



Published in final edited form as:

*J Math Imaging Vis.* 2017 October ; 59(2): 187–210. doi:10.1007/s10851-017-0726-4.

## A Transportation $L^p$ Distance for Signal Analysis

Matthew Thorpe<sup>1</sup>, Serim Park<sup>1</sup>, Soheil Kolouri<sup>3</sup>, Gustavo K. Rohde<sup>2</sup>, and Dejan Slepčev<sup>1</sup>

<sup>1</sup>Carnegie Mellon University, Pittsburgh, PA 15213, USA

<sup>2</sup>University of Virginia, Charlottesville, VA 22908, USA

<sup>3</sup>HRL Laboratories, Malibu, CA 90265, USA

### Abstract

Transport based distances, such as the Wasserstein distance and earth mover's distance, have been shown to be an effective tool in signal and image analysis. The success of transport based distances is in part due to their Lagrangian nature which allows it to capture the important variations in many signal classes. However these distances require the signal to be nonnegative and normalized. Furthermore, the signals are considered as measures and compared by redistributing (transporting) them, which does not directly take into account the signal intensity. Here we study a transport-based distance, called the  $TL^p$  distance, that combines Lagrangian and intensity modelling and is directly applicable to general, non-positive and multi-channelled signals. The distance can be computed by existing numerical methods. We give an overview of the basic properties of this distance and applications to classification, with multi-channelled non-positive one-dimensional signals and two-dimensional images, and color transfer.

### 1 Introduction

Enabled by advances in numerical implementation [4, 5, 16, 53, 71, 74], and their Lagrangian nature, transportation based distances for signal analysis are becoming increasingly popular in a large range of applications. Recent applications include astronomy [9, 18, 19], biomedical sciences [3, 25–27, 77, 81, 82, 88, 89], colour transfer [14, 17, 49, 62, 63], computer vision and graphics [7, 44, 60, 65, 68, 74, 75], imaging [36, 40, 64], information theory [78], machine learning [1, 15, 20, 34, 37, 48, 76], operational research [69] and signal processing [54, 58].

The success of transport based distances is due to the large number of applications that consider signals that are Lagrangian in nature (spatial rearrangements, i.e. transport, are a key factor when considering image differences). Many signals contain similar features for which transport based distances will outperform distances that only consider differences in intensity, such as the  $L^p$  distance. Optimal transport (OT) distances, for example the earth mover's distance or Wasserstein distance, are examples of transport distances. However these distances do not directly account for signal intensity. The  $L^p$  distance is the other extreme, this distance is based on intensity and does not take into account Lagrangian properties.

In this paper we develop the  $TL^p$  distance introduced in [21] which combines both Lagrangian and intensity based modeling. Our aim is to show that by including both transport and intensity within the distance we can better represent the similarities between

classes of data in many problems. For example, if a distance can naturally differentiate between classes, that is the within class distance is small compared to the between class separation, then the classification problem is made easier. This requires designing distances that can faithfully represent the structure within a given data set.

In the majority of the literature optimal transport distances interpret signals as either probability measures or as densities of probability measures. This places restrictions on the type of signals one can consider. Probability measures must be non-negative, integrate to unity and be real valued (i.e. cannot be applied to multi-channelled signals). In order to apply OT distances to a wider class of signals one has to use ad-hoc methods, which do not necessarily preserve metric properties, to transform the signal into a probability measure. This can often dampen the features, for example renormalization may reduce the intensity range of a signal. We do however note the works of Liero, Mielke and Savaré [42, 43], Chizat, Peyré, Schmitzer and Vialard [13, 14] and Kondratyev, Monsaingeon and Vorotnikov [38] who develop an optimal transport metric that is applicable to un-normalised positive measures. Similarly Pele and Wermen propose a variant of the earth movers distance that is applicable to un-normalised positive measures [59]. Whilst these are also promising avenues research there are still restrictive assumptions, such as signals must be non-negative and real valued, required in order to apply these distances.

Extensions to matrix valued optimal transport have been made in [11, 12, 52]. In [52] the authors propose a method for defining an optimal transport distance between matrix valued densities. As in the scalar valued case for a suitable class of matrix valued densities there exists a non-empty set of couplings. A matrix valued optimal transport distance is defined by minimising over the set of couplings a cost function that penalises both the transport of mass and rotations of the coupling. Similar ideas in [11, 12] use an analogue of the Brenier and Benamou formulation of optimal transport [4] to define a matrix valued optimal transport distance. Whilst these distances are applicable to matrix valued signals they still require the assumption of positivity (in this case positive definite). These distances are also specific to  $n \times n$  matrices which does not include vectors.

The  $TLP$  distance does not need the signal to be a probability measure and therefore the above restrictions do not apply. Rather, the  $TLP$  distance models the intensity directly. The applicability of the distance is sufficiently general as to include non-positive, multi-channelled and un-normalised signals on discrete or continuous domains.

Another property of OT distances, due to the lack of intensity modeling, is its insensitivity to high frequency perturbations. This is due to transport being on the order of the wavelength of the perturbation. Depending upon the application this can be an advantage or a disadvantage. For example in texture modeling one would want to be able to discriminate between a highly oscillating image and a constant image. On the other hand, the lack of sensitivity to high frequency noise, make the OT distance stable under such perturbations. Since the  $TLP$  distance directly models intensity then it inherits sensitivity to high frequency noise.

The aim of this paper is to develop the  $TL^p$  distance and demonstrate its applicability in a range of applications. We consider classification problems on data sets where we show that the  $TL^p$  distance better represents the underlying geometry, i.e. achieves a better between class to within class distance, than popular alternative distances.

We also consider the colour transfer problem in a context where spatial information, as well as intensity, is important. To apply standardised tests in applications such as medical imaging it is often necessary to normalise colour variation [33, 45, 73]. One solution is to match the means and variance of each colour channel (in some colour space e.g. RGB or LAB). However, by transferring the colour of one image onto the other it is possible to recolour an image with *exactly* the same colour profile.

A popular method is to use the OT distance on the histogram of images [14, 17, 49, 62, 63]. This allows one to take into account the intensity of pixels but includes no spatial information. The  $TL^p$  distance is able to include both spatial and intensity information.

Our methodology, therefore, has more in common with registration methods that aim to find a transformation that maximizes the similarity between two images where our measure of similarity includes both spatial and intensity information. One should compare our approach to [27] where the authors develop a numerical method for the Monge formulation of OT with the addition of an intensity term for image warping and registration. However, unlike the method presented in [27], the formulation presented here defines a metric.

## Paper Overview

The outline for this paper is the following. In the next section we review OT and give a formal definition of the  $TL^p$  distance followed by examples to illustrate its features and to compare with the OT and  $L^p$  distances. In Section 3 we give a more general definition and explain some of its key properties. In Section 4 we include applications of the  $TL^p$  distances. We first consider classification on synthetic and real- world signals and images. The data sets contain non-positive and un-normalised signals in either one or two dimensions. In addition one of the data sets is multi-channelled. A further application to the colour transfer problem is then given. Conclusions are given in Section 5.

## 2 Formal Definitions and Examples

### 2.1 Review of Optimal Transport and the $TL^p$ Distance

We begin by reviewing optimal transport in first the Kantorovich formulation and then the Monge formulation.

**The Kantorovich Formulation of Optimal Transport**—For measures  $\mu$  and  $\nu$  on  $\Omega \subset \mathbb{R}^d$  with the same mass and a continuous cost function  $c: \Omega \times \Omega \rightarrow [0, \infty)$  the Kantorovich formulation of OT is given by

$$\text{OT}(\mu, \nu) = \min_{\pi} \int_{\Omega \times \Omega} c(x, y) d\pi(x, y) \quad (1)$$

where the minimum is taken over probability measures  $\pi$  on  $\Omega \times \Omega$  such that the first marginal is  $\mu$  and the second marginal is  $\nu$ , i.e.  $\pi(A \times \Omega) = \mu(A)$  and  $\pi(\Omega \times B) = \nu(B)$  for all open sets  $A$  and  $B$ . We denote the set of such  $\pi$  by  $\Pi(\mu, \nu)$ . We call measures  $\pi \in \Pi(\mu, \nu)$  transport plans since  $\pi(A \times B)$  is the amount of mass in  $A$  that is transferred to  $B$ . Minimizers  $\pi^*$  of  $\text{OT}(\mu, \nu)$ , which we call optimal plans, exist when  $c$  is lower semi-continuous [84, Theorem 4.1]. When  $c$  is a metric  $\text{OT}(\mu, \nu)$  is also known as the earth mover's distance.

A common choice is  $c(x, y) = |x - y|^p = \sum_{i=1}^d |x_i - y_i|^p$  in which case we define

$d_{\text{OT}}(\mu, \nu) = \sqrt[p]{\text{OT}(\mu, \nu)}$ . When  $p = 2$  this is known as the Wasserstein distance. We will call  $d_{\text{OT}}$  the OT distance. With an abuse of notation we will sometimes write  $d_{\text{OT}}(f, g)$  when  $\mu$  and  $\nu$  have densities  $f$  and  $g$  respectively.

When  $\mu$  has a continuous density then the support of any optimal plan  $\pi^*$  is contained on the graph of a function  $T^*$ . In particular this implies  $\pi^*(A, B) = \mu(\{x : x \in A, T^*(x) \in B\})$  and furthermore that the optimal plan defines a mapping between  $\mu$  and  $\nu$ , see for example Figure 1a. This leads us to the Monge formulation of OT.

**The Monge Formulation of Optimal Transport**—An appealing property of optimal transport distances are their formulation in a Lagrangian setting. One can rewrite the optimal transport problem in the Monge formulation as

$$\text{OT}_M(\mu, \nu) = \inf_T \int_{\Omega \times \Omega} c(x, T(x)) d\mu(x) \quad (2)$$

where the infimum is taken over transport maps  $T: \Omega \rightarrow \Omega$  that rearrange  $\mu$  into  $\nu$ , i.e.  $\nu = T_{\#}\mu$  where we define the pushforward of  $\mu$  onto the range of  $T$  by  $T_{\#}\mu(A) = \mu(T^{-1}(A))$ , see Figure 1b. Historically the Monge formulation comes before the Kantorovich formulation; Monge formulated OT for the cost function  $c(x, y) = |x - y|$  in 1781 [47] and Kantorovich formulated OT (whilst being unaware of Monge's work) in 1942 [31]. In 1948 Kantorovich made the connection between his work and Monge's [32].

The Monge formulation is a non-convex optimization problem with nonlinear constraints. However when, for example,  $\mu$  and  $\nu$  have densities, then optimal transport maps  $T^*$  exist and give a natural interpolation between two measures. In particular when  $c(x, y) = |x - y|^p$  the map  $T_t(x) = (1 - t)x + tT^*(x)$  describes the path of particle  $x$  and furthermore the measure of  $\mu$  pushed forward by  $T_t$  is the geodesic (shortest path) between  $\mu$  and  $\nu$ . This property has had many uses in transport based morphometry applications such as biomedical

[3, 55, 80, 85], super-resolution [36] and has much in common with large deformation diffeomorphism techniques in shape analysis [23, 29].

**Optimal Transport in Signal and Image Processing**—To further motivate our development of the  $TL^p$  distance we point out some features of optimal transport important to signal and image processing. We refer to [35] and references therein for more details and a review of the subject.

Key to the success of OT is the ability to provide generative models which accurately represent various families of data distributions. The success and appeal of OT owes to (1) ability to capture well the signal variations due to spatial rearrangements (shifts, translations, transport), (2) that OT distances are theoretically well understood and have appealing features (for example the Wasserstein distance has a Riemannian structure and geodesics can be characterized), (3) efficiency and accuracy of numerical methods, (4) simplicity compared to other Lagrangian methods such as large deformation diffeomorphic metric mapping.

The Monge formulation of OT defines a mapping between images which has been used in, for example, *image registration* [24–27, 50, 82, 89] where one wishes to find a common geometric reference frame between two or more images. In addition to the properties listed above the success of OT is due to the fact that (5) the Monge problem is symmetric (i.e. if  $T$  is the optimal map from the first image to the second, then  $T^{-1}$  is the optimal map from the second image to the first) and (6) OT provides a landmark-free and parameter-free registration scheme.

We now introduce the  $TL^p$  distance in the simplest setting.

**The Transportation  $L^p$  Distance**—In this paper we use the  $TL^p$  distance (given in more generality in the next section), for functions  $f, g: \Omega \rightarrow \mathbb{R}^m$  defined by

$$d_{TL}^p(f, g) = \min_{\pi} \int_{\Omega \times \Omega} |x - y|^p + |f(x) - g(y)|^p d\pi(x, y)$$

where the minimum is taken over all probability measures  $\pi$  on  $\Omega \times \Omega$  such that both the marginals are the Lebesgue measure  $\mathcal{L}$  on  $\Omega$ , i.e.  $\pi \in (\mathcal{L}, \mathcal{L})$ . This can be understood in the following two ways.

The first is as an optimal transport distance of the Lebesgue measure with itself and cost  $c(x, y) = |x - y|^p + |f(x) - g(y)|^p$ . This observation allows one to apply existing numerical methods for OT where the effective dimension is  $d$  (recall that  $\Omega \subseteq \mathbb{R}^d$ ). For example, the Sinkhorn framework can be adapted to compute an entropy regularised approximation of the  $TL^p$  distance, see Appendix B. 2 for more details.

The second is as an OT distance between the Lebesgue measure raised onto the graphs of  $f$  and  $g$ . That is, given  $f, g: \Omega \rightarrow \mathbb{R}^m$  then we define the measures  $\tilde{\mu}, \tilde{\nu}$  on the graphs of  $f$  and  $g$  by

$$\tilde{\mu}(A \times B) = \mathcal{L}(\{x: x \in A, f(x) \in B\})$$

and

$$\tilde{\nu}(A \times B) = \mathcal{L}(\{y: y \in A, g(y) \in B\})$$

for any open sets  $A \subseteq \Omega$ ,  $B \subseteq \mathbb{R}^m$ . For example, given the function  $f: [0,1] \rightarrow [0,1]$  defined by  $f(x) = x$  and the Lebesgue measure on  $[0,1]$ , the pushforward of the measure  $\mu$  onto the graph of  $f$  is the measure  $\tilde{\mu}(C) = \frac{1}{\sqrt{2}} \text{Length}(\{x: (x, x) \in C\})$  for any (measurable)  $C \subset [0,1]^2$ ; it is intuitive that  $\tilde{\mu}(C)$  should be proportional to  $\text{Length}(\{x: (x, x) \in C\})$ , the constant of proportionality comes from  $\tilde{\mu}([0,1]^2) = \mu([0,1]) = 1$ . See also Figure 1b for an example where the measure  $\mu$  is Gaussian. The  $TL^p$  distance between  $f$  and  $g$  is then the OT distance between  $\tilde{\mu}$  and  $\tilde{\nu}$ .

