# Beta-divergence as a subclass of Bregman divergence

Romain Hennequin, Bertrand David, Roland Badeau

HAL Id: hal-00945202

https://hal.inria.fr/hal-00945202

Submitted on 25 Mar 2014

# Beta-divergence as a subclass of Bregman divergence

Romain Hennequin, *Student Member, IEEE*, Bertrand David, *Member, IEEE*,
and Roland Badeau, *Senior Member, IEEE*

*Abstract*—In this paper, the $\beta$-divergence is shown to be a particular case of Bregman divergence. This result makes it possible to straightforwardly apply theorems about Bregman divergences to $\beta$-divergences. This is of interest for numerous applications since these divergences are widely used, for instance in non-negative matrix factorization (NMF).

*Index Terms*—Beta-divergence, Bregman divergence, non-negative matrix factorization.

## I. Introduction

**D**IVERGENCES are distance-like functions, widely used to assess the similarity between two objects. For instance, Kullback-Liebler (KL) divergence [13] is used in information theory to compare two probability distributions, and the Itakura-Saito (IS) divergence is used as a measure of the perceptual difference between spectra [11]. Generalized classes of divergences, for instance Bregman divergences, are used in pattern classification and clustering [1]. In non-negative matrix factorization (NMF [14]), divergences are used as cost functions: NMF approximates an $F \times T$ non-negative matrix $\mathbf{V}$ with the product of two non-negative low-rank matrices:

$$\mathbf{V} \approx \mathbf{W}\mathbf{H},$$

where the size of $\mathbf{W}$ is $F \times R$ and the size of $\mathbf{H}$ is $R \times T$ (with $R < F$ and $R < T$).

This approximation is generally quantized with a cost function to be minimized with respect to $\mathbf{W}$ and $\mathbf{H}$. This cost function is often an element-wise divergence between $\mathbf{V}$ and $\mathbf{W}\mathbf{H}$ [14].

Numerous divergences are used as cost functions in NMF. Most common divergences probably are the Euclidean (EUC) distance, the KL divergence (see [14]) and the IS divergence (see [8]).

Several authors proposed generalized divergences which encompass these classical divergences:

- Csiszar's divergence [4], which is a generalization of Amari's $\alpha$-divergence [5]. Both these divergences encompass the KL divergence and its dual.
- Bregman divergence [3], [6], which encompasses the EUC distance, the KL divergence and the IS divergence.

- $\beta$-divergence, introduced in [7] and studied as a cost function for NMF in [12] which also encompasses the EUC distance, the KL divergence and the IS divergence.

NMF is widely used in numerous areas such as image processing [14], [10], text mining [17], email surveillance [2], spectroscopy [9] and audio processing [8], [19], [18].

In this paper, we give a proof that the $\beta$-divergence is actually a subclass of Bregman divergence. This result is supposed to be known in a certain community [15], [16], but we propose to give a full demonstration of it in the wide framework of element-wise divergences and we present application to illustrate its interest. This result indeed permits to immediately particularize properties derived for Bregman divergence to $\beta$-divergence.

## II. Divergence

In this section, we define the concept of divergence, element-wise divergence, and the particular case of Bregman divergence and $\beta$-divergence.

### A. Definition

Divergences are distance-like functions which measure the separation between two elements.

*Definition 2.1:* Let $\mathcal{S}$ be a set. A *divergence* on $S$ is a function $D : \mathcal{S} \times \mathcal{S} \to \mathbb{R}$ satisfying:

$$\forall (p,q) \in \mathcal{S} \times \mathcal{S} \quad D(p\|q) \geq 0, \text{ and } D(p\|q) = 0 \text{ iff } p = q.$$

As a distance, a divergence should be non-negative and separable. However, a divergence does not necessarily satisfy the triangle inequality and the symmetry axiom of a distance. In order to avoid the confusion with distances, the notation $D(p\|q)$ is often used instead of the classical distance notation $D(p,q)$.

