

# A scaled Bregman theorem with applications

## — Supplementary Material —

### Abstract

This is the Supplementary Material to Paper “A scaled Bregman theorem with applications” by R. Nock, A-K. Menon and C.-S. Ong. Theorems and Lemmata are numbered with letters (A, B, ...) to make a clear difference with the main file numbering.

### Table of contents

#### Supplementary material on proofs

(I) Proofs of results in main body	Pg 2
(II) Additional helper lemmata	Pg 6
(III) Working out examples of Table A1	Pg 15
(IV) Going deep: higher-order identities	Pg 25
(V) Additional application: perspective transform of exponential families	Pg 27
(VI) Additional application: computational information geometry	Pg 28
(VII) Review: binary density ratio estimation	Pg 30
(VIII) Comments on Theorem 1	Pg 31

#### Supplementary material on experiments

(IX) Multiclass density ratio experiments	Pg 33
(X) Adaptive filtering experiments	Pg 34

# I Proofs of results in main body

We present proofs of all results in the main body.

**Proof** [Proof of Theorem 1] Let  $\mathbf{J} : \mathcal{X} \rightarrow \mathcal{X}_g$  denote the Jacobian of  $h : \mathbf{x} \mapsto (1/g(\mathbf{x})) \cdot \mathbf{x}$ . By an elementary calculation,

$$g(\mathbf{x}) \cdot \mathbf{J} = \mathbf{I}_d - (1/g(\mathbf{x})) \cdot \mathbf{x} \nabla g(\mathbf{x})^\top,$$

which by the chain rule brings the following expression for the gradient of  $\check{\varphi}(\mathbf{y}) = g(\mathbf{y}) \cdot (\varphi \circ h)(\mathbf{y})$ :

$$\begin{aligned} \nabla \check{\varphi}(\mathbf{y}) &= \nabla g(\mathbf{y}) \cdot (\varphi \circ h)(\mathbf{y}) + g(\mathbf{y}) \cdot \nabla(\varphi \circ h)(\mathbf{y}) \\ &= \nabla g(\mathbf{y}) \cdot (\varphi \circ h)(\mathbf{y}) + g(\mathbf{y}) \cdot \mathbf{J}^\top \nabla \varphi(h(\mathbf{y})) \\ &= \nabla g(\mathbf{y}) \cdot (\varphi \circ h)(\mathbf{y}) + \nabla \varphi(h(\mathbf{y})) - (1/g(\mathbf{y})) \cdot \nabla g(\mathbf{y}) \mathbf{y}^\top \nabla \varphi(h(\mathbf{y})) \\ &= \nabla \varphi\left(\frac{1}{g(\mathbf{y})} \cdot \mathbf{y}\right) + \left(\varphi\left(\frac{1}{g(\mathbf{y})} \cdot \mathbf{y}\right) - \frac{1}{g(\mathbf{y})} \cdot \mathbf{y}^\top \nabla \varphi\left(\frac{1}{g(\mathbf{y})} \cdot \mathbf{y}\right)\right) \cdot \nabla g(\mathbf{y}). \end{aligned} \quad (1)$$

For simplicity, let  $\mathbf{u} = \mathbf{x}/g(\mathbf{x})$  and  $\mathbf{v} = \mathbf{y}/g(\mathbf{y})$ , so that  $\check{\varphi}(\mathbf{x}) = g(\mathbf{x}) \cdot \varphi(\mathbf{u})$  and  $\check{\varphi}(\mathbf{y}) = g(\mathbf{y}) \cdot \varphi(\mathbf{v})$ . The above then reads

$$\nabla \check{\varphi}(\mathbf{y}) = \nabla \varphi(\mathbf{v}) + (\varphi(\mathbf{v}) - \mathbf{v}^\top \nabla \varphi(\mathbf{v})) \cdot \nabla g(\mathbf{y}). \quad (2)$$

Now, the LHS of Equation (3) (main file) is

$$\begin{aligned} g(\mathbf{x}) \cdot D_\varphi\left(\frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{g(\mathbf{y})} \cdot \mathbf{y}\right) &= g(\mathbf{x}) \cdot D_\varphi(\mathbf{u} \parallel \mathbf{v}) \\ &= g(\mathbf{x}) \cdot \varphi(\mathbf{u}) - g(\mathbf{x}) \cdot \varphi(\mathbf{v}) - g(\mathbf{x}) \cdot \nabla \varphi(\mathbf{v})^\top (\mathbf{u} - \mathbf{v}) \\ &= \check{\varphi}(\mathbf{x}) - g(\mathbf{x}) \cdot \varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top (\mathbf{x} - g(\mathbf{x}) \cdot \mathbf{v}) \\ &= \check{\varphi}(\mathbf{x}) - g(\mathbf{x}) \cdot (\varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{x}, \end{aligned}$$

while the RHS is

$$\begin{aligned} D_{\check{\varphi}}(\mathbf{x} \parallel \mathbf{y}) &= \check{\varphi}(\mathbf{x}) - \check{\varphi}(\mathbf{y}) - \nabla \check{\varphi}(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &= \check{\varphi}(\mathbf{x}) - g(\mathbf{y}) \cdot \varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top (\mathbf{x} - \mathbf{y}) - (\varphi(\mathbf{v}) - \mathbf{v}^\top \nabla \varphi(\mathbf{v})) \cdot \nabla g(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}). \end{aligned}$$

Cancelling the common  $\check{\varphi}(\mathbf{x})$  and  $\nabla \varphi(\mathbf{v})^\top \mathbf{y}$  terms, the difference  $\Delta = \text{RHS} - \text{LHS}$  is

$$\begin{aligned} \Delta &= g(\mathbf{x}) \cdot (\varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{v}) - g(\mathbf{y}) \cdot \varphi(\mathbf{v}) + \nabla \varphi(\mathbf{v})^\top \mathbf{y} \\ &\quad - (\varphi(\mathbf{v}) - \mathbf{v}^\top \nabla \varphi(\mathbf{v})) \cdot \nabla g(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &= g(\mathbf{x}) \cdot (\varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{v}) - g(\mathbf{y}) \cdot \varphi(\mathbf{v}) + g(\mathbf{y}) \cdot \nabla \varphi(\mathbf{v})^\top \mathbf{v} \\ &\quad - (\varphi(\mathbf{v}) - \mathbf{v}^\top \nabla \varphi(\mathbf{v})) \cdot \nabla g(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &= g(\mathbf{x}) \cdot (\varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{v}) - g(\mathbf{y}) \cdot \left(\varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{v}\right) \\ &\quad - (\varphi(\mathbf{v}) - \mathbf{v}^\top \nabla \varphi(\mathbf{v})) \cdot \nabla g(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &= (\varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{v}) \cdot (g(\mathbf{x}) - g(\mathbf{y}) - \nabla g(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})) \\ &= (\varphi(\mathbf{v}) - \nabla \varphi(\mathbf{v})^\top \mathbf{v}) \cdot D_g(\mathbf{x} \parallel \mathbf{y}). \end{aligned}$$

Thus, the identity holds, if and only if either  $\varphi(\mathbf{v}) = \nabla\varphi(\mathbf{v})^\top \mathbf{v}$  for every  $\mathbf{v} \in \mathcal{X}_g$ , or  $D_g(\mathbf{x}||\mathbf{y}) = 0$ . The latter is true if and only if  $g$  is affine from Equation 2. The result follows.  $\blacksquare$

It is easy to check that Theorem 1 in fact holds for separable (matrix) trace divergences [Kulis et al., 2009] of the form

$$D_\varphi(\mathbf{x}||\mathbf{y}) \doteq \varphi(\mathbf{x}) - \varphi(\mathbf{y}) - \text{tr}(\nabla\varphi(\mathbf{y})^\top(\mathbf{x} - \mathbf{y})) , \quad (3)$$

with  $\varphi, g : \mathbf{S}(d) \rightarrow \mathbb{R}$  (for  $\mathbf{S}(d)$  the set of symmetric real matrices), with  $\varphi$  convex. In this case, the restricted positive homogeneity property becomes

$$\varphi(\mathbf{u}) = \text{tr}(\nabla\varphi(\mathbf{u})^\top \mathbf{u}) , \forall \mathbf{u} \in \mathcal{X}_g . \quad (4)$$

**Proof** [Proof of Lemma 2] Note that by construction,  $g(\mathbf{r}(\mathbf{x})) = \mathbb{P}(\mathbf{X} = \mathbf{x}) / ((1 - \pi_C) \cdot \mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = C))$ , and so

$$\begin{aligned} \left( \frac{1}{g(\mathbf{r}(\mathbf{x}))} \cdot \mathbf{r}(\mathbf{x}) \right)_c &= \frac{(1 - \pi_C) \cdot \mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = C)}{\mathbb{P}(\mathbf{X} = \mathbf{x})} \cdot \frac{\mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = c)}{\mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = C)} \\ &= \frac{(1 - \pi_C)}{\pi_c} \cdot \frac{\pi_c \mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = c)}{\mathbb{P}(\mathbf{X} = \mathbf{x})} \\ &= \eta(\mathbf{x}) . \end{aligned} \quad (5)$$

Furthermore,

$$\begin{aligned} \mathbb{P}(\mathbf{X} = \mathbf{x}) &= \sum_{c=1}^C \pi_c \mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = c) \\ &= (1 - \pi_C) \cdot \left( \frac{\pi_C}{1 - \pi_C} + \sum_{c < C} \frac{\pi_c}{1 - \pi_C} \cdot \frac{\mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = c)}{\mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = C)} \right) \cdot \mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = C) \\ &= (1 - \pi_C) \cdot g(\mathbf{r}(\mathbf{x})) \cdot \mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = C) . \end{aligned} \quad (6)$$

Now let

$$\hat{\mathbf{r}}(\mathbf{x}) = \frac{1}{\hat{\eta}_C(\mathbf{x})} \cdot \hat{\boldsymbol{\eta}}(\mathbf{x}) .$$

It then comes

$$\begin{aligned} &\mathbb{E}_M[D_\varphi(\boldsymbol{\eta}(\mathbf{X})||\hat{\boldsymbol{\eta}}(\mathbf{X}))] \\ &= (1 - \pi_C) \cdot \mathbb{E}_{P_C}[g(\mathbf{r}(\mathbf{x})) \cdot D_\varphi(\boldsymbol{\eta}(\mathbf{X})||\hat{\boldsymbol{\eta}}(\mathbf{X}))] \\ &= (1 - \pi_C) \cdot \mathbb{E}_{P_C}\left[g(\mathbf{r}(\mathbf{x})) \cdot D_\varphi\left(\frac{1}{g(\mathbf{r}(\mathbf{x}))} \cdot \mathbf{r}(\mathbf{x}) \parallel \hat{\boldsymbol{\eta}}(\mathbf{X})\right)\right] \\ &= (1 - \pi_C) \cdot \mathbb{E}_{P_C}\left[g(\mathbf{r}(\mathbf{x})) \cdot D_\varphi\left(\frac{1}{g(\mathbf{r}(\mathbf{x}))} \cdot \mathbf{r}(\mathbf{x}) \parallel \frac{1}{g(\hat{\mathbf{r}}(\mathbf{X}))} \cdot \hat{\mathbf{r}}(\mathbf{X})\right)\right] \\ &= (1 - \pi_C) \cdot \mathbb{E}_{P_C}[D_{\tilde{\varphi}}(\mathbf{r}(\mathbf{X})||\hat{\mathbf{r}}(\mathbf{X}))] , \end{aligned}$$

as claimed. ■

**Proof** [Proof of Lemma 3] For any  $\mathbf{x}$ ,  $\|\nabla \check{\varphi}_p(\mathbf{x})\|_q = W$  by Corollary B. Since  $\mathbf{w}_t = \nabla \check{\varphi}_p(\boldsymbol{\theta}_{t-1})$  for suitable  $\boldsymbol{\theta}_{t-1}$ , the result follows. The result for  $\|\nabla \check{\varphi}_q(\mathbf{w}_t)\|_p$  follows similarly by Corollary B.

Note that while  $\|\mathbf{w}_t\|_q = \|\nabla \varphi(\mathbf{w}_t)\|_p$  for the standard  $p$ -LMS update [Kivinen et al., 2006, Appendix I], these norms may vary with each iteration i.e.  $\mathbf{w}_t$  may not lie in the  $L_q$  ball. ■

**Proof** [Proof of Lemma 4] Similarly to the proof of Lemma E, a key to the proof of Lemma 4 relies on branching on Kivinen et al. [2006] through the use of Theorem 1. We first note that  $D_{\check{\varphi}_q}(\mathbf{u} \|\mathbf{w}_0) = W \cdot \|\mathbf{u}\|_q$  since  $\mathbf{w}_0 = \mathbf{0}$ , and  $D_{\check{\varphi}_q}(\mathbf{u} \|\mathbf{w}_{T+1}) \geq 0$ , and so

$$\begin{aligned}
W \cdot \|\mathbf{u}\|_q &\geq D_{\check{\varphi}_q}(\mathbf{u} \|\mathbf{w}_0) - D_{\check{\varphi}_q}(\mathbf{u} \|\mathbf{w}_{T+1}) \\
&= \sum_{t=1}^T \{D_{\check{\varphi}_q}(\mathbf{u} \|\mathbf{w}_{t-1}) - D_{\check{\varphi}_q}(\mathbf{u} \|\mathbf{w}_t)\} \text{ by telescoping property} \\
&= g_q(\mathbf{u}) \cdot \sum_{t=1}^T \left\{ D_{\varphi_q} \left( \frac{\mathbf{u}}{g_q(\mathbf{u})} \left\| \frac{\mathbf{w}_{t-1}}{g_q(\mathbf{w}_{t-1})} \right\| \right) - D_{\varphi_q} \left( \frac{\mathbf{u}}{g_q(\mathbf{u})} \left\| \frac{\mathbf{w}_t}{g_q(\mathbf{w}_t)} \right\| \right) \right\} \text{ by Theorem 1} \\
&= g_q(\mathbf{u}) \cdot \sum_{t=1}^T \left\{ D_{\varphi_q} \left( \frac{\mathbf{u}}{g_q(\mathbf{u})} \|\mathbf{w}_{t-1}\| \right) - D_{\varphi_q} \left( \frac{\mathbf{u}}{g_q(\mathbf{u})} \|\mathbf{w}_t\| \right) \right\} \text{ by Lemma J} . \tag{7}
\end{aligned}$$

Recall from Lemma I that

$$D_{\varphi_q} \left( \frac{\mathbf{u}}{g_q(\mathbf{u})} \|\mathbf{w}_{t-1}\| \right) - D_{\varphi_q} \left( \frac{\mathbf{u}}{g_q(\mathbf{u})} \|\mathbf{w}_t\| \right) \geq \frac{1}{4(p-1) \left( 2 + \frac{4M}{W} + \frac{2(Y+X_p W)}{(p-1)W^2} \right)} X_p^2 \cdot (s_t^2 - r_t^2)$$

where

$$\begin{aligned}
s_t &\doteq ((1/g_q(\mathbf{u})) \cdot \mathbf{u} - \mathbf{w}_{t-1})^\top \mathbf{x}_t \\
r_t &\doteq (1/g_q(\mathbf{u})) \cdot \mathbf{u}^\top \mathbf{x}_t - y_t.
\end{aligned}$$

Note that  $R_q(\mathbf{w}_{1:T}|\mathbf{u}) = \sum_{t=1}^T (s_t^2 - r_t^2)$  by definition. Summing the above for  $t = 1, 2, \dots, T$  and telescoping sums yields

$$\begin{aligned}
R_q(\mathbf{w}_{1:T}|\mathbf{u}) &\leq 4(p-1) \left( 2 + \frac{4M}{W} + \frac{2(Y+X_p W)}{(p-1)W^2} \right) X_p^2 W^2 \\
&= 4(p-1)X_p^2 W^2 + 16(p-1)MX_p^2 W + 8(Y+X_p W)X_p^2 \\
&\leq 4(p-1)X_p^2 W^2 + (16p-8)MX_p^2 W + 8YX_p^2 . \tag{8}
\end{aligned}$$

See Figure 1 for some geometric intuition about the updates. ■

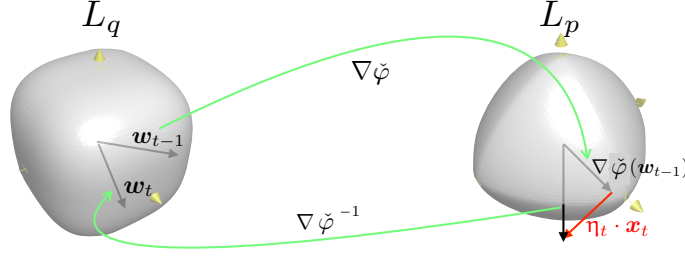


Figure 1: Illustration of the case  $W = 1$  for the  $\mathcal{B}_q(W)$ -update: all classifiers and image via  $\nabla\check{\varphi}$  belong to a ball of radius 1 (here,  $q = 3, p = 3/2$ ).

**Proof** [Proof of Lemma 5] We start by the sphere. Let  $\varphi(\mathbf{x}) \doteq (1/2) \cdot \|\mathbf{x}\|_2^2$ . Since a Bregman divergence is invariant to linear transformation, it comes from Table A1 that

$$D_\varphi\left(\frac{\mathbf{x}^S}{g_S(\mathbf{x}^S)} \parallel \frac{\mathbf{c}^S}{g_S(\mathbf{c}^S)}\right) = \frac{1}{g_S(\mathbf{c}^S)} \cdot D_{\check{\varphi}}(\mathbf{x} \parallel \mathbf{c}) = 1 - \cos D_G(\mathbf{x}, \mathbf{c}),$$

where we recall that  $D_G$  denotes the geodesic distance on the sphere (see Figure 1 and Appendix III). Equivalently,

$$\left\| \frac{1}{g_S(\mathbf{x}^S)} \cdot \mathbf{x}^S - \frac{1}{g_S(\mathbf{c}^S)} \cdot \mathbf{c}^S \right\|_2^2 = 1 - \cos D_G(\mathbf{x}, \mathbf{c}). \quad (9)$$

This equality allows us to use  $k$ -means++ using the LHS of (9) to compute the distribution that picks a center. The key to using the approximation property of  $k$ -means++ relies on the existence of a coordinate system on the sphere for which the cluster centroid is just the average of the cluster points (polar coordinates), an average that eventually has to be rescaled if the coordinate system is not that one [Dhillon and Modha, 2001, Endo and Miyamoto, 2015]. The existence of this coordinate system makes that the proof of Arthur and Vassilvitskii [2007] (and in particular the key Lemmata 3.2 and 3.3) can be carried out without modification to yield the same approximation ratio as that of Arthur and Vassilvitskii [2007] if the distortion at hand is the squared Euclidean distance, which turns out to be  $D_{\text{rec}}(\cdot : \cdot)$  from eq. (9).