Transport (i.e. matching) with respect to the  $TL^p$  distance is of the form  $(x, f(x)) \mapsto (y, g(y))$  and therefore has two components. We refer to horizontal transport as the transport  $x \mapsto y$  in  $\Omega$ , and vertical transport as the transport  $f(x) \mapsto g(y)$ .

Although the  $TL^p$  distance is a special case of OT we will, in order to make a clearer distinction between classical OT distances and the  $TL^p$  distance, assume that  $c(x, y) = |x - y|^p$  in (1).

In the next section we discuss the behaviour of the  $TL^p$  distance through three examples.

## 2.2 Examples Illustrating the Behaviour of $TL^p$

**No mass renormalization**—Unlike for OT, in the  $TL^p$  distance there is no need to assume that  $f$  and  $g$  are non-negative or that they have the same mass. If a signal is negative then a typical (ad-hoc) fix in OT is to add a constant to make the signal non-negative before computing the distance. How to choose this constant is often unclear unless a lower bound is known a-priori. Furthermore this may damage sensitivity to translations as the defining features of the signal become compressed. For example, considering the functions in Figure 2a, let  $g = f(\cdot - \ell)$  be the translation of  $f$ . OT will lose sensitivity when comparing

$$\hat{f} = \frac{f + \alpha}{f(f + \alpha)} \text{ and } \hat{g} = \frac{g + \alpha}{f(g + \alpha)}. \text{ In particular } d_{OT}(\hat{f}, \hat{g}) \text{ scales with the height of the}$$

renormalised function, which is of the order of  $\frac{1}{\alpha}$ , and the size of the shift:  $d_{OT}(\hat{f}, \hat{g}) \propto \frac{h_0 \ell}{\alpha}$

where  $h_0$  is the height of  $f$ . To ensure positivity one must choose a large  $\alpha$  but this also implies a small OT distance. Note also that both  $L^p$  and  $TL^p$  distances are invariant under adding a constant whereas OT is not.

Another approach to apply the OT distance to non-positive signals is to decompose each signal into the positive and negative components  $f = f^+ - f^-$ , where  $f^+ = \max\{0, f\}$  and  $f^- = \max\{0, -f\}$ , and compute the OT distance between each component. Whilst this method may be reasonable depending on the application it is not invariant under addition which

could produce some unnatural properties. For example, consider two signals, one slightly negative and one slightly positive. Then applying an (unbalanced) OT distance on the positive and negative components is equivalent to matching both signals to zero. Adding a small constant to the negative signal so that both signals are positive produces a qualitatively different result. Since both the  $TLP$  and  $LP$  distances are invariant under addition neither has this property. To mitigate this issue Bonneel, van de Panne, and Heidrich [8] consider signals decomposed into frequency bands. This also allows them to directly take into account the signal frequency. In some sense our approach is complimentary, as we seek a way to take into account the signal intensity.

**Sensitivity to High Frequency Perturbations**—The  $TLP$  distance inherits sensitivity to high frequency perturbations from the  $LP$  distance. For example, let  $g = f + A\xi$  where  $\xi$  is a high frequency perturbation with amplitude  $A$  and wavelength  $\omega$ . Suppose for simplicity that  $f$  is constant and  $\xi$  is a sinusoid (and that both signals are positive). The function  $g$  is the landscape consisting of piles of earth and trenches and  $f$  is the flat landscape. The OT distance between  $f$  and  $g$  measures the cost of moving piles of earth into trenches (in the most efficient manner). The two factors which determine the OT distance are the total amount of earth to be moved (which we assume fixed) and how far we move each piece of soil, which is determined by the wavelength. Hence the OT distance between  $f$  and  $g$  is on the order of the wavelength  $\omega$  of  $\xi$ , which is small, and independent of the amplitude  $A$ . On the other hand both the  $TLP$  distance and the  $LP$  distance are independent of the wavelength but scale linearly with amplitude, see Figure 2b. In particular OT is insensitive to high frequency noise regardless of the size of the amplitude whereas both  $TLP$  and  $LP$  distances scale linearly with the amplitude.

**Ability of the  $TLP$  Distance to Track Translations**—Another desirable property of both  $TLP$  and OT distances are their ability to keep track of translations for longer than the  $LP$  distance. Let  $f = A\chi_{[0,1]}$  be the indicator function of the set  $[0,1]$  on  $\mathbb{R}$  scaled by  $A > 1$  and  $g(x) = f(x - \ell)$  the translation of  $f$  by  $\ell$ . Once  $\ell > 1$  then the  $LP$  distance can no longer tell how far apart two humps are. On the other hand OT distances can track the hump indefinitely. In this example the  $TLP$  distance couples the graphs of  $f$  and  $g$  in one of three ways, see Figure 3. The first is when the transport is horizontal only in the graph (Figure 3 top left). In the second (top right) there is a mixture of horizontal and vertical transport. And in the third there is only vertical transport (bottom left), in which case the  $TLP$  distance coincides with the  $LP$  distance. One can calculate the range of the  $TLP$  distance which is on the order of  $A$ .

### 3 Definitions and Basic Properties of the $TLP$ Distance

In the previous section we defined the  $TLP$  distance for signals defined with respect to the Lebesgue measure. In this section we generalise to signals defined on a general class of measures. We let  $L^p(\mu)$  be the space of functions  $f$  such that  $\int_{\Omega} |f(x)|^p d\mu(x) < \infty$ . This is a Banach space with the usual norm.

We treat a signal as a pair  $(f, \mu)$  where  $\mu \in \mathcal{P}_p(\Omega)$  (the set of probability measures with finite  $p^{\text{th}}$  moment) and  $f: \Omega \rightarrow \mathbb{R}^m$  with  $f \in L^p(\mu)$ . The generality considered here allows us to

treat continuous and discrete signals simultaneously as well as allowing one to design the underlying measure in order to emphasise certain parts of the signal. We are also able to compare signals with different discretisations. However, unless otherwise stated,  $\mu = \nu$  is the Lebesgue measure. There is no assumption on the dimension  $m$  of the codomain. This allows us to consider multi-channelled signals.

The  $TL_\lambda^p$  distance for pairs  $(f, \mu) \in TLP$  where

$$TLP := \{(f, \mu) : f \in L^p(\mu), \mu \in \mathcal{P}_p(\Omega)\}$$

is defined by

$$d_{TL_\lambda^p}^p((f, \mu), (g, \nu)) = \min_{\pi \in \Pi(\mu, \nu)} \int_{\Omega \times \Omega} c_\lambda(x, y; f, g) d\pi \quad (3)$$

$$c_\lambda(x, y; f, g) = \frac{1}{\lambda} |x - y|^p + |f(x) - g(y)|^p \quad (4)$$

and  $\Pi(\mu, \nu)$  is the space of measures on  $\Omega \times \Omega$  such that the first marginal is  $\mu$  and the second marginal is  $\nu$ . Note that if  $f = g$  is constant then we recover the OT distance between the measures  $\mu$  and  $\nu$ . In the special cases, when  $\mu = \nu = \mathcal{L}$  are the Lebesgue measure, we write  $d_{TL_\lambda^p}(f, g) := d_{TL_\lambda^p}^p((f, \mathcal{L}), (g, \mathcal{L}))$  and, when  $\lambda = 1, \lambda = 1, d_{TL^p}(f, g) := d_{TL_1^p}(f, g)$ . The result of [21, Proposition 3.3] implies that  $d_{TL_\lambda^p}$  is a metric on  $TLP$ .

**Proposition 3.1.** [21] For any  $p \in [1, \infty]$  and  $\lambda > 0$ ,  $(TLP, d_{TL_\lambda^p})$  is a metric space.

When  $\mu = \nu = \mathcal{L}$  is the Lebesgue measure then an admissible plan is the identity plan:  $\pi (A \times B) = \mathcal{L} (A \cap B)$ . This implies that the  $TL_\lambda^p$  distance is bounded above by the  $L^p$  distance (for any  $\lambda$ ).

In fact the parameter  $\lambda$  controls how close the distance is to an  $L^p$  distance. As  $\lambda \rightarrow 0$  then the cost of horizontal transport:  $\frac{1}{\lambda} \int_{\Omega \times \Omega} |x - y|^p d\pi(x, y)$ , is very expensive which favours transport plans that are approximately the identity mapping. Hence  $d_{TL_0^p}(f, g) := \lim_{\lambda \rightarrow 0} d_{TL_\lambda^p}(f, g) = \|f - g\|_{L^p}$ . The following result, and the remainder of the results in this section, can be found in [79].



**Proposition 3.2.** [79] Let  $f, g \in L^p$  (with respect to the Lebesgue measure). The  $TL_\lambda^p$  distance is decreasing as a function of  $\lambda$  and

$$\lim_{\lambda \rightarrow 0} d_{TL_\lambda^p}(f, g) = \|f - g\|_{L^p}.$$

Moreover, if either  $f$  or  $g$  is Lipschitz then

$$d_{TL_\lambda^p}^p(f, g) \geq \begin{cases} \varepsilon^{p-1}(\lambda) \|f - g\|_{L^p}^p & \text{if } p > 1 \\ \|f - g\|_{L^p}^p & \text{if } p = 1 \text{ and } \lambda < \frac{1}{\kappa} \end{cases}$$

where  $\varepsilon(\lambda) = \frac{1}{1 + (\lambda\kappa)^{p-1}}$  and  $\kappa = (\min\{\text{Lip}(f), \text{Lip}(g)\})^p$ .

The above proposition implies that, when  $p = 1$ , if  $\frac{1}{\lambda}$  is chosen larger than the length scale given by the derivative then the  $TL_\lambda^1$  distance is exactly the  $L^1$  distance.

Recall that we can consider the  $TL_\lambda^p$  distance as an OT distance on the graphs of  $f$  and  $g$ :

$$d_{TL_\lambda^p}((f, \mu), (g, \nu)) = d_{\text{OT}}((\text{Id} \times f)_\# \mu, (\text{Id} \times g)_\# \nu). \quad (5)$$

When there exists a map  $T: \Omega \times \mathbb{R}^m \rightarrow \Omega \times \mathbb{R}^m$  realising the minimum of the Monge formulation of the RHS then we can understand the transport as a map  $(x, f(x)) \mapsto (y, g(y))$ . We recall that we refer to the transport  $x \mapsto y$  in the domain  $\Omega$  as horizontal transport and transport  $f(x) \mapsto g(y)$  in the codomain of  $f$  and  $g$  as vertical transport. We see that horizontal transport is cheap as  $\lambda \mapsto \infty$  and we only pay the cost of vertical transport. For example, if we consider  $f(x) = \chi_{[0,1]}$  and  $g(x) = \chi_{[1,2]}$  defined on the interval  $[0, 2]$  then the mapping  $T(x, y) = (T_1(x, y), T_2(x, y))$  where

$$T_1(x, y) = \begin{cases} x + 1 & \text{if } y = 1 \text{ and } x \in [0, 1] \\ x - 1 & \text{if } y = 0 \text{ and } x \in [0, 1] \end{cases}$$

and

$$T_2(x, y) = \begin{cases} 1 & \text{if } y = 1 \text{ and } x \in [0, 1] \\ 0 & \text{if } y = 0 \text{ and } x \in [0, 1] \end{cases}$$

defines a transport map on the support of  $(\text{Id} \times f)_{\#}\mathcal{L}$ . Furthermore, this implies

$$d_{TL_{\lambda}^p}^p(f, g) \leq \int_0^1 \frac{|x - T_1(x, 1)|^p}{\lambda} + |f(x) - g(T_1(x, 1))|^p dx + \int_1^2 \frac{|x - T_1(x, 0)|^p}{\lambda} + |f(x) - g(T_1(x, 0))|^p dx = \frac{2}{\lambda} \rightarrow 0$$

as  $\lambda \rightarrow \infty$ .

In this example  $d_{TL_{\infty}^p}(f, g) := \lim_{\lambda \rightarrow \infty} d_{TL_{\lambda}^p}(f, g) = 0$ . More generally the  $TL_{\infty}^p$  distance is an OT distance between the measures  $f_{\#}\mu$  and  $g_{\#}\nu$ .

**Proposition 3.3.** [79] Let  $\Omega \subseteq \mathbb{R}^d$ ,  $f, g : \Omega \rightarrow \mathbb{R}^m$  measurable functions and  $\mu, \nu \in \mathcal{P}_p(\Omega)$  where  $p \geq 1$ , then

$$\lim_{\lambda \rightarrow \infty} d_{TL_{\lambda}^p}((f, \mu), (g, \nu)) = d_{\text{OT}}(f_{\#}\mu, g_{\#}\nu)$$

where  $d_{\text{OT}}$  is the OT distance (on  $\mathcal{P}(\mathbb{R}^m)$ ) with cost  $c(x, y) = |x - y|^p$ .

As the example before the proposition showed,  $d_{TL_{\infty}^p}(f, g)$  is not a metric, however is non-negative, symmetric and the triangle inequality holds.

In [4, Section 2] the authors, using the fluid mechanics formulation of optimal transport, interpolate the optimal transport distance with quadratic cost with the  $L^2$  distance. The resulting interpolated distance can still be written in the fluid mechanics formulation which naturally gives rise to geodesics. By contrast the  $TL^p$  distance interpolates between  $L^p$  and the optimal transport distance of the push forward measures. This is well defined for any  $p \geq 1$  (unlike the previous method which requires  $p = 2$ ) however geodesics do not exist in the  $TL^p$  space. One must also treat the signals as probability measures in the approach of [4].

We observe that when  $\mu$  is a uniform measure (either in the discrete or continuous sense) the measure  $f_{\#}\mu$  is the histogram of  $f$ . The OT distance between histograms is a popular tool in histogram specification. Minimizers to the Monge formulation of  $d_{\text{OT}}(f_{\#}\mu, g_{\#}\nu)$  define a mapping between the histograms  $f_{\#}\mu$  and  $g_{\#}\nu$  [49, 62, 63]. However this mapping contains no spatial information. If instead one uses minimizers to the Monge formulation of the  $TL_{\lambda}^p$  distance (6) ( $\lambda < \infty$ ) then one can include spatial information in the histogram specification. We explore this further in Section 4.4 and apply the method to the colour transfer problem.

It is well known that there exists a minimizer (when  $c$  is lower semi-continuous) for OT. Since the  $TL_{\lambda}^p$  distance is closely related to an OT distance between measures in  $\mathbb{R}^{d+m}$  (i.e. measures supported on graphs) then there exists a minimizer to (3-4).

