

Supplementary Material for Hessian-Sum-Mixtures

Vassili Korotkine, Mitchell Cohen, and James Richard Forbes

Department of Mechanical Engineering, McGill University
817 Sherbrooke Street West, Montreal QC H3A 0C3

June 18, 2024

1 Practical Implementation

The Jacobians of the different mixtures as well as details of multipliers for practical implementation are covered in this section. The negative log likelihoods are proportional to, but not exactly equal to, the squared error terms because of the normalization constants required to make sure the square root argument does not become negative.

1.1 Max-Mixture

The log-likelihood is computed from the max-mixture approximation to the Gaussian mixture such that

$$-\log p_{\max}(\mathbf{y}|\mathbf{x}) = \frac{1}{2} \mathbf{e}_{\max}^T \mathbf{e}_{\max} \quad (1)$$

$$= -\log \alpha_{k^*} + \frac{1}{2} \mathbf{e}_{k^*}^T \mathbf{e}_{k^*}, \quad k^* = \arg \max_k \alpha_k \exp \left(-\frac{1}{2} \mathbf{e}_k^T \mathbf{e}_k \right) \quad (2)$$

$$\sim \frac{1}{2} \left\| \begin{bmatrix} \sqrt{2} \sqrt{\log c - \log \alpha_{k^*}} \\ \mathbf{e}_{k^*} \end{bmatrix} \right\|_2^2, \quad (3)$$

with the normalizing constant is given by $c = \max_k \alpha_k$ [1] to avoid the square root argument becoming negative, and the corresponding Jacobian given by

$$\frac{\partial \mathbf{e}_{\max}}{\partial \mathbf{x}} = \begin{bmatrix} \mathbf{0}^{1 \times n_e} \\ \frac{\partial \mathbf{e}_{k^*}}{\partial \mathbf{x}} \end{bmatrix}. \quad (4)$$

1.2 Sum-Mixture

The log-likelihood is computed from the full Gaussian Mixture, and in the sum-mixture case is given by the squared norm of a single scalar term,

$$-\log p_{\text{GM}}(\mathbf{y}|\mathbf{x}) \sim \frac{1}{2}e_{\text{SM}}^2, \quad (5)$$

$$\frac{1}{2}e_{\text{SM}}^2 = \frac{1}{2}\sqrt{2\left(\log c - \log \sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k\right)\right)}^2 \quad (6)$$

$$= \frac{1}{2}\sqrt{2\left(\log c + \frac{1}{2}\mathbf{e}_{k^*}^\top \mathbf{e}_{k^*} - \log \sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k + \frac{1}{2}\mathbf{e}_{k^*}^\top \mathbf{e}_{k^*}\right)\right)}^2, \quad (7)$$

where the normalization constant c is given by $\sum_{k=1}^K \alpha_k$, and the Jacobian is given by

$$\frac{\partial e_{\text{SM}}}{\partial \mathbf{x}} = \frac{1}{2e_{\text{SM}}} \frac{-2}{\sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k\right)} \sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k\right) (-\mathbf{e}_k) \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}} \quad (8)$$

$$= \frac{1}{e_{\text{SM}}} \frac{\sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k\right) \mathbf{e}_k^\top \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}}}{\sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k\right)}. \quad (9)$$

1.3 Max-Sum-Mixture

The algebraic value for the Max-Sum-Mixture is the same as for the sum-mixture, but the error term partitioning is different, such that

$$-\log p_{\text{GM}}(\mathbf{y}|\mathbf{x}) \sim \frac{1}{2}\mathbf{e}_{\text{MSM}}^\top \mathbf{e}_{\text{MSM}} \quad (10)$$

$$= \frac{1}{2}\mathbf{e}_{k^*}^\top \mathbf{e}_{k^*} + \frac{1}{2}\sqrt{2\left(\log c - \log \sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k + \frac{1}{2}\mathbf{e}_{k^*}^\top \mathbf{e}_{k^*}\right)\right)}^2 \quad (11)$$