The case of the hyperboloid follows the exact same path, but starts from the fact that Table A1 now brings

$$D_\varphi\left(\frac{\mathbf{x}^H}{g_H(\mathbf{x}^H)} \parallel \frac{\mathbf{c}^H}{g_H(\mathbf{c}^H)}\right) = \cosh D_G(\mathbf{y}, \mathbf{c}) - 1 = \left\| \frac{1}{g_H(\mathbf{x}^H)} \cdot \mathbf{x}^H - \frac{1}{g_H(\mathbf{c}^H)} \cdot \mathbf{c}^H \right\|_2^2.$$

To finish, in the same way as for the Sphere, we just need the existence of a coordinate system for which the centroid is an average of the cluster points, which can be obtained from hyperbolic barycentric coordinates [Ungar, 2014, Section 18]. ■

## II Additional helper lemmata

We begin with some helper lemmata that will be used in some of the proofs. In what follows, let

$$\begin{aligned}\varphi_q(\mathbf{w}) &= (1/2)(W^2 + \|\mathbf{w}\|_q^2) \\ \check{\varphi}_q(\mathbf{w}) &= W \cdot \|\mathbf{w}\|_q\end{aligned}$$

for some  $W > 0$  and  $p, q \in (1, \infty)$  such that  $1/p + 1/q = 1$ .

### II.1 Properties of $\varphi_q$ and $\check{\varphi}_q$

We use the following properties of  $\varphi_q, \check{\varphi}_q$ .

**Lemma A** *For any  $\mathbf{w}$ ,*

$$\begin{aligned}\nabla \varphi_q(\mathbf{w}) &= \|\mathbf{w}\|_q^{2-q} \cdot \text{sign}(\mathbf{w}) \otimes |\mathbf{w}|^{q-1} \\ \nabla \check{\varphi}_q(\mathbf{w}) &= W \cdot \|\mathbf{w}\|_q^{1-q} \cdot \text{sign}(\mathbf{w}) \otimes |\mathbf{w}|^{q-1},\end{aligned}$$

where  $\otimes$  denotes Hadamard product.

**Proof** The first identity was shown in [Kivinen et al., 2006, Example 1]. The second identity follows from a simple calculation. ■

This implies the following useful relations between the gradients of  $\varphi_q$  and  $\check{\varphi}_q$ .

**Corollary B** *For any  $\mathbf{w}$ ,*

$$\begin{aligned}\nabla \varphi_q(\mathbf{w}) &= (\|\mathbf{w}\|_q/W) \cdot \nabla \check{\varphi}_q(\mathbf{w}), \\ \|\nabla \check{\varphi}_q(\mathbf{w})\|_p &= W, \\ \|\nabla \varphi_q(\mathbf{w})\|_p &= \|\mathbf{w}\|_q.\end{aligned}$$

**Proof** [Proof of Corollary B] The proof follows by direct application of Lemma A and the definition of  $p, q$ . Note the third identity was shown in Kivinen et al. [2006, Appendix I]. ■

As a consequence, we conclude that the gradients of  $\varphi_q$  and  $\check{\varphi}_q$  coincide when considering vectors on the  $W$ -sphere.

**Lemma C** *For any  $\|\mathbf{w}\|_q = W$ ,*

$$\nabla \varphi_q(\mathbf{w}) = \nabla \check{\varphi}_q(\mathbf{w}).$$

**Proof** This follows from the relation between  $\nabla \varphi_q$  and  $\nabla \check{\varphi}_q$  from Lemma A. ■

Finally, we have the following result about the composition of gradients.

**Lemma D** *For any  $\mathbf{w}$ ,*

$$\nabla \varphi_q \circ \nabla \check{\varphi}_p(\mathbf{w}) = \nabla \check{\varphi}_q \circ \nabla \check{\varphi}_p(\mathbf{w}) = \frac{W}{\|\mathbf{w}\|_p} \cdot \mathbf{w}.$$

**Proof** For the first identity, applying Lemma A twice,

$$\begin{aligned}
\nabla \varphi_q \circ \nabla \check{\varphi}_p(\mathbf{w}) &= \frac{1}{\|\nabla \check{\varphi}_p(\mathbf{w})\|_q^{q-2}} \cdot \text{sign}(\nabla \check{\varphi}_p(\mathbf{w})) \otimes |\nabla \check{\varphi}_p(\mathbf{w})|^{q-1} \\
&= \frac{1}{W_q^{q-2}} \cdot \text{sign}(\mathbf{w}) \otimes \frac{W^{q-1}}{\|\mathbf{w}\|_p^{(p-1)(q-1)}} \cdot |\mathbf{w}|^{(p-1)(q-1)} \\
&= \frac{W}{\|\mathbf{w}\|_p} \cdot \mathbf{w} .
\end{aligned} \tag{10}$$

For the second identity, use Corollary B to conclude that

$$\begin{aligned}
\nabla \check{\varphi}_q \circ \nabla \check{\varphi}_p(\mathbf{w}) &= \frac{W}{\|\nabla \check{\varphi}_p(\mathbf{w})\|_q} \cdot \nabla \varphi_q(\nabla \check{\varphi}_p(\mathbf{w})) \\
&= W \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|_p},
\end{aligned}$$

as claimed. ■

## II.2 Bound on successive iterate divergence

The following Lemma extends [Kivinen et al., 2006, Appendix I] to  $\check{\varphi}_q$ .

**Lemma E** For any  $\mathbf{w}$  and  $\boldsymbol{\delta}$ ,

$$\begin{aligned}
D_{\check{\varphi}_q}(\mathbf{w} \|\nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta})) \\
\leq \frac{(p-1)\|\mathbf{w}\|_q W}{2} \cdot \left\| \frac{1}{\|\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}\|_p} \cdot (\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}) - \frac{1}{W} \cdot \nabla \check{\varphi}_q(\mathbf{w}) \right\|_p^2 .
\end{aligned} \tag{11}$$

**Proof** [Proof of Lemma E] In this proof,  $\circ$  denotes composition and  $\otimes$  is Hadamard product. The key step in the proof is the use of Theorem 1 to “branch” on the proof of [Kivinen et al., 2006, Appendix I] on the first following identity (letting  $\varphi_q(\mathbf{w}) \doteq (1/2) \cdot (W^2 + \|\mathbf{w}\|_q^2)$ ). We also make use of the dual

symmetry of Bregman divergences and we obtain third identity of:

$$\begin{aligned}
& D_{\check{\varphi}_q}(\mathbf{w} \|\nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta})) \\
&= \frac{\|\mathbf{w}\|_q}{W} \cdot D_{\varphi_q} \left( \frac{W}{\|\mathbf{w}\|_q} \cdot \mathbf{w} \left\| \frac{W}{\|\nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta})\|_q} \cdot \nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}) \right\| \right) \\
&= \frac{\|\mathbf{w}\|_q}{W} \cdot D_{\varphi_q} \left( \frac{W}{\|\mathbf{w}\|_q} \cdot \mathbf{w} \|\nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta})\| \right) \tag{12}
\end{aligned}$$

$$\begin{aligned}
&= \frac{\|\mathbf{w}\|_q}{W} \cdot D_{\varphi_p} \left( \nabla \varphi_q \circ \nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}) \left\| \nabla \varphi_q \left( \frac{W}{\|\mathbf{w}\|_q} \cdot \mathbf{w} \right) \right\| \right) \text{ by dual symmetry} \\
&= \frac{\|\mathbf{w}\|_q}{W} \cdot D_{\varphi_p} \left( \frac{W}{\|\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}\|_p} \cdot (\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}) \left\| \frac{W}{\|\mathbf{w}\|_q} \cdot \nabla \varphi_q(\mathbf{w}) \right\| \right) \tag{13}
\end{aligned}$$

$$= \frac{\|\mathbf{w}\|_q}{W} \cdot D_{\varphi_p} \left( \frac{W}{\|\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}\|_p} \cdot (\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}) \|\nabla \check{\varphi}_q(\mathbf{w})\| \right) \tag{14}$$

$$= \|\mathbf{w}\|_q W \cdot D_{\varphi_p} \left( \frac{1}{\|\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}\|_p} \cdot (\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}) \left\| \frac{1}{W} \cdot \nabla \check{\varphi}_q(\mathbf{w}) \right\| \right). \tag{15}$$

Equations (12) – (14) hold because of Corollary B. We now use Appendix I<sup>1</sup> in Kivinen et al. [2006] on Equation (15) and obtain

$$\begin{aligned}
& D_{\check{\varphi}_q}(\mathbf{w} \|\nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta})) \\
&\leq \frac{(p-1)\|\mathbf{w}\|_q W}{2} \cdot \left\| \frac{1}{\|\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}\|_p} \cdot (\nabla \check{\varphi}_q(\mathbf{w}) + \boldsymbol{\delta}) - \frac{1}{W} \cdot \nabla \check{\varphi}_q(\mathbf{w}) \right\|_p^2,
\end{aligned}$$

as claimed. ■

### II.3 Bound on successive iterate divergence to target

In what follows, we write the DN- $p$ -LMS updates as  $\mathbf{w}_t = \nabla \check{\varphi}_p(\boldsymbol{\theta}_t)$ , where

$$\boldsymbol{\theta}_t \doteq \nabla \check{\varphi}_q(\mathbf{w}_{t-1}) - \Delta_t$$

for  $\Delta_t = \eta_t \cdot (\mathbf{w}_{t-1}^\top \mathbf{x}_t - y_t) \cdot \mathbf{x}_t$ . Further, for notational ease, we write

$$\bar{\mathbf{u}} \doteq \frac{\mathbf{u}}{g_q(\mathbf{u})}$$

and

$$\bar{\boldsymbol{\theta}}_t \doteq \frac{\boldsymbol{\theta}_t}{\|\boldsymbol{\theta}_t\|_p}.$$

We have the following preliminary bound on the distance from iterates of DN- $p$ -LMS to the (normalised) target.

---

<sup>1</sup>This result is stated as a bound on  $D_{\varphi_q}(\mathbf{w} \|\nabla \varphi_q^{-1}(\nabla \varphi_q(\mathbf{w}) + \boldsymbol{\delta}))$ , which by the Bregman dual symmetry property is equivalent to a bound on  $D_{\varphi_p}(\nabla \varphi_q(\mathbf{w}) + \boldsymbol{\delta} \|\nabla \varphi_q(\mathbf{w}))$ .



**Lemma F** Fix any learning rate sequence  $\{\eta_t\}_{t=1}^T$ . Pick any  $\mathbf{u}$ , and consider iterates  $\{\mathbf{w}_t\}_{t=0}^T$  as per the update equation:

$$\mathbf{w}_t \doteq \nabla \check{\varphi}_p(\nabla \check{\varphi}_q(\mathbf{w}_{t-1}) - \eta_t \cdot \nabla \ell_t) . \quad (16)$$

Denote  $s_t \doteq (\bar{\mathbf{u}} - \mathbf{w}_{t-1})^\top \mathbf{x}_t$ ,  $r_t \doteq \bar{\mathbf{u}}^\top \mathbf{x}_t - y_t$ , and  $\alpha_t \doteq \frac{W}{\|\bar{\boldsymbol{\theta}}_t\|_p}$ . Suppose  $\|\mathbf{x}_t\|_p \leq X_p$ . Then,

$$D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_{t-1}) - D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_t) \geq Q + R + S + T ,$$

with

$$\begin{aligned} Q &\doteq \frac{\alpha_t}{2} \eta_t (s_t^2 - r_t^2), \\ R &\doteq (1 - \alpha_t) \cdot \underbrace{\left( W^2 - \bar{\mathbf{u}}^\top \nabla \check{\varphi}_q(\mathbf{w}_{t-1}) \right)}_{\in [0, 2W^2]}, \\ S &\doteq \frac{p-1}{2} \cdot \underbrace{\left( 2\alpha_t^2 \eta_t^2 (s_t - r_t)^2 X_p^2 - \|(s_t - r_t) \eta_t \alpha_t \cdot \mathbf{x}_t - (1 - \alpha_t) \cdot \nabla \check{\varphi}_q(\mathbf{w}_{t-1})\|_p^2 \right)}_{\geq -2(1-\alpha_t)^2 W^2}, \\ T &\doteq \frac{\alpha_t}{2} \eta_t (s_t - r_t)^2 (1 - 2(p-1) \eta_t \alpha_t X_p^2) . \end{aligned}$$

**Proof** [Proof of Lemma F] The Bregman triangle equality (also called the three points property) [Boissonnat et al., 2010, Property 5], [Cesa-Bianchi and Lugosi, 2006, Lemma 11.1] brings:

$$\begin{aligned} D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_{t-1}) - D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_t) &= (\bar{\mathbf{u}} - \mathbf{w}_{t-1})^\top (\nabla \varphi_q(\mathbf{w}_t) - \nabla \varphi_q(\mathbf{w}_{t-1})) - D_{\varphi_q}(\mathbf{w}_{t-1} \|\mathbf{w}_t) \\ &= (\bar{\mathbf{u}} - \mathbf{w}_{t-1})^\top (\nabla \check{\varphi}_q(\mathbf{w}_t) - \nabla \check{\varphi}_q(\mathbf{w}_{t-1})) - D_{\check{\varphi}_q}(\mathbf{w}_{t-1} \|\mathbf{w}_t) \text{ by Lemma 3 (main file) and Lemma C .} \end{aligned}$$

We now have

$$\nabla \check{\varphi}_q(\mathbf{w}_t) = \nabla \check{\varphi}_q \circ \nabla \check{\varphi}_p(\boldsymbol{\theta}_t) = W \cdot \bar{\boldsymbol{\theta}}_t$$

by Corollary B. We get

$$\begin{aligned} D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_{t-1}) - D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_t) &\geq (\bar{\mathbf{u}} - \mathbf{w}_{t-1})^\top (W \cdot \bar{\boldsymbol{\theta}}_t - \nabla \check{\varphi}_q(\mathbf{w}_{t-1})) - \frac{(p-1)W^2}{2} \cdot \left\| \bar{\boldsymbol{\theta}}_t - \frac{1}{W} \cdot \nabla \check{\varphi}_q(\mathbf{w}_{t-1}) \right\|_p^2 \\ &= (\bar{\mathbf{u}} - \mathbf{w}_{t-1})^\top (W \cdot \bar{\boldsymbol{\theta}}_t - \nabla \check{\varphi}_q(\mathbf{w}_{t-1})) - \frac{p-1}{2} \cdot \|W \cdot \bar{\boldsymbol{\theta}}_t - \nabla \check{\varphi}_q(\mathbf{w}_{t-1})\|_p^2 . \end{aligned}$$

Now, note that

$$\boldsymbol{\theta}_t = \nabla \check{\varphi}_q(\mathbf{w}_{t-1}) + \eta_t \cdot (s_t - r_t) \cdot \mathbf{x}_t .$$

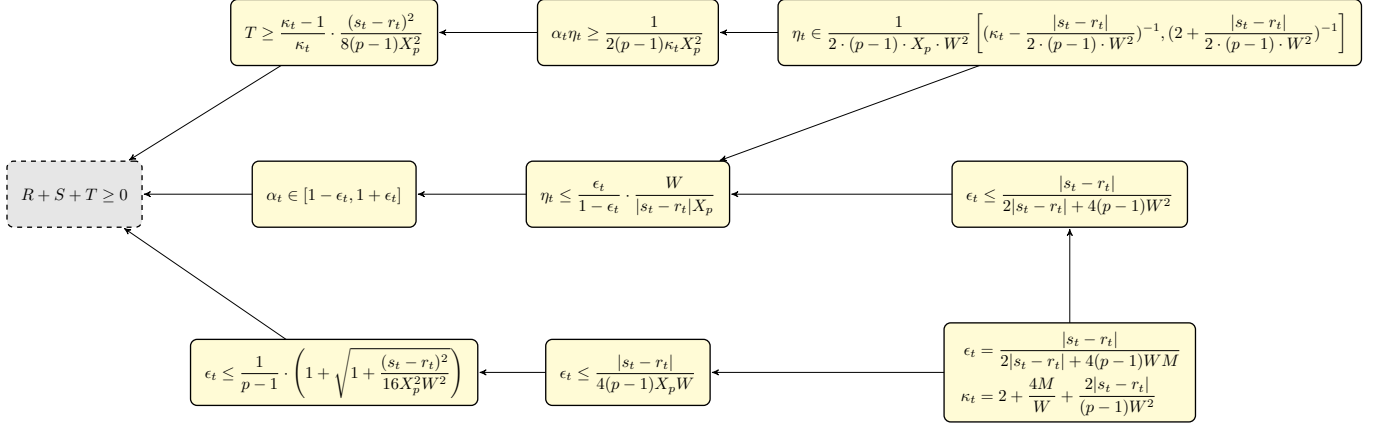


Figure 2: Schematic of proof of Lemma G. Arrows from equation  $A$  to  $B$  indicate that  $A \implies B$ .

We can thus rewrite the above as

$$\begin{aligned}
& D_{\varphi_q}(\bar{\mathbf{u}} \parallel \mathbf{w}_{t-1}) - D_{\varphi_q}(\bar{\mathbf{u}} \parallel \mathbf{w}_t) \\
& \geq s_t(s_t - r_t)\eta_t\alpha_t + (1 - \alpha_t) (\mathbf{w}_{t-1}^\top \nabla \check{\varphi}_q(\mathbf{w}_{t-1}) - \bar{\mathbf{u}}^\top \nabla \check{\varphi}_q(\mathbf{w}_{t-1})) \\
& \quad - \frac{p-1}{2} \cdot \|(s_t - r_t)\eta_t\alpha_t \cdot \mathbf{x}_t - (1 - \alpha_t) \cdot \nabla \check{\varphi}_q(\mathbf{w}_{t-1})\|_p^2 \\
& = s_t(s_t - r_t)\eta_t\alpha_t + (1 - \alpha_t) (W^2 - \bar{\mathbf{u}}^\top \nabla \check{\varphi}_q(\mathbf{w}_{t-1})) \\
& \quad - \frac{p-1}{2} \cdot \|(s_t - r_t)\eta_t\alpha_t \cdot \mathbf{x}_t - (1 - \alpha_t) \cdot \nabla \check{\varphi}_q(\mathbf{w}_{t-1})\|_p^2 \text{ by definition of } \nabla \check{\varphi}_q \\
& = Q + R + S + T,
\end{aligned}$$

as claimed. ■

We can show that the sum  $R + S + T \geq 0$ . This proof involves chaining together multiple simple inequalities. We give a high level overview in Figure 2.

**Lemma G** *Let  $R, S, T$  be as per Lemma F. Suppose we fix*

$$\eta_t = \gamma \cdot \frac{W}{4(p-1)MX_pW + |y_t - \mathbf{w}_{t-1}^\top \mathbf{x}_t|X_p}, \quad (17)$$

for any  $\gamma \in [1/2, 1]$ , and  $M \doteq \max\{W, X_p\}$ . Then,  $T + R + S \geq 0$ .

