

Advanced Probabilistic Machine Learning and Applications

Pablo Sánchez and Isabel Valera

1 Tutorial 4 (Solution Exercise 1)

Now, we are assuming the number of components is unbounded, $K \rightarrow \infty$, but it is constrained by the complexity of our dataset to a finite number K^+ . Additionally, we choose $\alpha_k = \frac{\alpha}{K}$ for convenience. We also choose to use a Dirichlet prior over the mixing coefficients of the form:

$$p(\pi|\alpha) = \text{Dir}(\pi | (\alpha/K, \dots, \alpha/K))$$

1. α : Concentration parameter related to the a priori knowledge about the number of clusters.
2. $H(\cdot)$: Prior distribution from which we are going to sample the likelihood parameters, i.e. base measure.
3. The DP provides a framework to have l samples from some distribution H with probability π_l .
4. The DP is the $\lim_{K \rightarrow \infty}$ of a Dirichlet distribution.
5. K^+ grows/decreases over the iterations. When a cluster k disappears we cannot recover it in the future nor its likelihood parameters.

Algorithm 1: Gibbs sampling with π collapsed

Since we have selected conjugate prior distribution for the mixing components π , we can marginalize them out.

Algorithm 1: π collapsed Gibbs sampling algorithm

```
Initialize  $K^+ = 1, \{z_i = 1\}_{i=1}^N$  and  $\theta_1 \sim p(\theta|\gamma)$ ;  
while not converged do  
  for  $n = 1, \dots, N$  do  
    Sample  $z_n \sim p(z_n|X, Z_{-n}, \Theta) = p(z_n|\mathbf{x}_n, Z_{-n}, \Theta)$ ;  
    if  $z_n = K^+ + 1$  then  
       $K^+ = K^+ + 1$ ;  
       $\theta_{K^+} \sim p(\theta|\mathbf{x}_n)$   
    end  
    If necessary, remove empty clusters;  
  end  
  for  $k = 1, \dots, K^+$  do  
    Sample  $\theta_k \sim p(\theta_k|X, Z)$ ;  
  end  
end
```

Posterior distribution over z_n : We first can write the posterior probability of the n -th sample belonging to cluster k is proportional to the joint distribution

$$p(z_n = k | \mathbf{x}_n, \mathbf{Z}_{-n}, \boldsymbol{\theta}) \propto p(z_n = k, \mathbf{x}_n | \mathbf{Z}_{-n}, \boldsymbol{\theta}) = p(\mathbf{x}_n | z_n = k, \mathbf{Z}_{-n}, \boldsymbol{\theta}) p(z_n = k | \mathbf{Z}_{-n})$$

which notice we must normalize the resulting distribution $\sum_k p(z_n = k | \mathbf{x}_n, \boldsymbol{\theta}, \mathbf{Z}_{-n}) = 1$. The prior term (given all previous cluster assignments) for a finite number of components was computed in Tutorial 2

$$p(z_n = k | \mathbf{Z}_{-n}) = \frac{\sum_{i \neq n} [z_i = k] + \alpha_k}{N - 1 + \sum_k \alpha_k} = \frac{m_k + \alpha_k}{N - 1 + \sum_k \alpha_k} = \frac{m_k + \alpha/K}{N - 1 + \alpha}$$

where m_k is the number of observations, except for n , assigned to component k . Notice the resulting distribution follows the scheme "rich get richer". The challenge in this Tutorial is that we are considering $K \rightarrow \infty$. The good point is that we still know that $\lim_{K \rightarrow \infty} \sum_{k=1}^K p(z_n = k | \mathbf{Z}_{-n}) = 1$. We can expand this limit in two terms: one with the clusters with observations and other with the clusters without any observation

$$\begin{aligned} \lim_{K \rightarrow \infty} \sum_{k=1}^K \frac{m_k + \alpha/K}{N - 1 + \alpha} &= \lim_{K \rightarrow \infty} \sum_{k=1}^{K^+} \frac{m_k + \alpha/K}{N - 1 + \alpha} + \sum_{k=K^++1}^K \frac{m_k + \alpha/K}{N - 1 + \alpha} \\ &= \sum_{k=1}^{K^+} \frac{m_k}{N - 1 + \alpha} + \lim_{K \rightarrow \infty} \sum_{k=K^++1}^K \frac{\alpha/K}{N - 1 + \alpha} = 1 \end{aligned}$$

so now we can compute the probability of the infinity number of empty clusters

$$\lim_{K \rightarrow \infty} \sum_{k=K^++1}^K \frac{\alpha/K}{N - 1 + \alpha} = 1 - \sum_{k=1}^{K^+} \frac{m_k}{N - 1 + \alpha} = 1 - \frac{N - 1}{N - 1 + \alpha}$$

