

---

# Generalizing Point Embeddings using the Wasserstein Space of Elliptical Distributions

---

**Boris Muzellec**  
CREST, ENSAE  
boris.muzellec@ensae.fr

**Marco Cuturi**  
Google Brain and CREST, ENSAE  
cuturi@google.com

## Abstract

Embedding complex objects as vectors in low dimensional spaces is a longstanding problem in machine learning. We propose in this work an extension of that approach, which consists in embedding objects as elliptical probability distributions, namely distributions whose densities have elliptical level sets. We endow these measures with the 2-Wasserstein metric, with two important benefits: (i) For such measures, the squared 2-Wasserstein metric has a closed form, equal to a weighted sum of the squared Euclidean distance between means and the squared Bures metric between covariance matrices. The latter is a Riemannian metric between positive semi-definite matrices, which turns out to be Euclidean on a suitable factor representation of such matrices, which is valid on the entire geodesic between these matrices. (ii) The 2-Wasserstein distance boils down to the usual Euclidean metric when comparing Diracs, and therefore provides a natural framework to extend point embeddings. We show that for these reasons Wasserstein elliptical embeddings are more intuitive and yield tools that are better behaved numerically than the alternative choice of Gaussian embeddings with the Kullback-Leibler divergence. In particular, and unlike previous work based on the KL geometry, we learn elliptical distributions that are not necessarily diagonal. We demonstrate the advantages of elliptical embeddings by using them for visualization, to compute embeddings of words, and to reflect entailment or hypernymy.

## 1 Introduction

One of the holy grails of machine learning is to compute meaningful low-dimensional embeddings for high-dimensional complex data. That ability has recently proved crucial to tackle more advanced tasks, such as for instance: inference on texts using word embeddings [Mikolov et al., 2013b, Pennington et al., 2014, Bojanowski et al., 2017], improved image understanding [Norouzi et al., 2014], representations for nodes in large graphs [Grover and Leskovec, 2016].

Such embeddings have been traditionally recovered by seeking *isometric* embeddings in lower dimensional Euclidean spaces, as studied in [Johnson and Lindenstrauss, 1984, Bourgain, 1985]. Given  $n$  input points  $x_1, \dots, x_n$ , one seeks as many embeddings  $\mathbf{y}_1, \dots, \mathbf{y}_n$  in a target space  $\mathcal{Y} = \mathbb{R}^d$  whose pairwise distances  $\|\mathbf{y}_i - \mathbf{y}_j\|_2$  do not depart too much from the original distances  $d_{\mathcal{X}}(x_i, x_j)$  in the input space. Note that when  $d$  is restricted to be 2 or 3, these embeddings  $(\mathbf{y}_i)_i$  provide a useful way to visualize the entire dataset. Starting with metric multidimensional scaling (mMDS) [De Leeuw, 1977, Borg and Groenen, 2005], several approaches have refined this intuition [Tenenbaum et al., 2000, Roweis and Saul, 2000, Hinton and Roweis, 2003, Maaten and Hinton, 2008]. More general criteria, such as reconstruction error [Hinton and Salakhutdinov, 2006, Kingma and Welling, 2014]; co-occurrence [Globerson et al., 2007]; or relational knowledge, be it in metric learning [Weinberger and Saul, 2009] or between words [Mikolov et al., 2013b] can be used to obtain vector embeddings. In such cases, distances  $\|\mathbf{y}_i - \mathbf{y}_j\|_2$  between embeddings, or alternatively their dot-products  $\langle \mathbf{y}_i, \mathbf{y}_j \rangle$

must comply with sophisticated desiderata. Naturally, more general and flexible approaches in which the embedding space  $\mathcal{Y}$  needs not be Euclidean can be considered, for instance in generalized MDS on the sphere [Maron et al., 2010], on surfaces [Bronstein et al., 2006], in spaces of trees [Bădoiu et al., 2007, Fakcharoenphol et al., 2003] or, more recently, computed in the Poincaré hyperbolic space [Nickel and Kiela, 2017].

**Probabilistic Embeddings.** Our work belongs to a recent trend, pioneered by Vilnis and McCallum, who proposed to embed data points as *probability measures* in  $\mathbb{R}^d$  [2015], and therefore generalize point embeddings. Indeed, point embeddings can be regarded as a very particular—and degenerate—case of probabilistic embedding, in which the uncertainty is infinitely concentrated on a single point (a Dirac). Probability measures can be more spread-out, or event multimodal, and provide therefore an opportunity for additional flexibility. Naturally, such an opportunity can only be exploited by defining a metric, divergence or dot-product on the space (or a subspace thereof) of probability measures. Vilnis and McCallum proposed to embed words as *Gaussians* endowed either with the Kullback-Leibler (KL) divergence or the expected likelihood kernel [Jebara et al., 2004]. The Kullback-Leibler and expected likelihood kernel on measures have, however, an important drawback: these geometries do not coincide with the usual Euclidean metric between point embeddings when the variances of these Gaussians collapse. Indeed, the KL divergence and the  $\ell_2$  distance between two Gaussians diverges to  $\infty$  or saturates when the variances of these Gaussians become small. To avoid numerical instabilities arising from this degeneracy, Vilnis and McCallum must restrict their work to diagonal covariance matrices. In a concurrent approach, Singh et al. represent words as distributions over their contexts in the optimal transport geometry [Singh et al., 2018].

**Contributions.** We propose in this work a new framework for probabilistic embeddings, in which point embeddings are seamlessly handled as a particular case. We consider arbitrary families of elliptical distributions, which subsume Gaussians, and also include uniform elliptical distributions, which are arguably easier to visualize because of their compact support. Our approach uses the 2-Wasserstein distance to compare elliptical distributions. The latter can handle degenerate measures, and both its value and its gradients admit closed forms [Gelbrich, 1990], either in their natural Riemannian formulation, as well as in a more amenable local Euclidean parameterization. We provide numerical tools to carry out the computation of elliptical embeddings in different scenarios, both to optimize them with respect to metric requirements (as is done in multidimensional scaling) or with respect to dot-products (as shown in our applications to word embeddings for entailment, similarity and hypernymy tasks) for which we introduce a proxy using a polarization identity.

**Notations**  $\mathcal{S}_{++}^d$  (resp.  $\mathcal{S}_+^d$ ) is the set of positive (resp. semi-)definite  $d \times d$  matrices. For two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  and a matrix  $\mathbf{M} \in \mathcal{S}_+^d$ , we write the Mahalanobis norm induced by  $\mathbf{M}$  as  $\|\mathbf{x} - \mathbf{c}\|_{\mathbf{M}}^2 = (\mathbf{x} - \mathbf{c})^T \mathbf{M} (\mathbf{x} - \mathbf{c})$  and  $|\mathbf{M}|$  for  $\det(\mathbf{M})$ . For  $V$  an affine subspace of dimension  $m$  of  $\mathbb{R}^d$ ,  $\lambda_V$  is the Lebesgue measure on that subspace.  $\mathbf{M}^\dagger$  is the pseudo inverse of  $\mathbf{M}$ .

## 2 The Geometry of Elliptical Distributions in the Wasserstein Space

We recall in this section basic facts about elliptical distributions in  $\mathbb{R}^d$ . We adopt a general formulation that can handle measures supported on subspaces of  $\mathbb{R}^d$  as well as Dirac (point) measures. That level of generality is needed to provide a seamless connection with usual vector embeddings, seen in the context of this paper as Dirac masses. We recall results from the literature showing that the squared 2-Wasserstein distance between two distributions from the same family of elliptical distributions is equal to the squared Euclidean distance between their means plus the squared Bures metric between their scale parameter scaled by a suitable constant.

**Elliptically Contoured Densities.** In their simplest form, elliptical distributions can be seen as generalizations of Gaussian multivariate densities in  $\mathbb{R}^d$ : their level sets describe concentric ellipsoids, shaped following a scale parameter  $\mathbf{C} \in \mathcal{S}_{++}^d$ , and centered around a mean parameter  $\mathbf{c} \in \mathbb{R}^d$  [Cambanis et al., 1981]. The density at a point  $\mathbf{x}$  of such distributions is  $f(\|\mathbf{x} - \mathbf{c}\|_{\mathbf{C}^{-1}}) / \sqrt{|\mathbf{C}|}$  where the generator function  $f$  is such that  $\int_{\mathbb{R}^d} f(\|\mathbf{x}\|^2) d\mathbf{x} = 1$ . Gaussians are recovered with  $f = g$ ,  $g(\cdot) \propto e^{-\cdot/2}$  while uniform distributions on full rank ellipsoids result from  $f = u$ ,  $u(\cdot) \propto \mathbf{1}_{\cdot \leq 1}$ .

Because the norm induced by  $\mathbf{C}^{-1}$  appears in formulas above, the scale parameter  $\mathbf{C}$  must have full rank for these definitions to be meaningful. Cases where  $\mathbf{C}$  does not have full rank can however

appear when a probability measure is supported on an affine subspace<sup>1</sup> of  $\mathbb{R}^d$ , such as lines in  $\mathbb{R}^2$ , or even possibly a space of null dimension when the measure is supported on a single point (a Dirac measure), in which case its scale parameter  $\mathbf{C}$  is  $\mathbf{0}$ . We provide in what follows a more general approach to handle these degenerate cases.

**Elliptical Distributions.** To lift this limitation, several reformulations of elliptical distributions have been proposed to handle degenerate scale matrices  $\mathbf{C}$  of rank  $\text{rk } \mathbf{C} < d$ . Gelbrich [1990, Theorem 2.4] defines elliptical distributions as measures with a density w.r.t the Lebesgue measure of dimension  $\text{rk } \mathbf{C}$ , in the affine space  $\mathbf{c} + \text{Im } \mathbf{C}$ , where the image of  $\mathbf{C}$  is  $\text{Im } \mathbf{C} \stackrel{\text{def}}{=} \{\mathbf{C}\mathbf{x}, \mathbf{x} \in \mathbb{R}^d\}$ . This approach is intuitive, in that it reduces to describing densities in their relevant subspace. A more elegant approach uses the parameterization provided by characteristic functions [Cambanis et al., 1981, Fang et al., 1990]. In a nutshell, recall that the characteristic function of a multivariate Gaussian is equal to  $\phi(\mathbf{t}) = e^{it^T \mathbf{c}} g(\mathbf{t}^T \mathbf{C} \mathbf{t})$  where, as in the paragraph above,  $g(\cdot) = e^{-\cdot/2}$ . A natural generalization to consider other elliptical distributions is therefore to consider for  $g$  other functions  $h$  of positive type [Ushakov, 1999, Theo.1.8.9], such as the indicator function  $u$  above, and still apply them to the same argument  $\mathbf{t}^T \mathbf{C} \mathbf{t}$ . Such functions are called *characteristic generators* and fully determine, along with a mean  $\mathbf{c}$  and a scale parameter  $\mathbf{C}$ , an elliptical measure. This parameterization does not require the scale parameter  $\mathbf{C}$  to be invertible, and therefore allows to define probability distributions that do not have necessarily a density w.r.t to the Lebesgue measure in  $\mathbb{R}^d$ . Both constructions are relatively complex, and we refer the interested reader to these references for a rigorous treatment.