**Proof** The triangle inequality and the fact that  $\|\nabla \check{\varphi}_q(\mathbf{w}_{t-1})\|_p = W$  brings

$$\begin{aligned}
\alpha_t & \in \left[ \frac{W}{\|\nabla \check{\varphi}_q(\mathbf{w}_{t-1})\|_p + \eta_t |s_t - r_t| \cdot \|\mathbf{x}_t\|_p}, \frac{W}{\|\nabla \check{\varphi}_q(\mathbf{w}_{t-1})\|_p - \eta_t |s_t - r_t| \cdot \|\mathbf{x}_t\|_p} \right] \\
& \subseteq \left[ \frac{W}{W + \eta_t |s_t - r_t| X_p}, \frac{W}{W - \eta_t |s_t - r_t| X_p} \right], \quad (18)
\end{aligned}$$

assuming that  $\eta_t$  is chosen so that

$$\eta_t \leq \frac{W}{|s_t - r_t| \cdot X_p}, \quad (19)$$

so that the right bound is non negative. To indeed ensure this, suppose that for some  $0 < \epsilon_t \leq 1/2$ , we fix

$$\eta_t \leq \frac{\epsilon_t}{1 - \epsilon_t} \cdot \frac{W}{|s_t - r_t| X_p}. \quad (20)$$

We would in addition obtain from Equation (18) that  $\alpha_t \in [1 - \epsilon_t, 1 + \epsilon_t]$ . Suppose  $\eta_t$  is also fixed to ensure

$$\eta_t \in \left[ \frac{W}{2(p-1)\kappa_t X_p W^2 - |s_t - r_t| X_p}, \frac{W}{4(p-1)X_p W^2 + |s_t - r_t| X_p} \right], \quad (21)$$

for some  $\kappa_t$  such that

$$\kappa_t \geq 2 + \frac{|s_t - r_t|}{(p-1)W^2}. \quad (22)$$

Notice that constraint on  $\kappa_t$  makes the interval non empty and its left bound strictly positive. Assuming (21) holds, we would have

$$\alpha_t \eta_t \in \left[ \frac{1}{2(p-1)\kappa_t X_p^2}, \frac{1}{4(p-1)X_p^2} \right]. \quad (23)$$

The left bound of (23) holds because

$$\begin{aligned} \alpha_t \eta_t &\geq \eta_t \cdot \frac{W}{W + \eta_t |s_t - r_t| X_p} \\ &\geq \frac{1}{2(p-1)\kappa_t X_p^2}. \end{aligned} \quad (24)$$

The first inequality holds because of (18) and the second one holds because of (21). The right bound of (23) holds because of (18), and so

$$\begin{aligned} \alpha_t \eta_t &\leq \eta_t \cdot \frac{W}{W - \eta_t |s_t - r_t| X_p} \\ &\leq \frac{1}{4(p-1)X_p^2}, \end{aligned} \quad (25)$$

where the last inequality is due to (21). Equation (23) makes that  $T(\eta_t \alpha_t)$  is at least its value when  $\alpha_t \eta_t$  attains the lower bound of (24), that is,

$$T(\eta_t \alpha_t) \geq \frac{\kappa_t - 1}{\kappa_t} \cdot \frac{(s_t - r_t)^2}{8(p-1)X_p^2}. \quad (26)$$

Now, to guarantee  $\alpha_t \in [1 - \epsilon_t, 1 + \epsilon_t]$ , it is sufficient that the right-hand side of inequality (20) belongs to interval (21) *and* we pick  $\eta_t$  within the interval [left bound (21), right-hand side (20)]. To guarantee that the right-hand side of inequality (20) falls in interval (21), we need first,

$$\frac{W}{2(p-1)\kappa_t X_p W^2 - |s_t - r_t| X_p} \leq \frac{\epsilon_t}{1 - \epsilon_t} \cdot \frac{W}{|s_t - r_t| X_p}, \quad (27)$$

that is,

$$\kappa_t \geq \frac{1}{\epsilon_t} \cdot \frac{|s_t - r_t|}{2(p-1)W^2}. \quad (28)$$

To guarantee that the right-hand side of inequality (20) falls in interval (21) we need then

$$\frac{W}{4(p-1)X_p W^2 + |s_t - r_t| X_p} \geq \frac{\epsilon_t}{1 - \epsilon_t} \cdot \frac{W}{|s_t - r_t| X_p}, \quad (29)$$

that is,

$$\epsilon_t \leq \frac{|s_t - r_t|}{2|s_t - r_t| + 4(p-1)W^2}. \quad (30)$$

To summarize, if we pick any strictly positive  $\epsilon_t$  following inequality (30) (note  $\epsilon_t < 1$ ) and

$$\kappa_t \doteq 2 + \frac{1}{\epsilon_t} \cdot \frac{|s_t - r_t|}{(p-1)W^2}, \quad (31)$$

then we shall have both  $\alpha_t \in [1 - \epsilon_t, 1 + \epsilon_t]$  and inequality (26) holds as well. In this case, we shall have

$$\begin{aligned} T + R + S &\geq \left(1 - \frac{1}{2 + \frac{1}{\epsilon_t} \cdot \frac{|s_t - r_t|}{(p-1)W^2}}\right) \cdot \frac{(s_t - r_t)^2}{8(p-1)X_p^2} - 2\epsilon_t W^2 - (p-1)\epsilon_t^2 W^2 \\ &\geq \left(1 - \frac{1}{2}\right) \cdot \frac{(s_t - r_t)^2}{8(p-1)X_p^2} - 2\epsilon_t W^2 - (p-1)\epsilon_t^2 W^2 \\ &= \frac{(s_t - r_t)^2}{16(p-1)X_p^2} - 2\epsilon_t W^2 - (p-1)\epsilon_t^2 W^2. \end{aligned} \quad (32)$$

To finish up, we want to solve for  $\epsilon_t$  the right-hand side such that it is non negative, and we find that  $\epsilon_t$  has to satisfy

$$\epsilon_t \leq \frac{1}{p-1} \cdot \left(1 + \sqrt{1 + \frac{(s_t - r_t)^2}{16X_p^2 W^2}}\right). \quad (33)$$

Since  $\sqrt{1+x} \geq \sqrt{x}$ , a sufficient condition is

$$\epsilon_t \leq \frac{|s_t - r_t|}{4(p-1)X_p W}. \quad (34)$$

To ensure this and inequality (30), it is sufficient that we fix

$$\epsilon_t \doteq \frac{|s_t - r_t|}{2|s_t - r_t| + 4(p-1)WM} , \quad (35)$$

where  $M \doteq \max\{W, X_p\}$ . With this expression for  $\epsilon_t$ , we get from (31),

$$\kappa_t \doteq 2 + \frac{4M}{W} + \frac{2|s_t - r_t|}{(p-1)W^2} . \quad (36)$$

For these choices, Lemma H implies that the given  $\eta_t$  is feasible. ■

**Lemma H** *Suppose  $\epsilon_t$  satisfies (35) and  $\kappa_t$  satisfies (36). Then, a sufficient condition for  $\eta_t$  to satisfy both (20) and (21) is*

$$\eta_t = \gamma \cdot \frac{W}{4(p-1)MX_pW + |y_t - \mathbf{w}_{t-1}^\top \mathbf{x}_t|X_p} ,$$

for any  $\gamma \in [1/2, 1]$ .

**Proof** [Proof of Lemma H] Notice the range of values authorized for  $\eta_t$ :

$$\begin{aligned} \eta_t &\in \left[ \frac{W}{2(p-1)\kappa_t X_p W^2 - |s_t - r_t|X_p}, \frac{\epsilon_t}{1 - \epsilon_t} \cdot \frac{W}{|s_t - r_t|X_p} \right] \\ &= \left[ \frac{W}{2(p-1) \left( 2 + \frac{4M}{W} + \frac{2|s_t - r_t|}{(p-1)W^2} \right) X_p W^2 - |s_t - r_t|X_p}, \frac{W}{4(p-1)MX_pW + |s_t - r_t|X_p} \right] \\ &= \left[ \frac{W}{2(2(p-1)W^2 + 4M(p-1)W + 2|s_t - r_t|)X_p - |s_t - r_t|X_p}, \frac{W}{4(p-1)MX_pW + |s_t - r_t|X_p} \right] \\ &= \left[ \frac{W}{4(p-1)X_pW^2 + 8(p-1)MX_pW + 3|s_t - r_t|X_p}, \frac{W}{4(p-1)MX_pW + |s_t - r_t|X_p} \right] \\ &\supset \left[ \frac{W}{8(p-1)MX_pW + 2|s_t - r_t|X_p}, \frac{W}{4(p-1)MX_pW + |s_t - r_t|X_p} \right] . \end{aligned} \quad (37)$$

A sufficient condition for  $\eta_t$  to fall in interval (37) is

$$\eta_t = \gamma \cdot \frac{W}{4(p-1)MX_pW + |y_t - \mathbf{w}_{t-1}^\top \mathbf{x}_t|X_p} ,$$

for any  $\gamma \in [1/2, 1]$ . ■

**Lemma I** Suppose we fix the learning rate as per (17). Pick any  $\mathbf{u}$ , and consider iterates  $\{\mathbf{w}_t\}_{t=0}^T$  as per Equation 10. Suppose  $\|\mathbf{x}_t\|_p \leq X_p$  and  $|y_t| \leq Y, \forall t \leq T$ . Then, for any  $t$ ,

$$D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_{t-1}) - D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_t) \geq \frac{1}{4(p-1) \left(2 + \frac{4M}{W} + \frac{2(Y+X_p W)}{(p-1)W^2}\right)} X_p^2 \cdot (s_t^2 - r_t^2)$$

where  $s_t \doteq (\bar{\mathbf{u}} - \mathbf{w}_{t-1})^\top \mathbf{x}_t$ ,  $r_t \doteq \bar{\mathbf{u}}^\top \mathbf{x}_t - y_t$ .

**Proof** [Proof of Lemma I] We start from the bound of Lemma F:

$$\begin{aligned} D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_{t-1}) - D_{\varphi_q}(\bar{\mathbf{u}} \|\mathbf{w}_t) &\geq Q + R + S + T \\ &\geq Q \text{ by Lemma G} \\ &= \frac{\alpha_t}{2} \eta_t (s_t^2 - r_t^2) \text{ by definition} \\ &\geq \frac{1}{4(p-1) \kappa_t X_p^2} \cdot (s_t^2 - r_t^2) \\ &\geq \frac{1}{4(p-1) \left(2 + \frac{4M}{W} + \frac{2 \max_t |y_t - \mathbf{w}_{t-1}^\top \mathbf{x}_t|}{(p-1)W^2}\right)} X_p^2 \cdot (s_t^2 - r_t^2) \\ &\geq \frac{1}{4(p-1) \left(2 + \frac{4M}{W} + \frac{2(Y+X_p W)}{(p-1)W^2}\right)} X_p^2 \cdot (s_t^2 - r_t^2) . \end{aligned} \quad (38)$$

The last constraint to check for this bound to be valid is our  $\epsilon_t$  in (35) has to be  $< 1/2$  from inequality (19), which trivially holds since  $4(p-1)WM \geq 0$ . We conclude by noting Lemma H provides a feasible value of  $\eta_t$ . ■

## II.4 Gauge normalisation

The following lemma about the gauge of  $\mathbf{x}$  will be useful.

**Lemma J** Let  $g_q(\mathbf{x}) = \|\mathbf{x}\|_q / W$  for some  $W > 0$ . Then, for the iterates  $\{\mathbf{w}_t\}$  as per Equation 9,  $g_q(\mathbf{w}_t) = 1$ .

**Proof** We have

$$\begin{aligned} g_q(\mathbf{w}_t) &= \frac{\|\mathbf{w}_t\|_q}{W} \\ &= \frac{W}{W} \text{ by Lemma 3} \\ &= 1, \forall t \geq 1 . \end{aligned} \quad (39)$$

■

### III Working out examples of Table A1

We fill in the details justifying each of the examples of Equation 3 provided in Table 2. We also provide the form of the corresponding divergences  $D_\varphi$  and distortions  $D_{\check{\varphi}}$  in the augmented Table A1.

	$\varphi$	$D_\varphi(\mathbf{x} \parallel \mathbf{y})$	$g$	$\check{\varphi}$	$D_{\check{\varphi}}(\mathbf{x} \parallel \mathbf{y})$
I	$\frac{1}{2} \cdot (1 + \ \mathbf{x}\ _2^2)$	$(1/2) \cdot \ \mathbf{x} - \mathbf{y}\ _2^2$	$\ \mathbf{x}\ _2$	$\ \mathbf{x}\ _2$	$\ \mathbf{x}\ _2 \cdot (1 - \cos \angle \mathbf{x}, \mathbf{y})$
II	$\frac{1}{2} \cdot (W + \ \mathbf{x}\ _q^2)$	$(1/2) \cdot (\ \mathbf{x}\ _q^2 - \ \mathbf{y}\ _q^2) - \sum_i \frac{(x_i - y_i) \cdot \text{sign}(y_i) \cdot  y_i ^{q-1}}{\ \mathbf{y}\ _q^{q-2}}$	$\frac{\ \mathbf{x}\ _q}{W}$	$W \cdot \ \mathbf{x}\ _q$	$W \cdot \ \mathbf{x}\ _q - W \cdot \sum_i \frac{x_i \cdot \text{sign}(y_i) \cdot  y_i ^{q-1}}{\ \mathbf{y}\ _q^{q-1}}$
III	$\frac{1}{2} \cdot (u^2 + \ \mathbf{x}^S\ _2^2)$	$(1/2) \cdot \ \mathbf{x}^S - \mathbf{y}^S\ _2^2$	$\frac{\ \mathbf{x}\ _2}{\sin \ \mathbf{x}\ _2}$	$\ \mathbf{x}^S\ _2$	$\frac{\ \mathbf{x}\ _2}{\sin \ \mathbf{x}\ _2} \cdot (1 - \cos D_G(\mathbf{x}, \mathbf{y}))$
IV	$\frac{1}{2} \cdot (u^2 + \ \mathbf{x}^H\ _2^2)$	$(1/2) \cdot \ \mathbf{x}^H - \mathbf{y}^H\ _2^2$	$-\frac{\ \mathbf{x}\ _2}{\sinh \ \mathbf{x}\ _2}$	$\ \mathbf{x}^H\ _2$	$-\frac{\ \mathbf{x}\ _2}{\sinh \ \mathbf{x}\ _2} \cdot (\cosh D_G(\mathbf{x}, \mathbf{y}) - 1)$
V	$\sum_i x_i \log x_i - x_i$	$\sum_i x_i \log \frac{x_i}{y_i} - \mathbf{1}^\top (\mathbf{x} - \mathbf{y})$	$\mathbf{1}^\top \mathbf{x}$	$\sum_i x_i \log x_i - \mathbf{1}^\top \mathbf{x} - (\mathbf{1}^\top \mathbf{x}) \log(\mathbf{1}^\top \mathbf{x})$	$\sum_i x_i \log \frac{x_i}{y_i} - d \cdot \mathbb{E}[\mathbf{X}] \cdot \log \frac{\mathbb{E}[\mathbf{X}]}{\mathbb{E}[\mathbf{Y}]}$
VI	$-d - \sum_i \log x_i$	$\sum_i \frac{x_i}{y_i} - \sum_i \log \frac{x_i}{y_i} - d$	$\prod_i x_i^{1/d}$	$-d \cdot \prod_i x_i^{1/d}$	$\sum_i \frac{x_i (\pi_{\mathbf{y}})^{1/d}}{y_i} - d(\pi_{\mathbf{x}})^{1/d}$
VII	$\text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x})$	$\text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x} \log \mathbf{y}) - \text{tr}(\mathbf{x}) + \text{tr}(\mathbf{y})$	$\text{tr}(\mathbf{x})$	$\text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x}) - \text{tr}(\mathbf{x}) \log \text{tr}(\mathbf{x})$	$\text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x} \log \mathbf{y}) - \text{tr}(\mathbf{x}) \cdot \log \frac{\text{tr}(\mathbf{X})}{\text{tr}(\mathbf{Y})}$
VIII	$-d - \log \det(\mathbf{x})$	$\text{tr}(\mathbf{x} \mathbf{y}^{-1}) - \log \det(\mathbf{x} \mathbf{y}^{-1}) - d$	$\det(\mathbf{x}^{1/d})$	$-d \cdot \det(\mathbf{x}^{1/d})$	$\det(\mathbf{y}^{1/d}) \text{tr}(\mathbf{x} \mathbf{y}^{-1}) - d \cdot \det(\mathbf{x}^{1/d})$

Table A1: Example of distortions (right columns) that can be “reverse engineered” as Bregman divergences involving a particular, non necessary linear  $g$ . Function  $\mathbf{x}^S \doteq f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^{d+1}$  is the (S)phere lifting map defined in (51), and  $\mathbf{x}^H$  is the (H)yperboloid lifting map defined in (62).  $D_G(\cdot, \cdot)$  is the geodesic distance between the exponential map of  $\mathbf{x}$  and  $\mathbf{y}$  on their respective manifold (sphere or hyperboloid). Related proofs are in Section III. Expectation  $\mathbb{E}[\mathbf{X}]$  is a shorthand for  $(1/d) \cdot \sum_i x_i$ .  $W \in \mathbb{R}_{+*}$  and  $u \in \mathbb{R}$  are constants.

**Row I** — for  $\mathcal{X} = \mathbb{R}^d$ , consider  $\varphi(\mathbf{x}) = (1 + \|\mathbf{x}\|_2^2)/2$  and  $g(\mathbf{x}) = \|\mathbf{x}\|_2$  (we project on the Euclidean sphere). It comes

$$\check{\varphi}(\mathbf{x}) = \|\mathbf{x}\|_2 \cdot \left( \frac{1 + \left\| \frac{1}{\|\mathbf{x}\|_2} \cdot \mathbf{x} \right\|_2^2}{2} \right) = \|\mathbf{x}\|_2. \quad (40)$$

$g$  is not linear (but it is homogeneous of degree 1), but we have

$$\varphi(\mathbf{x}) = 1 = \mathbf{x}^\top \nabla \varphi(\mathbf{x}), \forall \mathbf{x} : \|\mathbf{x}\|_2 = 1, \quad (41)$$

so  $\varphi$  is 1-homogeneous on the Euclidean sphere, and we can apply Theorem 1. We have

$$\begin{aligned} g(\mathbf{x}) \cdot D_\varphi \left( \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{g(\mathbf{y})} \cdot \mathbf{y} \right) &= \frac{\|\mathbf{x}\|_2}{2} \cdot \left\| \frac{1}{\|\mathbf{x}\|_2} \cdot \mathbf{x} - \frac{1}{\|\mathbf{y}\|_2} \cdot \mathbf{y} \right\|_2^2 \\ &= \|\mathbf{x}\|_2 \cdot \left( 1 - \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} \right) = \|\mathbf{x}\|_2 \cdot (1 - \cos(\mathbf{x}, \mathbf{y})) , \end{aligned} \quad (42)$$

and we also have

$$\begin{aligned} D_{\tilde{\varphi}}(\mathbf{x} \parallel \mathbf{y}) &= \|\mathbf{x}\|_2 - \|\mathbf{y}\|_2 - \frac{1}{\|\mathbf{y}\|_2} \cdot (\mathbf{x} - \mathbf{y})^\top \mathbf{y} \\ &= \|\mathbf{x}\|_2 - \|\mathbf{y}\|_2 - \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{y}\|_2} + \|\mathbf{y}\|_2 \end{aligned} \quad (43)$$

$$= \|\mathbf{x}\|_2 \cdot \left( 1 - \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} \right) = \|\mathbf{x}\|_2 \cdot (1 - \cos(\mathbf{x}, \mathbf{y})) , \quad (44)$$

which is equal to Equation (42), so we check that Theorem 1 applies in this case.  $D_{\tilde{\varphi}}$  has some interesting properties. One is a weak form of triangle inequality.

**Lemma K**  $D_{\tilde{\varphi}}(\mathbf{x} \parallel \mathbf{y}) + D_{\tilde{\varphi}}(\mathbf{y} \parallel \mathbf{z}) \leq D_{\tilde{\varphi}}(\mathbf{x} \parallel \mathbf{z})$ ,  $\forall \mathbf{x}, \mathbf{y}, \mathbf{z}$  such that  $\|\mathbf{y}\|_2 \leq \|\mathbf{x}\|_2$ .

