

Homework - 4

Pratik Mishra (14493)

Problem 1

Let z_d denote the topic assignment of all the words of the document d .

$$\therefore p(z_d | \alpha) = p(z_d | \theta_d) \cdot p(\theta_d | \alpha)$$

$$\begin{aligned} \therefore p(z_d | \theta_d) \cdot p(\theta_d | \alpha) &= \frac{\Gamma(\alpha)}{\Gamma(\alpha/k)^k} \cdot \prod_{k=1}^K \Gamma(d/k + \sum_{n=1}^{N_d} z_{dn}) \\ &= \frac{\Gamma(\alpha)}{\Gamma(\alpha/k)^k} \cdot \frac{\prod_{k=1}^K \Gamma(\alpha/k + N_d^{(k)})}{\Gamma(\alpha + N_d)} \end{aligned}$$

Let $z_{d,n}$ denote the topic assignment to every word in d document except the n^{th} word. So, since the documents are independent,

$$p(z_{d,n=k} | z_{-(d,n)}) = p(z_{d,n=k} | z_{d,n})$$

$$\therefore p(z_{d,n=k} | z_{d,n}, d) = \frac{p(z_d | \alpha)}{p(z_{d,n} | \alpha)}$$

$$p(z_{d,n} = k | z_{d,m}, \alpha) = \frac{N_d^{(k)}}{N_d - 1 + \alpha/k}$$

This makes intuitive sense as probability that the n^{th} word of a document belongs to the k^{th} topic depends on number of words in the document that belong to the same topic. Thus more words in a document of same topic implies higher probability of an unknown word to belong to that topic.

- let W_k denote the vocabulary assignment of all the words in all the documents that belong to the k^{th} topic.

$$\begin{aligned} p(W_k | \eta) &= p(W_k | \beta_k) \cdot p(\beta_k | \eta) \\ &= \frac{\Gamma(n)}{\Gamma(\eta/n)^n} \cdot \prod_{v=1}^V \frac{\Gamma(n/v + M_k^{(v)})}{\Gamma(n + M_k)} \end{aligned}$$

Here, $M_k^{(v)}$ denotes the total number of words in all the documents same as the v^{th} word in the vocabulary and M_k denotes the total ~~no. of~~ number

of words in all the documents assigned the k^{th} topic. ~~documents are partitioned into topics~~
~~whereas words are not~~

② ~~Vocabulary~~ ~~documents are partitioned into topics~~

$$\begin{aligned} & \therefore p(w_{d,n} = v | z_{d,n} = k, w_{k-(d,n)}, z_{-(d,n)}) \\ &= p(w_{d,n} = v | z_{d,n} = k, w_{k-(d,n)}) \\ & \therefore p(w_{d,n} = v | z_{d,n} = k, w_{k-(d,n)}) = \frac{p(w_k | n)}{p(w_{k-(d,n)} | n)} \\ &= \frac{M_k^{(v)} - 1 + n}{M_k + 1 + n} \end{aligned}$$

This makes intuitive sense because if more words of k^{th} topic are same as the v^{th} word in Vocabulary V , then probability of the unknown word to be the same would be higher.

PROBLEM 2

We know,

$$p(z_n = k | z_n) = \frac{m_k + \alpha / k - 1}{\alpha + N - 1}$$

$$\therefore p(x_n | z_n = k, w_{-n}, z_{-n})$$

$$= \int p(x_n | z_n = k, \mu_k, \Sigma_k) \cdot p(\mu_k, \Sigma_k | w_{-n}, z_{-n}) d\mu_k d\Sigma_k$$

$$= T(x | m_N, \frac{k_N + 1}{k_N (v_N - D + 1)} s_N, v_N - D + 1)$$

where N is the number of observations in cluster k and \bar{x} is the mean of those observations and,

$$m_N = \frac{k_0 m_0 + N \bar{x}}{k_N}$$

$$k_N = k_0 + N$$

$$v_N = v_0 + N$$

$$s_N = s_0 + \sum_{i=1}^N x_i x_i^T + k_0 m_0 m_0^T - k_N m_N m_N^T$$

so, when

$$K \rightarrow \infty,$$

$$p(z_n=k | z_n) = \frac{n_k - 1}{\alpha + N - 1}$$

where, k is the cluster already seen so far.