$$= \frac{1}{2}\left\|\left[\sqrt{2\left(\log c - \log \sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k + \frac{1}{2}\mathbf{e}_{k^*}^\top \mathbf{e}_{k^*}\right)\right)}\right]_{\mathbf{e}_{k^*}}\right\|_2^2 \quad (12)$$

$$= \left\|\begin{bmatrix} e_{\text{NL}} \\ \mathbf{e}_{k^*} \end{bmatrix}\right\|_2^2, \quad (13)$$

where $k^* = \arg \max \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k\right)$ is the dominant component and

$$e_{\text{NL}} = \sqrt{2\left(\log c - \log \sum_{k=1}^K \alpha_k \exp\left(-\frac{1}{2}\mathbf{e}_k^\top \mathbf{e}_k + \frac{1}{2}\mathbf{e}_{k^*}^\top \mathbf{e}_{k^*}\right)\right)}, \quad (14)$$

is defined similarly to e_{SM} in (6). The error Jacobian is given by

$$\frac{\partial \mathbf{e}_{\text{MSM}}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{e}_{k^*}}{\partial \mathbf{x}} \\ \frac{\partial e_{\text{NL}}}{\partial \mathbf{x}} \end{bmatrix} \quad (15)$$

and the Jacobian of the nonlinear term $\frac{\partial e_{\text{NL}}}{\partial \mathbf{x}}$ is given by

$$\frac{\partial e_{\text{NL}}}{\partial \mathbf{x}} = \frac{1}{e_{\text{NL}}} \frac{-\sum_{k=1}^K \alpha_k \exp(-\frac{1}{2} \mathbf{e}_k^T \mathbf{e}_k + \frac{1}{2} \mathbf{e}_{k^*}^T \mathbf{e}_{k^*}) \left(-\mathbf{e}_k^T \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}} + \mathbf{e}_{k^*}^T \frac{\partial \mathbf{e}_{k^*}}{\partial \mathbf{x}} \right)}{\sum_{k=1}^K \alpha_k \exp(-\frac{1}{2} \mathbf{e}_k^T \mathbf{e}_k + \frac{1}{2} \mathbf{e}_{k^*}^T \mathbf{e}_{k^*})}. \quad (16)$$

The normalization constant c is given by

$$c = K \max \alpha_k + \delta, \quad (17)$$

where δ is a damping constant [1] that controls the influence of the nonlinear term (14).

2 Robust Loss: Iterative Reweighted Least Squares

The robust loss formulation for a single factor may be written as

$$J = \rho(f(\mathbf{x})) \quad (1)$$

$$= \rho \left(\frac{1}{2} \mathbf{e}(\mathbf{x})^T \mathbf{e}(\mathbf{x}) \right), \quad (2)$$

where $\mathbf{e} : \mathbb{R}^{n_x} \rightarrow n_e$ is the error function, $\rho : \mathbb{R} \rightarrow \mathbb{R}$ is a robust loss function, and $f(\mathbf{x}) = \frac{1}{2} \mathbf{e}(\mathbf{x})^T \mathbf{e}(\mathbf{x})$ is a convenient intermediate quantity for the use of the chain rule.

The robustified Gauss-Newton least squares update is derived by considering Newton's method applied to (2), while assuming the Gauss-Newton Hessian approximation is valid for f ,

$$\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^T} \approx \frac{\partial \mathbf{e}}{\partial \mathbf{x}^T} \frac{\partial \mathbf{e}}{\partial \mathbf{x}}. \quad (3)$$

The Jacobian of the loss function (2) is then given by

$$\frac{\partial J}{\partial \mathbf{x}} = \frac{d\rho}{df} \frac{\partial f}{\partial \mathbf{x}}, \quad (4)$$

and the Hessian by

$$\frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^T} = \frac{d\rho}{df} \frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^T} + \frac{\partial f}{\partial \mathbf{x}^T} \frac{d^2 \rho}{df^2} \frac{\partial f}{\partial \mathbf{x}}. \quad (5)$$