**Proof**

$$\begin{aligned} D_{\tilde{\varphi}}(\mathbf{x} \parallel \mathbf{y}) + D_{\tilde{\varphi}}(\mathbf{y} \parallel \mathbf{z}) &= \|\mathbf{x}\|_2 \cdot (1 - \cos(\mathbf{x}, \mathbf{y})) + \|\mathbf{y}\|_2 \cdot (1 - \cos(\mathbf{y}, \mathbf{z})) \\ &= \|\mathbf{x}\|_2 \cdot ((1 - \cos(\mathbf{x}, \mathbf{y})) + (1 - \cos(\mathbf{y}, \mathbf{z}))) + (\|\mathbf{y}\|_2 - \|\mathbf{x}\|_2) \cdot (1 - \cos(\mathbf{y}, \mathbf{z})) \\ &\leq \|\mathbf{x}\|_2 \cdot (1 - \cos(\mathbf{x}, \mathbf{z})) + (\|\mathbf{y}\|_2 - \|\mathbf{x}\|_2) \cdot (1 - \cos(\mathbf{y}, \mathbf{z})) \\ &\leq D_{\tilde{\varphi}}(\mathbf{x} \parallel \mathbf{z}) + (\|\mathbf{y}\|_2 - \|\mathbf{x}\|_2) \cdot (1 - \cos(\mathbf{y}, \mathbf{z})) \\ &\leq D_{\tilde{\varphi}}(\mathbf{x} \parallel \mathbf{z}) , \end{aligned} \quad (45)$$

since  $\|\mathbf{y}\|_2 \leq \|\mathbf{x}\|_2$ . We have used the fact that  $(1 - \cos(\mathbf{x}, \mathbf{y}))$  is half the Euclidean distance between unit-normalized vectors. ■

Another good property is that  $D_{\tilde{\varphi}}(\mathbf{x} \parallel \boldsymbol{\mu})$  can be related to the log-likelihood of a von Mises-Fisher distribution with expected direction  $\boldsymbol{\mu}$ , which happens to be useful in text analysis [Reisinger et al., 2010].

**Row II** — Let  $\varphi(\mathbf{x}) \doteq (1/2) \cdot (u^2 + \|\mathbf{x}\|_q^2)$ , for  $q > 1$  [Kivinen et al., 2006]. We have

$$\varphi \left( \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \right) = \frac{u^2}{2} + \frac{1}{2} \cdot \left\| \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \right\|_q^2 = \frac{u^2}{2} + \frac{1}{2g^2(\mathbf{x})} \cdot \|\mathbf{x}\|_q^2 . \quad (46)$$



We also have

$$\begin{aligned}
\left(\frac{1}{g(\mathbf{x})} \cdot \mathbf{x}\right)^\top \nabla \varphi\left(\frac{1}{g(\mathbf{x})} \cdot \mathbf{x}\right) &= \frac{1}{g(\mathbf{x})} \cdot \sum_i \frac{x_i \cdot \text{sign}\left(\frac{1}{g(\mathbf{x})} \cdot x_i\right) \left|\frac{1}{g(\mathbf{x})} \cdot x_i\right|^{q-1}}{\left\|\frac{1}{g(\mathbf{x})} \cdot \mathbf{x}\right\|_q^{q-2}} \\
&= \sum_i \frac{\left|\frac{1}{g(\mathbf{x})} \cdot x_i\right|^q}{\left\|\frac{1}{g(\mathbf{x})} \cdot \mathbf{x}\right\|_q^{q-2}} \\
&= \frac{1}{g^2(\mathbf{x})} \cdot \|\mathbf{x}\|_q^2 .
\end{aligned} \tag{47}$$

To have the condition of Theorem 1 satisfied, we therefore need

$$\|\mathbf{x}\|_q = u g(\mathbf{x}) , \tag{48}$$

So we use  $g(\mathbf{x}) = \|\mathbf{x}\|_q/W$  and  $u = W$ , observing that  $\varphi$  is 1-homogeneous on the  $L_p$  sphere. We check that

$$\check{\varphi}(\mathbf{x}) = W \cdot \|\mathbf{x}\|_q . \tag{49}$$

and we obtain

$$D_{\check{\varphi}}(\mathbf{w} \|\mathbf{w}') = W \cdot \|\mathbf{w}\|_q - W \cdot \sum_i \frac{w_i \cdot \text{sign}(w'_i) \cdot |w'_i|^{q-1}}{\|\mathbf{w}'\|_q^{q-1}} . \tag{50}$$

**Row III —** As in Buss and Fillmore [2001], we assume  $\|\mathbf{x}\|_2 \leq \pi$ , or we renormalize or change the radius of the ball) We first lift the data points using the *Sphere* lifting map  $\mathbb{R}^d \ni \mathbf{x} \mapsto \mathbf{x}^S \in \mathbb{R}^{d+1}$ :

$$\mathbf{x}^S \doteq [x_1 \ x_2 \ \cdots \ x_d \ r_x \cot r_x]^\top , \tag{51}$$

where  $r_x \doteq \|\mathbf{x}\|_2$  is the Euclidean norm of  $\mathbf{x}$ . Notice that the last coordinate is a coordinate of the Hessian of the geodesic distance to the origin on the sphere [Buss and Fillmore, 2001]. We then let  $g(\mathbf{x}^S) \doteq r_x / \sin r_x$  (notice that  $g$  is computed using the first  $d$  coordinates). Finally, for  $\mathcal{X} = \mathbb{R}^{d+1}$  and  $u > 1$ , consider  $\varphi(\mathbf{x}^S) = (u^2 + \|\mathbf{x}^S\|_2^2)/2$ . The set of points for which  $\varphi(\mathbf{x}^S) = (\mathbf{x}^S)^\top \nabla \varphi(\mathbf{x}^S)$  is equivalently the subset  $\mathcal{X}_g \subseteq \mathbb{R}^{d+1}$  such that

$$\mathcal{X}_g \doteq \{\mathbf{x}^S : g^2(\mathbf{x}^S) = u^2\} . \tag{52}$$

So  $\varphi$  satisfies the restricted positive homogeneity of degree 1 on  $\mathcal{X}_g$  and we can apply Theorem 1. We first remark that:

$$\begin{aligned}
\|\mathbf{x}^S\|_2^2 &= r_x^2 + r_x^2 \cot^2 r_x \\
&= \frac{r_x^2}{\sin^2 r_x} = g^2(\mathbf{x}^S) ,
\end{aligned} \tag{53}$$

and

$$\check{\varphi}(\mathbf{x}^S) = \frac{r_x}{\sin r_x} \cdot \varphi\left(\frac{\sin r_x}{r_x} \cdot \mathbf{x}^S\right) = \frac{r_x}{\sin r_x} \cdot \left(\frac{\sin r_x}{r_x}\right)^2 \cdot \|\mathbf{x}^S\|_2^2 = \|\mathbf{x}^S\|_2, \quad (54)$$

and finally, because of the spherical law of cosines,

$$\sin r_x \sin r_y \cos(\mathbf{x}, \mathbf{y}) + \cos r_x \cos r_y = \cos D_G(\mathbf{x}, \mathbf{y}), \quad (55)$$

where we recall from eq. (16) that  $D_G(\mathbf{x}, \mathbf{y})$  is the geodesic distance between the image of the exponential maps of  $\mathbf{x}$  and  $\mathbf{y}$  on the sphere. We then derive

$$\begin{aligned} g(\mathbf{x}^S) \cdot D_\varphi\left(\frac{1}{g(\mathbf{x}^S)} \cdot \mathbf{x}^S \parallel \frac{1}{g(\mathbf{y}^S)} \cdot \mathbf{y}^S\right) \\ &= \frac{r_x}{2 \sin r_x} \cdot \left\| \frac{\sin r_x}{r_x} \cdot \mathbf{x}^S - \frac{\sin r_y}{r_y} \cdot \mathbf{y}^S \right\|_2^2 \\ &= \frac{r_x}{2 \sin r_x} \cdot \left( \frac{\sin^2 r_x}{\|\mathbf{x}\|_2^2} \cdot \|\mathbf{x}^S\|_2^2 + \frac{\sin^2 r_y}{\|\mathbf{y}\|_2^2} \cdot \|\mathbf{y}^S\|_2^2 - 2 \cdot \frac{\sin r_x}{r_x} \cdot \frac{\sin r_y}{r_y} \cdot (\mathbf{x}^S)^\top \mathbf{y}^S \right) \\ &= \frac{r_x}{\sin r_x} \cdot \left( 1 - \frac{\sin r_x}{r_x} \cdot \frac{\sin r_y}{r_y} \cdot (\mathbf{x}^S)^\top \mathbf{y}^S \right) \end{aligned} \quad (56)$$

$$= \frac{r_x}{\sin r_x} \cdot \left( 1 - \frac{\sin r_x}{r_x} \cdot \frac{\sin r_y}{r_y} \cdot (\mathbf{x}^\top \mathbf{y} + r_x r_y \cot r_x \cot r_y) \right) \quad (57)$$

$$\begin{aligned} &= \frac{r_x}{\sin r_x} \cdot (1 - \sin r_x \sin r_y \cdot (\cos(\mathbf{x}, \mathbf{y}) + \cot r_x \cot r_y)) \\ &= \frac{r_x}{\sin r_x} \cdot (1 - (\sin r_x \sin r_y \cos(\mathbf{x}, \mathbf{y}) + \cos r_x \cos r_y)) \\ &= \frac{r_x}{\sin r_x} \cdot (1 - \cos D_G(\mathbf{x}, \mathbf{y})) . \end{aligned} \quad (58)$$

In Equation (56), we use Equation (53), and we use Equation (55) in Equation (58). We also check

$$\begin{aligned} D_{\check{\varphi}}(\mathbf{x}^S \parallel \mathbf{y}^S) &= \|\mathbf{x}^S\|_2 - \|\mathbf{y}^S\|_2 - \frac{1}{\|\mathbf{y}^S\|_2} \cdot (\mathbf{x}^S - \mathbf{y}^S)^\top \mathbf{y}^S \\ &= \|\mathbf{x}^S\|_2 - \frac{1}{\|\mathbf{y}^S\|_2} \cdot (\mathbf{x}^S)^\top \mathbf{y}^S \\ &= \|\mathbf{x}^S\|_2 \cdot \left( 1 - \frac{(\mathbf{x}^S)^\top \mathbf{y}^S}{\|\mathbf{x}^S\|_2 \|\mathbf{y}^S\|_2} \right) \end{aligned} \quad (59)$$

$$= \frac{r_x}{\sin r_x} \cdot (1 - \cos D_G(\mathbf{x}, \mathbf{y})) . \quad (60)$$

To obtain (60), we use the fact that

$$\frac{(\mathbf{x}^S)^\top \mathbf{y}^S}{\|\mathbf{x}^S\|_2 \|\mathbf{y}^S\|_2} = \frac{\sin r_x}{r_x} \cdot \frac{\sin r_y}{r_y} \cdot (\mathbf{x}^\top \mathbf{y} + r_x r_y \cot r_x \cot r_y) , \quad (61)$$

and then plug it into Equation (60), which yields the identity between Equation (57) (and thus (58)) and (60). So Theorem 1 holds in this case as well. We also remark that  $(1/g(\mathbf{x}^S)) \cdot \mathbf{x}^S = \exp_0(\mathbf{x})$  is the exponential map for the sphere [Buss and Fillmore, 2001].

**Row IV** — In the same way as we did for row IV, we first create a lifting map, but this time *complex* valued, the Hyperboloid lifting map  $H: \mathbb{R}^d \ni \mathbf{x} \mapsto \mathbf{x}^H \in \mathbb{R}^d \times \mathbb{C}$ . With an abuse of notation, it is given by

$$\mathbf{x}^H \doteq [x_1 \ x_2 \ \cdots \ x_d \ ir_{\mathbf{x}} \coth r_{\mathbf{x}}]^\top, \quad (62)$$

and we let  $g(\mathbf{x}^H) \doteq -r_{\mathbf{x}}/\sinh r_{\mathbf{x}}$ , with  $\coth$  and  $\sinh$  defining respectively the hyperbolic cotangent and hyperbolic sine. We let  $0 \coth 0 = 0/\sinh 0 = 1$ . Notice that the complex number is pure imaginary and so  $H$  defines a  $d$  dimensional manifold that lives in  $\mathbb{R}^{d+1}$  assuming that the last coordinate is the imaginary axis. Let  $\exp_q(\mathbf{x}) \doteq (1/g(\mathbf{x}^H)) \cdot \mathbf{x}^H$ . Notice that

$$\begin{aligned} \|\exp_q(\mathbf{x})\|_2^2 &= \frac{\sinh^2 r_{\mathbf{x}}}{r_{\mathbf{x}}^2} \cdot (r_{\mathbf{x}}^2 + i^2 r_{\mathbf{x}}^2 \coth^2 r_{\mathbf{x}}) \\ &= \sinh^2 r_{\mathbf{x}} + i^2 \cosh^2 r_{\mathbf{x}} \\ &= \sinh^2 r_{\mathbf{x}} - \cosh^2 r_{\mathbf{x}} = -1, \end{aligned} \quad (63)$$

so  $\exp_q(\mathbf{x})$  defines a lifting map from  $\mathbb{R}^d$  to the hyperboloid model  $\mathbb{H}^d$  of hyperbolic geometry [Galperin, 1993]. In fact, it defines the exponential map for the plane  $T_q \mathbb{H}^d$  tangent to  $\mathbb{H}^d$  in point

$$\mathbf{q} \doteq [0 \ 0 \ \cdots \ 0 \ i] = \mathbf{0}^H.$$

To see this, remark that we can express the geodesic distance  $D_G$  with the hyperbolic metric between  $\mathbf{x}^H$  and  $\mathbf{y}^H$  as

$$D_G(\mathbf{x}^H, \mathbf{y}^H) \doteq \cosh^{-1}(-(\mathbf{x}^H)^\top \mathbf{y}^H), \quad (64)$$

where  $\cosh^{-1}$  is the inverse hyperbolic cosine. So, for any  $\mathbf{x} \in T_q \mathbb{H}^d$ , since  $r_{\mathbf{x}} = \|\mathbf{x} - \mathbf{0}\|_2$ , we have

$$\begin{aligned} D_G(\exp_q(\mathbf{x}), \mathbf{q}) &= \cosh^{-1}(-(\mathbf{x}^H)^\top \mathbf{0}^H) \\ &= \cosh^{-1}(-i^2 \cosh r_{\mathbf{x}}) \\ &= r_{\mathbf{x}} = \|\mathbf{x} - \mathbf{0}\|_2, \end{aligned} \quad (65)$$

and  $\exp_q(\mathbf{x})$  is indeed the exponential map for  $T_q \mathbb{H}^d$ . Now, remark that

$$\begin{aligned} \exp_q(\mathbf{x})^\top \exp_q(\mathbf{y}) &= \frac{\sinh r_{\mathbf{x}}}{r_{\mathbf{x}}} \cdot \frac{\sinh r_{\mathbf{y}}}{r_{\mathbf{y}}} \cdot (\mathbf{x}^H)^\top \mathbf{y}^H \\ &= \frac{\sinh r_{\mathbf{x}}}{r_{\mathbf{x}}} \cdot \frac{\sinh r_{\mathbf{y}}}{r_{\mathbf{y}}} \cdot (\mathbf{x}^\top \mathbf{y} + i^2 r_{\mathbf{x}} r_{\mathbf{y}} \coth r_{\mathbf{x}} \coth r_{\mathbf{y}}) \\ &= \sinh r_{\mathbf{x}} \sinh r_{\mathbf{y}} \cdot (\cos(\mathbf{x}, \mathbf{y}) - \coth r_{\mathbf{x}} \coth r_{\mathbf{y}}) \\ &= \sinh r_{\mathbf{x}} \sinh r_{\mathbf{y}} \cos(\mathbf{x}, \mathbf{y}) - \cosh r_{\mathbf{x}} \cosh r_{\mathbf{y}} \\ &= -\cosh D_G(\mathbf{x}^H, \mathbf{y}^H). \end{aligned} \quad (66)$$

Eq. (66) holds by the hyperbolic law of cosines. Now, we let  $\varphi(\mathbf{x}^H) = (u^2 + \|\mathbf{x}^H\|_2^2)/2$  and

$$\mathcal{X}_g \doteq \{\mathbf{x}^H : \|\mathbf{x}^H\|_2^2 = u^2\}. \quad (67)$$

We check that  $\varphi(\mathbf{x}^H) = u^2 = (\mathbf{x}^H)^\top \nabla \varphi(\mathbf{x}^H)$  for any  $\mathbf{x}^H \in \mathcal{X}_g$ , so we can apply Theorem 1. We then use eqs. (63) and (66) and derive

$$\begin{aligned}
& g(\mathbf{x}^H) \cdot D_\varphi \left( \frac{1}{g(\mathbf{x}^H)} \cdot \mathbf{x}^H \parallel \frac{1}{g(\mathbf{y}^H)} \cdot \mathbf{y}^H \right) \\
&= -\frac{r_{\mathbf{x}}}{2 \sinh r_{\mathbf{x}}} \cdot \|\exp_{\mathbf{q}}(\mathbf{x}) - \exp_{\mathbf{q}}(\mathbf{y})\|_2^2 \\
&= -\frac{r_{\mathbf{x}}}{2 \sinh r_{\mathbf{x}}} \cdot \left( \|\exp_{\mathbf{q}}(\mathbf{x})\|_2^2 + \|\exp_{\mathbf{q}}(\mathbf{y})\|_2^2 - 2 \exp_{\mathbf{q}}(\mathbf{x})^\top \exp_{\mathbf{q}}(\mathbf{y}) \right) \\
&= -\frac{r_{\mathbf{x}}}{\sinh r_{\mathbf{x}}} \cdot (\cosh D_G(\mathbf{x}^H, \mathbf{y}^H) - 1) .
\end{aligned} \tag{68}$$

Note that eq. (68) is a negative-valued and concave distortion.