**Proposition 3.4.** [79] Let  $\Omega \subseteq \mathbb{R}^d$  be open and bounded,  $f \in L^p(\mu)$ ,  $g \in L^p(\nu)$  where  $\mu, \nu \in \mathcal{P}(\Omega)$ ,  $\lambda \in [0, +\infty]$  and  $p \geq 1$ . Under these conditions there exists an optimal plan  $\pi \in \Pi(\mu, \nu)$  realising the minimum in  $d_{TL_\lambda^p}((f, \mu), (g, \nu))$ .

As in the OT case it is natural to set the  $TL_\lambda^p$  distance in the Monge formulation (2). We can write

$$d_{TL_\lambda^p}((f, \mu), (g, \nu)) = \inf_{T: T_\# \mu = \nu} \int_{\Omega} c\lambda(x, T(x); f, g) d\mu(x). \quad (6)$$

Minimizers to the above will not always exist. For example, consider when  $f = g$  then the  $TL_\lambda^p$  distance is the OT distance between  $\mu$  and  $\nu$ . If one chooses  $\mu = \frac{1}{3}\delta_{x_1} + \frac{1}{3}\delta_{x_2} + \frac{1}{3}\delta_{x_3}$  and  $\nu = \frac{1}{2}\delta_{y_1} + \frac{1}{2}\delta_{y_2}$  where all of  $x_i, y_j$  are distinct then there are no maps  $T: \{x_1, x_2, x_3\} \rightarrow \{y_1, y_2\}$  that pushforward  $\mu$  to  $\nu$ .

However, in terms of numerical implementation, an interesting and important case is when  $\mu$  and  $\nu$  are discrete measures (see also [83, pg 5, 14-15] for the following argument with the Monge OT problem). Let  $\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  and  $\nu = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$  then  $\pi = (\pi_{ij})_{i,j=1}^n \in \Pi(\mu, \nu)$  is a doubly stochastic matrix up to a factor of  $\frac{1}{n}$ , that is

$$\pi_{ij} \geq 0 \forall i, j, \sum_{i=1}^n \pi_{ij} = \frac{1}{n} \forall j \text{ and } \sum_{j=1}^n \pi_{ij} = \frac{1}{n} \forall i, \quad (7)$$

and the  $TL_\lambda^p$  distance can be written

$$d_{TL_\lambda^p}^p((f, \mu), (g, \nu)) = \min \sum_{i=1}^n \sum_{j=1}^n c\lambda(x_i, y_j; f, g) \pi_{ij} \quad (8)$$

where the minimum is taken over  $\pi$  satisfying (7). It is known (by Choquet's Theorem, e.g. [67, Theorem 32.3]) that the solution to this minimisation problem is an extremal point in the matrix set  $\Pi(\mu, \nu)$ . It is also known (by application of Birkhoff's Theorem, e.g. [6]) that extremal points in  $\Pi(\mu, \nu)$  are permutation matrices. This implies that there exists an optimal plan  $\pi^*$  that can be written as  $\pi_{ij}^* = \frac{1}{n} \delta_{j - \sigma(i)}$  for a permutation  $\sigma: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ .

Hence there exists an optimal plan to the Monge formulation of the  $TL_\lambda^p$  distance.

**Proposition 3.5.** For any  $f \in L^p(\mu)$  and  $g \in L^p(\nu)$  where  $\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  and  $\nu = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$  there exists a permutation  $\sigma: \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, n\}$  such that

$$d_{TL_\lambda^p}((f, \mu), (g, \nu)) = \frac{1}{n} \sum_{i=1}^n c_\lambda(x_i, x_{\sigma(i)}; f, g).$$

The above theorem implies that in the uniform discrete case there exists optimal plans (which are matrices) which will be sparse. In particular,  $\pi^*$  is an  $n \times n$  matrix with only  $n$  non-zero entries. This motivates the use of numerical methods that can take advantage of expected sparsity in the solution (e.g. iterative linear programming methods such as [53]).

### 4 $TL^p$ in Multivariate Signal and Image Processing

Written in the form (3) the  $TL_\lambda^p$  distance is an OT distance between the measures  $\mu$  and  $\nu$  with the cost function  $c$  given by (4) and which depends upon  $f$  and  $g$ . Hence, to compute  $TL_\lambda^p$  distances there are many algorithms for OT distances that we may apply, for example the multi-scale approaches of Schmitzer [71] and Oberman and Ruan [53], or the entropy regularized approaches of Cuturi [16] and Benamou, Carlier, Cuturi, Nenna and Peyre [5]. Our choice was the iterative linear programming method of Oberman and Ruan [53] for the multivariate signals which we find works well both in terms of accuracy and computation time. Our choice for the images was the entropy regularized solution due to Cuturi [5, 16]. Whilst this only produces an approximation of the  $TL^p$  distance we find it computationally efficient for 2D images. In particular this method regularizes the OT distance with  $\epsilon H(\pi)$  where  $H$  is entropy. We choose  $\epsilon$  as small as possible whilst avoiding numerical instability. In practice this corresponds to a choice of  $\epsilon \approx 0.005$ . For convenience we include a review of the numerical methods in Appendix B.

With respect to choosing  $\lambda$  there are two approaches we could take. The first is to compute the  $TL^p$  distance for a range of  $\lambda$  and then use cross-validation. There are two disadvantages to this approach: we would still have to know the range of  $\lambda$  and computing the  $TL_\lambda^p$  distance for multiple choices of  $\lambda$  would considerably increase computation time. The second approach, and the one we use for each example in this section, is to estimate  $\lambda$  by comparing length scales and desired behaviour. In particular we choose  $\lambda$  so that both horizontal and vertical transport make a contribution. For the applications in this section we want to stay away from the asymptotic regimes  $\lambda \approx 0$  and  $\lambda \gg$ . By balancing the vertical and horizontal length scale we can formally find an approximation of  $\lambda$  which in our results below works well. For example, if we expect a set of real valued time series to have range  $[f_{\min}, f_{\max}]$  and domain  $[t_{\min}, t_{\max}]$  then to balance the vertical and horizontal length scales we choose  $\lambda$  so that

$$\frac{|t_{\max} - t_{\min}|^p}{\lambda} \approx |f_{\max} - f_{\min}|^p.$$

We first consider two synthetic examples. Considering synthetic examples allows us to better demonstrate where the  $TLP$  distance will be successful. In particular synthetic examples can simplify the analysis and allow us to draw attention to features that may be obscured in real world applications.

The first synthetic example considers three classes where we can analytically compute the within class distances and between class separation. This allows us to compare how well we expect  $TLP$  distances to perform in a classification problem.

The second synthetic example uses simulated 2D data from one-hump and two-hump functions. We test how well the  $TLP$  distance recovers the classes and compare with OT and  $L^p$  distances.

Our first real world application is in classifying multivariate times series and 2D images. We choose a multivariate time series data set where we expect transport based methods to be successful but OT cannot be immediately applied. That is, OT must be applied to measures that are real valued so cannot be directly applied to multi-channelled signals. As a benchmark we find the OT distance on each channel separately and take the average over all channels. This would seem reasonable when channels are independent but is not a good assumption in the AUSLAN dataset where we expect temporally correlated signals.

Our chosen data set consists of sequences of sign language data (we define the data set in more detail shortly) which contains the position of both hands (parametrised by 22 variables) at each time. The  $TLP_\lambda^p$  distance can treat these signals as functions  $f: [0,1] \rightarrow \mathbb{R}^{22}$ . We expect to see certain features in the signals however these may be shifted based on the speed of the speaker. The second data set contains 2D images that must be normalised in order to apply the OT distance, this distorts some of the features leading to a poor performance.

We repeat the classification experiment on the AT&T Database of Faces. This is a database of ten 2D greyscale images of forty subjects. Note that if the images were in colour then one cannot immediately apply the OT distance.

The second real world application is histogram specification and colour transfer. Histogram specification or matching, where one defines a map  $T$  that matches one histogram with another, is widely used to define a colour transfer scheme. In particular let  $f: (x_i)_{i=1}^N \rightarrow \mathbb{R}^3$  represent a colour image by mapping pixels  $x_i$  to a colour  $f(x_i)$  (for example in RGB space), one defines a multidimensional histogram of colours on an image by

$\phi(c) = \frac{1}{N} \# \{x_i: f(x_i) = c\}$ . For colour images the histogram  $\phi$  is a measure on  $\mathbb{R}^3$ . For

notational clarity we will call  $\phi$  the colour histogram. One can equivalently define a histogram for grayscale images as a measure on  $\mathbb{R}$ .

Let  $\phi$  and  $\psi$  be two colour histograms for images  $f$  and  $g$  respectively. The OT map  $T$  defines a rearrangement of  $\phi$  onto  $\psi$ , that is  $\psi = T\#\phi$ . In colour transfer the map  $T$  is used to colour the image  $f$  using the palette of  $g$  by  $\hat{f}(x) = g(T(x))$ .

The histogram contains only intensity information and in particular there is no spatial dependence. Using the  $TL_\lambda^p$ -optimal map we define spatially correlated histogram specification and explain how this can be applied to the colour transfer problem. We demonstrate that  $TL_\lambda^p$  distance produces a visually more appealing solution than the OT solution when spatial information is important. We also observe that, for colour images, computing OT maps is a 3D problem (the domain and range of the transport maps is in colour space), whereas computing  $TL_\lambda^p$  maps is a 2D problem (the transport maps pixels to pixels). By Proposition 3.3 when  $\lambda$  is large we can approximate the OT distance between histograms by the  $TL_\lambda^p$  distance. For images that use the full spectrum of colours, i.e. the colour histograms are  $256 \times 256 \times 256$ , the size of the discretisation is  $256^3 \approx 16.8 \times 10^6$ . Hence the spatially correlated histogram specification method allows for a numerically efficient approximation of the OT induced histogram specification method when the size (in terms of number of pixels) of the images are less than  $4096 \times 4096$ .

#### 4.1 1D Class Separation for Synthetic Data

**Objective**—We compare the expected classification power of  $TL^p$ ,  $L^p$  and OT distances with three classes of 1D signals that differ by position (translations), shape (1 hump versus 2 hump) and frequency (hump versus chirp).

**Data Sets**—We consider data from three classes defined in Figure 4. The first class contains single hump function and the second class contains two hump functions. The third class consists of functions with one hump and one chirp, defined to be a high frequency perturbation of a hump. The classes are chosen to test the performance of the  $TL^2$  distance with  $L^2$  and OT distances with regards to identifying translations (where we expect the  $L^2$  distance to do poorly) with a class containing high frequency perturbations (where we expect the OT distance to do poorly).

**Methods**—For a distance to have good performance in classification and clustering problems it should be able to separate classes. To be able to quantify this we use the ratio of ‘between class separation’ to ‘class coverage radius’ that we define now.

Let  $\mathcal{E}_i^N = \{f_j^i\}_{j=1}^N$  be a sample of  $N$  functions from class  $\mathcal{E}_i$ . For a given radius  $r$  we let  $G_i(r)$  be the graph defined by connecting any two points in  $\mathcal{E}_i^N$  with distance less than  $r$ . The distance will be defined using the  $TL_\lambda^p$ ,  $L^2$  and OT distances. Let  $R_{TL_\lambda^p}(\mathcal{E}_i^N)$  be the smallest  $r$  such that  $G_i(r)$  is a connected graph using the  $TL_\lambda^p$  distance. Analogously we can define  $R_{L^p}$  and  $R_{OT}$ .

We define ‘between class separation’ as the Hausdorff distance between classes:

$$d_{H,\rho}(\mathcal{C}_i^N, \mathcal{C}_j^N) = \max \left\{ \sup_{f \in \mathcal{C}_i^N} \inf_{g \in \mathcal{C}_j^N} \rho(f, g), \sup_{g \in \mathcal{C}_j^N} \inf_{f \in \mathcal{C}_i^N} \rho(f, g) \right\}$$

where we will consider  $\rho$  to be one of the  $TL^2$ ,  $L^2$  or OT distances. Large values of  $d_{H,\rho}(\mathcal{C}_i^N, \mathcal{C}_j^N)$  imply that the classes  $\mathcal{C}_i^N$  and  $\mathcal{C}_j^N$  are well separated.

When  $R_\rho(\mathcal{C}_i^N) \leq d_{H,\rho}(\mathcal{C}_i^N, \mathcal{C}_j^N)$  then we say that the class  $\mathcal{C}_i^N$  is separable from class  $\mathcal{C}_j^N$  since for any  $f \in \mathcal{C}_i^N$  the nearest neighbour in  $(\mathcal{C}_i^N \cup \mathcal{C}_j^N) \setminus \{f\}$  is also in class  $\mathcal{C}_i^N$ . We define the pairwise property

$$\kappa_{ij}(\rho; N) = \frac{\mathbb{E}d_{H,\rho}(\mathcal{C}_i^N, \mathcal{C}_j^N)}{\max \{ \mathbb{E}R_\rho(\mathcal{C}_i^N), \mathbb{E}R_\rho(\mathcal{C}_j^N) \}}$$

where we take the expectation over sample classes  $\mathcal{C}_i^N$ . We will assume that the distribution over each class is uniform in the parameter  $\ell$ . When  $\kappa_{ij}(\rho; N) > 1$  then we expect classes  $\mathcal{C}_i^N$  and  $\mathcal{C}_j^N$  to be separable from each other.

As a performance metric we use the smallest value of  $N$  such that  $\kappa_{ij}(\rho; N) \geq 1$ . We let

$$N_{ij}^*(\rho) = \min \{ N : \kappa_{ij}(\rho; N) \geq 1 \}.$$

This measures how many data points we need in order to expect a good classification accuracy.

**Results**—We leave the calculation to the appendix but the conclusion is

$$\begin{aligned} N_{12}^*(TL^2) &< N_{12}^*(OT) < N_{12}^*(L^2) \\ N_{13}^*(TL^2) &< N_{13}^*(OT) < N_{13}^*(L^2) \\ N_{23}^*(TL^2) &< N_{23}^*(OT) < N_{23}^*(L^2). \end{aligned}$$

In each case the  $TL^2$  distance outperforms the  $L^2$  and OT distances.

In each class the  $L^2$  distance has a larger value of  $R$ . This implies a larger data set is needed to accurately cover each class. This is due to the Lagrangian nature of signals within each

class (translations) that is poorly represented by the  $L^2$  distance. The OT distance has the lowest (and therefore best) value of  $R$  in each class. Since each class is Lagrangian then the OT distance is very small between functions of the same class.

When considering between class separation the  $TL^2$  and  $L^2$  distances coincide and give a bigger (and better) between class distance than the OT distance. Since the class  $\mathcal{C}_3$  can be written as a high frequency perturbation of functions in the class  $\mathcal{C}_2$  then, in the OT distance, the functions from class  $\mathcal{C}_3$  approximate functions from the class  $\mathcal{C}_2$ . The distance  $d_{H,OT}(\mathcal{C}_2^N, \mathcal{C}_3^N)$  is therefore small so that one needs more data points in order to fully resolve these classes. We see a similar effect when considering  $d_{H,OT}$  for the other classes.