The intermediate function $f = \frac{1}{2} \mathbf{e}^T \mathbf{e}$ has Jacobian

$$\frac{\partial f}{\partial \mathbf{x}} = \mathbf{e}^T \frac{\partial \mathbf{e}}{\partial \mathbf{x}}, \quad (6)$$

and Hessian approximated by (3). The Hessian of the loss is thus given by

$$\frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^T} = \frac{d\rho}{df} \frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^T} + \frac{\partial f}{\partial \mathbf{x}^T} \frac{d^2 \rho}{df^2} \frac{\partial f}{\partial \mathbf{x}} \quad (7)$$

$$= \frac{d\rho}{df} \frac{\partial \mathbf{e}}{\partial \mathbf{x}^T} \frac{\partial \mathbf{e}}{\partial \mathbf{x}} + \left(2 \mathbf{e}^T \frac{\partial \mathbf{e}}{\partial \mathbf{x}} \right)^T \frac{d^2 \rho}{df^2} \left(2 \mathbf{e}^T \frac{\partial \mathbf{e}}{\partial \mathbf{x}} \right) \quad (8)$$

$$= \frac{d\rho}{df} \frac{\partial \mathbf{e}}{\partial \mathbf{x}^T} \frac{\partial \mathbf{e}}{\partial \mathbf{x}} + \frac{d^2 \rho}{df^2} \frac{\partial \mathbf{e}}{\partial \mathbf{x}^T} \mathbf{e} \mathbf{e}^T \frac{\partial \mathbf{e}}{\partial \mathbf{x}}. \quad (9)$$

Writing the Newton method update step with the Hessian approximation (9) and the Jacobian (6) yields

$$\left(\frac{d\rho}{df} \frac{\partial \mathbf{e}}{\partial \mathbf{x}^\top} \frac{\partial \mathbf{e}}{\partial \mathbf{x}} + \frac{d^2\rho}{df^2} \frac{\partial \mathbf{e}}{\partial \mathbf{x}^\top} \mathbf{e} \mathbf{e}^\top \frac{\partial \mathbf{e}}{\partial \mathbf{x}} \right) \delta \mathbf{x} = - \frac{d\rho}{df} \mathbf{e}^\top \frac{\partial \mathbf{e}}{\partial \mathbf{x}}, \quad (10)$$

Using only the robust loss function first order term is exactly the iteratively reweighted least squares approach. Some solvers, such as Ceres, also use the second order term, which is called the Triggs correction. Note that this approach does *not* correspond to defining an error term of the form

$$\tilde{e} = \sqrt{\rho(\mathbf{e}^\top \mathbf{e})}, \quad (11)$$

and using that in Gauss-Newton since that would be ill-conditioned, as the error term size reduces to one. Furthermore, for problems with many different factors, the robust loss ρ terms will all have different values.

3 Hessian Approximation for Sum-Mixtures

Consider an optimization problem consisting of a single factor of the form

$$J = \log c - \log \sum_{k=1}^K \alpha_k \exp \left(-\frac{1}{2} \mathbf{e}_k^\top \mathbf{e}_k \right) \quad (1)$$

$$= \sqrt{\log c - \log \sum_{k=1}^K \alpha_k \exp \left(-\frac{1}{2} \mathbf{e}_k^\top \mathbf{e}_k \right)}^2, \quad (2)$$

where the square root is present to make it similar to standard error evaluation for Gauss-Newton, and the normalization constant c is present for the square root argument to be positive.

Similarly to the robust loss case (2), the additional nonlinearity imposed by the Gaussian mixture needs to be taken into account for the Hessian approximation. Analogously to the robust loss case, $f_k = \frac{1}{2} \mathbf{e}_k^\top \mathbf{e}_k$ is assumed well-behaved such that

$$\frac{\partial^2 f_k}{\partial \mathbf{x} \partial \mathbf{x}^\top} \approx \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}^\top} \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}}, \quad (3)$$

and the definition for ρ follows as

$$\rho(f_1, \dots, f_K) = \log c - \log \sum_{k=1}^K \alpha_k \exp(-f_k). \quad (4)$$

Results from the robust loss case are unapplicable here due to the additional nonlinearity due to the LogSumExp term and the presence of multiple f_k terms instead of a single one.

