# CPSC532W: Homework 1

Name: Adam Jozefiak, Student No.: 27458158

1. We show that Gamma distribution is conjugate to the Poisson distribution. Let

$$\lambda \sim \text{Gamma}(\alpha, \beta)$$
$$x_1, \ldots, x_N \sim \text{Poisson}(\lambda)$$

Then by Bayes rule, the posterior distribution $p(\lambda|\mathbf{x})$ can be written as

$$p(\lambda|x_{1:N}) = \frac{p(x_{1:N}|\lambda)p(\lambda)}{p(x_{1:N})} \propto p(x_{1:N}|\lambda)p(\lambda)$$

where the last proportionallity follows by the fact that $p(x_{1:N}) = \int p(x_{1:N}|\lambda)p(\lambda)d\lambda$ is constant.

We note the probability density functions of $\text{Gamma}(\alpha, \beta)$ and $\text{Poisson}(\lambda)$:

$$p(\gamma) = \frac{\beta^\alpha}{\Gamma(\alpha)}\lambda^{\alpha-1}e^{-\beta\lambda}$$

$$p(x_{1:N}) = \prod_{i=1}^{N} \frac{\lambda^{x_i}e^{-\lambda}}{x_i!}$$

Then,

$$
\begin{aligned}
p(\lambda|x_{1:N}) \quad &\propto \quad p(x_{1:N}|\lambda)p(\lambda) && \text{as argued earlier} \\
&= \quad \left(\prod_{i=1}^{N} \frac{\lambda^{x_i}e^{-\lambda}}{x_i!}\right)\left(\frac{\beta^\alpha}{\Gamma(\alpha)}\lambda^{\alpha-1}e^{-\beta\lambda}\right) \\
&= \quad \frac{\beta^\alpha}{\prod_{i=1}^{N} x_i!}\lambda^{\sum_{i=1}^{N} x_i + \alpha - 1}e^{-(n+\beta)\lambda} \\
&\propto \quad \lambda^{\sum_{i=1}^{N} x_i + \alpha - 1}e^{-(n+\beta)\lambda} && \text{by } \lambda \text{ not appearing in } \frac{\beta^\alpha}{\prod_{i=1}^{N} x_i!}
\end{aligned}
$$

Therefore, up to a normalizing constant, $p(\lambda|x_{i:N}) \propto \lambda^{\sum_{i=1}^{N} x_i + \alpha - 1}e^{-(n+\beta)\lambda}$ and hence $p(\lambda|x_{1:N}) \sim \text{Gamma}(\sum_{i=1}^{N} x_i + \alpha, n + \beta)$. Therefore, the a prior Gamma distribution is a conjugate prior for a likelihood Poisson distribution.

2. Let $x$ and $x'$ be vectors of random variables. Then the Gibbs transition operator, for the change in the $k^{\text{th}}$ variable or index of $x$, is defined as $T_k(x, x') = p(x'_k|x_{-k}) \prod_{i\neq k} \mathbb{I}(x_i = x'_i)$. Then we can see that the Gibbs transition operator satisfies the detailed balance equattion by hte following argument

$$p(x)T(x, x')$$

$$= \quad p(x)p(x'_k|x_{-k}) \prod_{i \neq k} \mathbb{I}(x_i = x'_i)$$

$$= \quad p(x_k|x_{-k})p(x_{-k})p(x'_k|x_{-k}) \prod_{i \neq k} \mathbb{I}(x_i = x'_i) \quad \text{by the product rule}$$

$$= \quad p(x')p(x_k|x_{-k}) \prod_{i \neq k} \mathbb{I}(x_i = x'_i) \qquad\qquad \text{by the product rule on } p(x') = p(x'_k|x_{-k})p(x_{-k})$$

$$= \quad p(x')p(x_k|x'_{-k}) \prod_{i \neq k} \mathbb{I}(x_i = x'_i) \qquad\qquad \text{by } p(x_k|x_{-k}) \prod_{i \neq k} \mathbb{I}(x_i = x'_i) > 0 \Rightarrow x'_{-k} = x_{-k}$$

$$= \quad p(x')T(x', x)$$

This completes the proof of the Gibbs transition operator satisfying the detailed balance equation and therefore Gibbs transition operator can be interpreted as a Metropolis-Hastings transition operator that always accepts.

3. My solution to this problem is in Julia, based off the sample code. I found that the probability of it being cloudy given that the grass is wet is 57.58% by enumerating all possible world states and conditioning by counting which proportion are cloudy given the grass is wet. Using ancestral sampling and rejection, with 10000 successful samples, I found that the probability that it is cloudy given the grass is wet is 57.50% while 35.43% of samples were rejected. Using Gibbs sampling I found that the probability that it is cloudy given that the grass is wet is 58.66%. In my implementation, an index of 2 indicates a "true" assignment and an index 1 indicates a "false" assignment to a random variable.