**Rank Deficient Elliptical Distributions and their Variances.** For the purpose of this work, we will only require the following result: the variance of an elliptical measure is equal to its scale parameter  $\mathbf{C}$  multiplied by a scalar that only depends on its characteristic generator. Indeed, given a mean vector  $\mathbf{c} \in \mathbb{R}^d$ , a scale *semi*-definite matrix  $\mathbf{C} \in \mathcal{S}_+^d$  and a characteristic generator function  $h$ , we define  $\mu_{h,\mathbf{c},\mathbf{C}}$  to be the measure with characteristic function  $\mathbf{t} \mapsto e^{it^T \mathbf{c}} h(\mathbf{t}^T \mathbf{C} \mathbf{t})$ . In that case, one can show that the covariance matrix of  $\mu_{h,\mathbf{c},\mathbf{C}}$  is equal to its scale parameter  $\mathbf{C}$  times a constant  $\tau_h$  that only depends on  $h$ , namely

$$\text{var}(\mu_{h,\mathbf{c},\mathbf{C}}) = \tau_h \mathbf{C} . \quad (1)$$

For Gaussians, the scale parameter  $\mathbf{C}$  and its covariance matrix coincide, that is  $\tau_g = 1$ . For uniform elliptical distributions, one has  $\tau_u = 1/(d+2)$ : the covariance of a uniform distribution on the volume  $\{\mathbf{c} + \mathbf{C}\mathbf{x}, \mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\| = 1\}$ , such as those represented in Figure 1, is equal to  $\mathbf{C}/(d+2)$ .

**The 2-Wasserstein Bures Metric** A natural metric for elliptical distributions arises from optimal transport (OT) theory. We refer interested readers to [Santambrogio, 2015, Peyré and Cuturi, 2018] for exhaustive surveys on OT. Recall that for two arbitrary probability measures  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ , their squared 2-Wasserstein distance is equal to

$$W_2^2(\mu, \nu) \stackrel{\text{def}}{=} \inf_{X \sim \mu, Y \sim \nu} \mathbb{E} \|X - Y\|_2^2 .$$

This formula rarely has a closed form. However, in the footsteps of Dowson and Landau [1982] who proved it for Gaussians, Gelbrich [1990] showed that for  $\alpha \stackrel{\text{def}}{=} \mu_{h,\mathbf{a},\mathbf{A}}$  and  $\beta \stackrel{\text{def}}{=} \mu_{h,\mathbf{b},\mathbf{B}}$  in the *same* family  $\mathcal{P}_h = \{\mu_{h,\mathbf{c},\mathbf{C}}, \mathbf{c} \in \mathbb{R}^d, \mathbf{C} \in \mathcal{S}_+^d\}$ , one has

$$W_2^2(\alpha, \beta) = \|\mathbf{a} - \mathbf{b}\|_2^2 + \mathfrak{B}^2(\text{var } \alpha, \text{var } \beta) = \|\mathbf{a} - \mathbf{b}\|_2^2 + \tau_h \mathfrak{B}^2(\mathbf{A}, \mathbf{B}) , \quad (2)$$

<sup>1</sup>For instance, the random variable  $Y$  in  $\mathbb{R}^2$  obtained by duplicating the same normal random variable  $X$  in  $\mathbb{R}$ ,  $Y = [X, X]$ , is supported on a line in  $\mathbb{R}^2$  and has no density w.r.t the Lebesgue measure in  $\mathbb{R}^2$ .

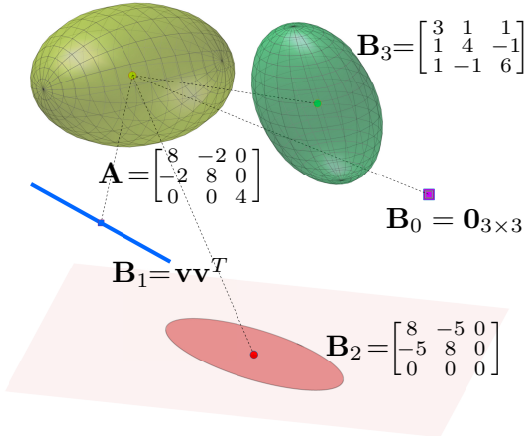


Figure 1: Five measures from the family of uniform elliptical distributions in  $\mathbb{R}^3$ . Each measure has a mean (location) and scale parameter. In this carefully selected example, the reference measure (with scale parameter  $\mathbf{A}$ ) is equidistant (according to the 2-Wasserstein metric) to the four remaining measures, whose scale parameters  $\mathbf{B}_0, \mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3$  have ranks equal to their indices (here,  $\mathbf{v} = [3, 7, -2]^T$ ).

where  $\mathfrak{B}^2$  is the (squared) Bures metric on  $\mathcal{S}_+^d$ , proposed in quantum information geometry [1969] and studied recently in [Bhatia et al., 2018, Malagò et al., 2018],

$$\mathfrak{B}^2(\mathbf{X}, \mathbf{Y}) \stackrel{\text{def}}{=} \text{Tr}(\mathbf{X} + \mathbf{Y} - 2(\mathbf{X}^{\frac{1}{2}} \mathbf{Y} \mathbf{X}^{\frac{1}{2}})^{\frac{1}{2}}) . \quad (3)$$

The factor  $\tau_h$  next to the rightmost term  $\mathfrak{B}^2$  in (2) arises from homogeneity of  $\mathfrak{B}^2$  in its arguments (3), which is leveraged using the identity in (1).

**A few remarks** (i) When both scale matrices  $\mathbf{A} = \text{diag } \mathbf{d}_\mathbf{A}$  and  $\mathbf{B} = \text{diag } \mathbf{d}_\mathbf{B}$  are diagonal,  $W_2^2(\alpha, \beta)$  is the sum of two terms: the usual squared Euclidean distance between their means, plus  $\tau_h$  times the squared *Hellinger* metric between the diagonals  $\mathbf{d}_\mathbf{A}, \mathbf{d}_\mathbf{B}$ :  $\mathfrak{H}^2(\mathbf{d}_\mathbf{A}, \mathbf{d}_\mathbf{B}) \stackrel{\text{def}}{=} \|\sqrt{\mathbf{d}_\mathbf{A}} - \sqrt{\mathbf{d}_\mathbf{B}}\|_2^2$ . (ii) The distance  $W_2$  between two Diracs  $\delta_{\mathbf{a}}, \delta_{\mathbf{b}}$  is equal to the usual distance between vectors  $\|\mathbf{a} - \mathbf{b}\|_2$ . (iii) The squared distance  $W_2^2$  between a Dirac  $\delta_{\mathbf{a}}$  and a measure  $\mu_{h, \mathbf{b}, \mathbf{B}}$  in  $\mathcal{P}_h$  reduces to  $\|\mathbf{a} - \mathbf{b}\|_2^2 + \tau_h \text{Tr} \mathbf{B}$ . The distance between a point and an ellipsoid distribution therefore always *increases* as the scale parameter of the latter increases. Although this point makes sense from the quadratic viewpoint of  $W_2^2$  (in which the quadratic contribution  $\|\mathbf{a} - \mathbf{x}\|_2^2$  of points  $\mathbf{x}$  in the ellipsoid that stand further away from  $\mathbf{a}$  than  $\mathbf{b}$  will dominate that brought by points  $\mathbf{x}$  that are closer, see Figure 3) this may be counterintuitive for applications to visualization, an issue that will be addressed in Section 4. (iv) The  $W_2$  distance between two elliptical distributions in the same family  $\mathcal{P}_h$  is always finite, no matter how degenerate they are. This is illustrated in Figure 1 in which a uniform measure  $\mu_{\mathbf{a}, \mathbf{A}}$  is shown to be exactly equidistant to four other uniform elliptical measures, some of which are degenerate. However, as can be hinted by the simple example of the Hellinger metric, that distance may not be differentiable for degenerate measures (in the same sense that  $(\sqrt{x} - \sqrt{y})^2$  is defined at  $x = 0$  but not differentiable w.r.t  $x$ ). (v) Although we focus in this paper on uniform elliptical distributions, notably because they are easier to plot and visualize, considering any other elliptical family simply amounts to changing the constant  $\tau_h$  next to the Bures metric in (2). Alternatively, increasing (or tuning) that parameter  $\tau_h$  simply amounts to considering elliptical distributions with increasingly heavier tails.

### 3 Optimizing over the Space of Elliptical Embeddings

Our goal in this paper is to use the set of elliptical distributions endowed with the  $W_2$  distance as an embedding space. To optimize objective functions involving  $W_2$  terms, we study in this section several parameterizations of the parameters of elliptical distributions. Location parameters only appear in the computation of  $W_2$  through their Euclidean metric, and offer therefore no particular challenge. Scale parameters are more tricky to handle since they are constrained to lie in  $\mathcal{S}_+^d$ . Rather than keeping track of scale parameters, we advocate optimizing directly on factors (square roots) of such parameters, which results in simple Euclidean (unconstrained) updates reviewed below.

**Geodesics for Elliptical Distributions** When  $\mathbf{A}$  and  $\mathbf{B}$  have full rank, the geodesic from  $\alpha$  to  $\beta$  is a curve of measures in the same family of elliptic distributions, characterized by location and scale parameters  $\mathbf{c}(t), \mathbf{C}(t)$ , where

$$\mathbf{c}(t) = (1 - t)\mathbf{a} + t\mathbf{b}; \quad \mathbf{C}(t) = ((1 - t)\mathbf{I} + t\mathbf{T}^{\mathbf{AB}}) \mathbf{A} ((1 - t)\mathbf{I} + t\mathbf{T}^{\mathbf{AB}}) , \quad (4)$$

and where the matrix  $\mathbf{T}^{\mathbf{AB}}$  is such that  $\mathbf{x} \rightarrow \mathbf{T}^{\mathbf{AB}}(\mathbf{x} - \mathbf{a}) + \mathbf{b}$  is the so-called Brenier optimal transportation map [1987] from  $\alpha$  to  $\beta$ , given in closed form as,

$$\mathbf{T}^{\mathbf{AB}} \stackrel{\text{def}}{=} \mathbf{A}^{-\frac{1}{2}} (\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}} , \quad (5)$$

and is the unique matrix such that  $\mathbf{B} = \mathbf{T}^{\mathbf{AB}} \mathbf{A} \mathbf{T}^{\mathbf{AB}}$  [Peyré and Cuturi, 2018, Remark 2.30]. When  $\mathbf{A}$  is degenerate, such a curve still exists as long as  $\text{Im } \mathbf{B} \subset \text{Im } \mathbf{A}$ , in which case the expression above is still valid using pseudo-inverse square roots  $\mathbf{A}^{\dagger/2}$  in place of the usual inverse square-root.