Finally, we have that

$$p(z_n = k | \mathbf{Z}_{-n}) = \frac{m_k}{n - 1 + \alpha} \quad p(z_n = K^+ + 1 | \mathbf{Z}_{-n}) = \frac{\alpha}{n - 1 + \alpha}$$

Then, the posterior of a non-empty component is of the form

$$\begin{aligned} p(z_n = k | \mathbf{x}_n, \mathbf{Z}_{-n}, \boldsymbol{\theta}) &\propto \frac{m_k}{n - 1 + \alpha} p(\mathbf{x}_n | z_n = k, \mathbf{Z}_{-n}, \boldsymbol{\theta}) \\ &\propto \frac{m_k}{n - 1 + \alpha} p(\mathbf{x}_n | \boldsymbol{\theta}_k) \end{aligned}$$

The likelihood term has the form

$$p(\mathbf{x}_n | z_n = k, \mathbf{Z}_{-n}, \boldsymbol{\theta}) = p(\mathbf{x}_n | z_n = k, \boldsymbol{\theta}) = p(\mathbf{x}_n | \boldsymbol{\theta}_k) = \prod_{j=1}^{W_n} \text{Cat}(x_{nj} | \boldsymbol{\theta}_k)$$

On the contrary, the posterior of an empty component is of the form

$$\begin{aligned} p(z_n = K_{\text{new}} | \mathbf{x}_n, \mathbf{Z}_{-n}, \boldsymbol{\theta}) &\propto \frac{\alpha}{n - 1 + \alpha} p(\mathbf{x}_n | z_n = K_{\text{new}}, \mathbf{Z}_{-n}, \boldsymbol{\theta}) \\ &= \frac{\alpha}{n - 1 + \alpha} p(\mathbf{x}_n | z_n = K_{\text{new}}) \\ &= \frac{\alpha}{n - 1 + \alpha} \int p(\mathbf{x}_n, \boldsymbol{\theta} | z_n = K_{\text{new}}) d\boldsymbol{\theta} \\ &= \frac{\alpha}{n - 1 + \alpha} \int p(\mathbf{x}_n | z_n = K_{\text{new}}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\gamma}) d\boldsymbol{\theta} \end{aligned} \tag{1}$$

where we need to integrate all the possible likelihood parameters, θ , that can be used in the new component K_{new} . Thus, we need to compute the integral

$$\begin{aligned}
\int p(\mathbf{x}_n | z_n = K_{\text{new}}, \theta) p(\theta | \gamma) d\theta &= \\
&= \int \prod_{j=1}^{W_n} \text{Cat}(x_{nj} | \theta) \text{Dir}(\theta | \gamma) d\theta \\
&= \int \prod_{j=1}^{W_n} \prod_{m=1}^{|I|} \theta_m^{[x_{nj}=m]} \frac{1}{B(\gamma)} \prod_{m=1}^{|I|} \theta_m^{\gamma_m-1} d\theta \\
&= \int \frac{1}{B(\gamma)} \prod_{m=1}^{|I|} \theta_m^{\gamma_m + c_{nm} - 1} d\theta \\
&= \frac{B(\gamma + \mathbf{c}_n)}{B(\gamma)} \\
&= \frac{\prod_{m=1}^{|I|} \prod_{i=0}^{c_{nm}-1} (\gamma_m + i)}{\prod_{j=0}^{W_n-1} (\sum_m \gamma_m + j)}
\end{aligned} \tag{2}$$

where $c_{nm} = \sum_{j=1}^{W_n} [x_{nj} = m]$. Note we can compute this a priori and save a bunch of computational resources.

Posterior distribution over θ_k :

$$\begin{aligned}
p(\theta_k | X, Z) &\propto p(\theta_k) p(X | \theta_k, Z) \\
&= \text{Dir}(\theta_k | \gamma) \prod_n p(\mathbf{x}_n | \theta_k)^{[z_n=k]} \\
&= \text{Dir}(\theta_k | \gamma) \prod_n \prod_j \prod_{m=1}^{W_n} \text{Cat}(x_{nj} | \theta_k)^{[z_n=k]} \\
&= \text{Dir}(\theta_k | \gamma) \prod_n \prod_j \prod_{m=1}^{W_n} \prod_{m=1}^{|I|} \theta_{km}^{[x_{nj}=m][z_n=k]} \\
&= \text{Dir}(\theta_k | \gamma) \prod_{m=1}^{|I|} \theta_{km}^{c_{km}}
\end{aligned} \tag{3}$$

$$= \text{Dir}(\theta_k | \gamma'_k) \tag{4}$$

where we have used $c_{km} = \sum_n [z_n = k] \sum_j [x_{nj} = m]$ in Equation 3 which is the number of occurrences of the m -th word in the cluster k ; and $\gamma'_{km} = \gamma_m + c_{km}$ in Equation 4. Finally, for a new component we will sample $\theta_{K_{\text{new}}}$

$$\begin{aligned}
p(\theta | \mathbf{x}_n) &\propto p(\theta) p(\mathbf{x}_n | \theta) \\
&= \text{Dir}(\theta | \gamma) p(\mathbf{x}_n | \theta) \\
&= \text{Dir}(\theta | \gamma_n)
\end{aligned}$$

where $\gamma_n = \gamma_m + \sum_j [x_{nj} = m]$. We choose to sample from the posterior instead of the prior because otherwise the sampled θ will probably not explain our sample correctly.