#### 4.2 2D Classification for Synthetic Data

**Objective**—We use simulated data to illustrate better separation of the  $TL^p$  distance compared to  $L^p$  and OT distances for 2D data from two classes of 1-hump and 2-hump functions.

**Data Sets**—The data set consists of two dimensional images simulated from the following classes

$$\mathbb{P} = \left\{ p[0, 1]^2: p(x) = \alpha\phi(x|\gamma, \sigma), \gamma \sim \text{unif}([0, 1]^2), \alpha \sim \text{unif}([0.5, 1]) \right\}$$

$$\mathbb{Q} = \left\{ [0, 1]^2: q(x) = \alpha\phi(x|\gamma_1, \sigma) - \alpha\phi(x|\gamma_2, \sigma), \gamma_1, \gamma_2 \text{ iid } \text{unif}([0, 1]^2), \alpha \sim \text{unif}([0.5, 1]) \right\}$$

where  $\phi(\cdot|\gamma, \sigma)$  is the multivariate normal pdf with mean  $\gamma \in \mathbb{R}^2$  and co-variance  $\sigma \in \mathbb{R}^{2 \times 2}$ . We choose  $\sigma = 0.01 \times \text{Id}$  where Id is the  $2 \times 2$  identity matrix. The first class,  $\mathbb{P}$ , are the set of multivariate Gaussians restricted to  $[0, 1]^2$  with mean uniformly sampled in  $[0, 1]^2$  and weighted by a uniformly sampled in  $[0.5, 1]$ . The second class,  $\mathbb{Q}$ , are the set of weighted differences between two Gaussian pdf's restricted to  $[0, 1]^2$  with means  $\gamma_1, \gamma_2$  sampled uniformly in  $[0, 1]^2$ . Note that the second class contains non-positive functions. See Figure 5 for examples from each class.

We simulate 25 from each set and denote the resulting set of functions by  $\mathcal{F} = \{f_i\}_{i=1}^N$  where  $N=50$ .

**Methods**—Let  $\left(\{f_i\}_{i=1}^N, D_{TL_\lambda^2}\right)$  be a finite dimensional metric space where  $D_{TL_\lambda^2}$  is the  $N \times N$  matrix containing all pairwise  $S_{L_\lambda^2}$  distances, i.e.  $D_{TL_\lambda^2}(i, j) = d_{TL_\lambda^2}(f_i, f_j)$ . Similarly for  $\left(\{f_i\}_{i=1}^N, D_{L^2}\right)$  and  $\left(\{f_i\}_{i=1}^N, D_{OT}\right)$  where the optimal transport distance is defined by  $d_{OT}(f, g) = \sqrt{OT(f, g)}$  and OT is given by (1) for  $\alpha(x, y) = |x - y|^2$ .



To apply the optimal transport distance we need to renormalise so that signals are all non-negative and integrate to the same value. We do this by applying the nonlinear transform  $\mathcal{N}(f) = \frac{f - \beta}{f - \beta}$  where  $\beta = \min_{f \in \mathcal{F}} \min_{x \in [0,1]} f(x)$ . Neither the  $L^2$  nor  $TL^2$  distances require normalisation.

We use non-metric multidimensional scaling (MDS) [39] to represent the graph in  $k$  dimensions. More precisely the aim is to approximate  $(\{f_i\}_{i=1}^N, D)$  by a metric space  $(\{x_i\}_{i=1}^N, D_{1,|2})$  embedded in  $\mathbb{R}^k$  ( $D_{1,|2}$  is the matrix of pairwise distances using the Euclidean distance, i.e.  $D_{1,|2}(i, j) = |x_i - x_j|_2$ ). This is done by minimising the stress  $S$  defined by

$$S_{TL^2_\lambda}^{(k)} = \frac{\sum_{i,j=1}^N \left( |x_i - x_j|_2^2 - F \left( D_{TL^2_\lambda}(i, j) \right) \right)^2}{\sum_{i,j=1}^N |x_i - x_j|_2^2}$$

over  $\{x_i\}_{i=1}^N \subset \mathbb{R}^k$  and monotonic transformations  $F: [0, \infty) \rightarrow [0, \infty)$ , with  $S_{L^2}$ ,  $S_{OT}$  defined analogously. The classical solution to finding the MDS projection (for Euclidean distances) is to use the  $k$  dominant eigenvectors of the matrix of squared distances, after double centring, as coordinates weighted by the square root of the eigenvalue. More precisely, define  $D^{(2)} = -\frac{1}{2} J \left[ |f_i - f_j|^2 \right]_{ij}$  where  $J = \text{Id} - \frac{1}{N} \mathbb{1}$  and  $\mathbb{1}$  is the  $N \times N$  matrix of ones. Let  $\Lambda_k$  be the matrix with the  $k$  largest eigenvalues of  $D^{(2)}$  on the diagonal and  $E_k$  to be the corresponding matrix of eigenvectors. Then  $X = E_k \Lambda_k^{-\frac{1}{2}}$  is the MDS projection. Increasing the dimension of the projected space  $k$  leads to a better approximation. In Figure 5 we show the projection in  $L^2$ ,  $TL^2$  and the OT distances for  $k = 2$  as well as the dependence of  $k$  on  $S$  for each choice of distance.

**Results**—Our results in Figure 5 show that the  $TL^2$  distance is the better distance for this problem. There is no separation in either  $L^2$  or OT distances whereas the  $TL^2$  distance completely separates the data. It should not therefore be surprising that the 1NN classifier with the  $TL^2$  distances outperforms the others. In fact, using 5 fold cross validation (CV) we get 100% accuracy with the  $TL^2$  distance, compared to 72% in the  $L^2$  distance and 86% in the OT distance. In addition we see that the stress  $S_p$  is much smaller and converges quickly to zero for the  $TL^2$  distance which indicates that the  $TL^2$  distance is, in this problem, more amenable to a low dimensional representation than either OT or  $L^2$  distances.

### 4.3 Classification with Real World Data Sets

**Objective**—We evaluate  $TL^2_\lambda$  as a distance to classify real world data sets where spatial and intensity information is expected to be important and compare with popular alternative

distances. We choose one dataset which is of the type multivariate time series and a second data set consisting of images.

**Data Sets**—We use two data sets. The first is the *AUS-LAN*[30, 41] data set which contains 95 classes (corresponding to different words) from a native AUSLAN speaker (Australian Sign Language) using 22 sensors on a CyberGlove (recording position of  $x$ ,  $y$ ,  $z$  axis, roll, yaw, pitch for left and right hand). Therefore signals are considered as functions from  $\{t_1, t_2, \dots, t_N\}$  to  $\mathbb{R}^{22}$ . There are 27 signals in each class which give a total of 2565 signals.

We make two pre-processing steps. The first is to truncate each signal so it is 44 frames in length. Empirically we find that the signal is constant after the 44th frame and therefore there is no loss of information in truncating the signal. The second pre-processing step is to normalise each channel independently. This is because some channels are orders of magnitude greater than others and would otherwise dominate each choice of distance.

The second data set we use is the AT&T Database of Faces [70]. The dataset consists of ten greyscale facial images from forty subjects, see Figure 6b for examples. There were 400 images in total. In order to reduce the computation time we reduced the size of the images from  $92 \times 112$  pixels to  $50 \times 50$  pixels.

**Methods**—For the multivariate time series we compare the performance of a INN classifier using the  $L^2$  and  $TL^2_\lambda$  distances as well as the state-of-the-art method dynamic time warping [22] and the OT distance average over each channel:

$$d_{\text{MOT}}(f, g) = \frac{1}{22} \sum_{i=1}^{22} d_{\text{OT}}(\hat{f}_i, \hat{g}_i)$$

where  $\hat{f}_i$  is the  $i^{\text{th}}$  channel of  $f$  after normalisation that is given by  $\hat{f}_i = \frac{f_i + c}{J(f_i + c)}$  and where  $c$  is chosen so that each signal is non-negative. We use the OT distance with cost  $c(x, y) = |x - y|^2$ .

There are three common variations of dynamic time warping. One can apply dynamic time warping directly to the signals  $f$  and  $g$  (denoted by DTW), to the derivative  $f'$  of the signals (denoted by DDTW) and to a weighted average of DTW and DDTW (denoted by WDTW). We define

$$\begin{aligned} d_{\text{DDTW}}(f, g) &= d_{\text{DTW}}(f', g') \\ d_{\text{WDTW}}(f, g) &= \alpha d_{\text{DTW}}(f, g) + (1 - \alpha) d_{\text{DDTW}}(f, g). \end{aligned}$$

The parameter  $\alpha$  is chosen by 5-fold 2nd depth cross validation. To be more precise the training data set is split into five partitions. One forms the testing data set (accounting for

20% of the data) and the other four form the training data set. To choose  $\alpha$  we further divide the training data set into five partitions (each accounting for 16% of the data set). For each  $\alpha = \frac{i}{100}$ , where  $i = 0, 1, \dots, 100$ , we compute how accurately one partition of the training set is classified using the remaining four parts. We then choose the value of  $\alpha$  which produces the best classification accuracy on the training data. This value of  $\alpha$  is then used to classify the testing data set.

The analogous distances for  $L^2$ ,  $TL_\lambda^2$  and multi-channelled OT are defined by

$$\begin{aligned} d_{DL^2}(f, g) &= d_{L^2}(f', g') \\ d_{DTL_\lambda^2}(f, g) &= d_{TL_\lambda^2}(f', g') \\ d_{DL^2}(f, g) &= d_{L^2}(f', g') \\ d_{DMOT}(f, g) &= d_{MOT}(f', g') \\ d_{WL^2}(f, g) &= \alpha d_{L^2}(f, g) + (1 - \alpha) d_{DL^2}(f, g) \\ d_{WTL_\lambda^2}(f, g) &= \alpha d_{TL_\lambda^2}(f, g) + (1 - \alpha) d_{DTL_\lambda^2}(f, g) \\ d_{WMOT}(f, g) &= \alpha d_{MOT}(f, g) + (1 - \alpha) d_{DMOT}(f, g). \end{aligned}$$

We do not have to choose the same value of  $\lambda$  in the  $TL_\lambda^2$  and  $DTL_\lambda^2$  distances however considering that signals are normalised, we will use the same value. Note that  $DL^2$ , DTW, DDTW, WDTW  $DTL_\lambda^2$  and DMOT are *not* metrics.

We remark that an alternative method for including derivatives in the  $TL_\lambda^p$  distance would be to extend the signal to include the derivative. We briefly assume that  $f$  is defined over a continuous domain. Let  $f: \mathbb{R} \rightarrow \mathbb{R}$ , and  $\tilde{f} = \left( f, \frac{df}{dx} \right)$ , then we define

$$d_{TW_\lambda^{1,p}}(f, g) = d_{TL_\lambda^p}(\tilde{f}, \tilde{g}).$$

We take our notation  $TW_\lambda^{k,p}$  from the Sobolev space notation where  $W^{k,p}$  is the Sobolev space with  $k$  weak derivatives integrable in  $L^p$ . There is no reason to limit this to one derivative, and we may define  $\tilde{f} = \left( f, \frac{df}{dx}, \dots, \frac{d^k f}{dx^k} \right)$  and

$$d_{TW_\lambda^{k,p}}(f, g) = d_{TL_\lambda^p}(\tilde{f}, \tilde{g}).$$

When the signals are discrete one should use a discrete approximation of the derivative. In order to be consistent with previous extensions of dynamic time warping we do not develop this approach here.

Dynamic time warping is only defined on time series so we are not able to apply it to the AT&T Database of Faces. We apply the optimal transport distance by normalising each image  $f \in \mathbb{R}^2 \rightarrow \{0, 1\}$   $\hat{f}(x) = \frac{f(x)}{\int_{[0,1]^2} f(y) dy}$ . There is no normalisation for either  $L^2$  or  $TL^2$  distances. We find the 1NN classifier using  $TL^2$ ,  $L^2$  and OT distances.

We will use  $\lambda = 1$  in AUSLAN and  $\lambda = 0.1$  in the AT&T Database of Faces for the  $TL^2_\lambda$  based distances. The underlying measure  $\mu$  is chosen to be the uniform measure defined on  $[0,1]$  or  $[0,1]^2$ .

**Results**—We considered two methods for comparing the performance of each distance. The first is the 1NN classification accuracy in each distance. We use the 1NN classification accuracy as a measure as to how well each distance captures the underlying geometry. A higher accuracy implies closest neighbours are more likely to belong to the same class.

The results are given in Table 1 where we report error rates using 5-fold cross-validation. In terms of the 1NN classifier for the AUSLAN data set we see that  $TL^2$  is better than  $L^2$  and is a modest improvement over dynamic time warping.

A rather surprising result is the difference between the MOT distance between the signals and the MOT distance between the derivative of signals. We believe this is most likely due to the length of the word being a good indicator of the word (this would also explain why the  $L^2$  distance has reasonable performance). We can see from the example signal in Figure 6a only the first part of the signal contains information (the word being spoken), the remainder of the signal is noise. Because we need to renormalise in order to apply the MOT distance then, similar to the example in Figure 2a, the difference between the first part of the signal (containing information) and the latter part of the signal (containing noise) is reduced.

On the other hand the derivative of the signal will place a lot of mass at the end of the word, with smaller masses in other places where the signal is changing. In particular MOT is now able to identify the length of the signal, leading to a big improvement in performance. Furthermore, some channels are likely to contain more information than others. The decoupling of channels in the MOT distance could be an advantage as simultaneously matching across all channels, as in the  $TL^P_\lambda$  distance, can mean the latter distance is corrupted by low quality channels. This artefact could be removed by weighting channels (this would require training the distance). Since we expect a temporally correlated distance to be a better model then, when weighting channels, we would expect to see an improved performance of the  $TL^P$  distance over the MOT distance.

Our results indicate that the  $TL^P_\lambda$  distance better represents the geometry of the dataset than any of the other (psuedo) distances. However a 1NN classifier should not be expected to

achieve the best classification results. We refer to [2] for a state-of-the-art neural network which produces a much better classification error than the 1NN method considered here (the smallest error rate in [2] for AUSLAN is 2.53%, this uses a training data set equal to  $\frac{4}{9}$ th of the total data set). We stress that the aim of this paper is to introduce a distance that better models signals where both spatial and intensity is important, not to define a new classification method.

For the AT&T Database of Faces the 1NN classifier using the OT distance performs worse with an error of 3.3%. The  $L^2$  distance does the second best with 2.5% and the  $TL^2$  distance is the best with 2%.