### B. Bregman divergence

*Definition 2.2:* Let $\mathcal{S}$ be a convex subset of a Hilbert space and $\Phi : \mathcal{S} \to \mathbb{R}$ a continuously differentiable strictly convex function. The *Bregman divergence* [3] $D_\Phi : \mathcal{S} \times \mathcal{S} \to \mathbb{R}_+$ (where $\mathbb{R}_+$ is the set of non-negative real numbers) is defined as:

$$D_\Phi(\mathbf{x}\|\mathbf{y}) = \Phi(\mathbf{x}) - \Phi(\mathbf{y}) - \langle \mathbf{x} - \mathbf{y}, \nabla\Phi(\mathbf{y}) \rangle,$$

where $\nabla\Phi(\mathbf{y})$ stands for the gradient of $\Phi$ evaluated at $\mathbf{y}$ and $\langle .,. \rangle$ is the standard Hermitian dot product.

The value of Bregman divergence $D_\Phi(\mathbf{x}||\mathbf{y})$ can be viewed as the difference between the function $\Phi(\mathbf{x})$ and its first order Taylor series at $\mathbf{y}$. Thus, adding an affine form to $\Phi$ does not change $D_\Phi$.

## III. ELEMENT-WISE DIVERGENCE

### A. Definition

In this section, $\mathcal{S} = \mathbb{R}_+^N$ or $\mathcal{S} = (\mathbb{R}_+\backslash\{0\})^N$. On such sets, one can define *element-wise divergences*: a divergence on $\mathbb{R}_+^N$ (resp. $(\mathbb{R}_+\backslash\{0\})^N$) is called element-wise if there exists a divergence $d$ on $\mathbb{R}_+$ (resp. $\mathbb{R}_+\backslash\{0\}$) such that:

$$\forall\mathbf{x} = (x_1, ..., x_n), \forall\mathbf{y} = (y_1, ..., y_n) \quad D(\mathbf{x}||\mathbf{y}) = \sum_{n=1}^{N} d(x_n|y_n).$$

### B. Element-wise Bregman divergence

Element-wise Bregman divergences are a subclass of Bregman divergences for which $\Phi$ is the sum of $N$ scalar, continuously differentiable and strictly convex element-wise functions:

$$\forall\mathbf{x} = (x_1, x_2, \ldots, x_N) \in \mathcal{S} \quad \Phi(\mathbf{x}) = \sum_{n=1}^{N} \phi(x_i).$$

Then $D_\Phi(\mathbf{x}||\mathbf{y}) = \sum_{i=1}^{N} d_\phi(x_i|y_i)$ where $d_\phi(x|y) = \phi(x) - \phi(y) - \phi'(y)(x-y)$ and thus, the divergence is element-wise. For element-wise Bregman divergences, we can equivalently denote the divergence $D_\Phi$ or $D_\phi$.

### C. β-divergence

*Definition 3.1:* Let $\beta \in \mathbb{R}$. The $\beta$-divergence on $\mathbb{R}_+\backslash\{0\}$ is defined by:

$$d_\beta(x|y) = \begin{cases} \frac{x}{y} - \log(\frac{x}{y}) - 1 & \beta = 0 \\ x(\log x - \log y) + (y - x) & \beta = 1 \\ \frac{x^\beta + (\beta-1)y^\beta - \beta x y^{\beta-1}}{\beta(\beta-1)} & \beta \in \mathbb{R}\backslash\{0,1\}. \end{cases}$$

One should notice that the previous definition of $\beta$-divergence is continuous with respect to $\beta$ in the sense that:

$$\forall\beta_0 \in \mathbb{R}, \quad \forall x, y \in \mathbb{R}_+\backslash\{0\} \quad d_{\beta_0}(x, y) = \lim_{\beta\to\beta_0} d_\beta(x, y),$$

particularly for $\beta_0 = 0$ and $\beta_0 = 1$.

