

Proofs for “Information Geometry and Minimum Description Length Networks”

An Approximation of $\ln N(\mathcal{B}, \boldsymbol{\alpha})$

As the value of $\ln N(\mathcal{B}, \boldsymbol{\alpha})$ does not depend on the choice of the coordinate system, we abuse notation and vary $\boldsymbol{\eta} = (\eta^{(1)}, \dots, \eta^{(\dim \mathcal{S})})$ to an *ideal coordinate system*, where $g(\boldsymbol{\eta})$ is everywhere identity. In theory, this is possible locally. However, for convenience, we assume that such a coordinate system exists globally. By definition,

$$\begin{aligned}
 \ln N(\mathcal{B}, \boldsymbol{\alpha}) &= \ln \left(\int_{\boldsymbol{\eta} \in \mathcal{S}} \sum_{i=1}^m \alpha_i \exp(-D(\boldsymbol{\eta} \parallel \boldsymbol{\eta}_i)) d\boldsymbol{\eta} \right) \\
 &= \ln \left(\sum_{i=1}^m \alpha_i \int_{\boldsymbol{\eta} \in \mathcal{S}} \exp(-D(\boldsymbol{\eta} \parallel \boldsymbol{\eta}_i)) d\boldsymbol{\eta} \right) \\
 &\approx \ln \left(\sum_{i=1}^m \alpha_i \int_{\eta^{(1)}} \cdots \int_{\eta^{(\dim \mathcal{S})}} \exp \left(-\frac{1}{2} (\boldsymbol{\eta} - \boldsymbol{\eta}_i)^T g(\boldsymbol{\eta}_i) (\boldsymbol{\eta} - \boldsymbol{\eta}_i) \right) \right. \\
 &\quad \left. \times \sqrt{|g(\boldsymbol{\eta})|} d\eta^{(1)} \cdots d\eta^{(\dim \mathcal{S})} \right) \\
 &= \ln \left(\sum_{i=1}^m \alpha_i \int_{\eta^{(1)}} \cdots \int_{\eta^{(\dim \mathcal{S})}} \exp \left(-\frac{1}{2} \|\boldsymbol{\eta} - \boldsymbol{\eta}_i\|_2^2 \right) d\eta^{(1)} \cdots d\eta^{(\dim \mathcal{S})} \right) \\
 &\approx \ln \left(\sum_{i=1}^m \alpha_i \exp \left(\frac{\dim \mathcal{S}}{2} \ln(2\pi) \right) \right) = \frac{\dim \mathcal{S}}{2} \ln(2\pi). \tag{S.1}
 \end{aligned}$$

The first “ \approx ” is by approximating D to a square distance, which is only accurate when $\boldsymbol{\eta}$ and $\boldsymbol{\eta}_i$ are close enough. The second “ \approx ” is by relaxing the domain of the integration from \mathcal{S} to $\mathfrak{R}^{\dim \mathcal{S}}$. This is a rough approximation for a general \mathcal{S} , to show the order of the term $\ln N(\mathcal{B}, \boldsymbol{\alpha})$, and to show its weak dependence to \mathcal{B} and $\boldsymbol{\alpha}$. More accurate approximations based on specific choices of \mathcal{S} can lead to better implementations of MDL networks and better criteria in accordance to MDL.

Proof of $E(\mathcal{N}, A) \leq \hat{E}(\mathcal{N}, A)$ (HARDN)

Proof. $\forall l, \forall i,$

$$\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})) \geq \max_j [\alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j}))]. \quad (\text{S.2})$$

As $-\ln(x)$ is monotonically decreasing,

$$\begin{aligned} E(\mathcal{N}, A) &= - \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})) \right) \\ &\leq - \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \ln \max_j [\alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j}))] \\ &= - \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \max_j [\ln \alpha_{l+1,j} - D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})] \\ &= \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \min_j [-\ln \alpha_{l+1,j} + D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})] = \hat{E}(\mathcal{N}, A). \quad (\text{S.3}) \end{aligned}$$

□

Proof of $E(\mathcal{N}, A) \leq \bar{E}(\mathcal{N}, A, B)$ (SOFTN)

Proof. Because of the convexity of $-\ln(x)$,

$$\begin{aligned} E(\mathcal{N}, A) &= - \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})) \right) \\ &= - \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \beta_{li}^j \cdot \frac{\alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j}))}{\beta_{li}^j} \right) \\ &\leq \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l+1}} \beta_{li}^j \left[-\ln \left(\frac{\alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j}))}{\beta_{li}^j} \right) \right] \\ &= \sum_{l=0}^{L-1} \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l+1}} \beta_{li}^j \left(\ln \frac{\beta_{li}^j}{\alpha_{l+1,j}} + D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j}) \right) = \bar{E}(\mathcal{N}, A, B). \quad (\text{S.4}) \end{aligned}$$