In the same spirit as Section 4.1 we define the performance metric  $\kappa_{ij}(\rho)$  as the ratio of distance between class  $i$  and class  $j$  and the maximum class coverage radius of class  $i$  and class  $j$ . For the distance between classes we use the Hausdorff distance (see Section 4.1) and for the class coverage radius we use the minimum radius  $r$  such that connecting any two data points in class  $i$  closer than  $r$  defines a connected graph. We plot the results in Figures 6c and 6d. The  $x$  axis represents pairs of classes where for visual clarity we have ordered the pairs so that the  $\kappa(L^2)$  is increasing. A large value of  $\kappa_{ij}$  indicates that it is easier to identify class  $i$  from class  $j$  whereas a small value indicates that identifying the two classes is a difficult problem.

For AUSLAN we see that the  $TL^2$  distance has, for the majority of pairs of classes, a larger value of  $\kappa_{ij}$  than the  $L^2$  distance and DTW and therefore better represents the class structure. The MOT distance does poorly, except in a few cases. We notice that all distances follow the trend that class separation is increasing with  $\kappa(L^2)$ .

For AT&T Database of Faces the  $L^2$  and  $TL^2$  distances perform very similarly. However both the  $TL^2$  and  $L^2$  distances are much more consistent than the OT distance, we can see that although between some classes the OT distance achieves the best results, with other classes the OT distance does extremely poorly (there are many more classes with a class separation close to 1).

#### 4.4 Histogram Specification and Colour Transfer with the $TL^p$ Distance

**Histogram specification and colour transfer**—Histogram specification concerns the problem of matching one histogram onto another. For a function  $f$  on a discrete domain  $X$  the histogram is given by  $f_{\#}\mu$  where  $\mu$  is the uniform discrete measure supported on  $N$  points. We do not make any assumption on the dimension of the codomain of  $f$  (so that  $f$  may be multivalued and the histogram may be multidimensional). This coincides with the definition given in the introduction to the section, that is

$$f_{\#}\mu(y) = \frac{1}{N} \# \{x \in X: f(x) = y\}.$$

Given two functions  $f: X \rightarrow \mathbb{R}^m$  and  $g: Y \rightarrow \mathbb{R}^m$ , with histograms  $\phi$  and  $\psi$  respectively, histogram specification is the problem of finding a map  $T: X \rightarrow Y$  such that  $\psi = T_{\#} \phi$ .

The colour transfer problem is the problem of colouring one image  $f$  with the palette of an exemplar image  $g$ . A common method used to solve this problem is to use histogram specification where  $T$  is the minimizer to Monge's optimal transport problem (2) between  $\phi$  and  $\psi$  [14, 17, 49, 62, 63]. Let our colour space be denoted by  $\mathcal{C}$  where for example if the colour space is 8 bit RGB then  $\mathcal{C} = \{0, 1, \dots, 255\}^3$ . The colour histogram then defines a measure over  $\mathcal{C}$ . If we consider two such histograms  $\phi$  and  $\psi$  corresponding to images  $f: X \rightarrow \mathcal{C}$  and  $g: Y \rightarrow \mathcal{C}$  respectively then a histogram specification is a map  $T: \mathcal{C} \rightarrow \mathcal{C}$  that satisfies  $\psi = T_{\#}\phi$ . The recoloured image  $f^{\circ} = g^{\circ} T$  has the same colour histogram as  $g$ . The solution  $f^{\circ}$  is a recolouring of  $f$  using the palette of  $g$ .

If we consider grayscale images then  $\mathcal{C} = [0, 1]$  and the optimal transport map (assuming it exists) is a monotonically increasing function. In particular this implies that if pixel  $x$  is lighter than pixel  $y$  (i.e.  $f(x) > f(y)$ ) then in the recoloured image  $f^{\circ} = T^{\circ} f$  pixel  $x$  is still lighter than pixel  $y$ . In this sense the OT solution preserves intensity ordering. But note that no spatial information is used to define  $T$ ; only the difference in intensity between pixels is used and not the distance between pixels.

**Spatially correlated histogram specification**—Let  $\phi$  and  $\psi$  be the histograms corresponding to images  $f: X \rightarrow \mathbb{R}^m$  and  $g: Y \rightarrow \mathbb{R}^m$  respectively. If we recall Proposition 3.3 then  $\lim_{\lambda \rightarrow \infty} d_{TL_{\lambda}^p}((f, \mu), (g, \nu)) = d_{OT}(f_{\#}\mu, g_{\#}\nu)$  (where  $\mu$  and  $\nu$  are the discrete uniform measures over the sets  $X$  and  $Y$ ). For  $\lambda < \infty$  the  $TL_{\lambda}^p$  distance includes spatial *and* intensity information. Hence the  $TL_{\lambda}^p$  distance provides a generalization of OT induced histogram specification.

Analogously to the OT induced histogram specification method we define the spatially correlated histogram specification to be histogram specification using the map  $T: X \rightarrow Y$  which is a minimizer to Monge's formulation of the  $TL_{\lambda}^p$  distance (6). When the images are of the same size then, by Proposition 3.5 such a map exists. The recoloured image  $\hat{f}$  of  $f$  is given by  $\hat{f} = g^{\circ} T$ . Furthermore when the images are of the same size the map  $T$  is a rearrangement of the pixels in  $X$  and therefore the histograms are invariant under  $T$ . In particular the histogram of  $\hat{f}$  is the same as the histogram of  $g$ .

Although we propose the spatially correlated histogram specification as a method to incorporate spatial structure we recall from the discussion at the start of Section 4 the value of the method as a numerically efficient approximation to OT induced histogram specification for colour images that are not too large. Motivated by Proposition 3.3 one expects that for large  $\lambda$  the  $TL_{\lambda}^p$  map is approximately the OT map between colour histograms. The OT problem is in the  $C$  space which, for colour images is 3 dimensional. However, the  $TL_{\lambda}^p$  problem is in the domain of the images  $\Omega$ , which is typically 2 dimensional. Hence one can use the  $TL_{\lambda}^p$  distance to approximate OT induced histogram specification in a lower dimensional space when  $O(n_c^3) = |\mathcal{C}| > |\Omega| = O(n_s^2)$  where  $n_c$  is the

size of discretisation in each colour channel and  $n_s$  is the size discretisation in each spatial dimension.

We briefly remark that histogram specification methods often include additional regularization terms. Such choices of regularization on the transport map include penalizing the gradients [17, 62, 63], sparsity [63], average transport [56] and rigidity [28]. One could apply any of the above regularizations to spatially correlated histogram specification.

**Examples**—First, let us consider the  $128 \times 128$  grayscale images in Figure 7. The objective is to combine the shading of the first image with the geometry of the second image. We are motivated by the scenario where one wishes to combine information about a scene obtained by two different measurements: one where intensities (dynamical range) are well resolved, but the spatial resolution (geometry) is not well captured, and another where dynamical range is poorly captured, but the geometry is well resolved. We furthermore allow that the scenes captured may be somewhat different. The desire is to combine the images to obtain a single image with both good geometry and intensity. The solution we propose is to use spatially correlated histogram specification to re-shade the image with low quality intensity.

The result, as given in Figure 7, produces what we consider to be the desired output. The shading has been transferred and the geometry has not been lost. One is not able to apply histogram specification (induced by the OT map) due to the lack of existence of an optimal transport map from the histogram of the original image  $\phi$  to the histogram exemplar image  $\psi$ . This is due to the histogram of the original image being a sum of two delta masses as in Figure 7d.

As a more challenging example we consider real world colour images. Images are  $128 \times 128$ . We compare our method with histogram specification using the OT mappings and the following state of the art methods for which code is freely available. Reinhard, Ashikhmin, Gooch and Shirley's renormalisation method (RAGS) [66] rescales the image so that the mean and standard deviation of the LAB channels match the exemplar image. Pitié and Kokaram (PK) [61] approximate colour histograms with a Gaussian and look for the best linear map between the two colour histograms. Essentially they look for couplings, as in the Monge formulation of optimal transport but they restrict to linear mappings. In general there may not be any linear mappings between two histograms, however it can be shown that when the histograms are Gaussians the set of mappings is not empty. The final method we compare to is the regularized transportation method due to Ferradans, Papadakis, Rabin, Peyré and Aujol (FPRPA) [17].

It is difficult to quantitatively assess the performance of each method objectively. Whether the output is satisfactory depends on content and artistic preferences. We refer to [86] for an objective quantitative measure of colour transfer results, however this explicitly marks against introducing colour artefacts. Whilst some colour artefacts are clearly undesirable introducing others, such as the northern lights in Figure 8 was the objective. Hence we have no way to quantify the performance of our method and instead rely on qualitative assessments.

In the first pair of images the exemplar image contains a few trees with the northern lights in the background, whilst the other image has a few trees with a mostly clear sky in the background. The challenge is to recreate the northern lights in the second image.

As one would expect, in Figure 8c we see that the histogram specification induced by OT loses the spatial structure. Indeed, it is hard to recognise the northern lights. Similar with each competing method in Figure 8g-8i, none of them successfully manage to reproduce the northern lights and the palm trees all pick up an unnatural reddish shade. The spatially correlated histogram specification solution does a much better job at preserving the ordering locally. As  $\lambda$  increases it becomes cheaper to match pixels that are further apart and therefore, for large  $\lambda$ , the matching does not preserve the local structure in the exemplar image.

In the second real world example we consider colour transfer between two images from Masson's trichrome staining procedure shown in Figures 9a and 9b. We manipulate the luminosity of the second images. The objective is to colour the second image using the palette of the first. In Figure 9 we compare the spatially correlated histogram specification method of  $TL_\lambda^p$  to the other methods.

Since we know the true image we may compare the colourisation with the true image. We report the  $L^2$  error computed by

$$\text{err}(f - g) = \frac{1}{N} \sqrt{\sum_{i=1}^3 \sum_j |f_i(x_j) - g_i(x_j)|^2}$$

where  $f = (f_1, f_2, f_3)$  and  $g = (g_1, g_2, g_3)$  are images in RGB space and  $N = 128^2$  is the number of pixels. The  $TL^p$  method (error 0.2885) gives a more accurate colourisation compared to the RAGS method (error 0.4040), the PK method (error 0.3568), the FPRPA method (error 0.4817), and the OT induced histogram specification method (error 0.4030). One can also see that the  $TL_1^2$  based method does not have the same artefacts as the other methods. In particular, (a) the darker band is still evident in RAGS and PK, (b) FPRPA fails to accurately recolour the white band on the right hand side, and (c) OT places too much white on the left hand side and not enough on the right hand side.

## 5 Conclusions

In this paper we have developed and applied a distance that directly accounts for the intensity of the signal within a Lagrangian framework. This differs from OT distances that do not directly measure intensity and the  $L^p$  distance which measures intensity only. Through applications we have shown the potential of this distance in signal analysis.

The distance is widely applicable, unlike in classical OT distances, such as the Wasserstein distance or the earth mover distance, the  $TL_\lambda^p$  distance does not require treating signals as measures. Treating a signal as a measure implies the following constraints: non-negativity,



normalised mass and single channelled. None of these assumptions are necessary for the  $TL_\lambda^p$  distance. Furthermore the distance is applicable to both discrete and continuous signals as well as allowing practitioners to emphasise features which in many cases should allow for a better representation of data sets, for example one could include derivatives.

Efficient existing methods, such as entropy regularized or multi-scale linear programming, for optimal transport are applicable to the  $TL_\lambda^p$  distance. In fact any numerical method for optimal transport that can cope with arbitrary cost functions is immediately available. This includes the entropy regularised approach of Cuturi [16]. However, there are more efficient methods that are specific to the OT distance with quadratic cost that are unavailable here, e.g. [74].

Via the representation as an OT distance between measures supported on graphs we expect many other results for OT distances to carry through to  $TL_\lambda^p$  distances. For example, one could extend the LOT method [85] for signal representation and analysis to the  $TL_\lambda^p$  distance. This would allow pairwise distances of a data set to be computed with numerical cost that is linear in the number of images. We leave the development for future work.

We considered a few examples where we expect (and then showed) that  $TL_\lambda^p$  will outperform OT and other distances. We expect the  $TL_\lambda^p$  distance to give a better performance than OT distances when intensity information is important. On the other hand, we do not expect the  $TL_\lambda^p$  distance to be robust to high frequency noise. In this case an OT distance would probably have superior performance.

The applications we considered were to classification and histogram specification in the context of colour transfer. For classification we chose data sets with a Lagrangian nature but were either multi-channelled or non-positive (so that in both cases one must apply ad-hoc methods in order to apply the OT distance). We showed the  $TL_\lambda^p$  distance better represented the underlying geometry. The 1NN classifier is a very simple method and we expect our results here could be significantly improved by, for example, replacing the  $L^2$  distance in the MDS projection approach of Weinberger and Chapelle [87] with the  $TL_\lambda^p$  distance. For the colour transfer problem we defined a spatially correlated histogram specification method which produced more visually appealing results when combining the colour of one image with the geometry of another.

Although the main motivation was to develop a distance which better represents Lagrangian data sets we also note that the  $TL_\lambda^p$  distance provides a numerically efficient (for images that are not too large) approximation for the OT induced histogram specification method by, for 2-dimensional images colour images, reducing the effective dimension of the problem from three for OT distances to two for the  $TL_\lambda^p$  distance. We also observe that the effective dimension of multi-channelled time signals is one. In particular the effective dimension is independent of the number of channels.

The applications we have considered are for demonstration on the performance of the  $TL_\lambda^p$  distance. A next step would be to consider a more detailed study of a specific problem. For example in the colour transfer application we could have considered regularization terms/constraints which would have improved the performance, e.g. [17, 28, 51, 56, 62, 63]. It was not the aim to propose a state-of-the-art method for each application, indeed each application would constitute a paper within its own right.

### Acknowledgments

Authors gratefully acknowledge funding from the NSF (CCF 1421502) and the NIH (GM090033, CA188938) in contributing to a portion of this work. DS also acknowledges funding by NSF (DMS-1516677). The authors are grateful to the Center for Nonlinear Analysis at CMU for its support. In addition the authors would like to thank the referees for their valuable comments that lead to significant improvements in the paper.