**Differentiability in Riemannian Parameterization** Scale parameters are restricted to lie on the cone  $\mathcal{S}_+^d$ . For such problems, it is well known that a direct gradient-and-project based optimization on scale parameters would prove too expensive. A natural remedy to this issue is to perform manifold optimization [Absil et al., 2009]. Indeed, as in any Riemannian manifold, the Riemannian gradient  $\text{grad}_x \frac{1}{2} d^2(x, y)$  is given by  $-\log_x y$  [Lee, 1997]. Using the expressions of the exp and log given in [Malagò et al., 2018], we can show that minimizing  $\frac{1}{2} \mathfrak{B}^2(\mathbf{A}, \mathbf{B})$  using Riemannian gradient descent corresponds to making updates of the form, with step length  $\eta$

$$\mathbf{A}' = ((1 - \eta)\mathbf{I} + \eta\mathbf{T}^{\mathbf{AB}}) \mathbf{A} ((1 - \eta)\mathbf{I} + \eta\mathbf{T}^{\mathbf{AB}}) . \quad (6)$$

When  $0 \leq \eta \leq 1$ , this corresponds to considering a new point  $\mathbf{A}'$  closer to  $\mathbf{B}$  along the Bures geodesic between  $\mathbf{A}$  and  $\mathbf{B}$ . When  $\eta$  is negative or larger than 1,  $\mathbf{A}'$  no longer lies on this geodesic but is guaranteed to remain PSD, as can be seen from (6). Figure 2 shows a  $W_2$  geodesic between two measures  $\mu_0$  and  $\mu_1$ , as well as its extrapolation following exactly the formula given in (4). That figure illustrates that  $\mu_t$  is not necessarily geodesic outside of the boundaries  $[0, 1]$  w.r.t. three relevant measures, because its metric derivative is smaller than 1 [Ambrosio et al., 2006, Theorem 1.1.2]. When negative steps are taken (for instance when the  $W_2^2$  distance needs to be increased), this lack of geodesicity has proved difficult to handle numerically for a simple reason: such updates may lead to degenerate scale parameters  $\mathbf{A}'$ , as illustrated around time  $t = 1.5$  of the curve in Figure 2. Another obvious drawback of Riemannian approaches is that they are not as well studied as simpler non-constrained Euclidean problems, for which a plethora of optimization techniques are available. This observations motivates an alternative Euclidean parameterization, detailed in the next paragraph.

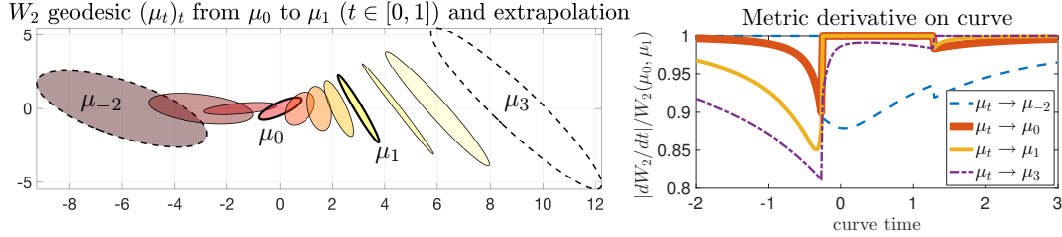


Figure 2: (left) Interpolation  $(\mu_t)_t$  between two measures  $\mu_0$  and  $\mu_1$  following the geodesic equation (4). The same formula can be used to interpolate on the left and right of times 0, 1. Displayed times are  $[-2, -1, -.5, 0, .25, .5, .75, 1, 1.5, 2, 3]$ . Note that geodesicity is not ensured outside of the boundaries  $[0, 1]$ . This is illustrated in the right plot displaying normalized metric derivatives of the curve  $\mu_t$  to four relevant points:  $\mu_0, \mu_1, \mu_{-2}, \mu_3$ . The curve  $\mu_t$  is not always locally geodesic, as can be seen by the fact that the metric derivative is strictly smaller than 1 in several cases.

**Differentiability in Euclidean Parameterization** A canonical way to handle a PSD constraint for  $\mathbf{A}$  is to rewrite it in factor form  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ . In the particular case of the Bures metric, we show that this simple parametrization comes without losing the geometric interest of manifold optimization, while benefiting from simpler additive updates. Indeed, one can (see supplementary material) that the gradient of the squared Bures metric has the following gradient:

$$\nabla_{\mathbf{L}} \frac{1}{2} \mathfrak{B}^2(\mathbf{A}, \mathbf{B}) = (\mathbf{I} - \mathbf{T}^{\mathbf{A}\mathbf{B}}) \mathbf{L}, \quad \text{with updates } \mathbf{L}' = ((1 - \eta)\mathbf{I} + \eta\mathbf{T}^{\mathbf{A}\mathbf{B}}) \mathbf{L}. \quad (7)$$

**Links between Euclidean and Riemannian Parameterization** The factor updates in (7) are exactly equivalent to the Riemannian ones (6) in the sense that  $\mathbf{A}' = \mathbf{L}'\mathbf{L}'^T$ . Therefore, by using a factor parameterization we carry out updates that stay on the Riemannian geodesic yet only require linear updates on  $\mathbf{L}$ , independently of the factor  $\mathbf{L}$  chosen to represent  $\mathbf{A}$  (given a factor  $\mathbf{L}$  of  $\mathbf{A}$ , any right-side multiplication of that matrix by a unitary matrix remains a factor of  $\mathbf{A}$ ).

When considering a general loss function  $\mathcal{L}$  that take as arguments squared Bures distances, one can also show that  $\mathcal{L}$  is geodesically convex w.r.t. to scale matrices  $\mathbf{A}$  if and only if it is convex in the usual sense with respect to  $\mathbf{L}$ , where  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ . Write now  $\mathbf{L}_{\mathbf{B}} = \mathbf{T}^{\mathbf{A}\mathbf{B}}\mathbf{L}$ . One can recover that  $\mathbf{L}_{\mathbf{B}}\mathbf{L}_{\mathbf{B}}^T = \mathbf{B}$ . Therefore, expanding the expression  $\mathfrak{B}^2$  for the right term below we obtain

$$\mathfrak{B}^2(\mathbf{A}, \mathbf{B}) = \mathfrak{B}^2(\mathbf{L}\mathbf{L}^T, \mathbf{L}_{\mathbf{B}}\mathbf{L}_{\mathbf{B}}^T) = \mathfrak{B}^2(\mathbf{L}\mathbf{L}^T, \mathbf{T}^{\mathbf{A}\mathbf{B}}\mathbf{L}(\mathbf{T}^{\mathbf{A}\mathbf{B}}\mathbf{L})^T) = \|\mathbf{L} - \mathbf{T}^{\mathbf{A}\mathbf{B}}\mathbf{L}\|_F^2$$

Indeed, the Bures distance simply reduces to the Frobenius distance between two factors of  $\mathbf{A}$  and  $\mathbf{B}$ . However these factors need to be carefully chosen: given  $\mathbf{L}$  for  $\mathbf{A}$ , the factor for  $\mathbf{B}$  must be computed according to an optimal transport map  $\mathbf{T}^{\mathbf{A}\mathbf{B}}$ .

**Polarization between Elliptical Distributions** Some of the applications we consider, such as the estimation of word embeddings, are inherently based on dot-products. By analogy with the polarization identity,  $\langle \mathbf{x}, \mathbf{y} \rangle = (\|\mathbf{x} - \mathbf{0}\|^2 + \|\mathbf{y} - \mathbf{0}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2)/2$ , we define a Wasserstein-Bures pseudo-dot-product, where  $\delta_{\mathbf{0}} = \mu_{\mathbf{0}, \mathbf{0}, d \times d}$  is the Dirac mass at  $\mathbf{0}$ ,

$$[\mu_{\mathbf{a}, \mathbf{A}} : \mu_{\mathbf{b}, \mathbf{B}}] \stackrel{\text{def}}{=} \frac{1}{2} (W_2^2(\mu_{\mathbf{a}, \mathbf{A}}, \delta_{\mathbf{0}}) + W_2^2(\mu_{\mathbf{b}, \mathbf{B}}, \delta_{\mathbf{0}}) - W_2^2(\mu_{\mathbf{a}, \mathbf{A}}, \mu_{\mathbf{b}, \mathbf{B}})) = \langle \mathbf{a}, \mathbf{b} \rangle + \text{Tr}(\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}$$

Note that  $[\cdot : \cdot]$  is not an actual inner product since the Bures metric is not Hilbertian, unless we restrict ourselves to diagonal covariance matrices, in which case it is the inner product between  $(\mathbf{a}, \sqrt{\mathbf{d}_A})$  and  $(\mathbf{b}, \sqrt{\mathbf{d}_B})$ . We use  $[\mu_{\mathbf{a},A} : \mu_{\mathbf{b},B}]$  as a similarity measure which has, however, some regularity: one can show that when  $\mathbf{a}, \mathbf{b}$  are constrained to have equal norms and  $\mathbf{A}$  and  $\mathbf{B}$  equal traces, then  $[\mu_{\mathbf{a},A} : \mu_{\mathbf{b},B}]$  is maximal when  $\mathbf{a} = \mathbf{b}$  and  $\mathbf{A} = \mathbf{B}$ . Differentiating all three terms in that sum, the gradient of this pseudo dot-product w.r.t.  $\mathbf{A}$  reduces to  $\nabla_{\mathbf{A}}[\mu_{\mathbf{a},A} : \mu_{\mathbf{b},B}] = \mathbf{T}^{\mathbf{A}\mathbf{B}}$ .