□

Proof of Theorem 3

Proof. Denote the true distribution with the components $\{\boldsymbol{\eta}_i^t\}$ and the weights $\{\alpha_i^t\}$ by $True(\mathbf{x})$. By eq. (7), $\forall \mathcal{N}, \forall A$, when $n \rightarrow \infty$,

$$\begin{aligned}
E(\mathcal{N}, A) &= -n \int True(\mathbf{x}) \ln \left(\sum_{j=1}^{n_1} \alpha_{1j} \exp(-D(\boldsymbol{\eta}(\mathbf{x}) \parallel \boldsymbol{\eta}_{1j})) \right) d\mathbf{x} \\
&\quad - \sum_{l=1}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})) \right) \\
&= -n \int True(\mathbf{x}) \ln \left(\sum_{j=1}^{n_1} \alpha_{1j} p(\mathbf{x} \mid \boldsymbol{\eta}_{1j}) \right) d\mathbf{x} + constant \\
&\quad - \sum_{l=1}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})) \right). \quad (\text{S.5})
\end{aligned}$$

We construct an MDL network \mathcal{N}^t , where \mathcal{L}_1^t is given by $\{\boldsymbol{\eta}_{1i}^t = \boldsymbol{\eta}_i^t\}$ with the weights $\{\alpha_{1i}^t = \alpha_i^t\}$. The rest of the cells $\{\boldsymbol{\eta}_{li}^t\}$ in higher levels, including their weights $\{\alpha_{li}^t\}$ are given by the sub-optimal solution which minimizes the above eq. (S.5) with \mathcal{L}_1^t and its weights fixed. Given that \mathcal{L}_0 is fixed by infinite samples corresponding to the truth, $\forall \mathcal{N}, \forall A$,

$$\begin{aligned}
E(\mathcal{N}, A) - E(\mathcal{N}^t, A^t) &= n \int True(\mathbf{x}) \ln \frac{True(\mathbf{x})}{\sum_{j=1}^{n_1} \alpha_{1j} p(\mathbf{x}_i \mid \boldsymbol{\eta}_{1j})} d\mathbf{x} \\
&\quad - \sum_{l=1}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j} \exp(-D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j})) \right) \\
&\quad + \sum_{l=1}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j}^t \exp(-D(\boldsymbol{\eta}_{li}^t \parallel \boldsymbol{\eta}_{l+1,j}^t)) \right). \quad (\text{S.6})
\end{aligned}$$

If $\{\boldsymbol{\eta}_{1j}\}$ in \mathcal{N} or $\{\alpha_{1j}\}$ in A does not correspond to $True(\mathbf{x})$, the first term on the right-hand-side of eq. (S.6) will go to $+\infty$ as $n \rightarrow \infty$. The second term is always non-negative, because of the non-negativity of D . Because of the sub-optimality discussed earlier, the third term is lower-bounded, as in

$$\sum_{l=1}^{L-1} \sum_{i=1}^{n_l} \ln \left(\sum_{j=1}^{n_{l+1}} \alpha_{l+1,j}^t \exp(-D(\boldsymbol{\eta}_{li}^t \parallel \boldsymbol{\eta}_{l+1,j}^t)) \right) \geq - \sum_{i=1}^{n_1} D(\boldsymbol{\eta}_i^t \parallel \tilde{\boldsymbol{\eta}}), \quad (\text{S.7})$$

where $\tilde{\boldsymbol{\eta}}$ can be any distribution, e.g., the right-handed Bregman centroid $\{\boldsymbol{\eta}_i^t\}$. The right-hand-side of eq. (S.7) is the negative cost of a simple structure (one cell in \mathcal{L}_2) to represent \mathcal{L}_1^t , which is upper-bounded by the sub-optimal negative cost on the left-hand-side. Integrating all the three terms on the right-hand-of eq. (S.6), $E(\mathcal{N}, A) > E(\mathcal{N}^t, A^t)$. Hence, in the optimal solution, \mathcal{L}_1 must be exactly $\{\boldsymbol{\eta}_i^t\}$ and the weights must be exactly $\{\alpha_i^t\}$. \square