From this divergence on $\mathbb{R}_+\backslash\{0\}$, one can define an element-wise $\beta$-divergence on $(\mathbb{R}_+\backslash\{0\})^N$:

$$D_\beta(\mathbf{x}||\mathbf{y}) = \sum_{n=1}^{N} d_\beta(x|y).$$

## IV. β-DIVERGENCE AS A BREGMAN DIVERGENCE

In this section, we show that the Bregman divergence encompasses the $\beta$-divergence in a natural way.

For $\beta \in \mathbb{R}$, let $\phi_\beta : \mathbb{R}_+\backslash\{0\} \to \mathbb{R}$ be the function defined as:

$$\forall x \in \mathbb{R}_+\backslash\{0\} \quad \phi_\beta(x) = \begin{cases} -\log x + x - 1 & \beta = 0 \\ x\log x - x + 1 & \beta = 1 \\ \frac{x^\beta}{\beta(\beta-1)} - \frac{x}{\beta-1} + \frac{1}{\beta} & \text{otherwise.} \end{cases}$$

As shown in Appendix A, this definition is continuous with respect to $\beta$ in the sense that:

$$\forall\beta_0 \in \mathbb{R}, \forall x \in \mathbb{R}_+\backslash\{0\} \quad \lim_{\beta\to\beta_0} \phi_\beta(x) = \phi_{\beta_0}(x).$$

For all $\beta \in \mathbb{R}$, $\phi_\beta$ is smooth on $\mathbb{R}_+\backslash\{0\}$ and its second derivative is:

$$\phi''_\beta(x) = x^{\beta-2}. \tag{1}$$

Thus $\phi_\beta$ is strictly convex and one can define the Bregman divergence $D_{\phi_\beta}$ associated to $\phi_\beta$:

$$D_{\phi_\beta}(\mathbf{x}||\mathbf{y}) = \sum_{n=1}^{N} \phi_\beta(x_n) - \phi_\beta(y_n) - \phi'_\beta(y_n)(x_n - y_n).$$

Straightforward calculations (see Appendix B) show that for all $\beta \in \mathbb{R}$, $D_{\phi_\beta} = D_\beta$ is a $\beta$-divergence. Thus the Bregman divergence encompasses $\beta$-divergence.

## V. APPLICATIONS

In this section, we present examples showing how our result can particularize properties of the Bregman divergence to the $\beta$-divergence, in order to illustrate its potential fields of application.

### A. Non-negative matrix factorization

The multiplicative update rule of $\mathbf{H}$ for minimizing a Bregman divergence $D_\phi$ cost function given in [6] is:

$$\mathbf{H} \leftarrow \mathbf{H}.\frac{\mathbf{W}^T(\phi''(\mathbf{WH}).\mathbf{V})}{\mathbf{W}^T(\phi''(\mathbf{WH}).(\mathbf{WH}))}.$$

The product ".", the fraction bar and $\phi''$ are element-wise operations on the corresponding matrices. We can directly derive the (already well-known [4]) update rule of $\mathbf{H}$ for a $\beta$-divergence $D_\beta$ cost function using (1):

$$\mathbf{H} \leftarrow \mathbf{H}.\frac{\mathbf{W}^T((\mathbf{WH})^{\cdot(\beta-2)}.\mathbf{V})}{\mathbf{W}^T((\mathbf{WH})^{\cdot(\beta-1)})}.$$

This illustrates the interest of deriving general properties about the Bregman divergence instead of the $\beta$-divergence.

## B. Right type centroid

The right type centroid is used in clustering as a "center" of a point cloud with respect to an asymmetric divergence: the right type centroid can thus be thought as an average typical point of a set.

*Definition 5.1:* Given a divergence $D$, the right type centroid of a finite set $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2..., \mathbf{x}_n\} \subset \mathcal{S}$ is defined as:

$$\mathbf{c}_{\text{right}}^D = \arg\min_{\mathbf{c}} \frac{1}{n} \sum_{i=1}^n D(\mathbf{x}_i \| \mathbf{c}).$$

*Theorem 5.1:* When $D = D_\beta$ is a $\beta$ divergence, $\mathbf{c}_{\text{right}}^{D_\beta}$ is unique, independent of $\beta$ and is equal to $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$.