**Computational Aspects** The computational bottleneck of gradient-based Bures optimization lies in the matrix square roots and inverse square roots operations that arise when instantiating transport maps  $\mathbf{T}$  as in (5). A naive method using eigenvector decomposition is far too time-consuming, and there is not yet, to the best of our knowledge, a straightforward way to perform it in batches on a GPU. We propose to use Newton-Schulz iterations (Algorithm 1, see [Higham, 2008, Ch. 6]) to approximate these root computations. These iterations producing both a root and an inverse root approximation, and, relying exclusively on matrix-matrix multiplications, stream efficiently on GPUs. Another problem lies in the fact that numerous roots and inverse-roots are required to form map  $\mathbf{T}$ . To solve this, we exploit an alternative formula for  $\mathbf{T}^{\mathbf{A}\mathbf{B}}$  (proof in the supplementary material):

$$\mathbf{T}^{\mathbf{A}\mathbf{B}} = \mathbf{A}^{-\frac{1}{2}}(\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}\mathbf{A}^{-\frac{1}{2}} = \mathbf{B}^{\frac{1}{2}}(\mathbf{B}^{\frac{1}{2}}\mathbf{A}\mathbf{B}^{\frac{1}{2}})^{-\frac{1}{2}}\mathbf{B}^{\frac{1}{2}}. \quad (8)$$

In a gradient update, both the loss and the gradient of the metric are needed. In our case, we can use the matrix roots computed during loss evaluation and leverage the identity above to compute on a budget the gradients with respect to either scale matrices  $\mathbf{A}$  and  $\mathbf{B}$ . Indeed, a naive computation of  $\nabla_{\mathbf{A}}\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$  and  $\nabla_{\mathbf{B}}\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$  would require the knowledge of 6 roots:

$$\mathbf{A}^{\frac{1}{2}}, \mathbf{B}^{\frac{1}{2}}, (\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}, (\mathbf{B}^{\frac{1}{2}}\mathbf{A}\mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}}, \mathbf{A}^{-\frac{1}{2}}, \text{ and } \mathbf{B}^{-\frac{1}{2}}$$

to compute the following transport maps

$$\mathbf{T}^{\mathbf{A}\mathbf{B}} = \mathbf{A}^{-\frac{1}{2}}(\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}\mathbf{A}^{-\frac{1}{2}}, \mathbf{T}^{\mathbf{B}\mathbf{A}} = \mathbf{B}^{-\frac{1}{2}}(\mathbf{B}^{\frac{1}{2}}\mathbf{A}\mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}}\mathbf{B}^{-\frac{1}{2}},$$

namely four matrix roots and two matrix inverse roots. We can avoid computing those six matrices using identity (8) and limit ourselves to two runs of Algorithm 1, to obtain the same quantities as

$$\{\mathbf{Y}_1 \stackrel{\text{def}}{=} \mathbf{A}^{\frac{1}{2}}, \mathbf{Z}_1 \stackrel{\text{def}}{=} \mathbf{A}^{-\frac{1}{2}}\}, \{\mathbf{Y}_2 \stackrel{\text{def}}{=} (\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}, \mathbf{Z}_2 \stackrel{\text{def}}{=} (\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}})^{-\frac{1}{2}}\} \\ \mathbf{T}^{\mathbf{A}\mathbf{B}} = \mathbf{Z}_1\mathbf{Y}_2\mathbf{Z}_1, \mathbf{T}^{\mathbf{B}\mathbf{A}} = \mathbf{Y}_1\mathbf{Z}_2\mathbf{Y}_1.$$

When computing the gradients of  $n \times m$  squared Wasserstein distances  $W_2^2(\alpha_i, \beta_j)$  in parallel, one only needs to run  $n$  Newton-Schulz algorithms (in parallel) to compute matrices  $(\mathbf{Y}_1^i, \mathbf{Z}_1^i)_{i \leq n}$ , and then  $n \times m$  Newton-Schulz algorithms to recover cross matrices  $\mathbf{Y}_2^{i,j}, \mathbf{Z}_2^{i,j}$ . On the other hand, using an automatic differentiation framework would require an additional backward computation of the same complexity as the forward pass evaluating computation of the roots and inverse roots, hence requiring roughly twice as many operations per batch.

**Avoiding Rank Deficiency at Optimization Time** Although  $\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$  is defined for rank deficient matrices  $\mathbf{A}$  and  $\mathbf{B}$ , it is not differentiable with respect to these matrices if they are rank deficient. Indeed, as mentioned earlier, this can be compared to the non-differentiability of the Hellinger metric,  $(\sqrt{x} - \sqrt{y})^2$  when  $x$  or  $y$  becomes 0, at which point it becomes *not* differentiable. If  $\text{Im } \mathbf{B} \not\subset \text{Im } \mathbf{A}$ , which is notably the case if  $\text{rk } \mathbf{B} > \text{rk } \mathbf{A}$ , then  $\nabla_{\mathbf{A}}\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$  no longer exists. However, even in that case,  $\nabla_{\mathbf{B}}\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$  exists iff  $\text{Im } \mathbf{A} \subset \text{Im } \mathbf{B}$ . Since it would be cumbersome to account for these subtleties in a large scale optimization setting, we propose to add a small common regularization term to all the factor products considered for our embeddings, and set  $\mathbf{A}_\varepsilon = \mathbf{L}\mathbf{L}^T + \varepsilon\mathbf{I}$  where  $\varepsilon > 0$  is a hyperparameter. This ensures that all matrices are full rank, and thus that all gradients exist. Most importantly, all our derivations still hold with this regularization, and can be shown to leave the method to compute the gradients w.r.t  $\mathbf{L}$  unchanged, namely remain equal to  $(\mathbf{I} - \mathbf{T}^{\mathbf{A}_\varepsilon\mathbf{B}})\mathbf{L}$ .

---

**Algorithm 1** Newton-Schulz

---

**Input:** PSD matrix  $\mathbf{A}$ ,  $\varepsilon > 0$

$\mathbf{Y} \leftarrow \frac{\mathbf{A}}{(1+\varepsilon)\|\mathbf{A}\|}, \mathbf{Z} \leftarrow \mathbf{I}$

**while** not converged **do**

$\mathbf{T} \leftarrow (3\mathbf{I} - \mathbf{Z}\mathbf{Y})/2$

$\mathbf{Y} \leftarrow \mathbf{Y}\mathbf{T}$

$\mathbf{Z} \leftarrow \mathbf{T}\mathbf{Z}$

**end while**

$\mathbf{Y} \leftarrow \sqrt{(1+\varepsilon)\|\mathbf{A}\|}\mathbf{Y}$

$\mathbf{Z} \leftarrow \frac{\mathbf{Z}}{\sqrt{(1+\varepsilon)\|\mathbf{A}\|}}$

**Output:** square root  $\mathbf{Y}$ , inverse square root  $\mathbf{Z}$

---

## 4 Experiments

We discuss in this section several applications of elliptical embeddings. We first consider a simple mMDS type visualization task, in which elliptical distributions in  $d = 2$  are used to embed isometrically points in high dimension. We argue that for such purposes, a more natural way to visualize ellipses is to use their precision matrices. This is due to the fact that the human eye somewhat acts in the opposite direction to the Bures metric, as discussed in Figure 3. We follow with more advanced experiments in which we consider the task of computing word embeddings on large corpora as a testing ground, and equal or improve on the state-of-the-art.

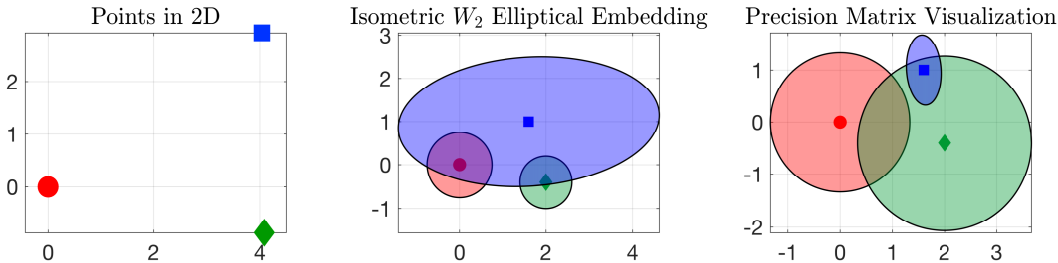


Figure 3: (left) three points on the plane. (middle) *isometric* elliptic embedding with the Bures metric: ellipses of a given color have the same respective distances as points on the left. Although the mechanics of optimal transport indicate that the blue ellipsoid is far from the two others, in agreement with the left plot, the human eye tends to focus on those areas that overlap (below the ellipsoid center) rather than those far away areas (north-east area) that contribute more significantly to the  $W_2$  distance. (right) the precision matrix visualization, obtained by considering ellipses with the same axes but inverted eigenvalues, agree better with intuition, since they emphasize that overlap and extension of the ellipse means on the contrary that those axis contribute less to the increase of the metric.

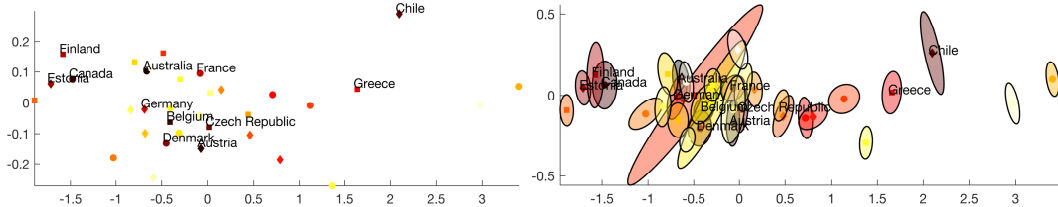


Figure 4: Toy experiment: visualization of a dataset of 10 PISA scores for 35 countries in the OECD. (left) MDS embeddings of these countries on the plane (right) elliptical embeddings on the plane using the precision visualization discussed in Figure 3. The normalized stress with standard MDS is 0.62. The stress with elliptical embeddings is close to  $5e - 3$  after 1000 gradient iterations, with random initializations for scale matrices (following a Standard Wishart with 4 degrees of freedom) and initial means located on the MDS solution.

**Visualizing Datasets Using Ellipsoids** Multidimensional scaling [De Leeuw, 1977] aims at embedding points  $\mathbf{x}_1, \dots, \mathbf{x}_n$  in a finite metric space in a lower dimensional one by minimizing the *stress*  $\sum_{ij} (\|\mathbf{x}_i - \mathbf{x}_j\| - \|\mathbf{y}_i - \mathbf{y}_j\|)^2$ . In our case, this translates to the minimization of  $\mathcal{L}_{\text{MDS}}(\mathbf{a}_1, \dots, \mathbf{a}_n, \mathbf{A}_1, \dots, \mathbf{A}_n) = \sum_{ij} (\|\mathbf{x}_i - \mathbf{x}_j\| - W_2(\mu_{\mathbf{a}_i, \mathbf{A}_i}, \mu_{\mathbf{a}_j, \mathbf{A}_j}))^2$ . This objective can be crudely minimized with a simple gradient descent approach operating on factors as advocated in Section 3, as illustrated in a toy example carried out using data from OECD’s PISA study<sup>2</sup>.