Nevertheless, a robustified Hessian approximation for the sum-mixture may be derived. The Jacobian is given by

$$\frac{\partial J}{\partial \mathbf{x}} = \sum_{k=1}^K \frac{\partial \rho}{\partial f_k} \frac{\partial f_k}{\partial \mathbf{x}}, \quad (5)$$

and the Hessian by

$$\frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^\top} = \sum_{k=1}^K \left(\frac{\partial \rho}{\partial f_k} \frac{\partial^2 f_k}{\partial \mathbf{x} \partial \mathbf{x}^\top} + \frac{\partial f_k}{\partial \mathbf{x}^\top} \sum_{j=1}^K \frac{\partial^2 \rho}{\partial f_j \partial f_k} \frac{\partial f_j}{\partial \mathbf{x}} \right). \quad (6)$$

The partial derivatives of ρ are given by

$$\frac{\partial \rho}{\partial f_k} = - \frac{-\alpha_k \exp(-f_k)}{\sum_{i=1}^K \alpha_i \exp(-f_i)} \quad (7)$$

$$= \frac{\alpha_k \exp(-f_k)}{\sum_{i=1}^K \alpha_i \exp(-f_i)} \quad (8)$$

$$\frac{\partial^2 \rho}{\partial f_j \partial f_k} = \frac{\delta_{jk}(-\alpha_k \exp(-f_k)) \sum_{i=1}^K \alpha_i \exp(-f_i) - \alpha_k \exp(-f_k)(-\alpha_j \exp(-f_j))}{(\sum_{i=1}^K \alpha_i \exp(-f_i))^2} \quad (9)$$

$$= \frac{-\delta_{jk}(\alpha_k \exp(-f_k)) \sum_{i=1}^K \alpha_i \exp(-f_i) + \alpha_k \alpha_j \exp(-f_k) \exp(-f_j)}{(\sum_{i=1}^K \alpha_i \exp(-f_i))^2} \quad (10)$$

and the Jacobian of $f_k = \frac{1}{2} \mathbf{e}_k^\top \mathbf{e}_k$ is given by

$$\frac{\partial f_k}{\partial \mathbf{x}} = \mathbf{e}_k^\top \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}}. \quad (11)$$

The Jacobian of the loss J is then given by (5),

$$\frac{\partial J}{\partial \mathbf{x}} = \sum_{k=1}^K \frac{\partial \rho}{\partial f_k} \frac{\partial f_k}{\partial \mathbf{x}} \quad (12)$$

$$= \sum_{k=1}^K \frac{\alpha_k \exp(-f_k)}{\sum_{i=1}^K \alpha_i \exp(-f_i)} \mathbf{e}_k^\top \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}} \quad (13)$$

$$= \sum_{k=1}^K \frac{\alpha_k \exp(-f_k)}{\sum_{i=1}^K \alpha_i \exp(-f_i)} \mathbf{e}_k^\top \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}}. \quad (14)$$

$$(15)$$

The Hessian of the loss is given by (6)

$$\frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^\top} = \sum_{k=1}^K \left(\frac{\partial \rho}{\partial f_k} \frac{\partial^2 f_k}{\partial \mathbf{x} \partial \mathbf{x}^\top} + \frac{\partial f_k}{\partial \mathbf{x}^\top} \sum_{j=1}^K \frac{\partial^2 \rho}{\partial f_j \partial f_k} \frac{\partial f_j}{\partial \mathbf{x}} \right) \quad (16)$$

$$= \sum_{k=1}^K \frac{\partial \rho}{\partial f_k} \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}^\top} \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}} + \left(\mathbf{e}_k^\top \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}} \right)^\top \sum_{j=1}^K \frac{\partial^2 \rho}{\partial f_j \partial f_k} \left(\mathbf{e}_j^\top \frac{\partial \mathbf{e}_j}{\partial \mathbf{x}} \right) \quad (17)$$

$$= \sum_{k=1}^K \frac{\partial \rho}{\partial f_k} \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}^\top} \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}} + \sum_{j=1}^K \frac{\partial^2 \rho}{\partial f_j \partial f_k} \frac{\partial \mathbf{e}_k}{\partial \mathbf{x}^\top} \mathbf{e}_k \mathbf{e}_j^\top \frac{\partial \mathbf{e}_j}{\partial \mathbf{x}}. \quad (18)$$