*Proof:* It was shown in [1] that, when $D = D_\Phi$ is a Bregman divergence, $\mathbf{c}_{\text{right}}^{D_\Phi}$ is unique, independent of $\Phi$ and is equal to $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$. As $\beta$-divergence is a subclass of Bregman divergence, the proof is straightforward. □

## VI. CONCLUSION

In this paper, we presented a proof that the general class of Bregman divergence encompasses the $\beta$-divergence in a natural way. This results permits to straightforwardly apply theorems about the Bregman divergence to the $\beta$-divergence. As the latter is widely used in methods such as NMF, which has applications in numerous areas (signal processing, clustering, data mining, spectroscopy), the field of application of this result is quite wide.

## RÉFÉRENCES

[1] Arindam Banerjee, Srujana Merugu, Inderjit S. Dhillon, and Joydeep Ghosh. Clustering with Bregman divergences. *Journal of Machine Learning Research*, 6:1705–1749, October 2005.
[2] Michael W. Berry and Murray Browne. Email surveillance using nonnegative matrix factorization. *Computational and Mathematical Organization Theory*, 11(3):249–264, February 2005.
[3] L. M. Bregman. The relaxation method of finding the common points of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):210–217, 1967.
[4] Andrzej Cichocki, Rafal Zdunek, and Sun-Ichi Amari. Csiszar's divergences for non-negative matrix factorization: Family of new algorithms. In *Conference on Independent Component Analysis and Blind Source Separation (ICA)*, pages 32 – 39, Charleston, SC, USA, March 2006.
[5] Andrzej Cichocki, Rafal Zdunek, Seungjin Choi, Robert J. Plemmons, and Shun-Ichi Amari. Non-negative tensor factorization using alpha and beta divergences. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 3, pages 1393 – 1396, Honolulu, Hawaii, USA, April 2007.
[6] Inderjit S. Dhillon and Suvrit Sra. Generalized nonnegative matrix approximations with Bregman divergences. In Y. Weiss, B. Schölkopf, and J. Platt, editors, *Neural Information Processing Systems conference (NIPS)*, pages 283–290. MIT Press, Cambridge, MA, December 2006.
[7] Shinto Eguchi and Yutaka Kano. Robustifying maximum likelihood estimation. Technical report, Institute of Statistical Mathematics, Tokyo, June 2001.
[8] Cédric Févotte, Nancy Bertin, and Jean-Louis Durrieu. Nonnegative matrix factorization with the Itakura-Saito divergence. With application to music analysis. *Neural Computation*, 11(3):793–830, March 2009.
[9] Cyril Gobinet, Eric Perrin, and Régis Huez. Application of non-negative matrix factorization to fluorescence spectroscopy. In *European Signal Processing Conference (EUSIPCO)*, Vienna, Austria, September 2004.
[10] Patrik O. Hoyer. Non-negative matrix factorization with sparseness constraints. *Journal of Machine Learning Research*, 5:1457–1469, November 2004.
[11] Fumitada Itakura and Shuzo Saito. Analysis synthesis telephony based on the maximum likelihood method. In *6th International Congress on Acoustics*, pages C–17–C–20, Tokyo, Japan, 1968.
[12] Raul Kompass. A generalized divergence measure for nonnegative matrix factorization. *Neural Computation*, 19(3):780–791, March 2007.
[13] Solomon Kullback and Richard A. Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86, March 1951.
[14] Daniel D. Lee and H. Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, October 1999.
[15] Frank Nielsen and Richard Nock. The dual voronoi diagrams with respect to representational bregman divergences. In *International Symposium on Voronoi Diagrams*, pages 71 – 78, Copenhagen, Denmark, June 2009.
[16] Families of Alpha-Beta, Gamma-Divergences: Flexible, and Robust Measures of Similarities. Andrzej cichocki and shun-ichi amari. *Entropy*, 12(6):1532–1568, June 2010.
[17] V. Paul Pauca, Farial Shahnaz, Michael W. Berry, and Robert J. Plemmons. Text mining using non-negative matrix factorizations. In *SIAM international conference on data mining*, pages 452–456, Lake Buena Vista, Florida, USA, January 2004.
[18] Jouni Paulus and Tuomas Virtanen. Drum transcription with non-negative spectrogram factorization. In *European Signal Processing Conference (EUSIPCO)*, Antalya, Turkey, September 2005.
[19] Tuomas Virtanen. Monaural sound source separation by nonnegative matrix factorization with temporal continuity. *IEEE Transactions on Audio, Speech and Language Processing*, 15(3):1066–1074, March 2007.

