

Supplementary Material

Estimation of Non-Normalized Mixture Models

Proof of Lemma 1 (Section 3.3)

Here, we consider the non-normalized mixture model

$$p(x | \theta, \pi) = \sum_{k=1}^K \pi_k \cdot p(x | \theta_k), \quad (34)$$

where the K components belong to the same parametric family:

$$p(x | \theta_k) = \frac{1}{Z(\theta_k)} \tilde{p}(x | \theta_k). \quad (35)$$

We also use the parametrization with (θ, c) defined by

$$p(x | \theta, c) = \sum_{k=1}^K p(x | \theta_k, c_k), \quad (36)$$

where

$$\log p(x | \theta_k, c_k) = \log \tilde{p}(x | \theta_k) + c_k. \quad (37)$$

Two parametrizations are connected by the transformation $c_k = \log \pi_k - \log Z(\theta_k)$.

From Theorem 1 of Gutmann and Hyvärinen (2012) with assumption (d),

$$(\theta, c) \in \arg \max_{\theta, c} J_{\text{NCE}}(\theta, c) \quad (38)$$

if and only if

$$p(x | \theta, c) = p(x | \theta^*, c^*), \quad (39)$$

which is rewritten as

$$\sum_{k=1}^K \pi_k p(x | \theta_k) = \sum_{k=1}^K \pi_k^* p(x | \theta_k^*). \quad (40)$$

Then, from assumptions (a)-(c), it leads to

$$\begin{aligned} & \{\pi_1 p(x | \theta_1), \dots, \pi_K p(x | \theta_K)\} \\ &= \{\pi_1^* p(x | \theta_1^*), \dots, \pi_K^* p(x | \theta_K^*)\}. \end{aligned} \quad (41)$$

Therefore, there exists $\sigma \in S_n$ such that

$$\pi_k p(x | \theta_k) = \pi_{\sigma(k)}^* p(x | \theta_{\sigma(k)}^*) \quad (42)$$

for $k = 1, \dots, K$, which is equivalent to

$$\pi_k = \pi_{\sigma(k)}^*, \quad \theta_k = \theta_{\sigma(k)}^* \quad (43)$$

for $k = 1, \dots, K$ by using assumption (a). Thus, we obtain (16).

Proof of Theorem 2 (Section 4.2)

Since $n(y_t)$ does not depend on θ and c , we can rewrite (3) and (4) as

$$(\hat{\theta}_{\text{NCE}}, \hat{c}_{\text{NCE}}) = \arg \max_{\theta, c} \tilde{J}_{\text{NCE}}(\theta, c), \quad (44)$$

where

$$\begin{aligned} \tilde{J}_{\text{NCE}}(\theta, c) &= \sum_{t=1}^N \log \frac{Np(x_t | \theta, c)}{Np(x_t | \theta, c) + Mn(x_t)} \\ &\quad + \sum_{t=1}^M \log \frac{1}{Np(y_t | \theta, c) + Mn(y_t)} \end{aligned} \quad (45)$$

Similarly, since $n_l(y_t^{(l)})$ does not depend on θ and c , we can rewrite (12) and (13) as

$$(\hat{\theta}_{\text{MNCE}}, \hat{c}_{\text{MNCE}}) = \arg \max_{\theta, c} \tilde{J}_{\text{MNCE}}(\theta, c), \quad (46)$$

where

$$\begin{aligned} &\tilde{J}_{\text{MNCE}}(\theta, c) \\ &= \sum_{t=1}^N \log \frac{Np(x_t | \theta, c)}{Np(x_t | \theta, c) + M_1 n_1(x_t) + \dots + M_L n_L(x_t)} \\ &\quad + \sum_{l=1}^L \sum_{t=1}^{M_l} \log \frac{1}{Np(y_t^{(l)} | \theta, c) + M_1 n_1(y_t^{(l)}) + \dots + M_L n_L(y_t^{(l)})}. \end{aligned} \quad (47)$$

Now, from (14), we obtain $\tilde{J}_{\text{NCE}}(\theta, c) = \tilde{J}_{\text{MNCE}}(\theta, c)$. Therefore,

$$(\hat{\theta}_{\text{MNCE}}, \hat{c}_{\text{MNCE}}) = (\hat{\theta}_{\text{NCE}}, \hat{c}_{\text{NCE}}). \quad (48)$$

Additional figure for Section 6

Figure 1 shows the histogram of the logit score of the posterior probability in the first cluster $\log p(z = 1 | x; \hat{\theta}, \hat{c}) - \log(1 - p(z = 1 | x; \hat{\theta}, \hat{c}))$. Compared to the proposed method, the Gaussian mixture models assign extremely large or small logit scores.

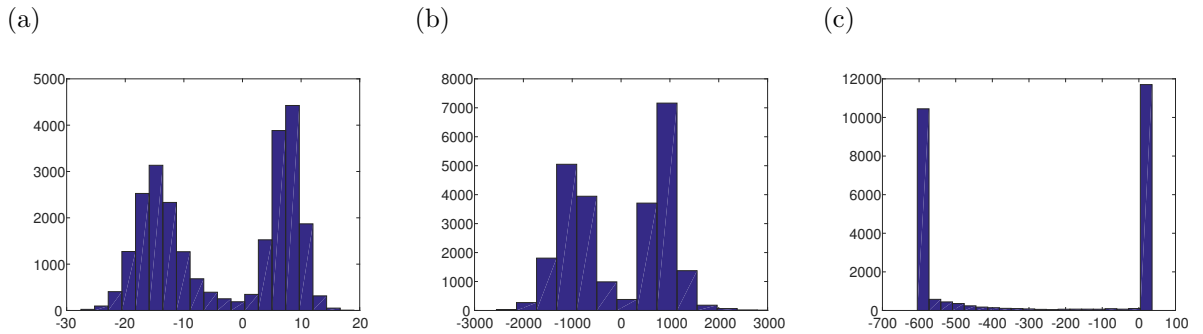


Figure 1: Histogram of the logit score of posterior probability. Please note different horizontal ranges in the three plots. (a) The proposed method. (b) Gaussian mixture model with diagonal covariance matrices. (c) Gaussian mixture model with isotropic covariance matrices.