4 Derivation of Normalization Constant

Here we derive the result related to the lower bound on ΔJ presented in Section IV of the paper. The solver error is given by

$$\mathbf{e}_{\text{solver}}^T = [\mathbf{e}_{\text{solver},1}^T \quad \sqrt{2(\gamma_{\text{HSM}} + \Delta J)}], \quad (1)$$

where

$$\mathbf{e}_{\text{solver},1}^T = \left[\sqrt{\frac{\partial \rho}{\partial f_1}} \mathbf{e}_1^T \quad \dots \quad \sqrt{\frac{\partial \rho}{\partial f_K}} \mathbf{e}_K^T \right], \quad (2)$$

and $\frac{\partial \rho}{\partial f_k}$ is defined such that

$$\frac{\partial \rho}{\partial f_k} = \frac{\alpha_k \exp\left(-\frac{1}{2} \mathbf{e}_k^T \mathbf{e}_k\right)}{\sum_{i=1}^{n_k} \alpha_i \exp\left(-\frac{1}{2} \mathbf{e}_i^T \mathbf{e}_i\right)}. \quad (3)$$

The desired error is given by the negative log-likelihood of the GMM as

$$J_{\text{GMM}}(\mathbf{x}) = -\log \sum_{k=1}^{n_k} \alpha_k \exp\left(-\frac{1}{2} \mathbf{e}_k(\mathbf{x})^T \mathbf{e}_k(\mathbf{x})\right). \quad (4)$$

By setting $\Delta J = J_{\text{GMM}} - \frac{1}{2} \mathbf{e}_{\text{solver},1}^T \mathbf{e}_{\text{solver},1}$, the evaluated cost becomes

$$\frac{1}{2} \mathbf{e}_{\text{solver}}^T \mathbf{e}_{\text{solver}} = \frac{1}{2} \mathbf{e}_{\text{solver},1}^T \mathbf{e}_{\text{solver},1} + \gamma_{\text{HSM}} + \Delta J - \frac{1}{2} \mathbf{e}_{\text{solver},1}^T \mathbf{e}_{\text{solver},1} \quad (5)$$

$$= \gamma_{\text{HSM}} + J_{\text{GMM}}. \quad (6)$$

The constant γ_{HSM} is set such that $\gamma_{\text{HSM}} + \Delta J \geq 0$, which requires a lower bound on ΔJ . The lower bound is obtained by first manipulating ΔJ as

$$\Delta J = J_{\text{GMM}} - \frac{1}{2} \mathbf{e}_{\text{solver},1}^T \mathbf{e}_{\text{solver},1} \quad (7)$$

$$= -\log \sum_{k=1}^{n_k} \alpha_k \exp\left(-\frac{1}{2} \mathbf{e}_k(\mathbf{x})^T \mathbf{e}_k(\mathbf{x})\right) - \frac{1}{2} \mathbf{e}_{\text{solver},1}^T \mathbf{e}_{\text{solver},1} \quad (8)$$

$$= -\log \sum_{k=1}^{n_k} \alpha_k \exp\left(-\frac{1}{2} \mathbf{e}_k(\mathbf{x})^T \mathbf{e}_k(\mathbf{x})\right) - \frac{1}{2} \sum_{k=1}^{n_k} \frac{\alpha_k \exp\left(-\frac{1}{2} \mathbf{e}_k^T \mathbf{e}_k\right)}{\sum_{i=1}^{n_k} \alpha_i \exp\left(-\frac{1}{2} \mathbf{e}_i^T \mathbf{e}_i\right)} \mathbf{e}_k^T \mathbf{e}_k. \quad (9)$$