## APPENDIX

### A. Continuity of $\phi_\beta$ with respect to $\beta$

With the little-o notation, one can write as $\beta \to 0$:

$$
\begin{aligned}
\phi_\beta(x) &= \frac{e^{\beta \log x}}{\beta(\beta - 1)} - \frac{x}{\beta - 1} + \frac{1}{\beta} \\
&= \frac{1 + \beta \log x + o(\beta)}{\beta(\beta - 1)} - \frac{x}{\beta - 1} + \frac{\beta - 1}{\beta(\beta - 1)} \\
&= \frac{1 + \log x - x}{\beta - 1} + o(1).
\end{aligned}
$$

Then:

$$\lim_{\beta \to 0} \phi_\beta(x) = -\log x + x - 1.$$

In a similar way, one can write as $\beta \to 1$:

$$
\begin{aligned}
\phi_\beta(x) &= \frac{x e^{(\beta-1) \log x}}{\beta(\beta - 1)} - \frac{\beta x}{\beta(\beta - 1)} + \frac{1}{\beta} \\
&= \frac{x(-\beta + 1 + (\beta - 1) \log x + o(\beta - 1))}{\beta(\beta - 1)} + \frac{1}{\beta} \\
&= \frac{-x + x \log x + 1}{\beta} + o(1).
\end{aligned}
$$

Then:

$$\lim_{\beta \to 1} \phi_\beta(x) = x \log x - x + 1.$$

### B. Equivalence between the Bregman divergence and the $\beta$-divergence

For $\beta \in \mathbb{R} \setminus \{0, 1\}$:

$$
\begin{aligned}
d_{\phi_\beta}(x|y) &= \frac{x^\beta}{\beta(\beta - 1)} - \frac{x}{\beta - 1} - \frac{y^\beta}{\beta(\beta - 1)} \\
&\quad + \frac{y}{\beta - 1} - \left( \frac{y^{\beta-1}}{\beta - 1} - \frac{1}{\beta - 1} \right)(x - y) \\
&= \frac{1}{\beta(\beta - 1)} (x^\beta + (\beta - 1)y^\beta - \beta x y^{\beta-1}) \\
&= d_\beta(x|y).
\end{aligned}
$$

It is straightforward to check that the equality $d_{\phi_\beta}(x|y) = d_\beta(x|y)$ also holds for $\beta \in \{0, 1\}$:

$$
\begin{aligned}
d_{\phi_0}(x|y) &= -\log x + x - (-\log y + y) - (-\frac{1}{y} + 1)(x - y) \\
&= -\log x + \log y + (x - y) + \frac{x}{y} - 1 - (x - y) \\
&= -\log \frac{x}{y} + \frac{x}{y} - 1 \\
&= d_0(x|y), \\
d_{\phi_1}(x|y) &= x \log x - x + 1 - (y \log y - y + 1) - \log y(x - y) \\
&= x(\log x - \log y) + (y - x) \\
&= d_1(x|y).
\end{aligned}
$$