Proof of Theorem 4

Proof. By the definition of $D(\boldsymbol{\eta}_1 \parallel \boldsymbol{\eta}_2)$ in section 2.3 as a Bregman divergence, $\forall \boldsymbol{\theta}(\boldsymbol{\eta})$, we have

$$\begin{aligned}
\text{gain}(\boldsymbol{\eta}) &= D(\boldsymbol{\eta}_1 \parallel \boldsymbol{\eta}_2) - D(\boldsymbol{\eta}_1 \parallel \boldsymbol{\eta}) - D(\boldsymbol{\eta} \parallel \boldsymbol{\eta}_2) \\
&= + \left(\psi^*(\boldsymbol{\eta}_1) - \psi^*(\boldsymbol{\eta}_2) - \boldsymbol{\theta}_2^T(\boldsymbol{\eta}_1 - \boldsymbol{\eta}_2) \right) \\
&\quad - \left(\psi^*(\boldsymbol{\eta}_1) - \psi^*(\boldsymbol{\eta}) - \boldsymbol{\theta}^T(\boldsymbol{\eta}_1 - \boldsymbol{\eta}) \right) \\
&\quad - \left(\psi^*(\boldsymbol{\eta}) - \psi^*(\boldsymbol{\eta}_2) - \boldsymbol{\theta}_2^T(\boldsymbol{\eta} - \boldsymbol{\eta}_2) \right) \\
&= (\boldsymbol{\theta}_2 - \boldsymbol{\theta})^T(\boldsymbol{\eta} - \boldsymbol{\eta}_1). \tag{S.8}
\end{aligned}$$

Let $\boldsymbol{\theta}_{lc} = (\boldsymbol{\theta}_1 + \boldsymbol{\theta}_2)/2$ be the left-handed Bregman centroid of $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$, then $\boldsymbol{\theta}_2 - \boldsymbol{\theta}_{lc} = \boldsymbol{\theta}_{lc} - \boldsymbol{\theta}_1$. Therefore,

$$\text{gain}(\boldsymbol{\eta}_{lc}) = (\boldsymbol{\theta}_2 - \boldsymbol{\theta}_{lc})^T(\boldsymbol{\eta}_{lc} - \boldsymbol{\eta}_1) = (\boldsymbol{\theta}_{lc} - \boldsymbol{\theta}_1)^T(\boldsymbol{\eta}_{lc} - \boldsymbol{\eta}_1). \tag{S.9}$$

On the other hand, $\forall \boldsymbol{\eta}_a, \boldsymbol{\eta}_b \in \mathcal{S}$, $\boldsymbol{\eta}_a \neq \boldsymbol{\eta}_b$,

$$\begin{aligned}
D(\boldsymbol{\eta}_a \parallel \boldsymbol{\eta}_b) + D(\boldsymbol{\eta}_b \parallel \boldsymbol{\eta}_a) &= + \left(\psi^*(\boldsymbol{\eta}_a) - \psi^*(\boldsymbol{\eta}_b) - \boldsymbol{\theta}_b^T(\boldsymbol{\eta}_a - \boldsymbol{\eta}_b) \right) \\
&\quad + \left(\psi^*(\boldsymbol{\eta}_b) - \psi^*(\boldsymbol{\eta}_a) - \boldsymbol{\theta}_a^T(\boldsymbol{\eta}_b - \boldsymbol{\eta}_a) \right) \\
&= (\boldsymbol{\theta}_a - \boldsymbol{\theta}_b)^T(\boldsymbol{\eta}_a - \boldsymbol{\eta}_b) > 0. \tag{S.10}
\end{aligned}$$

By eqs. (S.9) and (S.10),

$$\text{gain}(\boldsymbol{\eta}_{lc}) = D(\boldsymbol{\eta}_{lc} \parallel \boldsymbol{\eta}_1) + D(\boldsymbol{\eta}_1 \parallel \boldsymbol{\eta}_{lc}) > 0 \quad (\text{which proves } \textcircled{1}). \tag{S.11}$$

Similarly, we let $\boldsymbol{\eta}_{rc} = (\boldsymbol{\eta}_1 + \boldsymbol{\eta}_2)/2$ be the righted-handed Bregman centroid, then

$$\begin{aligned}
\text{gain}(\boldsymbol{\eta}_{rc}) &= (\boldsymbol{\theta}_2 - \boldsymbol{\theta}_{rc})^T(\boldsymbol{\eta}_{rc} - \boldsymbol{\eta}_1) = (\boldsymbol{\theta}_2 - \boldsymbol{\theta}_{rc})^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_{rc}) \\
&= D(\boldsymbol{\eta}_2 \parallel \boldsymbol{\eta}_{rc}) + D(\boldsymbol{\eta}_{rc} \parallel \boldsymbol{\eta}_2). \tag{S.12}
\end{aligned}$$

By eqs. (S.11) and (S.12), $\exists \boldsymbol{\eta} \in \mathcal{S}$ satisfying

$$\text{gain}(\boldsymbol{\eta}) \geq \max\{D(\boldsymbol{\eta}_{lc} \parallel \boldsymbol{\eta}_1) + D(\boldsymbol{\eta}_1 \parallel \boldsymbol{\eta}_{lc}), D(\boldsymbol{\eta}_2 \parallel \boldsymbol{\eta}_{rc}) + D(\boldsymbol{\eta}_{rc} \parallel \boldsymbol{\eta}_2)\}. \tag{S.13}$$

□