**Word Embeddings** The skipgram model [Mikolov et al., 2013a] computes word embeddings in a vector space by maximizing the log-probability of observing surrounding context words given an input central word. Vilnis and McCallum [2015] extended this approach to *diagonal* Gaussian embeddings using an energy whose overall principles we adopt here, adapted to elliptical distributions with *full* covariance matrices in the 2-Wasserstein space. For every word  $w$ , we consider an input (as a word) and an output (as a context) representation as an elliptical measure, denoted respectively  $\mu_w$  and  $\nu_w$ , both parameterized by a location vector and a scale parameter (stored in factor form).

<sup>2</sup><http://pisadataexplorer.oecd.org/ide/idepisa/>

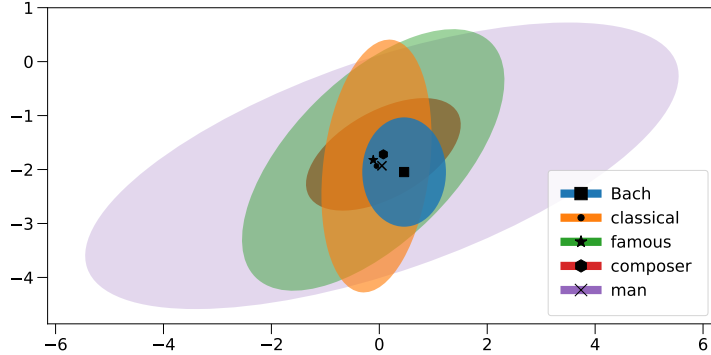


Figure 5: Precision matrix visualization of trained embeddings of a set of words on the plane spanned by the two principal eigenvectors of the covariance matrix of “Bach”.

Given a set  $\mathcal{R}$  of positive word/context pairs of words  $(w, c)$ , and for each input word a set  $N(w)$  of  $n$  negative contexts sampled randomly, we adapt Vilnis and McCallum’s loss function to the  $W_2^2$  distance to minimize the following hinge loss:

$$\sum_{(w,c) \in \mathcal{R}} \left[ M - [\mu_w : \nu_c] + \frac{1}{n} \sum_{c' \in N(w)} [\mu_w : \nu_{c'}] \right]_+$$

where  $M > 0$  is a margin parameter. We train our embeddings on the concatenated ukWaC and WaCkypedia corpora [Baroni et al., 2009], consisting of about 3 billion tokens, on which we keep only the tokens appearing more than 100 times in the text (for a total number of 261583 different words). We train our embeddings using adagrad [Duchi et al., 2011], sampling one negative context per positive context and, in order to prevent the norms of the embeddings to be too highly correlated with the corresponding word frequencies (see Figure in supplementary material), we use two distinct sets of embeddings for the input and context words.

We compare our full elliptical to diagonal Gaussian embeddings trained using the methods described in [Vilnis and McCallum, 2015] on a collection of similarity datasets by computing the Spearman rank correlation between the similarity scores provided in the data and the scores we compute based on our embeddings. Note that these results are obtained using context ( $\nu_w$ ) rather than input ( $\mu_w$ ) embeddings. For a fair comparison across methods, we set dimensions by ensuring that the number of free parameters remains the same: because of the symmetry in the covariance matrix, elliptical embeddings in dimension  $d$  have  $d + d(d+1)/2$  free parameters ( $d$  for the means,  $d(d+1)/2$  for the covariance matrices), as compared with  $2d$  for diagonal Gaussians. For elliptical embeddings, we use the common practice of using some form of normalized quantity (a cosine) rather than the direct dot product. We implement this here by computing the mean of two cosine terms, each corresponding separately to mean and covariance contributions:

$$\mathfrak{S}_{\mathfrak{B}}[\mu_{\mathbf{a}, \mathbf{A}}, \mu_{\mathbf{b}, \mathbf{B}}] := \frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\|\mathbf{a}\| \|\mathbf{b}\|} + \frac{\text{Tr}(\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}}{\sqrt{\text{Tr} \mathbf{A} \text{Tr} \mathbf{B}}}$$

Using this similarity measure rather than the Wasserstein-Bures dot product is motivated by the fact that the norms of the embeddings show some dependency with word frequencies (see figures in supplementary) and become dominant when comparing words with different frequencies scales. An alternative could have been obtained by normalizing the Wasserstein-Bures dot product in a more standard way that pools together means and covariances. However, as discussed in the supplementary material, this choice makes it harder to deal with the variations in scale of the means and covariances, therefore decreasing performance.

Table 1: Results for elliptical embeddings (evaluated using our cosine mixture) compared to diagonal Gaussian embeddings trained with the seomoz package (evaluated using expected likelihood cosine similarity as recommended by Vilnis and McCallum).

Dataset	W2G/45/C	Ell/12/CM
SimLex	<b>25.09</b>	24.09
WordSim	53.45	<b>66.02</b>
WordSim-R	61.70	<b>71.07</b>
WordSim-S	48.99	<b>60.58</b>
MEN	65.16	<b>65.58</b>
MC	59.48	<b>65.95</b>
RG	<b>69.77</b>	65.58
YP	<b>37.18</b>	25.14
MT-287	<b>61.72</b>	59.53
MT-771	<b>57.63</b>	56.78
RW	<b>40.14</b>	29.04



We also evaluate our embeddings on the Entailment dataset ([Baroni et al., 2012]), on which we obtain results roughly comparable to those of [Vilnis and McCallum, 2015]. Note that contrary to the similarity experiments, in this framework using the (unsymmetrical) KL divergence makes sense and possibly gives an advantage, as it is possible to choose the order of the arguments in the KL divergence between the entailing and entailed words.

**Hypernymy** In this experiment, we use the framework of [Nickel and Kiela, 2017] on hypernymy relationships to test our embeddings. A word A is said to be a *hypernym* of a word B if any B is a type of A, e.g. any *dog* is a type of *mammal*, thus constituting a tree-like structure on nouns. The WORDNET dataset [Miller, 1995] features a transitive closure of 743,241 hypernymy relations on 82,115 distinct nouns, which we consider as an undirected graph of relations  $\mathcal{R}$ . Similarly to the skipgram model, for each noun  $u$  we sample a fixed number  $n$  of negative examples and store them in set  $\mathcal{N}(u)$  to optimize the following loss:  $\sum_{(u,v) \in \mathcal{R}} \log \frac{e^{[\mu_u, \mu_v]}}{e^{[\mu_u, \mu_v]} + \sum_{v' \in \mathcal{N}(u)} e^{[\mu_u, \mu_{v'}]}}$ .

We train the model using SGD with only one set of embeddings. The embeddings are then evaluated on a link reconstruction task: we embed the full tree and rank the similarity of each positive hypernym pair  $(u, v)$  among all negative pairs  $(u, v')$  and compute the mean rank thus achieved as well as the mean average precision (MAP), using the Wasserstein-Bures dot product as the similarity measure. Elliptical embeddings consistently outperform Poincare embeddings for dimensions above a small threshold, as shown in Figure 6, which confirms our intuition that the addition of a notion of variance or uncertainty to point embeddings allows for a richer and more significant representation of words.

**Conclusion** We have proposed to use the space of elliptical distributions endowed with the  $W_2$  metric to embed complex objects. This latest iteration of probabilistic embeddings, in which a point an object is represented as a probability measure, can consider elliptical measures (including Gaussians) with arbitrary covariance matrices. Using the  $W_2$  metric we can provide a natural and seamless generalization of point embeddings in  $\mathbb{R}^d$ . Each embedding is described with a location  $\mathbf{c}$  and a scale  $\mathbf{C}$  parameter, the latter being represented in practice using a factor matrix  $\mathbf{L}$ , where  $\mathbf{C}$  is recovered as  $\mathbf{L}\mathbf{L}^T$ . The visualization part of work is still subject to open questions. One may seek a different method than that proposed here using precision matrices, and ask whether one can include more advanced constraints on these embeddings, such as inclusions or the presence (or absence) of intersections across ellipses. Handling multimodality using mixtures of Gaussians could be pursued. In that case a natural upper bound on the  $W_2$  distance can be computed by solving the OT problem between these mixtures of Gaussians using a simpler proxy: consider them as discrete measures putting Dirac masses in the space of Gaussians endowed with the  $W_2$  metric as a ground cost, and use the optimal cost of that proxy as an upper bound of their Wasserstein distance. Finally, note that the set of elliptical measures  $\mu_{\mathbf{c}, \mathbf{C}}$  endowed with the Bures metric can also be interpreted, given that  $\mathbf{C} = \mathbf{L}\mathbf{L}^T$ ,  $\mathbf{L} \in \mathbb{R}^{d \times k}$ , and writing  $\tilde{\mathbf{l}}_i = \mathbf{l}_i - \bar{\mathbf{l}}$  for the centered column vectors of  $\mathbf{L}$ , as a discrete point cloud  $(\mathbf{c} + \frac{1}{\sqrt{k}} \tilde{\mathbf{l}}_i)_i$  endowed with a  $W_2$  metric only looking at their first and second order moments. These  $k$  points, whose mean and covariance matrix match  $\mathbf{c}$  and  $\mathbf{C}$ , can therefore fully characterize the geometric properties of the distribution  $\mu_{\mathbf{c}, \mathbf{C}}$ , and may provide a simple form of multimodal embedding.

Table 2: Entailment benchmark: we evaluate our embeddings on the Entailment dataset using average precision (AP) and F1 scores. The threshold for F1 is chosen to be the best at test time.

Model	AP	F1
W2G/45/Cosine	0.70	0.74
W2G/45/KL	0.72	0.74
Ell/12/CM	0.70	0.73

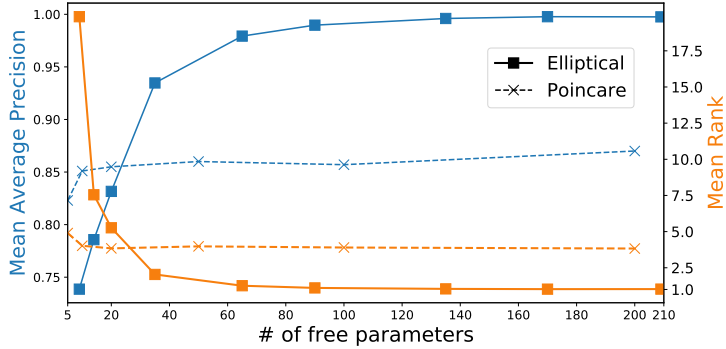


Figure 6: Reconstruction performance of our embeddings against Poincare embeddings (reported from [Nickel and Kiela, 2017], as we were not able to reproduce scores comparable to these values) evaluated by mean retrieved rank (lower=better) and MAP (higher=better).

