

---

## Supplementary Material

### 1 Online Posterior Update for the Inverse Noise Variance

To update  $q_{\text{cur}}(\tau)$ , we consider the blending distribution only in terms of the NN output  $f_o$  and  $\tau$ ,

$$\tilde{p}(f_o, \tau) \propto q_{\text{cur}}(f_o)q_{\text{cur}}(\tau)\mathcal{N}(y_{i_n}|f_o, \tau^{-1}) = \mathcal{N}(f_o|\alpha_n, \beta_n)\text{Gamma}(\tau|a, b)\mathcal{N}(y_{i_n}|f_o, \tau^{-1}). \quad (1)$$

Following the conditional expectation propagation (CEP) framework proposed by Wang and Zhe (2019), we first derive the conditional moments of  $\tau$  given  $f_o$  and then approximate the expectation of the conditional moments to obtain the moments and update the posterior of  $\tau$ . Specifically, from (1), we can easily derive the conditional blending distribution,

$$\tilde{p}(\tau|f_o) = \text{Gamma}(\tau|\hat{a}, \hat{b}) \quad (2)$$

where  $\hat{a} = a + \frac{1}{2}$  and  $\hat{b} = b + \frac{1}{2}(y_{i_n}^2 - 2f_o + f_o^2)$ . We can obtain the conditional moments of  $\tau$ ,

$$\mathbb{E}_{\tilde{p}(\tau|f_o)}[\tau] = \frac{\hat{a}}{\hat{b}}, \quad \mathbb{E}_{\tilde{p}(\tau|f_o)}[\log(\tau)] = \Psi(\hat{a}) - \log(\hat{b}).$$

where  $\Psi(\cdot)$  is the digamma function. Note that these moments are based on the sufficient statistics of Gamma distribution, which are standard for moment matching in ADF and EP framework. The true moments can therefore be calculated by taking the expectation of the conditional moments,

$$\begin{aligned} \mathbb{E}_{\tilde{p}}[\tau] &= \mathbb{E}_{\tilde{p}(f_o)}\mathbb{E}_{\tilde{p}(\tau|f_o)}[\tau] = \mathbb{E}_{\tilde{p}(f_o)}\left[\frac{\hat{a}}{\hat{b}}\right], \\ \mathbb{E}_{\tilde{p}}[\log(\tau)] &= \mathbb{E}_{\tilde{p}(f_o)}\mathbb{E}_{\tilde{p}(\tau|f_o)}[\log(\tau)] = \mathbb{E}_{\tilde{p}(f_o)}[\Psi(\hat{a}) - \log(\hat{b})]. \end{aligned}$$

However, the normalization constant for (1) is intractable and it is difficult to compute the marginal blending distribution  $\tilde{p}(f_o)$ . To overcome this problem, we approximate  $\tilde{p}(f_o)$  with the current posterior of  $f_o$ , namely  $q_{\text{cur}}(f_o)$ . This is reasonable, because  $\tilde{p}(f_o)$  is an integration of  $q(f_o)$  and one new data point; when we have processed many data points, adding one more data point is unlikely to significantly change the posterior. In other words, we can assume  $q(f_o)$  and  $\tilde{p}(f_o)$  are close in high density regions. Hence, we can approximate

$$\begin{aligned} \mathbb{E}_{\tilde{p}}[\tau] &\approx \mathbb{E}_{q_{\text{cur}}(f_o)}\left[\frac{\hat{a}}{\hat{b}}\right], \\ \mathbb{E}_{\tilde{p}}[\log(\tau)] &\approx \mathbb{E}_{q_{\text{cur}}(f_o)}[\Psi(\hat{a}) - \log(\hat{b})]. \end{aligned}$$

A second problem is that due to the nonlinearity of the conditional moments, even with  $q_{\text{cur}}(f_o)$  (which has a nice Gaussian form), we still cannot analytically compute the expectation. To address this issue, we further observe that the conditional moments are functions of  $f_o$  and  $f_o^2$ ,

$$\begin{aligned} h_1(f_o, f_o^2) &= \frac{\hat{a}}{\hat{b}} = \frac{a + \frac{1}{2}}{b + \frac{1}{2}(y_{i_n}^2 - 2f_o + f_o^2)}, \\ h_2(f_o, f_o^2) &= \Psi(\hat{a}) - \log(\hat{b}) = \Psi(a + \frac{1}{2}) - \log\left(b + \frac{1}{2}(y_{i_n}^2 - 2f_o + f_o^2)\right). \end{aligned}$$

Define  $\mathbf{f} = [f_o, f_o^2]^\top$ . We use a Taylor expansion at the mean of  $f_o$  and  $f_o^2$  to approximate the conditional moments,

$$\begin{aligned} h_1(f_o, f_o^2) &\approx h_1(\mathbb{E}_{q_{\text{cur}}}[f_o], \mathbb{E}_{q_{\text{cur}}}[f_o^2]) + (\mathbf{f} - \mathbb{E}_{q_{\text{cur}}}[\mathbf{f}])^\top \nabla h_1|_{\mathbf{f}=\mathbb{E}_{q_{\text{cur}}}[\mathbf{f}]}, \\ h_2(f_o, f_o^2) &\approx h_2(\mathbb{E}_{q_{\text{cur}}}[f_o], \mathbb{E}_{q_{\text{cur}}}[f_o^2]) + (\mathbf{f} - \mathbb{E}_{q_{\text{cur}}}[\mathbf{f}])^\top \nabla h_2|_{\mathbf{f}=\mathbb{E}_{q_{\text{cur}}}[\mathbf{f}]}. \end{aligned} \quad (3)$$