Summary

$$p(\boldsymbol{\theta}_k | \mathbf{X}, \mathbf{Z}) = \text{Dir}(\boldsymbol{\theta}_k | \boldsymbol{\gamma}'_k) \quad \gamma'_{km} = \gamma_m + \sum_n [z_n = k] \sum_j [x_{nj} = m] \quad (5)$$

$$p(\boldsymbol{\theta}_{K_{\text{new}}} | \mathbf{x}_n) = \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\gamma}_n) \quad \gamma_{nm} = \gamma_m + \sum_j [x_{nj} = m] \quad (6)$$

$$p(z_n = k | \mathbf{x}_n, \mathbf{Z}_{-n}, \boldsymbol{\theta}) \propto \frac{m_k}{n-1+\alpha} p(\mathbf{x}_n | \boldsymbol{\theta}_k) \quad m_k = \sum_{i \neq n} [z_i = k] \quad (7)$$

$$p(z_n = K_{\text{new}} | \mathbf{x}_n, \mathbf{Z}_{-n}, \boldsymbol{\theta}) \propto \frac{\alpha}{N-1+\alpha} \frac{B(\boldsymbol{\gamma} + \mathbf{c}_n)}{B(\boldsymbol{\gamma})} \quad c_{nm} = \sum_{j=1}^{W_n} [x_{nj} = m]$$

Algorithm 2: Gibbs sampling with π and $\boldsymbol{\theta}$ collapsed

Since we have selected the conjugate prior distribution for the likelihood parameters $\boldsymbol{\theta}$, we can marginalize them out.

Algorithm 2: $\pi, \boldsymbol{\theta}$ collapsed Gibbs sampling algorithm

Initialize $K^+ = 1$ and $\{z_i = 1\}_{i=1}^N$;

while not converged **do**

for $n = 1, \dots, N$ **do**

 Sample $z_n \sim p(z_n | \mathbf{X}, \mathbf{Z}_{-n})$;

if $z_n = K_{\text{new}}$ **then**

$K^+ = K^+ + 1$;

end

 If necessary, remove empty clusters;

end

end

Posterior distribution over z_n : Firstly, we can write the posterior probability of the n -th sample belonging to cluster k is proportional to the joint distribution

$$p(z_n = k | \mathbf{x}_n, \mathbf{X}_{-n}, \mathbf{Z}_{-n}) \propto p(z_n = k, \mathbf{x}_n | \mathbf{X}_{-n}, \mathbf{Z}_{-n}) = p(z_n = k | \mathbf{Z}_{-n}) p(\mathbf{x}_n | z_n = k, \mathbf{X}_{-n}, \mathbf{Z}_{-n})$$

which notice we must normalize the resulting distribution $\sum_k p(z_n = k | \mathbf{x}_n, \mathbf{X}_{-n}, \mathbf{Z}_{-n}) = 1$. We need to distinguish between existing ones. For new clusters it can be computed using Equation 1 in Exercise 2.

$$p(z_n = K_{\text{new}} | \mathbf{X}, \mathbf{Z}_{-n}) \propto \frac{\alpha}{n-1+\alpha} \frac{B(\boldsymbol{\gamma} + \mathbf{c}_n)}{B(\boldsymbol{\gamma})}$$

For existing clusters, the prior term was computed previously. The posterior predictive can be computed marginalizing the likelihood parameters

$$\begin{aligned}
p(\mathbf{x}_n | z_n = k, \mathbf{X}_{-n}, \mathbf{Z}_{-n}) &= \int p(\mathbf{x}_n, \boldsymbol{\theta}_k | z_n = k, \mathbf{X}_{-n}, \mathbf{Z}_{-n}) d\boldsymbol{\theta}_k \\
&= \int p(\mathbf{x}_n | z_n = k, \boldsymbol{\theta}_k) p(\boldsymbol{\theta}_k | \mathbf{X}_{-n}, \mathbf{Z}_{-n}) d\boldsymbol{\theta}_k \\
&= \int \prod_{j=1}^{W_n} \text{Cat}(x_{nj} | \boldsymbol{\theta}_k) \text{Dir}(\boldsymbol{\theta}_k | \boldsymbol{\gamma}_k'') d\boldsymbol{\theta}_k \tag{8}
\end{aligned}$$