**Row V** — for  $\mathcal{X} = \mathbb{R}_{+*}^d$ , consider  $\varphi(\mathbf{x}) = \sum_i x_i \log x_i - x_i$  and  $g(\mathbf{x}) = \mathbf{1}^\top \mathbf{x}$  (we normalize on the simplex). Since  $g$  is linear, we do not need to check for the homogeneity of  $\varphi$ , and we directly obtain:

$$\begin{aligned}
g(\mathbf{x}) \cdot D_\varphi \left( \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{g(\mathbf{y})} \cdot \mathbf{y} \right) &= \sum_i x_i \log x_i - (\mathbf{1}^\top \mathbf{x}) \log(\mathbf{1}^\top \mathbf{x}) - \mathbf{1}^\top \mathbf{x} \\
&\quad - \frac{\mathbf{1}^\top \mathbf{x}}{\mathbf{1}^\top \mathbf{y}} \cdot \sum_i y_i \log y_i - (\mathbf{1}^\top \mathbf{x}) \log(\mathbf{1}^\top \mathbf{y}) + \mathbf{1}^\top \mathbf{x} \\
&\quad - (\mathbf{1}^\top \mathbf{x}) \cdot \sum_i \left( \frac{x_i}{\mathbf{1}^\top \mathbf{x}} - \frac{y_i}{\mathbf{1}^\top \mathbf{y}} \right) \cdot \log \frac{y_i}{\mathbf{1}^\top \mathbf{y}} \\
&= \sum_i x_i \log \frac{x_i}{y_i} - (\mathbf{1}^\top \mathbf{x}) \cdot \log \frac{\mathbf{1}^\top \mathbf{x}}{\mathbf{1}^\top \mathbf{y}} .
\end{aligned} \tag{69}$$

Furthermore,

$$\check{\varphi}(\mathbf{x}) = \mathbf{1}^\top \mathbf{x} \cdot \left( \sum_i \frac{x_i}{\mathbf{1}^\top \mathbf{x}} \cdot \log \frac{x_i}{\mathbf{1}^\top \mathbf{x}} - 1 \right) = \sum_i x_i \log x_i - (\mathbf{1}^\top \mathbf{x}) \log(\mathbf{1}^\top \mathbf{x}) - \mathbf{1}^\top \mathbf{x} . \tag{70}$$

Noting that  $\check{\varphi}(\mathbf{x})$  is the sum of three terms, one of which is linear and can be removed for the divergence, so the divergence is just the sum of the two divergences with the two generators, which is found to be Equation (69) as well. Remark that while the KL divergence is convex in its both arguments,  $D_{\check{\varphi}}(\mathbf{x} \parallel \mathbf{y})$  may not be (jointly) convex. Indeed, its Hessian in  $\mathbf{y}$  equals:

$$\mathbf{H}_{\mathbf{y}}(D_{\check{\varphi}}) = \text{Diag}(\{x_i/y_i^2\}_i) - \frac{\mathbf{1}^\top \mathbf{x}}{(\mathbf{1}^\top \mathbf{y})^2} \cdot \mathbf{1} \mathbf{1}^\top , \tag{71}$$

which may be indefinite.

**Row VI** — for  $\mathcal{X} = \mathbb{R}_{+*}^d$ , consider  $\varphi(\mathbf{x}) = -d - \sum_i \log x_i$  and  $g(\mathbf{x}) = (\pi_{\mathbf{x}})^{1/d}$ , where we let  $\pi_{\mathbf{x}} \doteq \prod_i x_i$  (we normalize with the geometric average). It comes

$$\check{\varphi}(\mathbf{x}) = (\pi_{\mathbf{x}})^{1/d} \cdot \left( -d - \sum_i \log \frac{x_i}{(\pi_{\mathbf{x}})^{1/d}} \right) = -d \cdot (\pi_{\mathbf{x}})^{1/d} . \quad (72)$$

$g$  is not linear (but it is homogeneous of degree 1), and we have

$$\varphi(\mathbf{x}) = -d = \mathbf{x}^\top \nabla \varphi(\mathbf{x}) , \forall \mathbf{x} : \prod_i x_i = 1 , \quad (73)$$

so  $\varphi$  is 1-homogeneous on  $\mathcal{X}_g$ , and we can apply Theorem 1. We have

$$\begin{aligned} g(\mathbf{x}) \cdot D_\varphi \left( \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{g(\mathbf{y})} \cdot \mathbf{y} \right) &= (\pi_{\mathbf{x}})^{1/d} \cdot \sum_i \left( \frac{x_i (\pi_{\mathbf{y}})^{1/d}}{y_i (\pi_{\mathbf{x}})^{1/d}} - \log \frac{x_i (\pi_{\mathbf{y}})^{1/d}}{y_i (\pi_{\mathbf{x}})^{1/d}} \right) - d (\pi_{\mathbf{x}})^{1/d} \\ &= \sum_i \frac{x_i (\pi_{\mathbf{y}})^{1/d}}{y_i} - d (\pi_{\mathbf{x}})^{1/d} \log (\pi_{\mathbf{x}})^{1/d} - (\pi_{\mathbf{x}})^{1/d} \log \pi_{\mathbf{y}} + (\pi_{\mathbf{x}})^{1/d} \log \pi_{\mathbf{y}} \\ &\quad + d (\pi_{\mathbf{x}})^{1/d} \log (\pi_{\mathbf{x}})^{1/d} - d (\pi_{\mathbf{x}})^{1/d} \\ &= \sum_i \frac{x_i (\pi_{\mathbf{y}})^{1/d}}{y_i} - d (\pi_{\mathbf{x}})^{1/d} . \end{aligned} \quad (74)$$

We also have

$$\frac{\partial}{\partial x_i} \check{\varphi}(\mathbf{x}) = -(1/x_i) \cdot (\pi_{\mathbf{x}})^{1/d} , \quad (75)$$

and so

$$\begin{aligned} D_{\check{\varphi}}(\mathbf{x} \parallel \mathbf{y}) &= -d (\pi_{\mathbf{x}})^{1/d} + d (\pi_{\mathbf{y}})^{1/d} + \sum_i (x_i - y_i) \cdot \frac{(\pi_{\mathbf{y}})^{1/d}}{y_i} \\ &= -d (\pi_{\mathbf{x}})^{1/d} + d (\pi_{\mathbf{y}})^{1/d} + \sum_i \frac{x_i (\pi_{\mathbf{y}})^{1/d}}{y_i} - d (\pi_{\mathbf{y}})^{1/d} \\ &= \sum_i \frac{x_i (\pi_{\mathbf{y}})^{1/d}}{y_i} - d (\pi_{\mathbf{x}})^{1/d} , \end{aligned} \quad (76)$$

which is equal to Equation (74), so we check that Theorem 1 applies in this case.

**Row VII** — We use the following fact Kulis et al. [2009]. Let  $\mathbf{X} = \mathbf{U}\mathbf{L}\mathbf{U}^\top$  and  $\mathbf{Y} = \mathbf{V}\mathbf{T}\mathbf{V}^\top$  be the eigendecomposition of symmetric positive definite matrices  $\mathbf{X}$  and  $\mathbf{Y}$ , with  $\mathbf{L} \doteq \text{Diag}(\mathbf{l})$ ,  $\mathbf{T} \doteq \text{Diag}(\mathbf{t})$ , and  $\mathbf{U} \doteq [\mathbf{u}_1 | \mathbf{u}_2 | \dots | \mathbf{u}_d]$ ,  $\mathbf{V} \doteq [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_d]$  orthonormal; let  $\varphi = \text{tr}(\mathbf{X} \log \mathbf{X} - \mathbf{X})$ . Then we have

$$D_\varphi(\mathbf{X} \parallel \mathbf{Y}) = \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot D_{\varphi_2}(l_i \parallel t_j) , \quad (77)$$

with  $\varphi_2(x) = x \log x - x$ . We pick  $g(\mathbf{x}) = \text{tr}(\mathbf{x}) = \sum_i l_i$ , which brings from Equation (69)

$$\begin{aligned}
& g(\mathbf{x}) \cdot D_\varphi \left( \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{g(\mathbf{Y})} \cdot \mathbf{Y} \right) \\
&= \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \text{tr}(\mathbf{x}) \cdot D_{\varphi_2} \left( \frac{l_i}{\text{tr}(\mathbf{x})} \parallel \frac{t_j}{\text{tr}(\mathbf{Y})} \right) \\
&= \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \text{tr}(\mathbf{x}) \cdot \left( \frac{l_i}{\text{tr}(\mathbf{x})} \cdot \log \frac{l_i \cdot \text{tr}(\mathbf{Y})}{t_j \cdot \text{tr}(\mathbf{x})} - \frac{l_i}{\text{tr}(\mathbf{x})} + \frac{t_j}{\text{tr}(\mathbf{Y})} \right) \\
&= \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \left( l_i \log \frac{l_i}{t_j} - l_i + t_j \right) + \log \left( \frac{\text{tr}(\mathbf{Y})}{\text{tr}(\mathbf{x})} \right) \cdot \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot l_i \\
&\quad + \frac{\text{tr}(\mathbf{x})}{\text{tr}(\mathbf{Y})} \cdot \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot t_j - \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot t_j .
\end{aligned} \tag{78}$$

Because  $\mathbf{U}, \mathbf{V}$  are orthonormal, we also get  $\sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot l_i = \sum_i l_i \sum_j \cos^2(\mathbf{u}_i, \mathbf{v}_j) = \sum_i l_i = \text{tr}(\mathbf{x})$  and  $\sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot t_j = \text{tr}(\mathbf{y})$ , and so Equation (78) becomes

$$\begin{aligned}
& g(\mathbf{x}) \cdot D_\varphi \left( \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{g(\mathbf{Y})} \cdot \mathbf{Y} \right) \\
&= \text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x} \log \mathbf{Y}) - \text{tr}(\mathbf{x}) + \text{tr}(\mathbf{Y}) + \text{tr}(\mathbf{x}) \cdot \log \left( \frac{\text{tr}(\mathbf{Y})}{\text{tr}(\mathbf{x})} \right) + \text{tr}(\mathbf{x}) - \text{tr}(\mathbf{Y}) \\
&= \text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x} \log \mathbf{Y}) - \text{tr}(\mathbf{x}) \cdot \log \left( \frac{\text{tr}(\mathbf{x})}{\text{tr}(\mathbf{Y})} \right) .
\end{aligned} \tag{79}$$

We also check that

$$\begin{aligned}
\check{\varphi}(\mathbf{x}) &= \text{tr}(\mathbf{x}) \cdot \text{tr} \left( \frac{1}{\text{tr}(\mathbf{x})} \cdot \mathbf{x} \log \left( \frac{1}{\text{tr}(\mathbf{x})} \cdot \mathbf{x} \right) - \frac{1}{\text{tr}(\mathbf{x})} \cdot \mathbf{x} \right) \\
&= \text{tr} \left( \mathbf{x} \log \left( \frac{1}{\text{tr}(\mathbf{x})} \cdot \mathbf{x} \right) \right) - \text{tr}(\mathbf{x}) ,
\end{aligned} \tag{80}$$

and

$$\begin{aligned}
\mathbf{x} \log \left( \frac{1}{\text{tr}(\mathbf{x})} \cdot \mathbf{x} \right) &= \mathbf{U} \mathbf{L} \mathbf{U}^\top \mathbf{U} \log \left( \frac{1}{\mathbf{1}^\top \mathbf{L}} \cdot \mathbf{L} \right) \mathbf{U}^\top \\
&= \mathbf{U} \mathbf{L} \log \left( \frac{1}{\mathbf{1}^\top \mathbf{L}} \cdot \mathbf{L} \right) \mathbf{U}^\top
\end{aligned} \tag{81}$$

$$= \mathbf{U} \mathbf{L} \log \mathbf{L} \mathbf{U}^\top - \log \text{tr}(\mathbf{x}) \cdot \mathbf{U} \mathbf{L} \mathbf{U}^\top , \tag{82}$$

so that  $\check{\varphi}(\mathbf{x}) = \text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x}) - \text{tr}(\mathbf{x}) \cdot \log \text{tr}(\mathbf{x})$ . Let  $\varphi_3(\mathbf{x}) \doteq \text{tr}(\mathbf{x}) \cdot \log \text{tr}(\mathbf{x})$ . We have  $\nabla \varphi_3(\mathbf{x}) = (1 + \log \text{tr}(\mathbf{x})) \cdot \mathbf{I}$ . Since a (Bregman) divergence involving a sum of generators is the sum of (Bregman)

divergences, we get

$$\begin{aligned}
D_{\check{\varphi}}(\mathbf{x} \parallel \mathbf{Y}) &= \text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x} \log \mathbf{Y} - \mathbf{x} + \mathbf{Y}) - \text{tr}(\mathbf{x}) \cdot \log \text{tr}(\mathbf{x}) + \text{tr}(\mathbf{Y}) \cdot \log \text{tr}(\mathbf{Y}) \\
&\quad + (1 + \log \text{tr}(\mathbf{Y})) \cdot \text{tr}(\mathbf{x} - \mathbf{Y}) \\
&= \text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x} \log \mathbf{Y}) - \text{tr}(\mathbf{x}) \cdot \log \text{tr}(\mathbf{x}) + \text{tr}(\mathbf{x}) \cdot \log \text{tr}(\mathbf{Y}) \\
&= \text{tr}(\mathbf{x} \log \mathbf{x} - \mathbf{x} \log \mathbf{Y}) - \text{tr}(\mathbf{x}) \cdot \log \left( \frac{\text{tr}(\mathbf{x})}{\text{tr}(\mathbf{Y})} \right), \tag{83}
\end{aligned}$$

which is Equation (79).

**Row VIII —** We have the same property as for Row V, but this time with  $\varphi_2 = -d - \log x$  [Kulis et al., 2009]. We check that whenever  $\det(\mathbf{x}) = 1$ , we have

$$\begin{aligned}
\varphi(\mathbf{x}) = -d - \log \det(\mathbf{x}) &= -d \\
&= -\det(\mathbf{x}) \text{tr}(\mathbf{I}) \\
&= \text{tr}(\det(\mathbf{x}) \mathbf{x}^{-1} \mathbf{x}) = \text{tr}(\nabla \varphi(\mathbf{x})^\top \mathbf{x}). \tag{84}
\end{aligned}$$

For  $g(\mathbf{x}) \doteq \det \mathbf{x}^{1/d}$ , we get:

$$\begin{aligned}
\check{\varphi}(\mathbf{x}) &= \det \mathbf{x}^{1/d} \cdot \left( -d - \log \det \left( \frac{1}{\det \mathbf{x}^{1/d}} \cdot \mathbf{x} \right) \right) \\
&= \det \mathbf{x}^{1/d} \cdot \left( -d - \log \frac{1}{\det \mathbf{x}} \cdot \det \mathbf{x} \right) = -d \cdot \det \mathbf{x}^{1/d}, \tag{85}
\end{aligned}$$

and furthermore

$$\begin{aligned}
\nabla \check{\varphi}(\mathbf{x}) &= -d \cdot \nabla(\det \mathbf{x}^{1/d})(\mathbf{x}) \\
&= -\det(\mathbf{x}^{1/d}) \cdot \mathbf{x}^{-1} \tag{86}
\end{aligned}$$

So,

$$\begin{aligned}
D_{\check{\varphi}}(\mathbf{x} \parallel \mathbf{Y}) &= -d \cdot \det \mathbf{x}^{1/d} + d \cdot \det \mathbf{Y}^{1/d} + \text{tr}(\det(\mathbf{Y}^{1/d}) \cdot \mathbf{Y}^{-1}(\mathbf{x} - \mathbf{Y})) \\
&= -d \cdot \det \mathbf{x}^{1/d} + d \cdot \det \mathbf{Y}^{1/d} + \det(\mathbf{Y}^{1/d}) \text{tr}(\mathbf{x} \mathbf{Y}^{-1}) - d \cdot \det \mathbf{Y}^{1/d} \\
&= \det(\mathbf{Y}^{1/d}) \text{tr}(\mathbf{x} \mathbf{Y}^{-1}) - d \cdot \det \mathbf{x}^{1/d}. \tag{87}
\end{aligned}$$

We check that it is equal to:

$$\begin{aligned}
g(\mathbf{x}) \cdot D_{\varphi} \left( \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{g(\mathbf{Y})} \cdot \mathbf{Y} \right) \\
= \det \mathbf{x}^{1/d} \cdot \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \left( \frac{l_i \det \mathbf{Y}^{1/d}}{t_j \det \mathbf{x}^{1/d}} - \log \frac{l_i \det \mathbf{Y}^{1/d}}{t_j \det \mathbf{x}^{1/d}} - d \right). \tag{88}
\end{aligned}$$

To check it, we use the fact that, since  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal,

$$\begin{aligned}
& \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \log \frac{l_i \det \mathbf{Y}^{1/d}}{t_j \det \mathbf{X}^{1/d}} \\
&= \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \log l_i - \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \log \det \mathbf{X}^{1/d} \\
&\quad + \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \log \det \mathbf{Y}^{1/d} - \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \log t_j \\
&= \underbrace{\sum_i \log l_i - d \cdot \log \det \mathbf{X}^{1/d}}_{=0} + \underbrace{d \cdot \log \det \mathbf{Y}^{1/d} - \sum_j \log t_j}_{=0} = 0 , \tag{89}
\end{aligned}$$

which yields

$$\begin{aligned}
g(\mathbf{X}) \cdot D_\varphi \left( \frac{1}{g(\mathbf{X})} \cdot \mathbf{X} \parallel \frac{1}{g(\mathbf{Y})} \cdot \mathbf{Y} \right) &= \det \mathbf{X}^{1/d} \cdot \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \left( \frac{l_i \det \mathbf{Y}^{1/d}}{t_j \det \mathbf{X}^{1/d}} - d \right) \\
&= \det \mathbf{Y}^{1/d} \cdot \sum_{i,j} (\mathbf{u}_i^\top \mathbf{v}_j)^2 \cdot \frac{l_i}{t_j} - d \cdot \det \mathbf{X}^{1/d} \\
&= \det \mathbf{Y}^{1/d} \cdot \text{tr}(\mathbf{X} \mathbf{Y}^{-1}) - d \cdot \det \mathbf{X}^{1/d} , \tag{90}
\end{aligned}$$

which is equal to Equation (87).



## IV Going deep: higher-order identities

We can generalize Theorem 1 to higher order identities. For this, consider  $k > 0$  an integer, and let  $g_1, g_2, \dots, g_k : \mathcal{X} \rightarrow \mathbb{R}_*$  be a sequence of differentiable functions. For any  $\ell, \ell' \in [k]_*$  such that  $\ell \leq \ell'$ , we let  $\tilde{g}_{\ell, \ell'}$  be defined recursively as:

$$\tilde{g}_{\ell, \ell'}(\mathbf{x}) \doteq \begin{cases} \tilde{g}_{\ell-1, \ell'}(\mathbf{x}) \cdot g_{\ell'-(\ell-1)}\left(\frac{1}{\tilde{g}_{\ell-1, \ell'}(\mathbf{x})} \cdot \mathbf{x}\right) & \text{if } 1 < \ell \leq \ell' , \\ g_{\ell'}(\mathbf{x}) & \text{if } \ell = 1 , \end{cases} \quad (91)$$

and, for any  $\ell \in [k]$ ,

$$\check{\varphi}^{(\ell)}(\mathbf{x}) \doteq \begin{cases} g_{\ell}(\mathbf{x}) \cdot \check{\varphi}^{(\ell-1)}\left(\frac{1}{g_{\ell}(\mathbf{x})} \cdot \mathbf{x}\right) & \text{if } 0 < \ell \leq k , \\ \varphi(\mathbf{x}) & \text{if } \ell = 0 . \end{cases} \quad (92)$$

Notice that even when all  $g$ . are affine, this does not guarantee that some  $\tilde{g}_{\ell, \ell'}$  for  $\ell \neq 1$  is going to be affine. However, if for example  $g_{\ell'}$  is affine and all “preceeding”  $g_{\ell}$  ( $\ell \leq \ell'$ ) are homogeneous of degree 1, then all  $\tilde{g}_{\ell, \ell'}$  ( $\forall \ell \leq \ell'$ ) are affine. The following result can be seen as extension of the functional composition rules known for generalized perspective transforms of functions [Maréchal, 2005b] to composition rules for perspective transforms of divergences (See Section VIII).