Changing variables to $f_k = \frac{1}{2} \mathbf{e}_k(\mathbf{x})^T \mathbf{e}_k(\mathbf{x})$ yields

$$\Delta J = -\log \sum_{k=1}^{n_k} \alpha_k \exp(-f_k) - \frac{1}{2} \sum_{k=1}^{n_k} \frac{\alpha_k \exp(-f_k)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} 2f_k \quad (10)$$

$$= -\log \sum_{k=1}^{n_k} \alpha_k \exp(-f_k) - \sum_{k=1}^{n_k} \frac{\alpha_k \exp(-f_k)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} f_k, \quad (11)$$

with $\alpha_k > 0, f_k \geq 0$. Then, exponentiating and taking logarithm of the second term yields

$$\Delta J = -\log \sum_{k=1}^{n_k} \alpha_k \exp(-f_k) - \sum_{k=1}^{n_k} \frac{\alpha_k \exp(-f_k)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} f_k \quad (12)$$

$$= -\log \left(\sum_{k=1}^{n_k} \alpha_k \exp(-f_k) \right) - \log \left(\exp \sum_{k=1}^{n_k} \frac{\alpha_k \exp(-f_k)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} f_k \right) \quad (13)$$

$$= -\log \sum_{k=1}^{n_k} \left(\alpha_k \exp(-f_k) \exp \sum_{j=1}^{n_k} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} f_j \right) \quad (14)$$

$$= -\log \sum_{k=1}^{n_k} \left(\alpha_k \exp \left(-f_k + \sum_{j=1}^{n_k} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} f_j \right) \right) \quad (15)$$

$$= -\log \sum_{k=1}^{n_k} \left(\alpha_k \exp \left(\sum_{j=1}^{n_k} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} (f_j - f_k) \right) \right). \quad (16)$$

The overall expression in (16) is decreasing as a function of the exponent argument S_k ,

$$S_k = \sum_{j=1}^{n_k} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} (f_j - f_k). \quad (17)$$

Therefore, an upper bound on S_k is required. In the worst case scenario, $f_j \geq f_k$ for all j in the summation. This can be clarified as follows. Defining $A_{f_j \geq f_k}$ as the set of indices j with $f_j \geq f_k$, as well as the set $A_{f_j < f_k}$ with $f_j < f_k$, and splitting S_k up into sums over these two sets yields

$$S_k = \sum_{A_{f_j \geq f_k}} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} (f_j - f_k) + \sum_{A_{f_j < f_k}} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} (f_j - f_k) \quad (18)$$

$$\leq \sum_{A_{f_j \geq f_k}} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} (f_j - f_k), \quad (19)$$

since $f_j - f_k < 0$ in the second term of (18). Thus, proceeding with $f_j - f_k \geq 0$ allows to write

$$S_k = \sum_{j=1}^{n_k} \frac{\alpha_j \exp(-f_j)}{\sum_{i=1}^{n_k} \alpha_i \exp(-f_i)} (f_j - f_k) \quad (20)$$

$$\leq \sum_{j=1}^{n_k} \frac{\alpha_j \exp(-f_j)}{\alpha_j \exp(-f_j) + \alpha_k \exp(-f_k)} (f_j - f_k) \quad (21)$$

$$= \sum_{j=1}^{n_k} \frac{1}{1 + \frac{\alpha_k}{\alpha_j} \exp(f_j - f_k)} (f_j - f_k) \quad (22)$$

$$\leq \sum_{j=1}^{n_k} \frac{1}{\frac{\alpha_k}{\alpha_j} \exp(f_j - f_k)} (f_j - f_k) \quad (23)$$

$$= \frac{1}{\alpha_k} \sum_{j=1}^{n_k} \alpha_j \exp(-(f_j - f_k)) (f_j - f_k). \quad (24)$$