We take expectation over the Taylor expansion, and obtain a closed-form result,

$$\mathbb{E}_{\tilde{p}}[\tau] = \mathbb{E}_{\tilde{p}}[h_1] \approx \frac{a^*}{b^*}, \quad \mathbb{E}_{\tilde{p}}[\log \tau] = \mathbb{E}_{\tilde{p}}[h_2] \approx \Psi(a^*) - \log(b^*) \quad (4)$$

where

$$a^* = a + \frac{1}{2}, \quad b^* = b + \frac{1}{2}((y_{i_n} - \alpha_n)^2 + \beta_n).$$

Finally, from these moments, we can obtain the updated posterior,  $q(\tau) = \text{Gamma}(\tau|a^*, b^*)$ .

## 2 The Updates for Spike-and-Slab Prior Approximation

In our streaming posterior inference, after we execute ADF to process all the entries in the newly received batch, we use standard EP to update the spike-and-slab prior approximation. In this way, we can refine the approximation quality so as to effectively sparsify and condense the neural network to prevent overfitting. Specifically, for each weight  $w_{mjt}$ , we first divide the posterior by the prior approximation to obtain the calibrated (or context) distribution,

$$q^\backslash(w_{mjt}, s_{mjt}) \propto \frac{q_{\text{cur}}(w_{mjt}, s_{mjt})}{A(w_{mjt}, s_{mjt})} = \text{Bern}(s_{mjt}|\rho_0)\mathcal{N}(w_{mjt}|\mu_{mjt}^\backslash, v_{mjt}^\backslash)$$

where  $A(w_{mjt}, s_{mjt}) = \text{Bern}(s_{mjt}|c(\rho_{mjt}))\mathcal{N}(w_{mjt}|\mu_{mjt}^0, v_{mjt}^0)$  (see (13) in the main paper). Because both  $q_{\text{cur}}$  and  $A$  belong to the exponential family, this can be easily done by subtracting the natural parameters. Note that  $\text{Bern}(s_{mjt}|\rho_0)$  is the prior of  $s_{mjt}$  (see (4) and (5) in the main paper) — this comes from the fact that the (approximate) posterior of  $s_{mjt}$  is proportional to the product of its prior and the approximation term in  $A$ .

Next, we combine the calibrated distribution and the exact prior to obtain a tilted distribution (which is similar to the blending distribution in the streaming case),

$$\tilde{p}(w_{mjt}, s_{mjt}) \propto q^\backslash(w_{mjt}, s_{mjt})(s_{mjt}\mathcal{N}(w_{mjt}|0, \sigma_0^2) + (1 - s_{mjt})\delta(w_{mjt})). \quad (5)$$

We then project  $\tilde{p}$  to the exponential family to obtain the updated posterior,

$$q^*(w_{mjt}, s_{mjt}) = \text{Bern}(s_{mjt}|c(\rho_{mjt}^*))\mathcal{N}(w_{mjt}|\mu_{mjt}^*, v_{mjt}^*),$$

where  $c(\cdot)$  is the sigmoid function,

$$\rho_{mjt}^* = \log\left(\frac{\mathcal{N}(\mu_{mjt}^\backslash|0, \sigma_0^2 + v_{mjt}^\backslash)}{\mathcal{N}(\mu_{mjt}^\backslash|0, v_{mjt}^\backslash)}\right), \quad (6)$$

$$\mu_{mjt}^* = c(\hat{\rho}_{mjt})\hat{\mu}_{mjt}, \quad (7)$$

$$v_{mjt}^* = c(\hat{\rho}_{mjt})(\hat{v}_{mjt} + (1 - c(\hat{\rho}_{mjt}))\hat{\mu}_{mjt}^2), \quad (8)$$

$$(9)$$

and

$$\hat{\rho}_{mjt} = \rho_{mjt}^* + c^{-1}(\rho_0),$$

$$\hat{v}_{mjt} = ((v_{mjt}^\backslash)^{-1} + \sigma_0^{-2})^{-1},$$

$$\hat{\mu}_{mjt} = \hat{v}_{mjt} \frac{\mu_{mjt}^\backslash}{v_{mjt}^\backslash}.$$

Finally, we can update the prior approximation term via dividing the updated posterior by the calibrated distribution,  $A^*(w_{mjt}, s_{mjt}) \propto q^*(w_{mjt}, s_{mjt})/q^\backslash(w_{mjt}, s_{mjt})$ . Now, we replace the current prior approximation by  $A^*$  and set  $q_{\text{cur}} = q^*$ , to prepare the streaming inference for the next batch. Therefore, the learned posterior weights and selection probabilities are consistent, and they effectively deactivate many weights to adjust the complexity of the network.

## References

Wang, Z. and Zhe, S. (2019). Conditional expectation propagation. In *UAI*, page 6.