Below is my code for computing the probability by enumerating all possible world states and conditioning by counting which proportion are cloudy given the grass is wet:

```julia
45   ## condition and marginalize:
46   p_C_given_W = 0.0
47   p_C_and_W = 0.0
48   p_W =  0.0
49   for c in 1:2
50       for s in 1:2
51           for r in 1:2
52               global p_W += p[c,s,r,2]
53           end
54       end
55   end
56   for s in 1:2
57       for r in 1:2
58           global p_C_and_W += p[2,s,r,2]
59       end
60   end
61   p_C_given_W = p_C_and_W / p_W
62
63   println("There is a ", p_C_given_W*100, "% chance it is cloudly given the grass is wet.")
```

Below is my code for computing the probability by ancestral sampling and rejection:

```
65  ##2. ancestral sampling and rejection:
66  num_samples = 10000
67  samples = zeros(num_samples)
68  rejections = 0
69  i = 1
70  while i <= num_samples
71      C = 0
72      S = 0
73      R = 0
74      W = 0
75      p = rand()
76      if p > p_C(1)
77          C = 2
78      else
79          C = 1
80      end
81      p = rand()
82      if p > p_S_given_C(1,C)
83          S = 2
84      else
85          S = 1
86      end
87      p = rand()
88      if p > p_R_given_C(1,C)
89          R = 2
90      else
91          R = 1
92      end
93      p = rand()
94      if p > p_W_given_S_R(1,S,R)
95          W = 2
96      else
97          W = 1
98      end
99      if W == 2
100         global samples[i] = C-1
101         global i += 1
102     else
103         global rejections += 1
104     end
105 end
106
107
108 println("The chance of it being cloudy given the grass is wet is ", (sum(samples)./num_samples)*100, "%.")
109 println(100*rejections/(num_samples+rejections), "% of the total samples were rejected.")
```

Below is my code for computing the probability by Gibbs sampling:

```julia
139  ##gibbs sampling
140  # num_samples = 10000
141  num_samples = 10000
142  samples = zeros(num_samples)
143  state = ones(Int64,4)
144  #c,s,r,w, set w = True
145  state[4] = 2
146  i = 1
147  while i <= num_samples
148
149      # Sample r
150      prob = rand()
151      if prob > p_R_given_C_S_W[state[1],state[2],1,state[4]]
152          state[3] = 2
153      else
154          state[3] = 1
155      end

157      # Sample s
158      prob = rand()
159      if prob > p_S_given_C_R_W[state[1],1,state[3],state[4]]
160          state[2] = 2
161      else
162          state[2] = 1
163      end
164
165      # Sample c
166      prob = rand()
167      if prob > p_C_given_S_R[1,state[2],state[3]]
168          state[1] = 2
169      else
170          state[1] = 1
171      end
172      samples[i] = state[1]
173      global i += 1
174
175  end
176
177
178  println("The chance of it being cloudy given the grass is wet is, ",100*sum(samples .- 1)/num_samples ,"%.")
```

4.
- We first derive the updates for Metropolis Hastings sampling of the $\hat{t}$ block, that is we want to sample from $p(\hat{t}|\hat{x}, t, x, w, \sigma^2, \alpha)$.

  Let the proposal distribution be $\hat{t}'|\hat{t} \sim$ Normal$(\hat{t}, \sigma^2)$. Then $\hat{t}'|\hat{t}$ and $\hat{t}|\hat{t}'$ are symmetric, i.e $q(\hat{t}'|\hat{t}) = q(\hat{t}|\hat{t}')$, and so $\frac{q(\hat{t}|\hat{t}')}{q(\hat{t}'|\hat{t})} = 1$. Then after sampling $\hat{t}'$ from $q(\hat{t}'|\hat{t})$ the accepting probability is:

4

$$A(\hat{t}'|\hat{t})$$

$$= \min(1, \frac{p(\hat{t}'|\hat{x},t,x,w,\sigma^2,\alpha)q(\hat{t}|\hat{t}')}{p(\hat{t}|\hat{x},t,x,w,\sigma^2,\alpha)q(\hat{t}'|\hat{t})})$$

$$= \min(1, \frac{p(\hat{t}'|\hat{x},t,x,w,\sigma^2,\alpha)}{p(\hat{t}|\hat{x},t,x,w,\sigma^2,\alpha)}) \qquad \text{by the symmetry of the proposal distribution}$$

$$= \min(1, \frac{p(\hat{t}',\hat{x},t,x,w,\sigma^2,\alpha)p(\hat{x},t,x,w,\sigma^2,\alpha)}{p(\hat{t},\hat{x},t,x,w,\sigma^2,\alpha)p(\hat{x},t,x,w,\sigma^2,\alpha)}) \qquad \text{by the product rule}$$