### A Performance of $TL_\lambda p$ in Classification Problems with Simple and Oscillatory Signals

We compare the performance of  $TL_\lambda^2$ ,  $L^2$  and OT distances with respect to classification/clustering for the three classes  $\{\mathcal{C}_j\}_{j=1,2,3}$  of signals defined in Figure 4. We test how each distance performs by finding the smallest number of data points such that the classes  $\mathcal{C}_i^N = \{f_i\}_{i=1}^N \subset \mathcal{C}_i$  are separable. For sufficiently large  $N$  the approximation  $d_{H,\rho}(\mathcal{C}_i^N, \mathcal{C}_j^N) \approx d_{H,\rho}(\mathcal{C}_i, \mathcal{C}_j)$  is used to simplify the computation. Similarly, as a proxy for  $\mathbb{E}R\rho(\mathcal{C}_i^N)$  we use  $R\rho(\widehat{\mathcal{C}}^N)$  where

$$\widehat{\mathcal{C}}_i^N = \left\{ f_\ell : \ell = \ell_{\min}^i + \frac{n-1}{N-1}(\ell_{\min}^i - \ell_{\max}^i), n \in \{1, 2, \dots, N\} \right\}$$

is the uniform sample from class  $\mathcal{C}_i$  (recall that class  $\mathcal{C}_i$  is parameterized by  $\ell \in [\ell_{\min}^i, \ell_{\max}^i]$ ) and with an abuse of notation we use the subscript of  $\widehat{f}$  to denote the dependence of  $\widehat{\mathcal{C}}$ .

It follows that the class separation distances and class coverage radius are approximated by

$$\begin{aligned}
 d_{H,L^2}^2(\mathcal{C}_1^N, \mathcal{C}_2^N) &\approx \frac{\alpha}{2} R_{L^2}^2(\mathcal{C}_1^N) \approx \frac{2}{N} \\
 d_{H,L^2}^2(\mathcal{C}_1^N, \mathcal{C}_3^N) &\approx \frac{3\alpha}{4} R_{L^2}^2(\mathcal{C}_2^N) \approx \frac{1}{N} \\
 d_{H,L^2}^2(\mathcal{C}_2^N, \mathcal{C}_3^N) &\approx \frac{\alpha}{4} R_{L^2}^2(\mathcal{C}_3^N) \approx \frac{2\alpha}{N\gamma} \\
 d_{H,OT}^2(\mathcal{C}_1^N, \mathcal{C}_2^N) &\approx \frac{\beta^2\alpha}{4} R_{OT}^2(\mathcal{C}_1^N) \approx \frac{\alpha}{N^2} \\
 d_{H,OT}^2(\mathcal{C}_1^N, \mathcal{C}_3^N) &\approx \frac{\beta^2\alpha}{4} R_{OT}^2(\mathcal{C}_2^N) \approx \frac{\alpha}{N^2} \\
 d_{H,OT}^2(\mathcal{C}_2^N, \mathcal{C}_3^N) &\approx \frac{\beta^2\alpha}{4} R_{OT}^2(\mathcal{C}_3^N) \approx \frac{\alpha}{N^2} \\
 d_{H,TL_\lambda^2}^2(\mathcal{C}_1^N, \mathcal{C}_2^N) &\approx \frac{\alpha}{2} R_{TL_\lambda^2}^2(\mathcal{C}_1^N) \approx \frac{\alpha^2}{N} \\
 d_{H,TL_\lambda^2}^2(\mathcal{C}_1^N, \mathcal{C}_3^N) &\approx \frac{3\alpha}{4} R_{TL_\lambda^2}^2(\mathcal{C}_2^N) \approx \frac{4\alpha^2}{N} \\
 d_{H,TL_\lambda^2}^2(\mathcal{C}_2^N, \mathcal{C}_3^N) &\approx \frac{\alpha}{4} R_{TL_\lambda^2}^2(\mathcal{C}_3^N) \approx \frac{\alpha^2}{N}.
 \end{aligned}$$

We have

$$\begin{aligned}
 \kappa_{12}^2(L^2; N) &\approx \frac{\alpha N}{4}, \kappa_{13}^2(L^2; N) \approx \frac{3\gamma N}{8}, \\
 \kappa_{12}^2(OT; N) &\approx \frac{\beta^2 N}{4}, \kappa_{13}^2(OT; N) \approx \frac{\beta^2 N^2}{4}, \\
 \kappa_{12}^2(TL_\lambda^2; N) &\approx \frac{N}{8\alpha}, \kappa_{13}^2(TL_\lambda^2; N) \approx \frac{3N}{4\alpha}, \\
 \kappa_{23}^2(L^2; N) &\approx \frac{\gamma N}{8}, \\
 \kappa_{23}^2(OT; N) &\approx \frac{\gamma^2 N^2}{8}, \\
 \kappa_{23}^2(TL_\lambda^2; N) &\approx \frac{N}{16\alpha}.
 \end{aligned}$$

Finally we can compute  $N^*$ ,

$$\begin{aligned}
 N_{12}^*(L^2) &\approx \frac{4}{\alpha}, N_{13}^*(L^2) \approx \frac{8}{3\gamma}, N_{23}^*(L^2) \approx \frac{8}{\gamma} \\
 N_{12}^*(OT) &\approx \frac{2}{\beta}, N_{13}^*(OT) \approx \frac{2}{\beta}, N_{23}^*(OT) \approx \frac{\sqrt{8}}{\gamma} \\
 N_{12}^*(TL^2) &\approx \frac{\alpha}{8}, N_{13}^*(TL^2) \approx \frac{4\alpha}{3}, N_{23}^*(TL^2) \approx 16\alpha
 \end{aligned}$$

which for  $\beta > \frac{\alpha}{2}, \beta > \frac{3\gamma}{4}$  and  $\gamma < \frac{\sqrt{2\alpha}}{8}$  implies the ordering given Section 4.1.

## B Numerical Methods

In principle any numerical method for computing OT distances capable of dealing with an arbitrary cost function can be adapted to compute  $TL_\lambda^p$  the distance. Here we describe two numerical methods we used in Section 4.

### B.1 Iterative Linear Programming

Here we describe the iterative linear programming method of Oberman and Ruan [53] which we abbreviate OR. Although this method is not guaranteed to find the minimum in (3) we find it works well in practice and is easier to implement than, for example, methods due to Schmitzer [71] that provably minimize (3) but require a more advanced refinement procedure. See also [46] and references therein for a multiscale descent approach.

The linear programming problem restricted to a subset  $M \subseteq \Omega_h^2$  is

$$\begin{aligned} &\text{minimize: } \sum_{(i,j) \in M} c_\lambda(x_i, x_j; f_h, g_h) \pi_{ij} \text{ over } \pi \\ &\text{subject to } \sum_{i: (i,j) \in M} \pi_{ij} = q_j, \quad \sum_{j: (i,j) \in M} \pi_{ij} = p_i \end{aligned} \quad (\text{LP}_h)$$

where  $c_\lambda$  is given by (4). When  $\mathcal{M} = \Omega_h^2$  then the  $TL_\lambda^p$  distance between  $(f_h, \mu_h)$  and  $(g_h, \nu_h)$  is the minimum to the above linear programme. Furthermore if  $\pi_h$  is the minimizer in the  $TL_\lambda^p$  distance then it is also the solution to the linear programme in  $(\text{LP}_h)$  for any  $\mathcal{M}$  containing the support of  $\pi_h$ . That is if one already knows (or can reasonably estimate) the set of nodes  $\mathcal{M}$  for which the optimal plan is non-zero then one need only consider the linear programme on  $\mathcal{M}$ . This is advantageous when  $\mathcal{M}$  is a much smaller set. Motivated by Proposition 3.5 we expect to be able to write the optimal plan as a map. This implies whilst  $\pi_h$  has  $n^2$  unknowns we only expect  $n$  of them to be non-zero.

The method proposed by OR is given in Algorithm 1. An initial discretisation scale  $h_0$  is given and an estimate  $\pi_{h_0}$  found for the linear programme  $(\text{LP}_h)$  with  $\mathcal{M} = \Omega_{h_0}^2$ . One then iteratively finds  $\mathcal{M}_r \subseteq \Omega_{h_r}^2$ , where  $h_r = \frac{h_r - 1}{2}$ , to be the set of nodes defined by the following refinement procedure. Find the set of nodes for which  $\pi_{h_{r-1}}$  is non-zero, add the neighbouring nodes and then project onto the refined grid  $\Omega_{h_r}^2$ . The optimal plan  $\pi_{h_r}$  on  $\Omega_{h_r}^2$  is then estimated by solving the linear programme  $(\text{LP}_h)$  with  $\mathcal{M} = \mathcal{M}_r$ .

The grid  $\Omega_{h_r}$  will scale as  $(2^{rd} h_0^{-1})^2$ . If the linear programme is run  $N$  times then at the  $r^{\text{th}}$  step the linear programme has on the order of  $2^{rd} h_0^{-1}$  variables. In particular on the last (and

most expensive) step the number of variables is  $O(2^{Nd}h_0^{-1})$ . This compares to size  $(2^{Nd}h_0^{-1})^2$  if the linear programme was run on the final grid without this refinement procedure.

---

**Algorithm 1** An Iterative Linear Programming Approach [53]

---

**Input:** functions  $f, g \in L^p(\Omega)$ , measures  $\mu, \nu \in \mathcal{P}(\Omega)$  and parameters  $h_0, N$ .

- 1: Set  $r = 0$ .
- 2: **repeat**
- 3:     Define  $\mathcal{S}_r = \Omega_{h_r}^2$  where  $\Omega_{h_r}$  is the square grid lattice with distances between neighbouring points  $h_r$  and discretise functions  $f, g$  and measures  $\mu, \nu$  on  $\Omega_{h_r}$ .
- 4:     **if**  $r = 0$  **then**
- 5:         Solve  $(LP_h)$  on  $\mathcal{S}_0$  and call the output  $\pi_{h_0}$ .
- 6:     **else**
- 7:         Find the set of nodes on  $\mathcal{S}_{r-1}$  for which  $\pi_{h_{r-1}}$  is non-zero and call the set  $\mathcal{K}_{r-1}$ .
- 8:         To  $\mathcal{K}_{r-1}$  add all neighbouring nodes and call this set  $\mathcal{N}_{r-1}$ .
- 9:         Define  $\mathcal{M}_r$  to be the set of nodes on  $\mathcal{S}_r$  that are children of nodes in  $\mathcal{N}_{r-1}$ .
- 10:         Solve  $(LP_h)$  restricted to  $\mathcal{M}_r$  and call the optimal plan  $\pi_{h_r}$ .
- 11:     **end if**
- 12:     Set  $h_{r+1} = \frac{h_r}{2}$  and  $r \mapsto r+1$ .
- 13: **until**  $r = N$

**Output:** The optimal plan  $\pi_{h_{N-j}}$  for  $(LP_h)$ .

---

## B.2 Entropic Regularisation

Cuturi, in the context of computing OT distances, proposed regularizing the minimization in (3) with entropy [16]. This was further developed by Benamou, Carlier, Cuturi, Nenna and Peyré [5], abbreviated to BCCNP, which is the method we describe here. Instead of considering the distance  $TL_\lambda^p$  we consider

$$S_\varepsilon = \inf_{\pi \in \Pi(\mu, \nu)} \left\{ \sum_{i=1}^n \sum_{j=1}^n c_\lambda(x_i, x_j; f, g) \pi_{ij} - \varepsilon H(\pi) \right\}$$

where  $H(\pi) = -\sum_{i=1}^n \sum_{j=1}^n \pi_{ij} \log \pi_{ij}$  is the entropy. In the OT case the distance  $S_\varepsilon$  is also known as the Sinkhorn distance. It is a short calculation to show

$$S_\varepsilon = \varepsilon \inf_{\pi \in \Pi(\mu, \nu)} \{ \text{KL}(\pi | \kappa) \}$$

where  $\mathcal{K}_{ij} = \exp\left(-\frac{c_\lambda(x_i, x_j; f, g)}{\varepsilon}\right)$  (the exponential is taken pointwise) and KL is the Kullback-Leibler divergence defined by

$$\text{KL}(\pi | \mathcal{K}) = \sum_{i=1}^n \sum_{j=1}^n \pi_{ij} \log\left(\frac{\pi_{ij}}{\mathcal{K}_{ij}}\right)$$

It can be shown that the optimal choice of  $\pi$  for  $S_\varepsilon$  can be written in the form  $\pi^* = \text{diag}(u) \mathcal{K} \text{diag}(v)$  where  $u, v \in \mathbb{R}^n$  are limits, as  $r \rightarrow \infty$ , of the sequence

$$v^{(0)} = \mathbb{1}, u^{(r)} = \frac{p}{\mathcal{K}v^{(r)}}, v^{(r+1)} = \frac{q}{\mathcal{K}^T u^{(r)}}$$

and  $p = (p_1, \dots, p_n), q = (q_1, \dots, q_n)$  (multiplication is the usual matrix-vector multiplication, division is pointwise and T denotes the matrix transpose). The algorithm given in 2 is a special case of iterative Bregman projections and also known as the Sinkhorn algorithm.

The stopping condition proposed in [16] is to let  $\pi^{(r)} = \text{diag}(u^{(r)}) \mathcal{K} \text{diag}(v^{(r)})$  then stop when

$$\left| \frac{\sum_{i,j=1}^n \mathcal{K}_{ij} \pi_{ij}^{(r)} - \varepsilon H(\pi^{(r)})}{\sum_{i,j=1}^n \mathcal{K}_{ij} \pi_{ij}^{(r-1)} - \varepsilon H(\pi^{(r-1)})} - 1 \right| < 10^{-4}.$$

Note that although as  $\varepsilon \rightarrow 0$  we will recover the unregularised  $TL_\lambda^p$  distance we also suffer numerical instability as  $\mathcal{K} \rightarrow 0$  exponentially in  $\varepsilon$ . These instabilities have been addressed in, for example, [14, 72].

For optimal transport with quadratic cost  $c(x, y) = |x - y|^2$  the Sinkhorn algorithm can be more efficiently implemented using Gaussian convolutions [74]. The two numerical methods described so far use the formulation of  $TL_\lambda^2$  given by (3-4) which interprets the  $TL_\lambda^2$  as an OT distance between measures  $\mu$  and  $\nu$  for a (non-quadratic) cost function  $c_\lambda(\cdot, \cdot; f, g)$ , hence one cannot make use of previous OT methods such as [74].

However, we also recall that we can define the  $TL_\lambda^2$  distance as the optimal transport distance between measures  $(f \times \text{Id})_\# \mu$  and  $(g \times \text{Id})_\# \nu$ , see (5), in which case the entropy regularized approach can be implemented using Gaussian convolutions in dimension  $d + m$  (when  $p = 2$ ), where  $f: \Omega \subseteq \alpha^d \rightarrow \alpha^m$ . Although this means that the numerical method is based in a higher dimension we note the success of the bilateral grid method for bilateral filters that are also based on computing a Gaussian filter in a higher dimension [10, 57]. For colour images,

where  $m = 3$  this approach may not be efficient however for  $m = 1$  these ideas have the potential for an improved algorithm.

---

**Algorithm 2** An Entropy Regularised Approach [5, 16]

---

**Input:** discrete functions  $f = (f_1, \dots, f_n)$ ,  $g = (g_1, \dots, g_n)$ , discrete measures  $\mu = \sum_{i=1}^n p_i \delta x_i$ ,  $\nu = \sum_{j=1}^n q_j \delta x_j$ , the parameter  $\epsilon$  and a stopping condition.