## References

- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- L. Ambrosio, N. Gigli, and G. Savaré. *Gradient flows in metric spaces and in the space of probability measures*. Springer, 2006.
- Marco Baroni, Silvia Bernardini, Adriano Ferraresi, and Eros Zanchetta. The wacky wide web: a collection of very large linguistically processed web-crawled corpora. *Language Resources and Evaluation*, 43(3): 209–226, September 2009.
- Marco Baroni, Raffaella Bernardi, Ngoc-Quynh Do, and Chung-chieh Shan. Entailment above the word level in distributional semantics. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 23–32. ACL, 2012.
- Rajendra Bhatia, Tanvi Jain, and Yongdo Lim. On the Bures-Wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 2018.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- Ingwer Borg and Patrick JF Groenen. *Modern multidimensional scaling: Theory and applications*. Springer Science & Business Media, 2005.
- Jean Bourgain. On Lipschitz embedding of finite metric spaces in Hilbert space. *Israel Journal of Mathematics*, 52(1):46–52, 1985.
- Yann Brenier. Décomposition polaire et réarrangement monotone des champs de vecteurs. *CR Acad. Sci. Paris Sér. I Math*, 305(19):805–808, 1987.
- Alexander M Bronstein, Michael M Bronstein, and Ron Kimmel. Generalized multidimensional scaling: a framework for isometry-invariant partial surface matching. *Proceedings of the National Academy of Sciences*, 103(5):1168–1172, 2006.
- Elia Bruni, Nam Khanh Tran, and Marco Baroni. Multimodal distributional semantics. *J. Artif. Int. Res.*, 49(1): 1–47, January 2014.
- Mihai Bădoiu, Piotr Indyk, and Anastasios Sidiropoulos. Approximation algorithms for embedding general metrics into trees. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 512–521. Society for Industrial and Applied Mathematics, 2007.
- Donald Bures. An extension of Kakutani’s theorem on infinite product measures to the tensor product of semifinite  $w^*$ -algebras. *Transactions of the American Mathematical Society*, 135:199–212, 1969.
- Stamatis Cambanis, Steel Huang, and Gordon Simons. On the theory of elliptically contoured distributions. *Journal of Multivariate Analysis*, 11(3):368 – 385, 1981.
- Jan De Leeuw. Applications of convex analysis to multidimensional scaling. In *Recent Developments in Statistics*, 1977.
- DC Dowson and BV Landau. The Fréchet distance between multivariate normal distributions. *Journal of multivariate analysis*, 12(3):450–455, 1982.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- Jittat Fakcharoenphol, Satish Rao, and Kunal Talwar. A tight bound on approximating arbitrary metrics by tree metrics. In *Proceedings of the thirty-fifth annual ACM symposium on Theory of computing*, pages 448–455. ACM, 2003.
- KT Fang, S Kotz, and KW Ng. *Symmetric Multivariate and Related Distributions*. Chapman and Hall/CRC, 1990.
- Lev Finkelstein, Evgeniy Gabrilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppín. Placing search in context: the concept revisited. *ACM Trans. Inf. Syst.*, 20(1):116–131, 2002.
- Matthias Gelbrich. On a formula for the l2 Wasserstein metric between measures on Euclidean and Hilbert spaces. *Mathematische Nachrichten*, 147(1):185–203, 1990.

- Amir Globerson, Gal Chechik, Fernando Pereira, and Naftali Tishby. Euclidean embedding of co-occurrence data. *Journal of Machine Learning Research*, 8(Oct):2265–2295, 2007.
- Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864. ACM, 2016.
- Guy Halawi, Gideon Dror, Evgeniy Gabrilovich, and Yehuda Koren. Large-scale learning of word relatedness with constraints. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '12, pages 1406–1414, New York, NY, USA, 2012. ACM.
- Nicholas J. Higham. *Functions of Matrices: Theory and Computation (Other Titles in Applied Mathematics)*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2008.
- Felix Hill, Roi Reichart, and Anna Korhonen. Simlex-999: Evaluating semantic models with genuine similarity estimation. *Comput. Linguist.*, 41(4):665–695, December 2015.
- Geoffrey E Hinton and Sam T Roweis. Stochastic neighbor embedding. In *Advances in Neural Information Processing Systems*, pages 857–864, 2003.
- Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- Tony Jebara, Risi Kondor, and Andrew Howard. Probability product kernels. *Journal of Machine Learning Research*, 5:819–844, 2004.
- William B. Johnson and Joram Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. In *Conference in modern analysis and probability (New Haven, Conn., 1982)*, volume 26 of *Contemp. Math.*, pages 189–206. Amer. Math. Soc., Providence, RI, 1984.
- Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *Proceedings of the International Conference on Learning Representations*, 2014.
- J.M. Lee. *Riemannian Manifolds: An Introduction to Curvature*. Graduate Texts in Mathematics. Springer New York, 1997.
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.
- Luigi Malagò, Luigi Montrucchio, and Giovanni Pistone. Wasserstein-Riemannian geometry of positive-definite matrices. *arXiv preprint arXiv:1801.09269*, 2018.
- Yariv Maron, Michael Lamar, and Elie Bienenstock. Sphere embedding: An application to part-of-speech induction. In *Advances in Neural Information Processing Systems*, pages 1567–1575, 2010.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *ICLR Workshop*, 2013a.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119, 2013b.
- George A. Miller. Wordnet: A lexical database for english. *Commun. ACM*, 38(11):39–41, November 1995.
- George A. Miller and Walter G. Charles. Contextual correlates of semantic similarity. *Language and Cognitive Processes*, 6(1):1–28, 1991.
- Maximillian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6341–6350. Curran Associates, Inc., 2017.
- Mohammad Norouzi, Tomas Mikolov, Samy Bengio, Yoram Singer, Jonathon Shlens, Andrea Frome, Greg Corrado, and Jeffrey Dean. Zero-shot learning by convex combination of semantic embeddings. In *Proceedings of the International Conference on Learning Representations*, 2014.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.

- Gabriel Peyré and Marco Cuturi. Computational optimal transport. *arXiv preprint arXiv:1803.00567*, 2018.
- Kira Radinsky, Eugene Agichtein, Evgeniy Gabrilovich, and Shaul Markovitch. A word at a time: Computing word relatedness using temporal semantic analysis. In *Proceedings of the 20th International Conference on World Wide Web, WWW '11*, pages 337–346, New York, NY, USA, 2011. ACM.
- Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000.
- Herbert Rubenstein and John B. Goodenough. Contextual correlates of synonymy. *Commun. ACM*, 8(10): 627–633, October 1965.
- Filippo Santambrogio. *Optimal Transport for Applied Mathematicians*. Birkhauser, 2015.
- Sidak Pal Singh, Andreas Hug, Aymeric Dieuleveut, and Martin Jaggi. Context mover’s distance & barycenters: Optimal transport of contexts for building representations. *arXiv preprint arXiv:1808.09663*, 2018.
- Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- Minh thang Luong, Richard Socher, and Christopher D. Manning. Better word representations with recursive neural networks for morphology. In *Proceedings of the Thirteenth Annual Conference on Natural Language Learning*. Tomas Mikolov, Wen-tau, 2013.
- Nikolai G Ushakov. *Selected topics in characteristic functions*. Walter de Gruyter, 1999.
- Luke Vilnis and Andrew McCallum. Word representations via Gaussian embedding. *Proceedings of the International Conference on Learning Representations*, 2015. arXiv preprint arXiv:1412.6623.
- K.Q. Weinberger and L.K. Saul. Distance metric learning for large margin nearest neighbor classification. *The Journal of Machine Learning Research*, 10:207–244, 2009.
- Dongqiang Yang and David M. W. Powers. Measuring semantic similarity in the taxonomy of wordnet. In *Proceedings of the Twenty-eighth Australasian Conference on Computer Science - Volume 38, ACSC '05*, pages 315–322, Darlinghurst, Australia, Australia, 2005. Australian Computer Society, Inc.

# Supplementary Material

## Equivalent formulations of $\mathbf{T}^{\mathbf{AB}}$

$\mathbf{T}^{\mathbf{AB}}$  is defined as the unique PSD matrix verifying  $\mathbf{T}^{\mathbf{AB}}\mathbf{A}\mathbf{T}^{\mathbf{AB}} = \mathbf{B}$ . Using this definition, we derive two equivalent formulations for  $\mathbf{T}^{\mathbf{AB}}$  :

$$\begin{aligned}\mathbf{T}^{\mathbf{AB}} &= \mathbf{A}^{-\frac{1}{2}} \left( \mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} \right)^{\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}} \\ &= \mathbf{B}^{\frac{1}{2}} \left( \mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} \right)^{-\frac{1}{2}} \mathbf{B}^{\frac{1}{2}}\end{aligned}$$

The first is derived as in [Malagò et al., 2018]:

$$\begin{aligned}\mathbf{T}^{\mathbf{AB}}\mathbf{A}\mathbf{T}^{\mathbf{AB}} &= \mathbf{B} \\ \mathbf{A}^{\frac{1}{2}}\mathbf{T}^{\mathbf{AB}}\mathbf{A}^{\frac{1}{2}}\mathbf{A}^{\frac{1}{2}}\mathbf{T}^{\mathbf{AB}}\mathbf{A}^{\frac{1}{2}} &= \mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}} \\ \mathbf{A}^{\frac{1}{2}}\mathbf{T}^{\mathbf{AB}}\mathbf{A}^{\frac{1}{2}} &= \left( \mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}} \right)^{\frac{1}{2}} \\ \mathbf{T}^{\mathbf{AB}} &= \mathbf{A}^{-\frac{1}{2}} \left( \mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}} \right)^{\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}}\end{aligned}$$