$$= \min(1, \frac{p(\hat{t}',\hat{x},t,x,w,\sigma^2,\alpha)}{p(\hat{t},\hat{x},t,x,w,\sigma^2,\alpha)})$$

$$= \min(1, \frac{\prod_{i=1}^{N} p(t_i|x,w,\sigma^2)p(w|\alpha)p(\hat{t}'|\hat{x},w,\sigma^2)}{\prod_{i=1}^{N} p(t_i|x,w,\sigma^2)p(w|\alpha)p(\hat{t}|\hat{x},w,\sigma^2)})$$

$$= \min(1, \frac{p(\hat{t}'|\hat{x},w,\sigma^2)}{p(\hat{t}|\hat{x},w,\sigma^2)})$$

Therefore, in the Metropolis Hastings sampling block for $\hat{t}$, we sample $\hat{t}'$ from $\text{Normal}(\hat{t}, \sigma^2)$ and we accept $\hat{t}'$ with probability $\min(1, \frac{p(\hat{t}'|\hat{x},w,\sigma^2)}{p(\hat{t}|\hat{x},w,\sigma^2)})$.

Next, we derive the updates for Metropolis Hastings sampling of the $w$ block, that is we want to sample from $p(w|\hat{t}, \hat{x}, t, x, \sigma^2, \alpha)$. We can treat the $\hat{t}$ and $\hat{x}$ as $t$ and $x$ and so we will derive the updates for Metropolis Hastings sampling of $p(w|t, x, \sigma^2, \alpha)$.

Let the proposal distribution be $w'|w \sim \text{Normal}(w, \alpha I)$. Then $w'|w$ and $w|w'$ are symmetric, i.e $q(w'|w) = q(w|w')$, and so $\frac{q(w|w')}{q(w'|w)} = 1$. Then after sampling $w'$ from $q(w'|w)$ the accepting probability is:

$$A(w'|w)$$

$$= \min(1, \frac{p(w'|t,x,\sigma^2,\alpha)q(w|w')}{p(w|t,x,\sigma^2,\alpha)q(w'|w)})$$

$$= \min(1, \frac{p(w'|t,x,\sigma^2,\alpha)}{p(w|t,x,\sigma^2,\alpha)}) \qquad \text{by the symmetry of the proposal distribution}$$

$$= \min(1, \frac{p(w',t,x,\sigma^2,\alpha)p(t,x,\sigma^2,\alpha)}{p(w,t,x,\sigma^2,\alpha)p(t,x,\sigma^2,\alpha)}) \qquad \text{by the product rule}$$

$$= \min(1, \frac{p(w',t,x,\sigma^2,\alpha)}{p(w,t,x,\sigma^2,\alpha)})$$

$$= \min(1, \frac{\prod_{i=1}^{N} p(t_i|x,w,\sigma^2)p(w|\alpha)}{\prod_{i=1}^{N} p(t_i|x,w',\sigma^2)p(w'|\alpha)})$$

Therefore, in the Metropolis Hastings sampling block for $w$, we sample $w'$ from $\text{Normal}(w, \alpha I)$ and we accept $w'$ with probability $\min(1, \frac{\prod_{i=1}^{N} p(t_i|x,w,\sigma^2)p(w|\alpha)}{\prod_{i=1}^{N} p(t_i|x,w',\sigma^2)p(w'|\alpha)})$.

- To perform pure Gibbs sampling for $w$ and $\hat{t}$ we would like to first sample $w$ from $p(w|\hat{t}, \hat{x}, t, x, \sigma^2, \alpha)$. As argued earlier, we can treat $\hat{t}$ and $\hat{x}$ as being part of $t$ and $x$ and so we might as well sample $w$ from $p(w|t, x, \sigma^2, \alpha)$. Then we have that

$$p(w|t, x, \sigma^2, \alpha)$$

$$= \frac{p(w, t, x, \sigma^2, \alpha)}{p(t, x, \sigma^2, \alpha)} \qquad \text{by the product rule}$$

$$\propto p(w, t, x, \sigma^2, \alpha)$$

$$= \prod_{i=1}^{N} p(t_i|x_i, w\sigma^2)p(w|\alpha)$$

Hence, in Gibbs sampling, we can sample $w$ proportionally to the joint distribution or specifically the distribution $\prod_{i=1}^{N} p(t_i|x_i, w\sigma^2)p(w|\alpha)$. As will be shown in the next bullet point, this distribution is the posterior of $w$ and it is

$$w|t, x, \sigma^2, \alpha \sim \text{Normal}([\frac{1}{\alpha}I + \frac{1}{\sigma^2}xx^T]^{-1}\frac{1}{\sigma^2}xt, [\frac{1}{\alpha}I + \frac{1}{\sigma^2}xx^T]^{-1})$$