1:

$$\text{Set } r = 0, \mathcal{K} = \left( \exp \left( - \frac{c(x_i, x_j; f, g)}{\epsilon} \right) \right)_{ij} \text{ and } u^{(0)} = \mathbb{1} \in \mathbb{R}^n.$$

2:

**repeat**

3:

$$\text{Let } r \rightarrow r+1, v^{(r)} = \frac{q}{\mathcal{K} T_u^{(r-1)}} \text{ and } u^{(r)} = \frac{p}{\mathcal{K} u^{(r)}} \text{ where } \underline{p} = (p_1, \dots, p_n), \underline{q} = (q_1, \dots, q_n)$$

4:

**until** Stopping condition has been reached

5:

Set  $\pi = \text{diag}(u^{(r)}) \mathcal{K} \text{diag}(v^{(r)})$ .

**Output:** An estimate  $\pi$  on the optimal plan for  $S_\epsilon$  where the accuracy is determined by the stopping condition.

---

## References

1. Åström F, Petra S, Schmitzer B, Schnörr C. Image labeling by assignment. arXiv. 2016 1603.05285.
2. Aswolinskiy W, Reinhart RF, Steil J. chapter Time Series Classification in Reservoir- and Model-Space: A Comparison. Springer International Publishing; 2016. Artificial Neural Networks in Pattern Recognition: 7th IAPR TC3 Workshop, AN- NPR 2016, Ulm, Germany, September 28-30, 2016, Proceedings; 197–208.
3. Basu S, Kolouri S, Rohde GK. Detecting and visualizing cell phenotype differences from microscopy images using transport-based morphometry. Proceedings of the National Academy of Sciences. 2014; 111(9):3448–3453.
4. Benamou JD, Brenier Y. A computational fluid mechanics solution to the Monge-Kantorovich mass transfer problem. Numerische Mathematik. 2000; 84(3):375–393.
5. Benamou JD, Carlier G, Cuturi M, Nenna L, Peyré G. Iterative bregman projections for regularized transportation problems. SIAM Journal on Scientific Computing. 2015; 37(2):A1111–A1138.
6. Birkhoff G. Three observations on linear algebra. Univ Nac Tacum an Rev Ser A. 1946; 5:147–151.
7. Bonneel N, Rabin J, Peyré G, Pfister H. Sliced and Radon Wasserstein barycenters of measures. Journal of Mathematical Imaging and Vision. 2015; 51(1):22–45.
8. Bonneel N, Van De Panne M, Paris S, Heidrich W. Displacement interpolation using Lagrangian mass transport. ACM Transactions on Graphics (TOG). 2011; 30:158.
9. Brenier Y, Frisch U, Hénon M, Loeper G, Matarrese S, Mohayaee R, Sobolevski A. Reconstruction of the early universe as a convex optimization problem. Monthly Notices of the Royal Astronomical Society. 2003; 346(2):501–524.
10. Chen J, Paris S, Durand F. Real-time edge-aware image processing with the bilateral grid. ACM SIGGRAPH 2007 Papers, SIGGRAPH '07 ACM. 2007
11. Chen Y, Georgiou TT, Tannenbaum A. Matrix optimal mass transport: a quantum mechanical approach. arXiv. 2016 1610.03041.
12. Chen Y, Georgiou TT, Tannenbaum A. Interpolation of matrices and matrix-valued measures: The unbalanced case. arXiv. 2017 1612.05914.
13. Chizat L, Peyré G, Schmitzer B, Vialard FX. An interpolating distance between optimal transport and Fisher-Rao metrics. Foundations of Computational Mathematics. 2016:1–44.

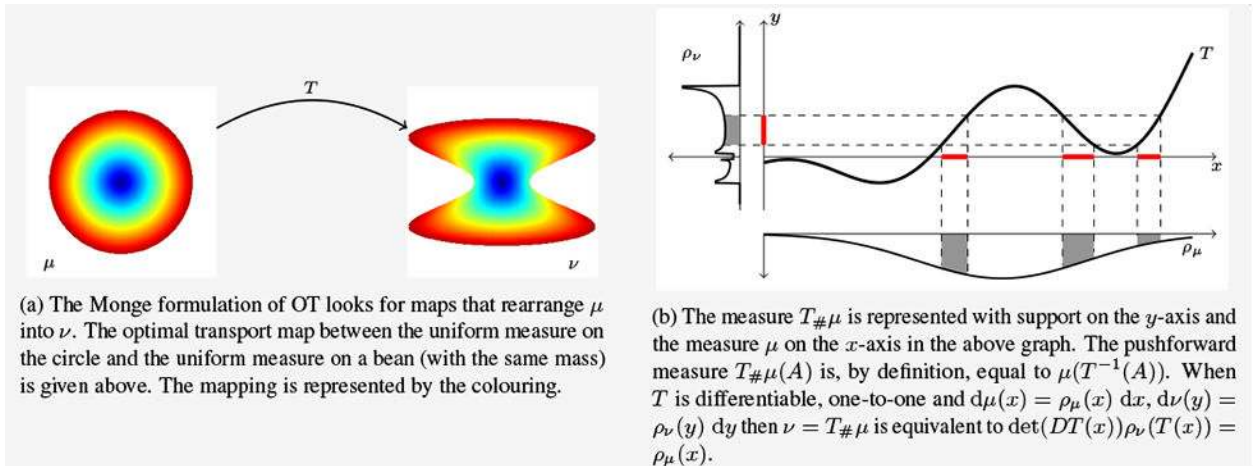
14. Chizat L, Peyré G, Schmitzer B, Vialard FX. Scaling algorithms for unbalanced optimal transport problems. arXiv. 2016 1607.05816.
15. Courty N, Flamary R, Tuia D. Proceedings, Part I chapter Domain Adaptation with Regularized Optimal Transport. Springer Berlin Heidelberg; 2014. Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2014, Nancy, France, September 15-19, 2014; 274–289.
16. Cuturi M. Sinkhorn distances: Lightspeed computation of optimal transport. Advances in Neural Information Processing Systems. 2013:2292–2300.
17. Ferradans S, Papadakis N, Rabin J, Peyré G, Aujol JF. Proceedings chapter Regularized Discrete Optimal Transport. Springer Berlin Heidelberg; 2013. Scale Space and Variational Methods in Computer Vision: 4th International Conference, SSVN 2013, Schloss Seggau, Leibnitz, Austria, June 2-6, 2013; 428–439.
18. Frisch U, Matarrese S, Mohayaee R, Sobolevski A. A reconstruction of the initial conditions of the universe by optimal mass transportation. Nature. 2002; 417(6886):260–262. [PubMed: 12015595]
19. Frisch U, Sobolevskii A. Application of optimal transportation theory to the reconstruction of the early universe. Journal of Mathematical Sciences (New York). 2004; 133(1):303–309.
20. Frogner C, Zhang C, Mobahi H, Araya-Polo M, Poggio TA. Learning with a Wasserstein loss. Advances in Neural Information Processing Systems (NIPS). 2015; 28
21. García Trillos N, Slepčev D. Continuum limit of Total Variation on point clouds. Archive for Rational Mechanics and Analysis. 2015:1–49.
22. Górecki T, Łuczak M. Multivariate time series classification with parametric derivative dynamic time warping. Expert Systems with Applications. 2015; 42(5):2305–2312.
23. Grenander U, Miller MI. Computational anatomy: An emerging discipline. Q Appl Math. 1998; LVI(4):617–694.
24. Haber E, Rehman T, Tannenbaum A. An efficient numerical method for the solution of the  $L_2$  optimal mass transfer problem. SIAM Journal on Scientific Computing. 2010; 32(1):197–211. [PubMed: 21278828]
25. Haker S, Tannenbaum A. On the Monge-Kantorovich problem and image warping. IMA Volumes in Mathematics and its Applications. 2003; 133:65–86.
26. Haker S, Tannenbaum A, Kikinis R. chapter Mass Preserving Mappings and Image Registration. Springer Berlin Heidelberg; 2001. Medical Image Computing and Computer-Assisted Intervention - MICCAI2001: 4th International Conference Utrecht, The Netherlands, October 1417, 2001 Proceedings; 120–127.
27. Haker S, Zhu L, Tannenbaum A, Angenent S. Optimal mass transport for registration and warping. International Journal of Computer Vision. 2004; 60(3):225–240.
28. Hug R, Maitre E, Papadakis N. Multi-physics optimal transportation and image interpolation. ESAIM: Mathematical Modelling and Numerical Analysis. 2015; 49(6):1671–1692.
29. Joshi SC, Miller MI. Landmark matching via large deformation diffeomorphisms. IEEE Transactions on Image Processing. 2000; 9(8):1357–1370. [PubMed: 18262973]
30. Kadous MW. PhD thesis. The University of New South Wales; 2002. Temporal classification: Extending the classification paradigm to multivariate time series.
31. Kantorovich LV. On the translocation of masses. Dokl Akad Nauk SSSR. 1942; 37:199–201.
32. Kantorovich LV. A problem of Monge. Uspekhi Mat Nauk. 1948; 3(24):225–226.
33. Khan AM, Rajpoot N, Treanor D, Magee D. A nonlinear mapping approach to stain normalization in digital histopathology images using image-specific color deconvolution. IEEE Transactions on Biomedical Engineering. 2014; 61(6):1729–1738. [PubMed: 24845283]
34. Kolouri S, Park S, Rohde GK. The Radon cumulative distribution transform and its application to image classification. Image Processing, IEEE Transactions on. 2016; 25(2):920–934.
35. Kolouri S, Park S, Thorpe M, Slepčev D, Rohde GK. Transport-based analysis, modeling, and learning from signal and data distributions. arXiv. 2016 1609.04767.
36. Kolouri S, Rohde GK. Transport-based single frame super resolution of very low resolution face images; Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on; 2015. 4876–4884.



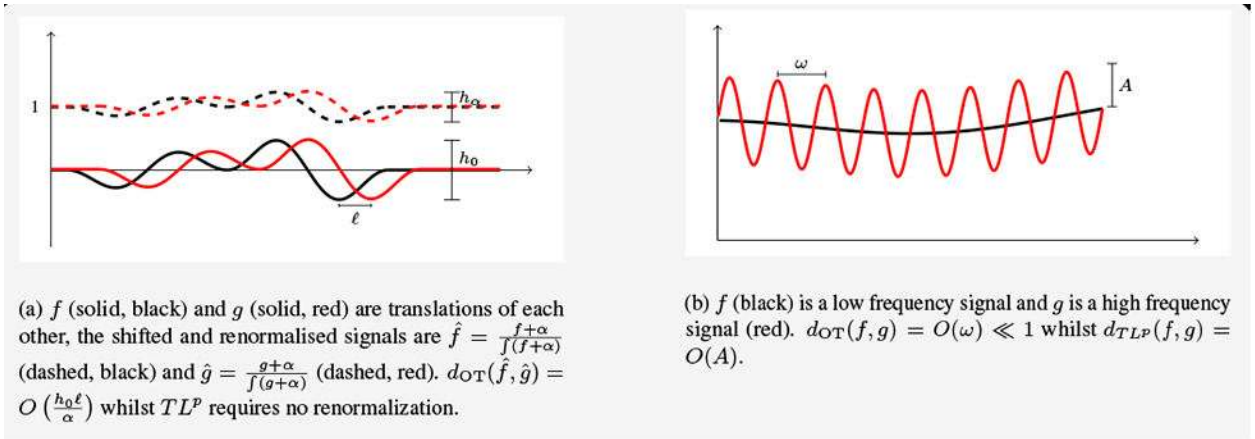
37. Kolouri S, Tosun AB, Ozolek JA, Rohde GK. A continuous linear optimal transport approach for pattern analysis in image datasets. *Pattern Recognition*. 2016; 51:453–462. [PubMed: 26858466]
38. Kondratyev S, Monsaingeon L, Vorotnikov D. A new optimal transport distance on the space of finite radon measures. *Advances in Differential Equations*. 2016; 21(11-12):1117–1164.
39. Kruskal JB. Nonmetric multidimensional scaling: A numerical method. *Psychometrika*, 1964; 29(2):115–129.
40. Lellmann J, Lorenz DA, Schönlieb C, Valkonen T. Imaging with Kantorovich-Rubinstein discrepancy. *SIAM Journal on Imaging Sciences*. 2014; 7(4):2833–2859.
41. Lichman M. UCI machine learning repository. 2013
42. Liero M, Mielke A, Savaré G. Optimal transport in competition with reaction: the hellinger-kantorovich distance and geodesic curves. arXiv. 2015 1509.00068.
43. Liero M, Mielke A, Savaré G. Optimal entropy-transport problems and a new Hellinger-Kantorovich distance between positive measures. arXiv. 2016 1508.07941.
44. Lipman Y, Daubechies I. Conformal Wasserstein distances: Comparing surfaces in polynomial time. *Advances in Mathematics*. 2011; 227(3):1047–1077.
45. Magee D, Treanor D, Crellin D, Shires M, Smith K, Mohee K, Quirke P. Colour normalisation in digital histopathology images. *Proc Optical Tissue Image analysis in Microscopy, Histopathology and Endoscopy (MICCAI Workshop)*. 2009; 100
46. Mérigot Q. A multiscale approach to optimal transport. *Computer Graphics Forum*. 2011; 30(5): 1583–1592.
47. Monge G. Mémoire sur la théorie des déblais et des remblais. *Histoire de l'Académie Royale des Sciences*. 1781:666–704.
48. Montavon G, Müller KR, Cuturi M. Wasserstein training of Boltzmann machines. arXiv. 2015 1507.01972.
49. Morovic J, Sun PL. Accurate 3d image colour histogram transformation. *Pattern Recognition Letters*. 2003; 24(11):1725–1735.
50. Museyko O, Stiglmayr M, Klamroth K, Leugering G. On the application of the Monge-Kantorovich problem to image registration. *SIAM Journal on Imaging Sciences*. 2009; 2(4):1068–1097.
51. Nikolova M, Wen YW, Chan R. Exact histogram specification for digital images using a variational approach. *Journal of Mathematical Imaging and Vision*. 2013; 46(3):309–325.
52. Ning L, Georgiou TT, Tannenbaum A. Matrix-valued monge-kantorovich optimal mass transport; 52nd IEEE Conference on Decision and Control; 2013. 3906–3911.
53. Oberman AM, Ruan Y. An efficient linear programming method for optimal transportation. arXiv. 2015 1509.03668.
54. Oudre L, Jakubowicz J, Bianchi P, Simon C. Classification of periodic activities using the Wasserstein distance. *IEEE Transactions on Biomedical Engineering*. 2012; 59(6):1610–1619. [PubMed: 22434794]
55. Ozolek JA, Tosun AB, Wang W, Chen C, Kolouri S, Basu S, Huang H, Rohde GK. Accurate diagnosis of thyroid follicular lesions from nuclear morphology using supervised learning. *Medical Image Analysis*. 2014; 18(5):772–780. [PubMed: 24835183]
56. Papadakis N, Bugeau A, Caselles V. Image editing with spatiograms transfer. *IEEE Transactions on Image Processing*. 2012; 21(5):2513–2522. [PubMed: 22249712]
57. Paris S, Durand F. A fast approximation of the bilateral filter using a signal processing approach. *International Journal of Computer Vision*. 2009; 81(1):24–52.
58. Park S, Kolouri S, Kundu S, Rohde G. The cumulative distribution transform and linear pattern classification. arXiv. 2015 1507.05936.
59. Pele O, Werman M. A linear time histogram metric for improved sift matching; European conference on computer vision; 2008. 495–508.
60. Pele O, Werman M. Fast and robust Earth Mover's Distances; 2009 IEEE 12th International Conference on Computer Vision; 2009. 460–467.
61. Pitié F, Kokaram A. The linear Monge-Kantorovich linear colour mapping for example-based colour transfer; 4th European Conference on Visual Media Production; 2007. 1–9.