**Corollary L** *For any  $k \in \mathbb{N}_*$ , let  $\varphi : \mathcal{X} \rightarrow \mathbb{R}$  be convex differentiable, and  $g_{\ell} : \mathcal{X} \rightarrow \mathbb{R}_*$  ( $\ell \in [k]$ ) a sequence of  $k$  differentiable functions. Then the following relationship holds, for any  $\ell, \ell' \in [k]_*$  with  $\ell \leq \ell'$ :*

$$\tilde{g}_{\ell, \ell'}(\mathbf{x}) \cdot D_{\check{\varphi}^{(\ell'-\ell)}}\left(\frac{1}{\tilde{g}_{\ell, \ell'}(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{\tilde{g}_{\ell, \ell'}(\mathbf{y})} \cdot \mathbf{y}\right) = D_{\check{\varphi}^{(\ell')}}(\mathbf{x} \parallel \mathbf{y}) , \forall \mathbf{x}, \mathbf{y} \in \mathcal{X} , \quad (93)$$

with  $\tilde{g}_{\ell, \ell'}$  defined as in Equation (91) and  $\check{\varphi}^{(\ell')}$  defined as in Equation (92), if and only if at least one of the two following conditions hold:

- (i)  $\tilde{g}_{\ell, \ell'}$  is affine on  $\mathcal{X}$ ;
- (ii)  $\check{\varphi}^{(\ell'-\ell)}$  is positive homogeneous of degree 1 on  $\mathcal{X}_{\ell, \ell'} \doteq \{(1/\tilde{g}_{\ell, \ell'}(\mathbf{x})) \cdot \mathbf{x} : \mathbf{x} \in \mathcal{X}\}$ .

We check that whenever  $\varphi$  is convex and all  $g$ . are non-negative, then all  $\check{\varphi}^{(\ell)}$  are convex ( $\forall \ell \in [k]$ ). To prove this, we choose  $\ell' = \ell$  and rewrite Equation (3), which brings, since  $\check{\varphi}^{(\ell'-\ell)} = \check{\varphi}^{(0)} = \varphi$ ,

$$\tilde{g}_{\ell, \ell}(\mathbf{x}) \cdot D_{\varphi}\left(\frac{1}{\tilde{g}_{\ell, \ell}(\mathbf{x})} \cdot \mathbf{x} \parallel \frac{1}{\tilde{g}_{\ell, \ell}(\mathbf{y})} \cdot \mathbf{y}\right) = D_{\check{\varphi}^{(\ell)}}(\mathbf{x} \parallel \mathbf{y}) , \forall \mathbf{x}, \mathbf{y} \in \mathcal{X} . \quad (94)$$

Since  $\varphi$  is convex, a sufficient condition to prove our result is to show that  $\tilde{g}_{\ell, \ell}$  is non-negative — which will prove that the right hand side of (94) is non-negative, and therefore  $\check{\varphi}^{(\ell)}$  is convex —. This can easily be proven by induction from the expression of  $\tilde{g}_{\ell, \ell'}$  in (91) and the fact that all  $g$ . are non-negative.

One interesting candidate for simplification is when all  $g_\ell$  are the same affine function, say  $g_\ell(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + b, \forall \ell \in [k]$ . In this case, we have indeed:

$$\begin{aligned}\tilde{g}_{\ell,\ell'}(\mathbf{x}) &= \mathbf{a}^\top \mathbf{x} + b \cdot \tilde{g}_{\ell-1,\ell'}(\mathbf{x}) \\ &= b^\ell + \mathbf{a}^\top \mathbf{x} \cdot \sum_{j=1}^{\ell-1} b^j, \end{aligned} \quad (95)$$

$$\check{\varphi}^{(\ell')}(\mathbf{x}) = \left( b^{\ell'} + \mathbf{a}^\top \mathbf{x} \cdot \sum_{j=1}^{\ell'-1} b^j \right) \cdot \varphi \left( \frac{1}{b^{\ell'} + \mathbf{a}^\top \mathbf{x} \cdot \sum_{j=1}^{\ell'-1} b^j} \cdot \mathbf{x} \right). \quad (96)$$

**Proof** [Proof of Corollary L] To check eq. (4), we first remark ( $\ell'$  being fixed) that it holds for  $\ell = 1$  (this is eq. (4)), and then proceed by an induction from the induction base hypothesis that, for some  $\ell \leq \ell'$ ,

$$\check{\varphi}^{(\ell')}(\mathbf{x}) = \tilde{g}_{\ell,\ell'}(\mathbf{x}) \cdot \check{\varphi}^{(\ell'-\ell)} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right). \quad (97)$$

We now have

$$\begin{aligned}\check{\varphi}^{(\ell')}(\mathbf{x}) &= \frac{\tilde{g}_{\ell+1,\ell'}(\mathbf{x})}{g_{\ell'-\ell} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right)} \cdot \check{\varphi}^{(\ell'-\ell)} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right) \end{aligned} \quad (98)$$

$$= \frac{\tilde{g}_{\ell+1,\ell'}(\mathbf{x})}{g_{\ell'-\ell} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right)} \cdot g_{\ell'-\ell} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right) \cdot \check{\varphi}^{(\ell'-(\ell+1))} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x}) g_{\ell'-\ell} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right)} \cdot \mathbf{x} \right) \quad (99)$$

$$\begin{aligned} &= \tilde{g}_{\ell+1,\ell'}(\mathbf{x}) \cdot \check{\varphi}^{(\ell'-(\ell+1))} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x}) g_{\ell'-\ell} \left( \frac{1}{\tilde{g}_{\ell,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right)} \cdot \mathbf{x} \right) \\ &= \tilde{g}_{\ell+1,\ell'}(\mathbf{x}) \cdot \check{\varphi}^{(\ell'-(\ell+1))} \left( \frac{1}{\tilde{g}_{\ell+1,\ell'}(\mathbf{x})} \cdot \mathbf{x} \right). \end{aligned} \quad (100)$$

Eq. (98) comes from eq. (97) and the definition of  $\tilde{g}_\ell$  in (91), eq. (99) comes from the definition of  $\check{\varphi}^{(\ell'-\ell)}$  in (92), eq. (100) is a second use of the definition of  $\tilde{g}_\ell$  in (91). ■

Notice the eventual high non-linearities introduced by the composition in eqs (91,92), which justifies the "deep" characterization.

## V Additional application: perspective transform of exponential families

Let  $\varphi$  be the cumulant function of a regular  $\varphi$ -exponential family with pdf  $p_\varphi(\cdot|\boldsymbol{\theta})$ , where  $\boldsymbol{\theta} \in \mathcal{X}$  is its natural parameter. Let  $\Omega(\cdot)$  be a norm on  $\mathcal{X}$ . Let  $\boldsymbol{\theta}_\Omega$  be the image of  $\boldsymbol{\theta} \in \mathcal{X}$  by the application from  $\mathcal{X}$  onto the  $\Omega$ -ball of unit norm defined by  $\boldsymbol{x} \mapsto (1/\Omega(\boldsymbol{x})) \cdot \boldsymbol{x}$ . For any two  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathcal{X}$ , let

$$\text{KL}_\varphi(\boldsymbol{\theta}||\boldsymbol{\theta}') \doteq \int p_\varphi(\boldsymbol{x}|\boldsymbol{\theta}) \log \frac{p_\varphi(\boldsymbol{x}|\boldsymbol{\theta})}{p_\varphi(\boldsymbol{x}|\boldsymbol{\theta}')} d\boldsymbol{x} \quad (101)$$

be the KL divergence between the two densities  $p_\varphi(\cdot|\boldsymbol{\theta})$  and  $p_\varphi(\cdot|\boldsymbol{\theta}')$ . A function is called regular if  $\Phi \doteq \{\boldsymbol{\theta} : \varphi(\boldsymbol{\theta}) \ll \infty\}$  is open.

**Lemma M** (*perspective transform of exponential families*) *For any convex regular  $\varphi$  which is restricted positive 1-homogeneous on  $\mathcal{X}_\Omega$ , the KL-divergence between two members of the same  $\varphi$ -exponential family satisfies:*

$$\Omega(\boldsymbol{\theta}') \cdot \text{KL}_\varphi(\boldsymbol{\theta}_\Omega||\boldsymbol{\theta}'_\Omega) = D_{\check{\varphi}}(\boldsymbol{\theta}'||\boldsymbol{\theta}) = \text{KL}_{\check{\varphi}}(\boldsymbol{\theta}||\boldsymbol{\theta}') . \quad (102)$$

**Proof** We know that  $\text{KL}(\boldsymbol{\theta}||\boldsymbol{\theta}') = D_\varphi(\boldsymbol{\theta}'||\boldsymbol{\theta})$  [Boissonnat et al., 2010]. Hence,

$$\begin{aligned} D_{\check{\varphi}}(\boldsymbol{\theta}'||\boldsymbol{\theta}) &= \Omega(\boldsymbol{\theta}') \cdot D_\varphi\left(\frac{1}{\Omega(\boldsymbol{\theta}')} \cdot \boldsymbol{\theta}' || \frac{1}{\Omega(\boldsymbol{\theta})} \cdot \boldsymbol{\theta}\right) \\ &= \Omega(\boldsymbol{\theta}') \cdot D_\varphi(\boldsymbol{\theta}'_\Omega||\boldsymbol{\theta}_\Omega) \\ &= \Omega(\boldsymbol{\theta}') \cdot \text{KL}_\varphi(\boldsymbol{\theta}_\Omega||\boldsymbol{\theta}'_\Omega) , \end{aligned} \quad (103)$$

as claimed. To prove the rest of the Lemma, we remark that  $\check{\Phi}$  is open because  $\Phi$  is open and  $\check{\varphi}$  is convex because  $\Omega(\boldsymbol{\theta}) \geq 0$  and  $\varphi$  is convex, so  $\check{\varphi}$  is convex regular and defines the cumulant of a regular exponential family for which  $D_{\check{\varphi}}(\boldsymbol{\theta}'||\boldsymbol{\theta}) = \text{KL}_{\check{\varphi}}(\boldsymbol{\theta}||\boldsymbol{\theta}')$  [Banerjee et al., 2005]. ■

The interest in Lemma M is to provide an integral-free expression of the KL-divergence when natural parameters are scaled by non-trivial transformations (left inequality). Furthermore, the equality

$$\Omega(\boldsymbol{\theta}') \cdot \text{KL}_\varphi\left(\frac{1}{\Omega(\boldsymbol{\theta})} \cdot \boldsymbol{\theta} || \frac{1}{\Omega(\boldsymbol{\theta}')} \cdot \boldsymbol{\theta}'\right) = \text{KL}_{\check{\varphi}}(\boldsymbol{\theta}||\boldsymbol{\theta}') \quad (104)$$

states a valid generalized perspective transform equality because  $\Omega$  is proper convex [Maréchal, 2005a,b]. Notice however that the rescaling of the KL divergence on the left uses its *right* parameter, unlike in Theorem 1. To summarize the content of Lemma M, we first "polarize" (left / right) the perspective transform of a Bregman divergence as a reference of which parameter is used to rescale the divergence. We also define as the *perspective transform of an exponential family* as the new distribution whose cumulant is the perspective transform of the cumulant (we do not change the natural parameter). We can then summarize eq. (104) by:

*"the right perspective transform of the KL divergence between two distributions of the same exponential family is the KL divergence between the perspective transform of the distributions"*

Finally, it is out of the scope of this paper, but Lemma M can also be extended to generalized exponential families [Fongillo and Reid, 2013].

## VI Additional application: computational information geometry

Two important objects of central importance in (computational) geometry are balls and Voronoi diagrams induced by a distortion, with which we can characterize the topological and computational aspects of major structures (Voronoi diagrams, triangulations, nearest neighbor topologies, etc.) [Boissonnat et al., 2010].

### VI.1 Bregman balls

Since a Bregman divergence is not necessarily symmetric, there are two types of (dual) balls that can be defined, the first or second types, where the variable  $\mathbf{x}$  is respectively placed in the left or right position. The first type Bregman balls are convex while the second type are not necessarily convex. A (closed) Bregman ball of the second type (with center  $\mathbf{c}$  and "radius"  $r$ ) is defined as:

$$B'(\mathbf{c}, r | \mathcal{X}, \varphi) \doteq \{ \mathbf{x} \in \mathcal{X} : D_\varphi(\mathbf{c} \| \mathbf{x}) \leq r \} . \quad (105)$$

It turns out that any divergence  $D_{\check{\varphi}}$  induces a ball of the second type, which is not necessarily analytically a Bregman ball (when  $\check{\varphi}$  is not convex), *but* turns out to define the *same* ball as a Bregman ball over properly scaled arguments (notice that the scaling is the same for both arguments, the ball's center and radius).

**Theorem N** *Let  $(\varphi, g, \check{\varphi})$  satisfy the conditions of Theorem 1, with  $g$  non negative and  $g(\mathbf{c}) \neq 0$ . Then*

$$B'(\mathbf{c}, r | \check{\varphi}, \mathcal{X}) = B' \left( \frac{1}{g(\mathbf{c})} \cdot \mathbf{c}, \frac{r}{g(\mathbf{c})} \middle| \varphi, \mathcal{X}_g \right) . \quad (106)$$

**Proof** From Theorem 1, we have

$$D_{\check{\varphi}}(\mathbf{c} \| \mathbf{x}) \leq r \quad (107)$$

iff

$$D_\varphi \left( \frac{1}{g(\mathbf{c})} \cdot \mathbf{c} \middle\| \frac{1}{g(\mathbf{c})} \cdot \mathbf{x} \right) \leq \frac{1}{g(\mathbf{c})} \cdot r . \quad (108)$$

Hence,

$$\begin{aligned} B'(\mathbf{c}, r | \check{\varphi}, \mathcal{X}) &= \{ \mathbf{x} \in \mathcal{X} : D_{\check{\varphi}}(\mathbf{c} \| \mathbf{x}) \leq r \} \\ &= \left\{ \mathbf{x} \in \mathcal{X} : D_\varphi \left( \frac{1}{g(\mathbf{c})} \cdot \mathbf{c} \middle\| \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} \right) \leq \frac{r}{g(\mathbf{c})} \right\} \\ &= B' \left( \frac{1}{g(\mathbf{c})} \cdot \mathbf{c}, \frac{r}{g(\mathbf{c})} \middle| \varphi, \mathcal{X}_g \right) , \end{aligned} \quad (109)$$

as claimed. ■

In other words and to be a little bit more specific,

*”any  $\mathbf{x}$  belongs to the ball of the second type induced by  $D_{\tilde{\varphi}}$  over  $\mathcal{X}$  **iff**  $(1/g(\mathbf{x})) \cdot \mathbf{x} (\in \mathcal{X}_g)$  belongs to the Bregman ball of the second type induced by  $D_{\varphi}$  over  $\mathcal{X}_g$  (obtained by scaling both the center and radius by  $g(\mathbf{c})$ ).”*

This property is not true for balls of the first type. What Theorem N says is that the topology induced by  $D_{\tilde{\varphi}}$  over  $\mathcal{X}$  is just *no different* from that induced by  $D_{\varphi}$  over  $\mathcal{X}_g$ .

## VI.2 Bregman Voronoi diagrams

Let us now investigate Bregman Voronoi diagrams. In the same way as there exists two types of Bregman balls, we can define two types of Bregman Voronoi diagrams that depend on the equation of the *Bregman bisector* [Boissonnat et al., 2010]. Of particular interest is the Bregman bisector of the *first type*:

$$BB_{\varphi}(\mathbf{x}, \mathbf{y}|\mathcal{X}) = \{ \mathbf{z} \in \mathcal{X} : D_{\varphi}(\mathbf{z}|\mathbf{x}) = D_{\varphi}(\mathbf{z}|\mathbf{y}) \} . \quad (110)$$

Let us define  $\mathbf{x}, \mathbf{y}$  as the Bregman bisector parameters. It turns out that any divergence  $D_{\tilde{\varphi}}$  induces a bisector of the first type which is not necessarily analytically a Bregman bisector (when  $\tilde{\varphi}$  is not convex), *but* turns out to define the *same* bisector as a Bregman bisector over transformed coordinates.

**Theorem O** *Let  $(\varphi, g, \tilde{\varphi})$  satisfy the conditions of Theorem 1. Then*

$$BB_{\tilde{\varphi}}(\mathbf{x}, \mathbf{y}|\mathcal{X}) = BB_{\varphi}(\mathbf{x}, \mathbf{y}|\mathcal{X}_g) . \quad (111)$$

(proof similar to Theorem N) Again, we get more precisely

*”any  $\mathbf{x}$  belongs to a Bregman bisector of the first type induced by  $D_{\tilde{\varphi}}$  over  $\mathcal{X}$  **iff**  $(1/g(\mathbf{x})) \cdot \mathbf{x} (\in \mathcal{X}_g)$  belongs to the corresponding Bregman bisector of the first type induced by  $D_{\varphi}$  over  $\mathcal{X}_g$  (obtained by scaling both bisector parameters by  $g(\cdot)$ ).”*

This property is not true for Bregman bisectors of the second type (obtained by permuting  $\mathbf{z}$  with the Bregman bisector parameters in eq . (110)).

## VI.3 Consequences

Theorems N, O have several important algorithmic consequences, some of which are listed now:

- the Voronoi diagram (resp. Delaunay triangulation) of the first type associated to  $\tilde{\varphi}$  can be constructed via the Voronoi diagram (resp. Delaunay triangulation) of the first type associated to  $\varphi$  [Boissonnat et al., 2010];
- range search using ball trees on  $D_{\tilde{\varphi}}$  can be efficiently implemented using Bregman divergence  $D_{\varphi}$  on  $\mathcal{X}_g$  [Cayton, 2009];
- the minimum enclosing ball problem, the one-class clustering problem (an important problem in machine learning), with balls of the second type on  $D_{\tilde{\varphi}}$  can be solved via the minimum Bregman enclosing ball problem on  $D_{\varphi}$  [Nock and Nielsen, 2005].

## VII Review: binary density ratio estimation

For completeness, we quickly review the central result of Menon and Ong [2016, Proposition 3]. Let  $(P, Q, \pi)$  be densities giving  $\mathbb{P}(X|Y = 1)$ ,  $\mathbb{P}(X = \mathbf{x}|Y = -1)$ ,  $\mathbb{P}(Y = 1)$  respectively, and  $M$  giving  $\mathbb{P}(X = \mathbf{x})$  accordingly. Let  $r(\mathbf{x}) \doteq \mathbb{P}(X = \mathbf{x}|Y = 1)/\mathbb{P}(X = \mathbf{x}|Y = -1)$  be the density ratio of the class-conditional densities, and  $\eta(\mathbf{x}) \doteq \mathbb{P}[Y = 1|X = \mathbf{x}]$  be the class-probability function. Then, we have the following, which extends [Menon and Ong, 2016, Proposition 6] for the case  $\pi \neq \frac{1}{2}$ .