Next, we need to sample $\hat{t}$ which requires sampling from $p(\hat{t}|\hat{x}, t, x, w, \sigma^2, \alpha)$. We can see that

$$p(\hat{t}|\hat{x}, t, x, w, \sigma^2, \alpha)$$

$$= \frac{p(\hat{t}, \hat{x}, t, x, w, \sigma^2, \alpha)}{p(\hat{x}, t, x, w, \sigma^2, \alpha)} \qquad \text{by the product rule}$$

$$\propto p(\hat{t}, \hat{x}, t, x, w, \sigma^2, \alpha)$$

$$= \prod_{i=1}^{N} p(t_i|x_i, w, \sigma^2)p(w|\alpha)p(\hat{t}|\hat{x}, w, \sigma^2)$$

$$\propto p(\hat{t}|\hat{x}, w, \sigma^2)$$

Hence, Gibbs sampling, we can sample $\hat{t}$ proportionally to the distriubtion $p(\hat{t}|\hat{x}, w, \sigma^2)$ which has distribution $\text{Normal}(w^T\hat{x}, \sigma^2)$.
- We compute the analytic form of the posterior predictive: $p(\hat{t}|\hat{x}, t, x, \sigma^2, \alpha)$. We begin by noting the following

$$p(\hat{t}|\hat{x}, t, x, \sigma^2, \alpha)$$

$$= \int p(\hat{t}|\hat{x}, t, x, \sigma^2, w, \alpha)ntp(\hat{t}, w|\hat{x}, t, x, \sigma^2, \alpha)dw \quad \text{by the sum rule}$$

$$= \int p(\hat{t}|\hat{x}, t, x, \sigma^2, w, \alpha)p(w|\hat{x}, t, x, \sigma^2, \alpha)dw \qquad \text{by the product rule}$$

$$= \int p(\hat{t}|\hat{x}, \sigma^2, w)p(w|t, x, \sigma^2, \alpha)dw$$

Where the last line above follows by the fact that $p(\hat{t}|\hat{x}, t, x, \sigma^2, w, \alpha) = p(\hat{t}|\hat{x}, \sigma^2, w)$ and $p(w|\hat{x}, t, x, \sigma^2, \alpha) = p(w|t, x, \sigma^2, \alpha)$, the latter following by the fact that $w$ can only have a dependency on $\hat{x}$ if $\hat{t}$ is known.

Next, we will argue that both $\hat{t}|\hat{x}, \sigma^2, w$ and $w|t, x, \sigma^2, \alpha$ are Gaussian random variables.
Firstly, $\hat{t}|\hat{x}, \sigma^2, w \sim \text{Normal}(w^T\hat{x}, \sigma^2)$ by the likelihood given in the question.
As for $w|t, x, \sigma^2, \alpha$ we note that

$$w|\alpha \sim \text{Normal}(0, \alpha I)$$

$$t|w, x, \sigma^2 \sim \text{Normal}(w^T x, \sigma^2)$$

Then we note that

$$p(w|t, x, \sigma^2, \alpha) \propto p(t|w, x, \sigma^2, \alpha)p(w, x|\sigma^2, \alpha) \propto p(t|w, x, \sigma^2)p(w|\alpha)$$

Then by 2.116, page 93, in Pattern Recognition and MAchine Learning by Bishop it follows that

$$w|t, x, \sigma^2, \alpha \sim \text{Normal}([\frac{1}{\alpha}I + \frac{1}{\sigma^2}xx^T]^{-1}\frac{1}{\sigma^2}xt, [\frac{1}{\alpha}I + \frac{1}{\sigma^2}xx^T]^{-1})$$

Then by the earlier observations that

$$p(\hat{t}|\hat{x}, t, x, \sigma^2, \alpha) \quad = \quad \int p(\hat{t}|\hat{x}, \sigma^2, w)p(w|t, x, \sigma^2, \alpha)dw$$

and

$$\hat{t}|\hat{x}, \sigma^2, w \quad \sim \quad \text{Normal}(w^T\hat{x}, \sigma^2)$$

and by linear combination rules of Gaussian random variables we have that

$$\hat{t}|\hat{x}, t, x, \sigma^2, \alpha \sim \text{Normal}([(\frac{1}{\alpha}I + \frac{1}{\sigma^2}xx^T)^{-1}\frac{1}{\sigma^2}xt]^T\hat{x}, \hat{x}^T(\frac{1}{\alpha}I + \frac{1}{\sigma^2}xx^T)^{-1}\hat{x} + \sigma^2)$$

This completes the derivation of the analytic for mof hte posterior predictive.

5.