Using the change of variables $t = f_j - f_k$, with $t > 0$, makes explicit that

$$\exp(-(f_j - f_k))(f_j - f_k) = \exp(-t)(t), \quad (25)$$

which has a single maximum at $t = 1$ with $\exp(-t)(t) = \exp(-1)$. Therefore,

$$S_k \leq \frac{1}{\alpha_k} \sum_{j=1}^{n_k} \alpha_j \exp(-1) \quad (26)$$

$$\leq \frac{1}{\alpha_k} \sum_{j=1}^{n_k} \alpha_j. \quad (27)$$

Therefore,

$$\Delta J = -\log \sum_{k=1}^{n_k} (\alpha_k \exp(S_k)) \geq -\log \sum_{k=1}^{n_k} \left(\alpha_k \exp \left(\frac{1}{\alpha_k} \sum_{j=1}^{n_k} \alpha_j \right) \right), \quad (28)$$

and the normalization constant γ_{HSM} may be set to

$$\gamma_{\text{HSM}} = \log \sum_{k=1}^{n_k} \left(\alpha_k \exp \left(\frac{1}{\alpha_k} \sum_{j=1}^{n_k} \alpha_j \right) \right). \quad (29)$$

5 Identities

5.1 Jacobian of vector-scalar product

For $a(\mathbf{x})$ a scalar function of \mathbf{x} and $\mathbf{v}(\mathbf{x})$ a column vector function of \mathbf{x} , the Jacobian of their product is obtained as

$$\left. \frac{\partial a\mathbf{v}}{\partial \mathbf{x}} \right|_{i,j} = \frac{\partial}{\partial x_j} a v_i \quad (1)$$

$$= \frac{\partial a}{\partial x_j} v_i + a \frac{\partial v_i}{\partial x_j} \quad (2)$$

$$\frac{\partial a\mathbf{v}}{\partial \mathbf{x}} = \mathbf{v} \frac{\partial a}{\partial \mathbf{x}} + a \frac{\partial \mathbf{v}}{\partial \mathbf{x}} \quad (3)$$

5.2 Hessian of a function through Jacobian of gradient

The Hessian $\frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^\top}$ of a function $J(x)$ is given by a matrix with entries

$$\left. \frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^\top} \right|_{i,j} = \frac{\partial J}{\partial x_i \partial x_j}. \quad (4)$$

The gradient $\frac{\partial J}{\partial \mathbf{x}}^\top$ of J has entries given by

$$\left. \frac{\partial J}{\partial \mathbf{x}} \right|_i = \frac{\partial J}{\partial x_i}, \quad (5)$$

and the Jacobian of the gradient is given by

$$\left. \frac{\partial}{\partial \mathbf{x}} \frac{\partial J^\top}{\partial \mathbf{x}} \right|_{i,j} = \frac{\partial J}{\partial x_j \partial x_i}. \quad (6)$$

Therefore,

$$\frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^\top} = \left(\frac{\partial}{\partial \mathbf{x}} \frac{\partial J^\top}{\partial \mathbf{x}} \right)^\top. \quad (7)$$

5.3 2D wedge operator multiplication flip

Consider a general vector $\mathbf{v} \in \mathbb{R}^2$ such that

$$\delta \xi^{\phi^\wedge} \mathbf{v} = \begin{bmatrix} 0 & -\delta \xi^\phi \\ \delta \xi^\phi & 0 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \quad (8)$$

$$= \begin{bmatrix} -\delta \xi^\phi v_2 \\ \delta \xi^\phi v_1 \end{bmatrix} \quad (9)$$

$$= \begin{bmatrix} -v_2 \\ v_1 \end{bmatrix} \delta \xi^\phi \quad (10)$$

$$= \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \mathbf{v} \delta \xi^\phi. \quad (11)$$

References

- [1] T. Pfeifer, S. Lange, and P. Protzel, “Advancing Mixture Models for Least Squares Optimization,” *IEEE Robotics and Automation Letters (RA-L)*, vol. 6, no. 2, pp. 3941–3948, 2021.