PROBLEM 3

$$p(z|\alpha) = \int_0^1 \int_0^1 \dots \int_0^1 \left(\prod_{n=1}^N p(z_{nk} | \pi_n) \right) p(\pi_k | \alpha) d\pi_k d\pi_1 d\pi_2 \dots$$

$$= \int_0^1 \int_0^1 \dots \int_0^1 \left(\prod_{n=1}^N \pi_k^{z_{nk}} (1 - \pi_k)^{1 - z_{nk}} \right) \frac{\pi_k^{\alpha/k - 1} \Gamma(\alpha/k + 1)}{\Gamma(\alpha/k)} d\pi_k d\pi_1 d\pi_2 \dots$$

$$= \int_0^1 \int_0^1 \dots \int_0^1 \frac{\alpha^k}{k^k} \prod_{n=1}^N \pi_k^{\sum_{n=1}^N z_{nk} + \alpha/k - 1} (1 - \pi_k)^{N - \sum_{n=1}^N z_{nk}} d\pi_k d\pi_1 d\pi_2 \dots$$

$$= \frac{\alpha^k}{k^k} \prod_{n=1}^k \left(\frac{\Gamma(\sum_n z_{nk} + \alpha/k) \Gamma(N+1 - \sum_n z_{nk})}{\Gamma(N+1 + \alpha/k)} \right)$$

$$p(z|\alpha) = \frac{\alpha^k}{(k \times \Gamma(N+1 + \alpha/k))^k} \cdot \prod_{n=1}^k \left(\Gamma(\sum_n z_{nk} + \alpha/k) \times \Gamma(N+1 - \sum_n z_{nk}) \right)$$

Let Z_{nk} be the entries of the n^{th} column

$$P(Z_{nk} | Z_{-nk}) = \frac{P(Z_n | \alpha)}{P(Z_{-nk} | \alpha)}$$

Total $= \sum_{k=1}^K$

$$= \frac{1}{N + \frac{\alpha}{K}} \times \frac{\prod_{n=1}^N (Z_{nk} + \alpha/K) \prod_{n=1}^N (N+1 - \sum_{m \neq n} Z_{mk})}{\prod_{n=1}^N (\sum_{m \neq n} Z_{mk} - Z_{nk} + \alpha/K) \prod_{n=1}^N (N - \sum_{m \neq n} Z_{mk} + Z_{nk})}$$

$$= \frac{1}{N + \frac{\alpha}{K}} \times \left(\sum_{n=1}^N Z_{nk} + \alpha/K - 1 \right)^{Z_{nk}} \left(\frac{N - \sum_{n=1}^N Z_{nk}}{N + \alpha/K} \right)^{1-Z_{nk}}$$

$$E \left[\sum_{n=1}^N Z_{nk} \right] = \sum_{n=1}^N E[Z_{nk}]$$

$$= \sum_{n=1}^N 1 \cdot P(Z_{nk} = 1)$$

$$= \sum_{n=1}^N \int p(Z_{nk}=1 | \pi_k) \cdot p(\pi_k | \alpha) d\pi_k$$

$$= \sum_{n=1}^N \int \pi_k \cdot \frac{\alpha}{K} \cdot \pi_k^{\alpha/K - 1} d\pi_k$$

$$= \sum_{n=1}^N \frac{\alpha}{K} \cdot \frac{1}{\alpha/K + 1} = \frac{Nd}{\alpha + K}$$

Total number of 1's

$$= \sum_{k=1}^K \left(E \left[\sum_{n=1}^N z_{nk} \right] \right) = \sum_{k=1}^K \frac{N\alpha}{\alpha+k}$$

$$= \frac{NK\alpha}{\alpha+k}$$