We then adapt this derivation to obtain a second formulation of  $\mathbf{T}^{\mathbf{AB}}$ :

$$\begin{aligned}\mathbf{T}^{\mathbf{AB}}\mathbf{A}\mathbf{T}^{\mathbf{AB}} &= \mathbf{B} \\ (\mathbf{T}^{\mathbf{AB}})^{-1}\mathbf{B}(\mathbf{T}^{\mathbf{AB}})^{-1} &= \mathbf{A} \\ \mathbf{B}^{\frac{1}{2}}(\mathbf{T}^{\mathbf{AB}})^{-1}\mathbf{B}^{\frac{1}{2}}\mathbf{B}^{\frac{1}{2}}(\mathbf{T}^{\mathbf{AB}})^{-1}\mathbf{B}^{\frac{1}{2}} &= \mathbf{B}^{\frac{1}{2}}\mathbf{A}\mathbf{B}^{\frac{1}{2}} \\ \mathbf{B}^{\frac{1}{2}}(\mathbf{T}^{\mathbf{AB}})^{-1}\mathbf{B}^{\frac{1}{2}} &= \left( \mathbf{B}^{\frac{1}{2}}\mathbf{A}\mathbf{B}^{\frac{1}{2}} \right)^{\frac{1}{2}} \\ (\mathbf{T}^{\mathbf{AB}})^{-1} &= \mathbf{B}^{-\frac{1}{2}} \left( \mathbf{B}^{\frac{1}{2}}\mathbf{A}\mathbf{B}^{\frac{1}{2}} \right)^{\frac{1}{2}} \mathbf{B}^{-\frac{1}{2}} \\ \mathbf{T}^{\mathbf{AB}} &= \mathbf{B}^{\frac{1}{2}} \left( \mathbf{B}^{\frac{1}{2}}\mathbf{A}\mathbf{B}^{\frac{1}{2}} \right)^{-\frac{1}{2}} \mathbf{B}^{\frac{1}{2}}\end{aligned}$$

## Derivation of the Riemannian gradient updates

From [Malagò et al., 2018], we have that the exp and log maps of the Riemannian Bures metric are given by:

$$\begin{aligned}\exp_{\mathbf{C}}(\mathbf{V}) &= (\mathcal{L}_{\mathbf{C}}(\mathbf{V}) + \mathbf{I}) \mathbf{C} (\mathcal{L}_{\mathbf{C}}(\mathbf{V}) + \mathbf{I}) \\ \log_{\mathbf{C}}(\mathbf{B}) &= (\mathbf{T}^{\mathbf{CB}} - \mathbf{I}) \mathbf{C} + \mathbf{C} (\mathbf{T}^{\mathbf{CB}} - \mathbf{I})\end{aligned}$$

where  $\mathcal{L}_{\mathbf{C}}(\mathbf{V})$  is the solution of *Lyapunov* equation  $\mathcal{L}_{\mathbf{C}}(\mathbf{V})\mathbf{C} + \mathbf{C}\mathcal{L}_{\mathbf{C}}(\mathbf{V}) = \mathbf{V}$ . One can show that the  $\mathcal{L}_{\mathbf{C}}$  operator is linear, and that the following identity holds:  $\mathcal{L}_{\mathbf{C}}(\mathbf{XC} + \mathbf{CX})$ . In particular,  $\mathcal{L}_{\mathbf{C}}(\log_{\mathbf{C}} \mathbf{B}) = \mathbf{T}^{\mathbf{CB}} - \mathbf{I}$ .

From this, since  $\text{grad}_{\mathbf{A}} \frac{1}{2} \mathfrak{B}^2(\mathbf{A}, \mathbf{B}) = -\log_{\mathbf{A}} \mathbf{B}$ , the Riemannian gradient update is given by

$$\begin{aligned} \mathbf{A}_{t+1} &= \exp_{\mathbf{A}_t}(\eta_t \log_{\mathbf{A}_t} \mathbf{B}) \\ &= (\eta_t \mathcal{L}_{\mathbf{A}_t}(\log_{\mathbf{A}_t} \mathbf{B}) + \mathbf{I}) \mathbf{A}_t (\eta_t \mathcal{L}_{\mathbf{A}_t}(\log_{\mathbf{A}_t} \mathbf{B}) + \mathbf{I}) \\ &= ((1 - \eta_t)\mathbf{I} + \eta_t \mathbf{T}^{\mathbf{A}_t \mathbf{B}}) \mathbf{A}_t ((1 - \eta_t)\mathbf{I} + \eta_t \mathbf{T}^{\mathbf{A}_t \mathbf{B}}) \end{aligned}$$

## Derivation of the Euclidean gradient

**Notations:**  $\otimes$  is the Kronecker product of matrices. Recall that

$$\begin{aligned} [\mathbf{B}^\top \otimes \mathbf{A}] \text{vec}(\mathbf{X}) &= \text{vec}(\mathbf{A}\mathbf{X}\mathbf{B}) \\ [\mathbf{A} \otimes \mathbf{B}][\mathbf{C} \otimes \mathbf{D}] &= [\mathbf{A}\mathbf{C} \otimes \mathbf{B}\mathbf{D}] \end{aligned}$$

In the following, we will often omit the  $\text{vec}(\cdot)$  and treat matrices as vectors when the context makes it clear. We will make use of the following identities:

$$\begin{aligned} \partial_{\mathbf{X}} f \circ g(\mathbf{X}) &= \partial_{\mathbf{X}} f(g(\mathbf{X})) \partial_{\mathbf{X}} g(\mathbf{X}) \\ \partial_{\mathbf{X}} (fg)(\mathbf{X}) &= [g(\mathbf{X})^\top \otimes \mathbf{I}] \partial_{\mathbf{X}} f(\mathbf{X}) + [\mathbf{I} \otimes g(\mathbf{X})] \partial_{\mathbf{X}} g(\mathbf{X}) \end{aligned}$$

and

$$\partial_{\mathbf{X}} \mathbf{X}^{\frac{1}{2}} = [\mathbf{X}^{\frac{1}{2}} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{X}^{\frac{1}{2}}]^{-1}$$

$$\text{Let } f(\mathbf{A}, \mathbf{B}) = \text{Tr}(\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}}.$$

Let us differentiate  $f$  w.r.t  $\mathbf{A}$ :

$$\begin{aligned} \nabla_{\mathbf{A}} f(\mathbf{A}, \mathbf{B}) &= \left[ \partial_{\mathbf{A}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}} \right]^\top \mathbf{I} \\ &= \left[ \left[ \mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} \right]^{\frac{1}{2}} \otimes \mathbf{I} + \mathbf{I} \otimes \left[ \mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} \right]^{\frac{1}{2}} \right]^{-1} \partial_{\mathbf{A}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}} \right]^\top \mathbf{I} \\ &= \left[ \mathbf{B}^{\frac{1}{2}} \otimes \mathbf{B}^{\frac{1}{2}} \right] \left[ \left[ \mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} \right]^{\frac{1}{2}} \otimes \mathbf{I} + \mathbf{I} \otimes \left[ \mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} \right]^{\frac{1}{2}} \right]^{-1} \mathbf{I} \\ &= \left[ \mathbf{B}^{\frac{1}{2}} \otimes \mathbf{B}^{\frac{1}{2}} \right] \frac{1}{2} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{-\frac{1}{2}} \\ &= \frac{1}{2} \mathbf{B}^{\frac{1}{2}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{-\frac{1}{2}} \mathbf{B}^{\frac{1}{2}} \end{aligned}$$

$$\text{Therefore } \nabla_{\mathbf{A}} f(\mathbf{A}, \mathbf{B}) = \frac{1}{2} \mathbf{T}^{\mathbf{A}\mathbf{B}}$$

Let now  $\mathbf{A} = \mathbf{L}\mathbf{L}^\top$ , let us differentiate w.r.t  $\mathbf{L}$ :

$$\begin{aligned} \nabla_{\mathbf{L}} f(\mathbf{L}\mathbf{L}^\top, \mathbf{B}) &= \left[ \partial_{\mathbf{L}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}} \right]^\top \mathbf{I} \\ &= \partial_{\mathbf{L}} \mathbf{A}^\top \left[ \partial_{\mathbf{A}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{\frac{1}{2}} \right]^\top \mathbf{I} \\ &= [\mathbf{L}^\top \otimes \mathbf{I}] [\mathbf{I} + \mathbf{T}_{n,n}] \frac{1}{2} \mathbf{B}^{\frac{1}{2}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{-\frac{1}{2}} \mathbf{B}^{\frac{1}{2}} \\ &= \mathbf{B}^{\frac{1}{2}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{-\frac{1}{2}} \mathbf{B}^{\frac{1}{2}} \mathbf{L} \end{aligned}$$

where  $\mathbf{T}_{n,n}$  is the transposition tensor, such that  $\forall \mathbf{X} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{T}_{n,n} \text{vec}(\mathbf{X}) = \text{vec}(\mathbf{X}^\top)$ .

$$\text{Therefore } \nabla_{\mathbf{L}} f(\mathbf{L}\mathbf{L}^\top, \mathbf{B}) = \mathbf{T}^{\mathbf{A}\mathbf{B}} \mathbf{L}.$$

Using the same calculations, one can see that if  $\mathbf{A} = \mathbf{L}\mathbf{L}^\top + \varepsilon \mathbf{I}$ , then we still have

$$\nabla_{\mathbf{L}} f(\mathbf{L}\mathbf{L}^\top + \varepsilon \mathbf{I}, \mathbf{B}) = \mathbf{T}^{\mathbf{A}\mathbf{B}} \mathbf{L}$$

since  $\partial_{\mathbf{L}} [\mathbf{L}\mathbf{L}^\top + \varepsilon \mathbf{I}] = \partial_{\mathbf{L}} [\mathbf{L}\mathbf{L}^\top]$

## Model Hyperparameters and Training Details

**Word Embeddings** We train our embeddings on the concatenated ukWaC and WaCkypedia corpora [Baroni et al., 2009], consisting of about 3 billion tokens, on which we keep only the tokens appearing more than 100 times in the text after lowercasing and removal of all punctuation (for a total number of 261583 different words). We optimize 5 epoches using adagrad [Duchi et al., 2011] with  $\epsilon = 10^{-8}$  with a learning rate of 0.01. We use a window size of 10 (i.e. positive examples consist of the first 5 preceding and first 5 succeeding words), set the margin to 10, sample one negative context per positive context and, in order to prevent the norms of the embeddings to be too highly correlated with the corresponding word frequencies (see Figure 7), we use two distinct sets of embeddings for the input and context words. In order to use as much parallelization as possible, we use batches of size 10000, but believe that smaller batches would lead to improved performances. We limit matrix square root approximations to 6 Newton-Schulz iterations and add  $0.01\mathbf{I}$  to the covariances to ensure non-singularity.