**Lemma P** *Given a class-probability estimator  $\hat{\eta}: \mathcal{X} \rightarrow [0, 1]$ , let the density ratio estimator  $\hat{r}$  be*

$$\hat{r}(\mathbf{x}) = \frac{1 - \pi}{\pi} \cdot \frac{\hat{\eta}(\mathbf{x})}{1 - \hat{\eta}(\mathbf{x})} . \quad (112)$$

*Then for any convex differentiable  $\varphi: [0, 1] \rightarrow \mathbb{R}$ ,*

$$\mathbb{E}_{X \sim M} [D_\varphi(\eta(X) \parallel \hat{\eta}(X))] = \pi \cdot \mathbb{E}_{X \sim Q} [D_{\check{\varphi}}(r(X) \parallel \hat{r}(X))] . \quad (113)$$

*where  $\check{\varphi}$  is as per Equation 4 with  $g(z) \doteq \frac{1-\pi}{\pi} + z$  .*

**Proof** [Proof of Lemma P] Note that

$$\begin{aligned} \frac{1}{g(r(\mathbf{x}))} \cdot r(\mathbf{x}) &= \frac{\pi \mathbb{P}(X = \mathbf{x}|Y = -1)}{\mathbb{P}(X = \mathbf{x})} \cdot \frac{\mathbb{P}(X = \mathbf{x}|Y = 1)}{\mathbb{P}(X = \mathbf{x}|Y = -1)} \\ &= \frac{\pi \mathbb{P}(X = \mathbf{x}|Y = 1)}{\mathbb{P}(X = \mathbf{x})} \\ &= \eta(\mathbf{x}) , \end{aligned} \quad (114)$$

and furthermore

$$\begin{aligned} \mathbb{P}(X = \mathbf{x}) &= (1 - \pi) \mathbb{P}(X = \mathbf{x}|Y = -1) + \pi \mathbb{P}(X = \mathbf{x}|Y = 1) \\ &= \pi \cdot \left( \frac{1 - \pi}{\pi} + \frac{\mathbb{P}(X = \mathbf{x}|Y = 1)}{\mathbb{P}(X = \mathbf{x}|Y = -1)} \right) \cdot \mathbb{P}(X = \mathbf{x}|Y = -1) \\ &= \pi \cdot g(r(\mathbf{x})) \cdot \mathbb{P}(X = \mathbf{x}|Y = -1) . \end{aligned} \quad (115)$$

So,

$$\mathbb{E}_{X \sim M} [D_\varphi(\eta(X) \parallel \hat{\eta}(X))] = \pi \cdot \mathbb{E}_{X \sim Q} [g(r(X)) \cdot D_\varphi(\eta(X) \parallel \hat{\eta}(X))] \quad (116)$$

$$= \pi \cdot \mathbb{E}_{X \sim Q} \left[ g(r(X)) \cdot D_\varphi \left( \frac{1}{g(r(X))} \cdot r(X) \parallel \hat{\eta}(X) \right) \right] \quad (117)$$

$$= \pi \cdot \mathbb{E}_{X \sim Q} \left[ g(r(X)) \cdot D_\varphi \left( \frac{1}{g(r(X))} \cdot r(X) \parallel \frac{1}{g(\hat{r}(X))} \cdot \hat{r}(X) \right) \right] \quad (118)$$

$$= \pi \cdot \mathbb{E}_{X \sim Q} [D_{\check{\varphi}}(r(X) \parallel \hat{r}(X))] , \quad (119)$$

as claimed. Equation (116) comes from (115), Equation (117) comes from (114), Equation (118) comes from (112) and the definition of  $g$ . Equation (119) comes from Theorem 1, noting that  $g$  is linear. ■

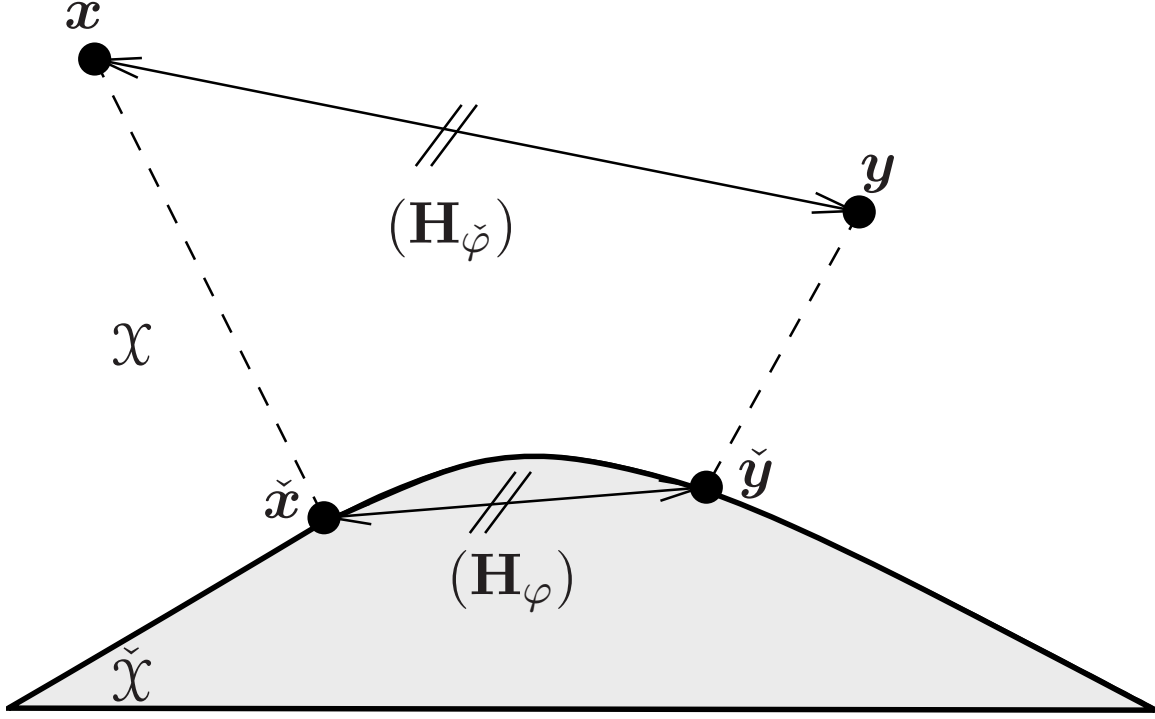


Figure 3: A depiction of the adaptive isometry that Theorem 1 provides. To simplify the picture as much as possible, we have used the shorthands  $\tilde{x} \doteq (1/g(x)) \cdot x$ ,  $\tilde{X} \doteq X_g$ . Double bars mean same metric length with respect to the (square root of) the Hessians in parenthesis.

## VIII Comments on Theorem 1

### VIII.1 Theorem 1 vs scaled isometries

Theorem 1 states in fact an isometry under some conditions, but an adaptive one in the sense that metrics involved rely on all parameters, and in particular on the points involved in the divergences (See Figure 3). Indeed, a simple Taylor expansion of the equation (2) (main file) shows that any such Bregman distortion with a twice differentiable generator can be expressed as:

$$D_\varphi(x\|y) = \frac{1}{2} \cdot (x - y)^\top \mathbf{H}_\varphi(x - y) , \quad (120)$$

for *some* value of the Hessian  $\mathbf{H}_\varphi$  depending on  $x, y$  (see for example [Kivinen et al., 2006, Appendix I], [Amari and Nagaoka, 2000]). Hence, under the constraint that both  $\varphi$  and  $\tilde{\varphi}$  are twice differentiable, eq. (3) becomes

$$g(x) \cdot \left( \frac{1}{g(x)} \cdot x - \frac{1}{g(y)} \cdot y \right)^\top \mathbf{H}_\varphi \left( \frac{1}{g(x)} \cdot x - \frac{1}{g(y)} \cdot y \right) = (x - y)^\top \mathbf{H}_{\tilde{\varphi}}(x - y) . \quad (121)$$

Notice that eq. (121) holds even when  $\mathbf{H}_{\check{\varphi}}$  is indefinite. Assuming  $g$  non-negative (which, by the way, enforces the convexity of  $\check{\varphi}$  and prevents  $\mathbf{H}_{\check{\varphi}}$  from being indefinite), we get by taking square roots,

$$\sqrt{g(\mathbf{x})} \cdot \left\| \frac{1}{g(\mathbf{x})} \cdot \mathbf{x} - \frac{1}{g(\mathbf{y})} \cdot \mathbf{y} \right\|_{\mathbf{H}_{\check{\varphi}}} = \|\mathbf{x} - \mathbf{y}\|_{\mathbf{H}_{\check{\varphi}}} , \quad (122)$$

which is a scaled isometry relationship between  $\mathcal{X}_g$  (left) and  $\mathcal{X}$  (right), but again the metrics involved depend on the arguments. Nevertheless, eq. (122) displays a sophisticated relationship between distances in  $\mathcal{X}_g$  and in  $\mathcal{X}$  which may prove useful in itself. With this in mind and keeping into account the restrictions on  $g$ , Theorem 1 states via eq. (122) that

*”distances on  $\mathcal{X}$  with metric  $\mathbf{H}_{\check{\varphi}}^{1/2}$  equal scaled distances after mapping  $\mathbf{x} \mapsto (1/g(\mathbf{x})) \cdot \mathbf{x} (\in \mathcal{X}_g)$  with metric  $\mathbf{H}_{\varphi}^{1/2}$ ”.*

## VIII.2 Theorem 1 vs generalized perspective transforms

Perspective transforms, also defined as epi-multiplication [Bauschke et al., 2008], are well known objects in convex analysis, and used in machine learning in particular to design and analyse loss functions [Reid and Williamson, 2011]. [Maréchal, 2005a,b] has defined a generalized notion of perspective transforms which coincidentally happens to define  $\check{\varphi}$  when assumptions are made about  $g$ . More precisely, Maréchal’s generalized perspective transform of functions  $\varphi$  and  $g$  is defined as:

$$(\varphi \triangle g)(\mathbf{x}, \mathbf{y}) \doteq \begin{cases} g(\mathbf{y}) \cdot \varphi\left(\frac{1}{g(\mathbf{y})} \cdot \mathbf{x}\right) & \text{if } g(\mathbf{y}) \in (0, +\infty) \\ \varphi 0^+(\mathbf{x}) & \text{if } g(\mathbf{y}) = 0 \\ \infty & \text{if } g(\mathbf{y}) = +\infty \end{cases} \quad (123)$$

where  $\varphi 0^+$  is the recession function of  $\varphi$ . For the definition to be valid, both  $\varphi$  and  $g$  have to be proper convex and  $g$  has to be positive. In this case, one remarks that

$$\check{\varphi}(\mathbf{x}) = (\varphi \triangle g)(\mathbf{x}, \mathbf{x}) , \quad (124)$$

but of course this holds only when significant restrictions are put on  $g$ . This does not prevent very interesting cases for the application of Theorem 1, as witnessed by the application to exponential families given in Section V, as well as several examples in Table A1 (rows I, II, III, V, VII). In this case, Theorem 1 gives an indication of how to define the generalized perspective transform of a Bregman divergence, which would be just the left hand side of eq. (4) in Theorem 1. To our knowledge, the use perspective transforms in machine learning has been limited to the definition of Csiszar’s duals for loss function and divergences [Reid and Williamson, 2011], and they have not been used to define the perspective of a divergence. Our results indicate that such objects would not be just mathematical curiosities, but could eventually be the ground for new methods to deal with popular problems. This is out of the scope of this paper, but if we resort to perspectives, then Theorem 1 can be roughly summarized by the property that

*”the perspective transform of the divergence equals the divergence of the perspective transform”.*



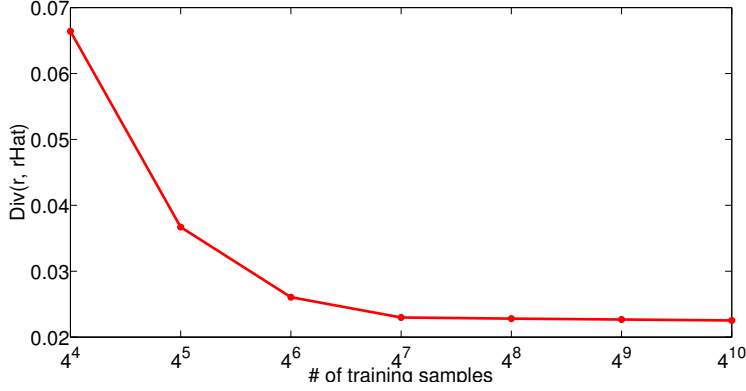


Figure 4: Density ratio estimate divergence  $\mathbb{E}_{\mathbf{X} \sim P_C} [D_{\tilde{\varphi}}(r(\mathbf{X}), \hat{r}(\mathbf{X}))]$  as a function of # of training samples.

## IX Additional experiments: Multiclass density ratio experiments

We consider a synthetic multiclass density ratio estimation problem. We fix  $\mathcal{X} = \mathbb{R}^2$ , and consider  $C = 3$  classes. We consider a distribution where the class-conditionals  $\Pr(\mathbf{X}|\mathbf{Y} = c)$  are multivariate Gaussians with means  $\boldsymbol{\mu}_c$  and covariance  $\sigma_c^2 \cdot \text{Id}$ . As the class-conditionals have a closed form, we can explicitly compute  $\boldsymbol{\eta}$ , as well the density ratio  $\mathbf{r}$  to the reference class  $c^* = C$ .

For fixed class prior  $\boldsymbol{\pi} = \Pr(\mathbf{Y} = c)$ , we draw  $N_{\text{Tr}}$  samples from  $\Pr(\mathbf{X}, \mathbf{Y})$ . From this, we estimate the class-probability  $\hat{\boldsymbol{\eta}}$  using multiclass logistic regression. This can be seen as minimising  $\mathbb{E}_{\mathbf{X} \sim M} [D_{\varphi}(\boldsymbol{\eta}(\mathbf{X}) \| \hat{\boldsymbol{\eta}}(\mathbf{X}))]$  where  $\varphi(z) = \sum_i z_i \log z_i$  is the generator for the KL-divergence.

We then use Equation 6 (main file) to estimate the density ratios  $\hat{\mathbf{r}}$  from  $\hat{\boldsymbol{\eta}}$ . On a fresh sample of  $N_{\text{te}}$  instances from  $\Pr(\mathbf{X}, \mathbf{Y})$ , we estimate the right hand side of Lemma 2, viz.  $\mathbb{E}_{\mathbf{X} \sim P_C} [D_{\tilde{\varphi}}(r(\mathbf{X}) \| \hat{r}(\mathbf{X}))]$ , where  $\tilde{\varphi}$  uses the  $g$  as specified in Lemma 2. From the result of Lemma 2, we expect this divergence to be small when  $\hat{\boldsymbol{\eta}}$  is a good estimator of  $\boldsymbol{\eta}$ .

We perform the above for sample sizes  $N \in \{4^4, 4^5, \dots, 4^{10}\}$ , with  $N_{\text{Tr}} = 0.8N$  and  $N_{\text{te}} = 0.2N$ . For each sample size, we perform  $T = 25$  trials, where in each trial we randomly draw  $\boldsymbol{\pi}$  uniformly over  $(1/C)\mathbf{1} + (1 - 1/C) \cdot [0, 1]^C$ ,  $\boldsymbol{\mu}_c$  from  $0.1 \cdot \mathcal{N}(\mathbf{0}, 1)$ , and  $\sigma_c$  uniformly from  $[0.5, 1]$ . Figure 4 summarises the mean divergence across the  $T$  trials for each sample size. We see that, as expected, with more training samples the divergence decreases in a monotone fashion.

## **X Additional experiments: Adaptive filtering experiments**

Tables A2 – A7 present *in extenso* the experiments of  $p$ -LMS vs DN- $p$ -LMS, as a function of  $(p, q)$ , whether target  $\mathbf{u}$  is sparse or not, and the misestimation factor  $\rho$  for  $X_p$ . We refer to Kivinen et al. [2006] for the formal definitions used for sparse / dense targets as well as for the experimental setting, which we have reproduced with the sole difference that the signal changes periodically each 1 000 iterations.

## References

- S.-I. Amari and H. Nagaoka. *Methods of Information Geometry*. Oxford University Press, 2000.
- D. Arthur and S. Vassilvitskii.  $k$ -means++ : the advantages of careful seeding. In *19<sup>th</sup> SODA*, 2007.
- A. Banerjee, S. Merugu, I. Dhillon, and J. Ghosh. Clustering with Bregman divergences. *JMLR*, 6: 1705–1749, 2005.
- H.-H. Bauschke, R. Goebel, Y. Lucet, and X. Wang. The proximal average: Basic theory. *SIAM J. Opt.*, 19:766–785, 2008.
- J.-D. Boissonnat, F. Nielsen, and R. Nock. Bregman Voronoi diagrams. *DCG*, 44(2):281–307, 2010.
- S.-R. Buss and J.-P. Fillmore. Spherical averages and applications to spherical splines and interpolation. *ACM Transactions on Graphics*, 20:95–126, 2001.
- L. Cayton. Efficient bregman range search. In *NIPS\*22*, pages 243–251, 2009.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning and Games*. Cambridge University Press, 2006.
- I. Dhillon and D.-S. Modha. Concept decompositions for large sparse text data using clustering. *MLJ*, 42:143–175, 2001.
- Y. Endo and S. Miyamoto. Spherical  $k$ -means++ clustering. In *Proc. of the 12<sup>th</sup> MDAI*, pages 103–114, 2015.
- R.-M. Fongillo and M. Reid. Convex foundations for generalized maxent models. In *MaxEnt’13*, 2013.
- G.-A. Galperin. A concept of the mass center of a system of material points in the constant curvature spaces. *Communications in Mathematical Physics*, 154:63–84, 1993.
- J. Kivinen, M. Warmuth, and B. Hassibi. The  $p$ -norm generalization of the LMS algorithm for adaptive filtering. *IEEE Trans. SP*, 54:1782–1793, 2006.
- B. Kulis, M.-A. Sustik, and I.-S. Dhillon. Low-rank kernel learning with Bregman matrix divergences. *JMLR*, 10:341–376, 2009.
- P. Maréchal. On a functional operation generating convex functions, part 1: duality. *J. of Optimization Theory and Applications*, 126:175–189, 2005a.
- P. Maréchal. On a functional operation generating convex functions, part 2: algebraic properties. *J. of Optimization Theory and Applications*, 126:375–366, 2005b.
- A.-K. Menon and C.-S. Ong. Linking losses for class-probability and density ratio estimation. In *ICML*, 2016.

- R. Nock and F. Nielsen. Fitting the Smallest Enclosing Bregman Ball. In *Proc. of the 16<sup>th</sup> European Conference on Machine Learning*, pages 649–656. Springer-Verlag, 2005.
- M.-D. Reid and R.-C. Williamson. Information, divergence and risk for binary experiments. *JMLR*, 12:731–877, 2011.
- J. Reisinger, A. Waters, B. Silverthorn, and R.-J. Mooney. Spherical topic models. In *27<sup>th</sup> ICML*, pages 903–910, 2010.
- A.-A. Ungar. *Mathematics Without Boundaries: Surveys in Interdisciplinary Research*, chapter An Introduction to Hyperbolic Barycentric Coordinates and their Applications, pages 577–648. Springer New York, 2014.