$$\begin{aligned}
&= \int \prod_{j=1}^{W_n} \prod_{m=1}^{|I|} \boldsymbol{\theta}_{km}^{[x_{nj}=m]} C \prod_{m=1}^{|I|} \boldsymbol{\theta}_{km}^{\gamma_{km}''-1} d\boldsymbol{\theta}_k \\
&= \frac{1}{B(\boldsymbol{\gamma}_k'')} \int \prod_{m=1}^{|I|} \boldsymbol{\theta}_{km}^{c_{nm} + \gamma_{km}''-1} d\boldsymbol{\theta}_k \tag{9} \\
&= \frac{B(\boldsymbol{\gamma}_k''(n) + \mathbf{c}_n)}{B(\boldsymbol{\gamma}_k''(n))}
\end{aligned}$$

where we have used the result in Exercise 1 to get $\gamma_{km}''(n) = \gamma_m + \sum_{i \neq n} [\mathbf{z}_i = k] \sum_j [\mathbf{x}_{ij} = m]$ in Equation 8; the quantity $c_{nm} = \sum_j [x_{nj} = m]$ in Equation 9 represents the number of occurrences of the m -th word in document n ; in steps 13 we use $\sum_m c_{nm} = W_n$. We can further develop the ratio between the two Beta functions

$$\begin{aligned}
\frac{B(\boldsymbol{\gamma}_k'' + \mathbf{c}_n)}{B(\boldsymbol{\gamma}_k'')} &= \frac{\prod_{m=1}^{|I|} \Gamma(\gamma_{km}'' + c_{nm})}{\Gamma(\sum_m \gamma_{km}'' + c_{nm})} \frac{\Gamma(\sum_m \gamma_{km}'')}{\prod_{m=1}^{|I|} \Gamma(\gamma_{km}'')} \\
&= \frac{\prod_{m=1}^{|I|} \prod_{i=0}^{c_{nm}-1} (\gamma_{km}'' + i) \Gamma(\gamma_{km}'')}{\Gamma(\sum_m \gamma_{km}'' + W_n)} \frac{\Gamma(\sum_m \gamma_{km}'')}{\prod_{m=1}^{|I|} \Gamma(\gamma_{km}'')} \\
&= \frac{\prod_{m=1}^{|I|} \prod_{i=0}^{c_{nm}-1} (\gamma_{km}'' + i) \Gamma(\gamma_{km}'')}{\prod_{j=0}^{W_n-1} (\sum_m \gamma_{km}'' + j) \Gamma(\sum_m \gamma_{km}'')} \frac{\Gamma(\sum_m \gamma_{km}'')}{\prod_{m=1}^{|I|} \Gamma(\gamma_{km}'')} \\
&= \frac{\prod_{m=1}^{|I|} \prod_{i=0}^{c_{nm}-1} (\gamma_{km}'' + i)}{\prod_{j=0}^{W_n-1} (\sum_m \gamma_{km}'' + j)} \tag{10}
\end{aligned}$$

We compute the log posterior predictive to avoid numerical instabilities

$$\begin{aligned}
\log p(\mathbf{x}_n | z_n = k, \mathbf{X}_{-n}, \mathbf{Z}_{-n}) &= \log \prod_{m=1}^{|I|} \prod_{i=0}^{c_{nm}-1} (\gamma_{km}'' + i) - \log \prod_{j=0}^{W_n-1} \left(\sum_m \gamma_{km}'' + j \right) \\
&= \sum_{m=1}^{|I|} \sum_{i=0}^{c_{nm}-1} \log(\gamma_{km}'' + i) - \sum_{j=0}^{W_n-1} \log \left(\sum_m \gamma_{km}'' + j \right)
\end{aligned}$$

Summary

$$p(z_n = k | \mathbf{x}_n, \mathbf{X}_{-n}, \mathbf{Z}_{-n}) \propto \frac{m_k}{N-1+\alpha} \frac{B(\boldsymbol{\gamma}_k''(n) + \mathbf{c}_n)}{B(\boldsymbol{\gamma}_k''(n))} \quad \gamma_{km}''(n) = \gamma_m + \sum_{i \neq n} [\mathbf{z}_i = k] \sum_j [\mathbf{x}_{ij} = m]$$

$$p(z_n = K_{\text{new}} | \mathbf{x}_n, \mathbf{X}_{-n}, \mathbf{Z}_{-n}) \propto \frac{\alpha}{N-1+\alpha} \frac{B(\boldsymbol{\gamma} + \mathbf{c}_n)}{B(\boldsymbol{\gamma})}$$