To generate batches, we use the same sampling tricks as in [Mikolov et al., 2013b], namely sub-sampling the frequent terms (using a threshold of  $10^{-5}$  as recommended for large datasets) and smoothing the negative distribution by using probabilities  $\{f_i^{3/4}/Z\}$  where  $f_i$  is the frequency of word  $i$  for sampling negative contexts  $\{c'_i\}$ .

We then evaluate our embeddings on the following datasets: Simlex [Hill et al., 2015], WordSim [Finkelstein et al., 2002], MEN [Bruni et al., 2014], MC [Miller and Charles, 1991], RG [Rubenstein and Goodenough, 1965], YP [Yang and Powers, 2005], MTurk [Radinsky et al., 2011] [Halawi et al., 2012], RW [thang Luong et al., 2013], using the context embeddings and the Wasserstein-Bures cosine as a similarity measure.

**Hypernymy** We train our embeddings on the transitive closure of the WORDNET dataset [Miller, 1995] which features 743,241 hypernymy relations on 82,115 distinct nouns. For disambiguation, note that if  $(u, v)$  is a hypernymy relation with  $u \neq v$ ,  $(v, u)$  is in general *not* a positive relation, but  $(u, u)$  is as a noun is always its own hypernym.

We perform our optimization using SGD with batches of 1000 relations, a learning rate 0.02 for dimensions 3 and 4 and 0.01 for higher dimensions, sample 50 negative examples per positive relation, use 6 square root iterations and add  $0.01\mathbf{I}$  to the covariances. Contrary to the skipgram experiment, we use a single set of embeddings and use the Wasserstein-Bures dot product as a similarity measure.

### Wasserstein-Bures Cosine

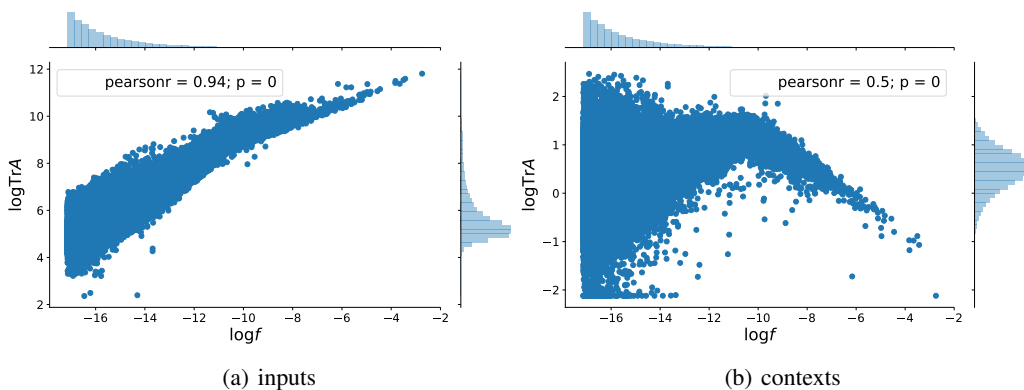


Figure 7: log-log plot of the traces of the embeddings’ covariances vs. word frequency: the sizes of the input embeddings follow a power law, whereas context embeddings give less importance to very frequent words and emphasize on medium frequency words.

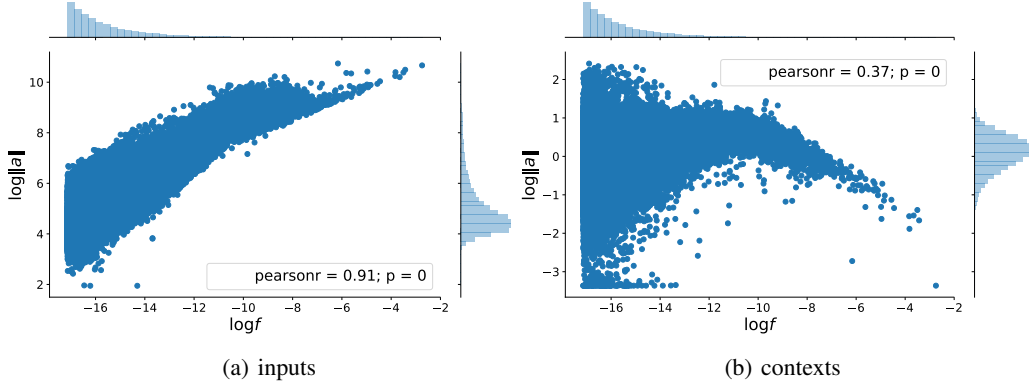


Figure 8: log-log plot of the norms of the embeddings’ means vs. word frequency: the sizes of the input embeddings follow a power law, whereas context embeddings give less importance to very frequent words and emphasize on medium frequency words.

As discussed in section 4, a natural choice of similarity measure would be the Wasserstein-Bures cosine, obtained by normalizing the Wasserstein-Bures dot product with the means’ norms and covariances’ root traces jointly:

$$\text{cos}_{\mathfrak{B}}[\rho_{\mathbf{a},\mathbf{A}}, \rho_{\mathbf{b},\mathbf{B}}] := \frac{\langle \mathbf{a}, \mathbf{b} \rangle + \text{Tr}[\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}}]^{\frac{1}{2}}}{(\|\mathbf{a}\|^2 + \text{Tr} \mathbf{A})^{\frac{1}{2}} (\|\mathbf{b}\|^2 + \text{Tr} \mathbf{B})^{\frac{1}{2}}}$$

However, we have found that in some applications (and notably in our skipgram experiments) such a joint normalization can result in either the means or the covariances to have a negligible contribution if the scales of the parameters differ too much. To circumvent this problem, we introduce another similarity measure, which is a mixture of two cosine terms:

$$\mathfrak{S}_{\mathfrak{B}}[\rho_{\mathbf{a},\mathbf{A}}, \rho_{\mathbf{b},\mathbf{B}}] := \frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\|\mathbf{a}\| \|\mathbf{b}\|} + \frac{\text{Tr}[\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}}]^{\frac{1}{2}}}{\sqrt{\text{Tr} \mathbf{A} \text{Tr} \mathbf{B}}}$$

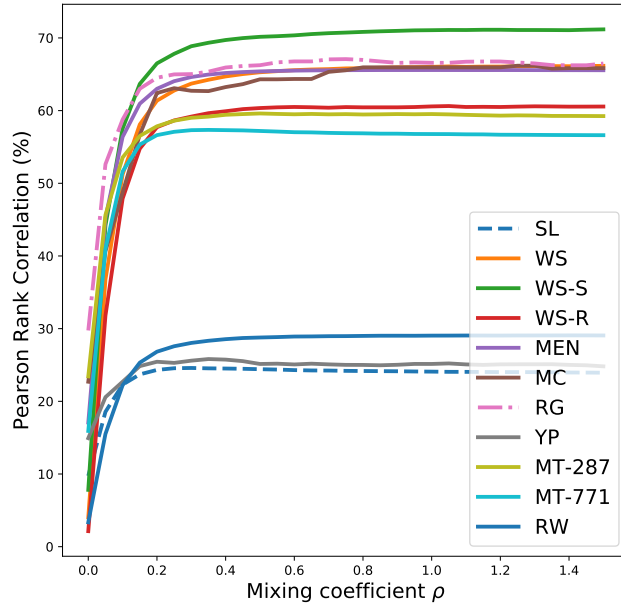
This latter similarity measure allows to gather information from the means and the covariances independently. Note that while the term corresponding to the covariances is obtained in a cosine-like normalization, it takes values between 0 and 1 as it only involve traces of PSD matrices, whereas the means term is a regular Euclidean cosine and therefore takes values between -1 and 1. We compare the behaviors of these two measures on the word similarity evaluation task by introducing a mixing coefficient  $\rho$ , and defining

$$\text{cos}_{\mathfrak{B}}[\rho_{\mathbf{a},\mathbf{A}}, \rho_{\mathbf{b},\mathbf{B}}; \rho] := \frac{\langle \mathbf{a}, \mathbf{b} \rangle + \rho \text{Tr}[\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}}]^{\frac{1}{2}}}{(\|\mathbf{a}\|^2 + \rho \text{Tr} \mathbf{A})^{\frac{1}{2}} (\|\mathbf{b}\|^2 + \rho \text{Tr} \mathbf{B})^{\frac{1}{2}}}$$

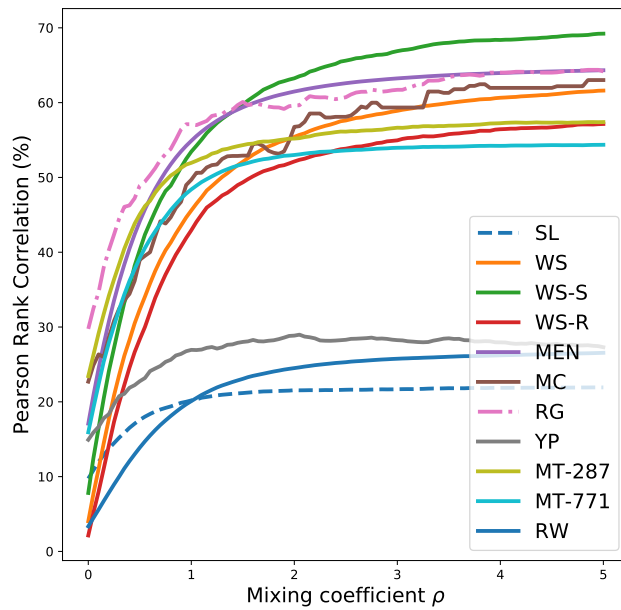
$$\mathfrak{S}_{\mathfrak{B}}[\rho_{\mathbf{a},\mathbf{A}}, \rho_{\mathbf{b},\mathbf{B}}; \rho] := \frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\|\mathbf{a}\| \|\mathbf{b}\|} + \rho \frac{\text{Tr}[\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}}]^{\frac{1}{2}}}{\sqrt{\text{Tr} \mathbf{A} \text{Tr} \mathbf{B}}}$$

As can be seen from figure 9, the Wasserstein-Bures cosine is less well behaved and makes it difficult to find an optimal mixing value. On the other hand, the mixture of cosines similarity measure varies more smoothly and seems to reach a performance maximum around  $\rho = 1$ , and achieves better performance than the Wasserstein-Bures cosine on most datasets.





(a)  $\mathcal{S}_{23}$



(b)  $\cos_{23}$

Figure 9: Pearson rank correlation scores on similarity benchmarks as a function of the mixing coefficient:  $\mathcal{S}_{23}$  smoothly attains a maximum in performance around  $\rho = 1$ , whereas  $\cos_{23}$  has a not so smooth behavior.