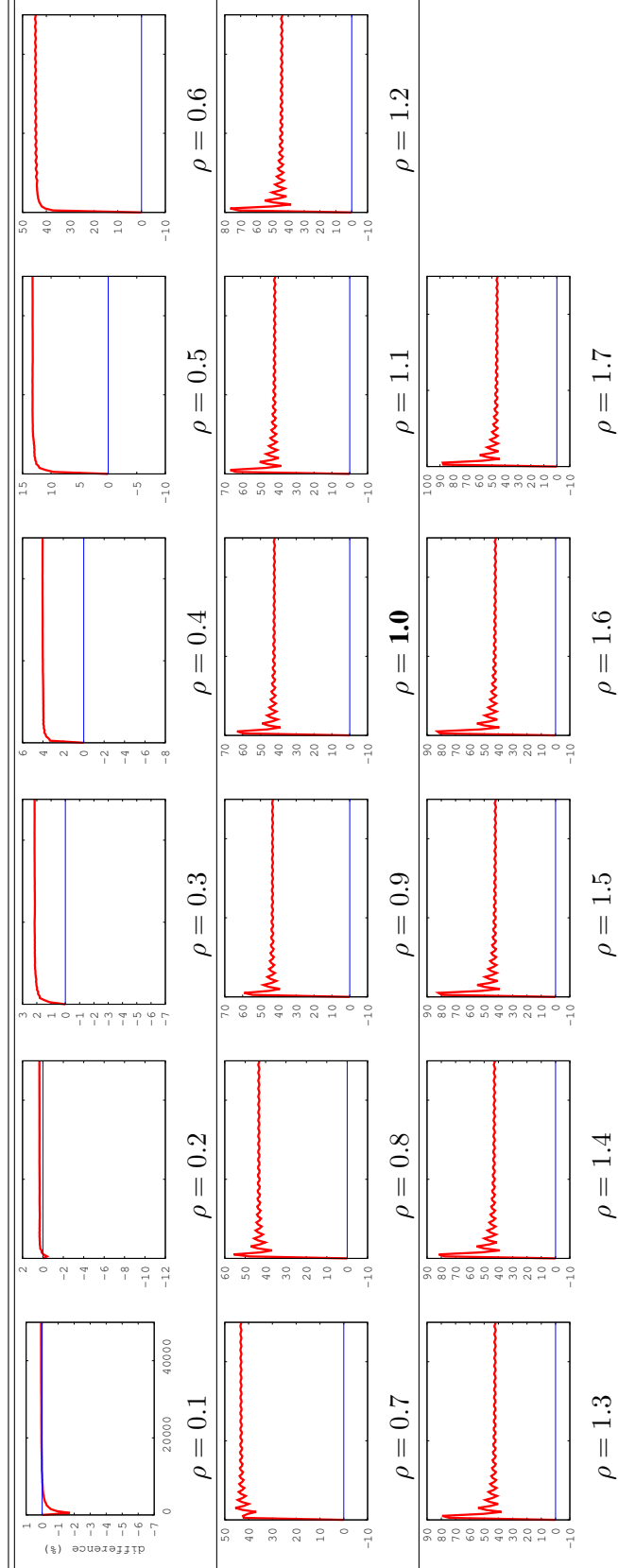


Table A2: Error( $p$ -LMS) - Error(DN- $p$ -LMS) as a function of  $t$  ( $\in \{1, 2, \dots, 50000\}$ ),  $\mathbf{u} = \text{dense}$ ,  $(p, q) = (1.17, 6.9)$ .

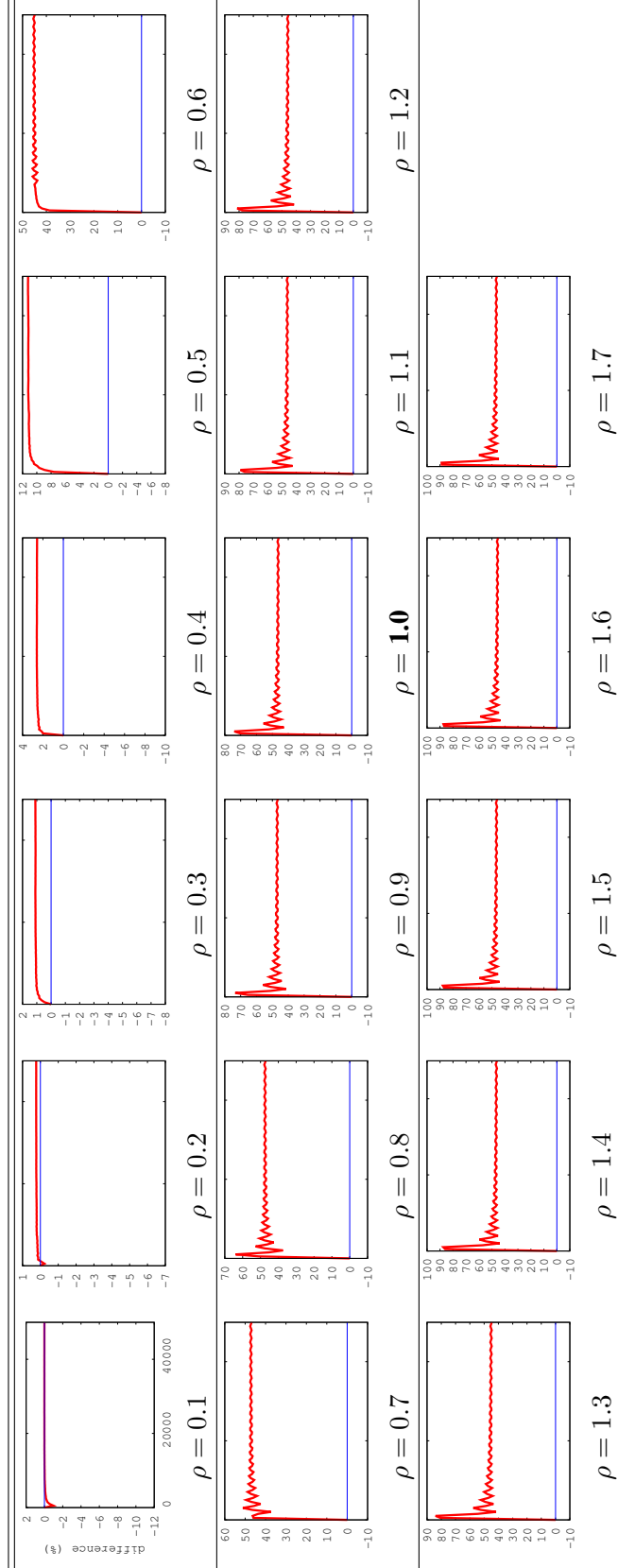


Table A3: Error( $p$ -LMS) - Error(DN- $p$ -LMS) as a function of  $t$  ( $\in \{1, 2, \dots, 50000\}$ ),  $\mathbf{u} = \text{sparse}$ ,  $(p, q) = (1.17, 6.9)$ .

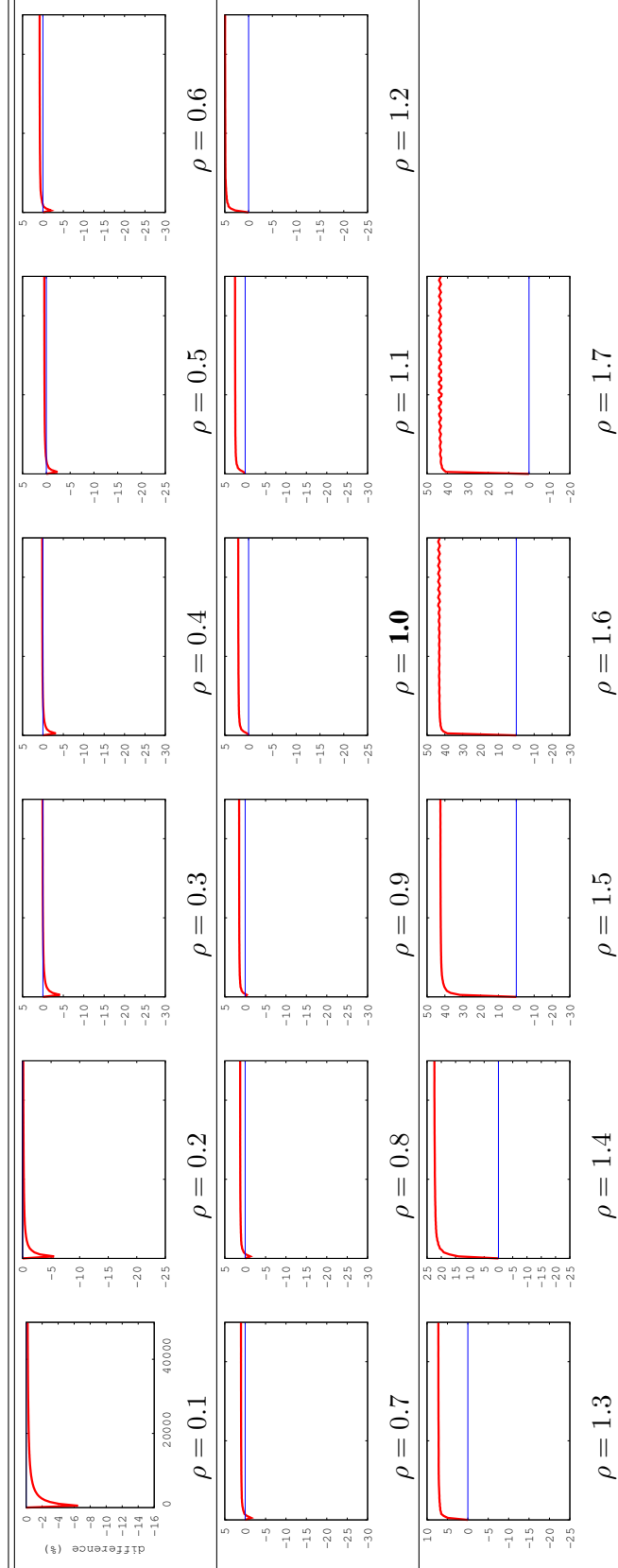


Table A4: Error( $p$ -LMS) - Error(DN- $p$ -LMS) as a function of  $t$  ( $\in \{1, 2, \dots, 50000\}$ ),  $\mathbf{u} = \text{dense}$ ,  $(p, q) = (2.0, 2.0)$ .

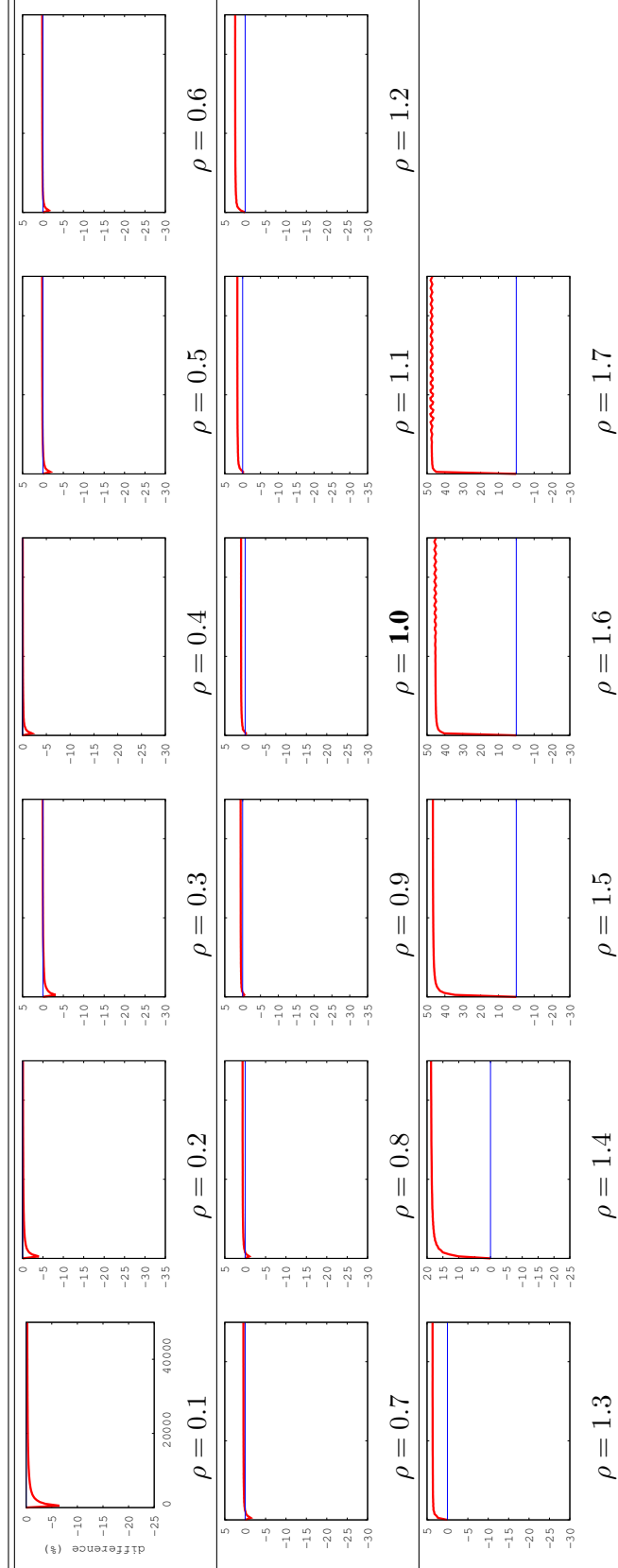


Table A5: Error( $p$ -LMS) - Error(DN- $p$ -LMS) as a function of  $t$  ( $\in \{1, 2, \dots, 50000\}$ ),  $\mathbf{u} = \text{sparse}$ ,  $(p, q) = (2.0, 2.0)$ .



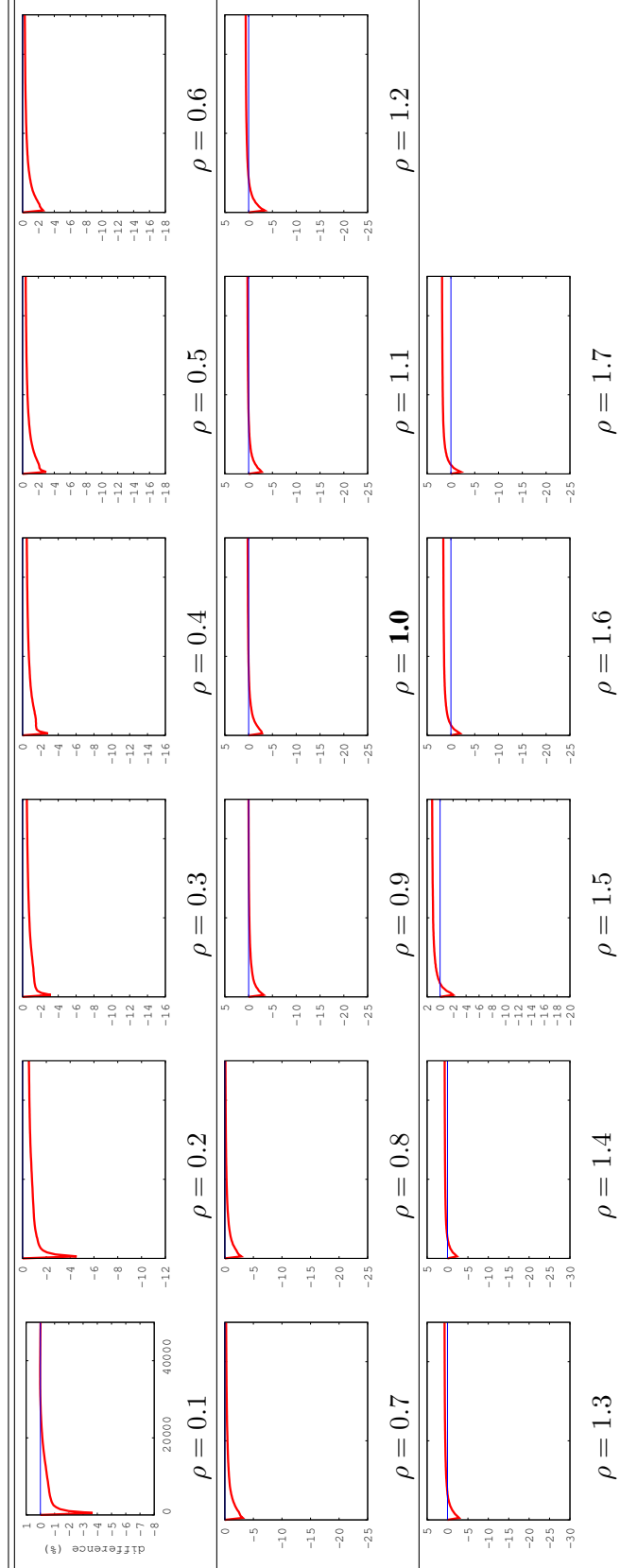


Table A6: Error( $p$ -LMS) - Error(DN- $p$ -LMS) as a function of  $t \in \{1, 2, \dots, 50000\}$ ,  $\mathbf{u} = \text{dense}$ ,  $(p, q) = (6.9, 1.17)$ .

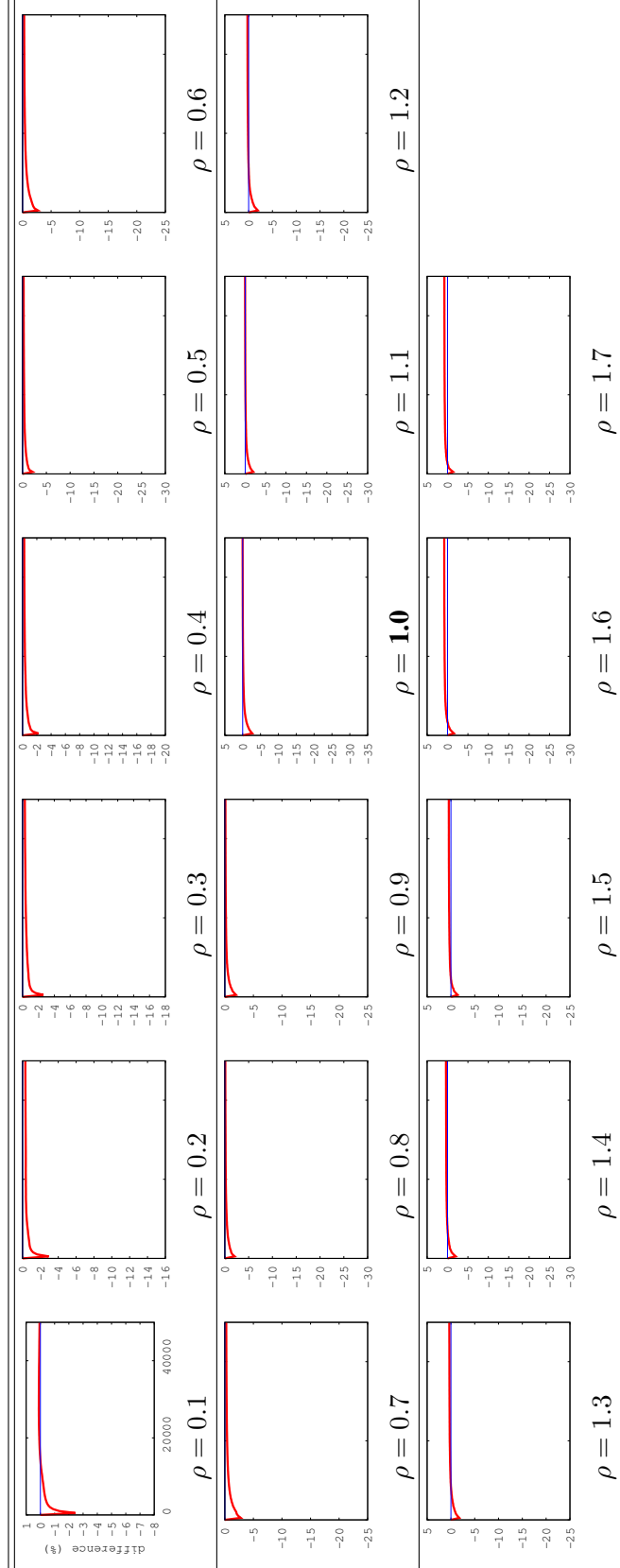


Table A7: Error( $p$ -LMS) - Error(DN- $p$ -LMS) as a function of  $t$  ( $\in \{1, 2, \dots, 50000\}$ ),  $\mathbf{u} = \text{sparse}$ ,  $(p, q) = (6.9, 1.17)$ .