0=1$  simulate the Markov chain with distribution assigned as above for t=1000 consecutive times # State space S = 1:3# Transition probability matrix P < - matrix(c(0, 0.5, 0.5,5/8, 1/8, 1/4, 2/3, 1/3, 0), nrow = 3,byrow = T)[,1] [,2] [,3] ## [1,] 0.0000000 0.5000000 0.50 ## [2,] 0.6250000 0.1250000 0.25 ## [3,] 0.6666667 0.3333333 0.00 set.seed(1234) sim\_Markov\_chain <- function(space, P, x0, N\_steps){</pre> # Vector to store the simulated values sim\_chain <- rep(NA, N\_steps+1)</pre> sim\_chain[1] <- x0 # Simulation for(t in 1:N\_steps){ sim\_chain[t+1] <- sample(space, size=1, prob=P[sim\_chain[t], ])</pre> return(sim\_chain) # Number of steps N\_steps <- 1000 # Initial state x\_0 <- 1 # Simulated chain sim\_chain <- sim\_Markov\_chain(S, P, x\_0, N\_steps)</pre> plot(sim\_chain, type = 'b') 3.0 2.5 2.0 <del>ر</del>ن 0. 200 400 600 800 1000 Index compute the empirical relative frequency of the two states in your simulation sim\_chain Freq 0.3986014 0.3066933 0.2947053 0.3 sim\_chain ₽ 0.2-0.1 sim\_chain repeat the simulation for 500 times and record only the final state at time  $t=1000\,$  for each of the 500 simulated chains. Compute the relative frequency of the 500 final states. What distribution are you approximating in this way? Try to formalize the difference between this point and the previous point. M <- 500 t <- 1000 # Simulate M chains final\_states <- sapply(1:M, function(x) sim\_Markov\_chain(S, P, x\_0, N\_steps)[t + 1])</pre> final\_states Freq 0.380 0.328 0.292 final\_states final\_states In the previous step we were approximating the stationary distribution, ie the distribution we get when  $t o +\infty$  , while in this step we're approximating the distribution for a fixed  $t=1000\,.$ compute the theoretical stationary distribution  $\pi$  and explain how you have obtained it. Is it well approximated by the simulated empirical relative frequencies computed in (b) and (c)? Remembering that the stationary distribution  $\pi(\cdot)$  satisfy the following conditions:  $egin{aligned} \bullet & \pi_j \geq 0 \ orall \ j \in S \ egin{aligned} \bullet & \sum_{j \in S} \pi_j = 1 \ egin{aligned} \bullet & \pi P = \pi \end{aligned}$ we can compute it as the normalized eigenvector relative to the eigenvalue  $1. \,$ stationary\_vect <- eigen(t(P))\$vectors[, 1]</pre> stationary\_distr <- stationary\_vect / sum(stationary\_vect)</pre> stationary\_distr 0.3917526 0.3298969 0.2783505 0.4node node As we can see the stationary distribution  $\pi(\cdot)$  is well approximated by the simulated empirical relative frequencies computed in (b) and (c). what happens if we start at t=0 from state  $X_0=2$  instead of  $X_0=1$ ? # Initial state x\_0 <- 2 # Simulated chain sim\_chain\_2 <- sim\_Markov\_chain(S, P, x\_0, N\_steps)</pre> sim\_chain\_2 Freq 0.3836164 0.3476523 0.2687313 0.3 sim\_chain\_2 ₽ 0.2 -L 3 sim\_chain\_2 As we can see the results we get starting from  $X_0=2$  are not much different from the ones we got before in the case of  $X_0=1$  . The reason for this is because as  $t o +\infty$  the empirical relative frequencies distribution converges to the stationary distribution.

2.a

sim\_chain

2.b

1

2

3

0.4 -

0.0 -

2.c

2

3

0.3 -

₽ 0.2 -

0.1-

0.0 -

2.d/e

node

2

0.3 -

0.1 -

0.0 -

2.f

3

0.1 -

0.0 -