62. Rabin J, Ferradans S, Papadakis N. Adaptive color transfer with relaxed optimal transport; 2014 IEEE International Conference on Image Processing (ICIP); 2014. 4852–4856.
63. Rabin J, Papadakis N. chapter Non-convex Relaxation of Optimal Transport for Color Transfer Between Images. Springer International Publishing; 2015. Geometric Science of Information: Second International Conference, GSI2015, Palaiseau, France, October 28-30, 2015, Proceedings; 87–95.
64. Rabin J, Peyré G. Wasserstein regularization of imaging problem; Image Processing (ICIP), 2011 18th IEEE International Conference on; 2011. 1521–1544.
65. Rabin J, Peyré G, Cohen LD. chapter Geodesic Shape Retrieval via Optimal Mass Transport. Springer Berlin Heidelberg; 2010. Computer Vision - ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part V; 771–784.
66. Reinhard E, Adhikhmin M, Gooch B, Shirley P. Color transfer between images. IEEE Computer Graphics and Applications. 2001; 21(5):34–41.
67. Rockafellar RT. Convex Analysis. Princeton University Press; 1970.
68. Rubner Y, Tomasi C, Guibas LJ. The Earth Mover's Distance as a metric for image retrieval. International Journal of Computer Vision. 2000; 40(2):99–121.
69. Russell EJ. Letters to the editor-extension of Dantzig's algorithm to finding an initial near-optimal basis for the transportation problem. Operations Research. 1969; 17(1):187–191.
70. Samaria FS, Harter AC. Parameterisation of a stochastic model for human face identification. Proceedings of 1994 IEEE Workshop on Applications of Computer Vision. 1994:138–142.
71. Schmitzer B. chapter A sparse algorithm for dense optimal transport. Springer International Publishing; 2015. Scale Space and Variational Methods in Computer Vision: 5th International Conference, SSVN 2015, Lge-Cap Ferret, France, May 31 - June 4, 2015, Proceedings; 629–641.
72. Schmitzer B. Stabilized sparse scaling algorithms for entropy regularized transport problems. arXiv. 2016 1610.06519.
73. Shinohara RT, Sweeney EM, Goldsmith J, Shiee N, Mateen FJ, Calabresi PA, Jarso S, Pham DL, Reich DS, Crainiceanu CM. Statistical normalization techniques for magnetic resonance imaging. Neuroimage: Clinical. 2014
74. Solomon J, de Goes F, Peyré G, Cuturi M, Butscher A, Nguyen A, Du T, Guibas L. Convolutional Wasserstein distances: Efficient optimal transportation on geometric domains. ACM Trans Graph. 2015; 34(4):66–66. 1–66.
75. Solomon J, Rustamov R, Guibas L, Butscher A. Earth mover's distances on discrete surfaces. ACM Transactions on Graphics (TOG). 2104; 33(4):67.
76. Solomon J, Rustamov R, Leonidas G, Butscher A. Wasserstein propagation for semi-supervised learning. In: Jebara T, Xing EP, editors Proceedings of the 31st International Conference on Machine Learning (ICML-14). JMLR Workshop and Conference Proceedings; 2014. 306–214.
77. Su Z, Zeng W, Wang Y, Lu ZL, Gu X. Shape classification using Wasserstein distance for brain morphometry analysis. Information Processing in Medical Imaging. 2015; 24:411–423. [PubMed: 26221691]
78. Tannenbaum E, Georgiou T, Tannenbaum A. Signals and control aspects of optimal mass transport and the Boltzmann entropy; 49th IEEE Conference on Decision and Control (CDC); 2010. 1885–1890.
79. Thorpe M, Slepčev D. Transportation  $L^p$  distances: Properties and extensions. Preparation. 2016
80. Tosun AB, Yergiyev O, Kolouri S, Silverman JF, Rohde GK. Novel computer-aided diagnosis of mesothelioma using nuclear structure of mesothelial cells in effusion cytology specimens. Proc SPIE. 2014; 9041:90410Z–90410Z-6.
81. Tosun AB, Yergiyev O, Kolouri S, Silverman JF, Rohde GK. Detection of malignant mesothelioma using nuclear structure of mesothelial cells in effusion cytology specimens. Cytometry Part A. 2015; 87(4):326–333.
82. ur Rehman T, Haber E, Pryor G, Melonakos J, Tannenbaum A. 3D nonrigid registration via optimal mass transport on the gpu. Medical image analysis. 2009; 13(6):931–940. [PubMed: 19135403]
83. Villani C. Graduate Studies in Mathematics. American Mathematical Society; 2003. Topics in Optimal Transportation.

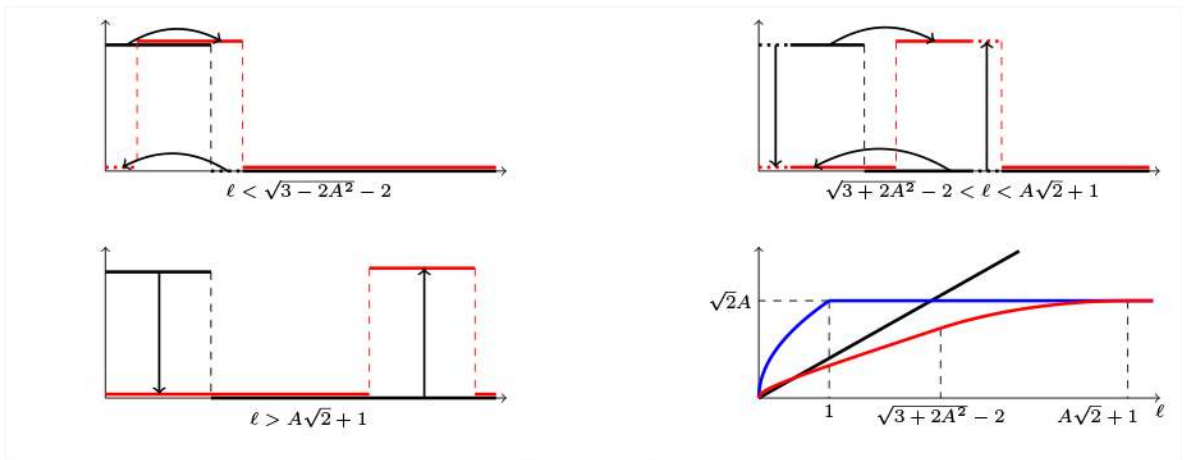
84. Villani C. Optimal transport: old and new. Springer Science & Business Media; 2008.
85. Wang W, Slepčev D, Basu S, Ozolek JA, Rohde GK. A linear optimal transportation framework for quantifying and visualizing variations in sets of images. *International Journal of Computer Vision*. 2012; 101(2):254–269.
86. Wang Z, Bovik AC, Sheikh HR, Simoncelli EP. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*. 2003; 13(4):600–612.
87. Weinberger KQ, Chapelle O. Large margin taxonomy embedding for document categorization. *Advances in Neural Information Processing Systems*. 2009:1737–1744.
88. Zhu L, Haker S, Tannenbaum A. Flattening maps for the visualization of multibranched vessels. *Medical Imaging, IEEE Transactions on*. 2005; 24(2):191–198.
89. Zhu L, Yang Y, Haker S, Tannenbaum A. An image morphing technique based on optimal mass preserving mapping. *Image Processing, IEEE Transactions on*. 2007; 16(6):1481–1495.



**Figure 1.** On the left an example optimal transport map for OT, on the right an illustration of the pushforward measure.



**Figure 2.**  
A Comparison of  $TL^P$  with OT.



**Figure 3.**  $TL^2$  transport between  $f(x) = A\chi_{[0,1]}$  (black) and  $g(x) = f(x - \ell)$  (red) and the  $TL^2$  distance (red),  $L^2$  distance (blue) and OT (black) between  $f$  and  $g$  (bottom right).

(C<sub>1</sub>) **One hump functions:** of the form

$$f = \chi_{[\ell, \ell + \alpha]}$$

where  $\ell \in [0, 1 - \alpha]$ .

(C<sub>2</sub>) **Two hump functions:** of the form

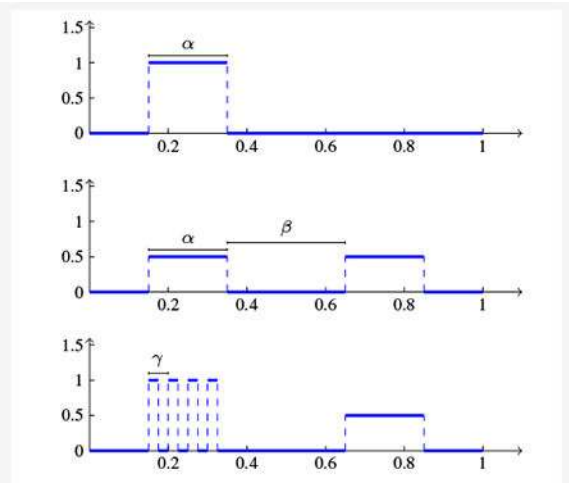
$$f = \frac{1}{2} (\chi_{[\ell, \ell + \alpha]} + \chi_{[\ell + \beta + \alpha, \ell + \beta + 2\alpha]})$$

where  $\ell \in [0, 1 - \beta - 2\alpha]$ .

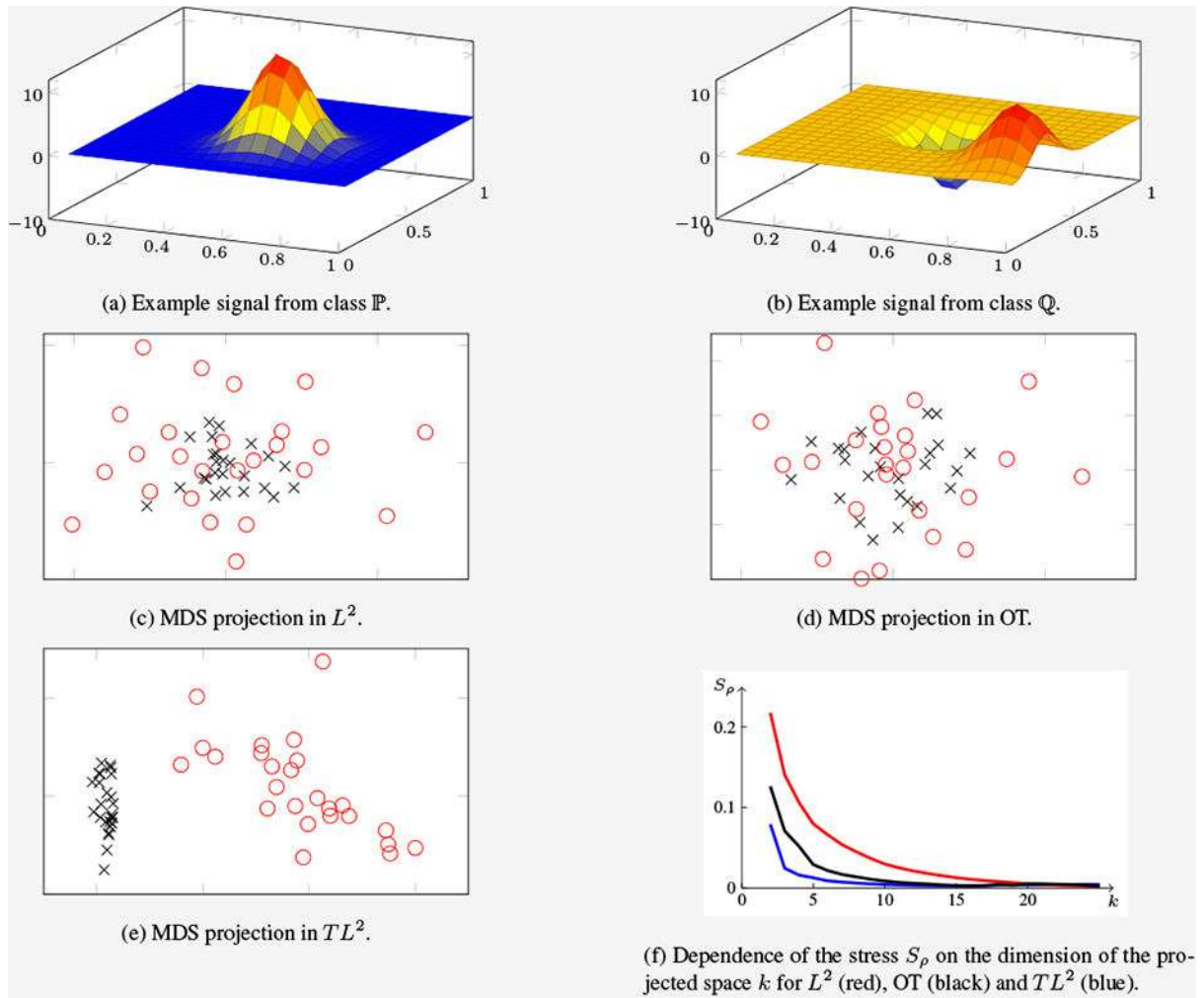
(C<sub>3</sub>) **One hump, one chirp functions:** of the form

$$f = \sum_{i=0}^{\frac{\alpha}{\gamma} - 1} \chi_{[\ell + i\gamma, \ell + \frac{(2i+1)\gamma}{2}]} + \frac{1}{2} \chi_{[\ell + \beta + \alpha, \ell + \beta + 2\alpha]}$$

where  $\ell \in [0, 1 - \beta - 2\alpha]$ .

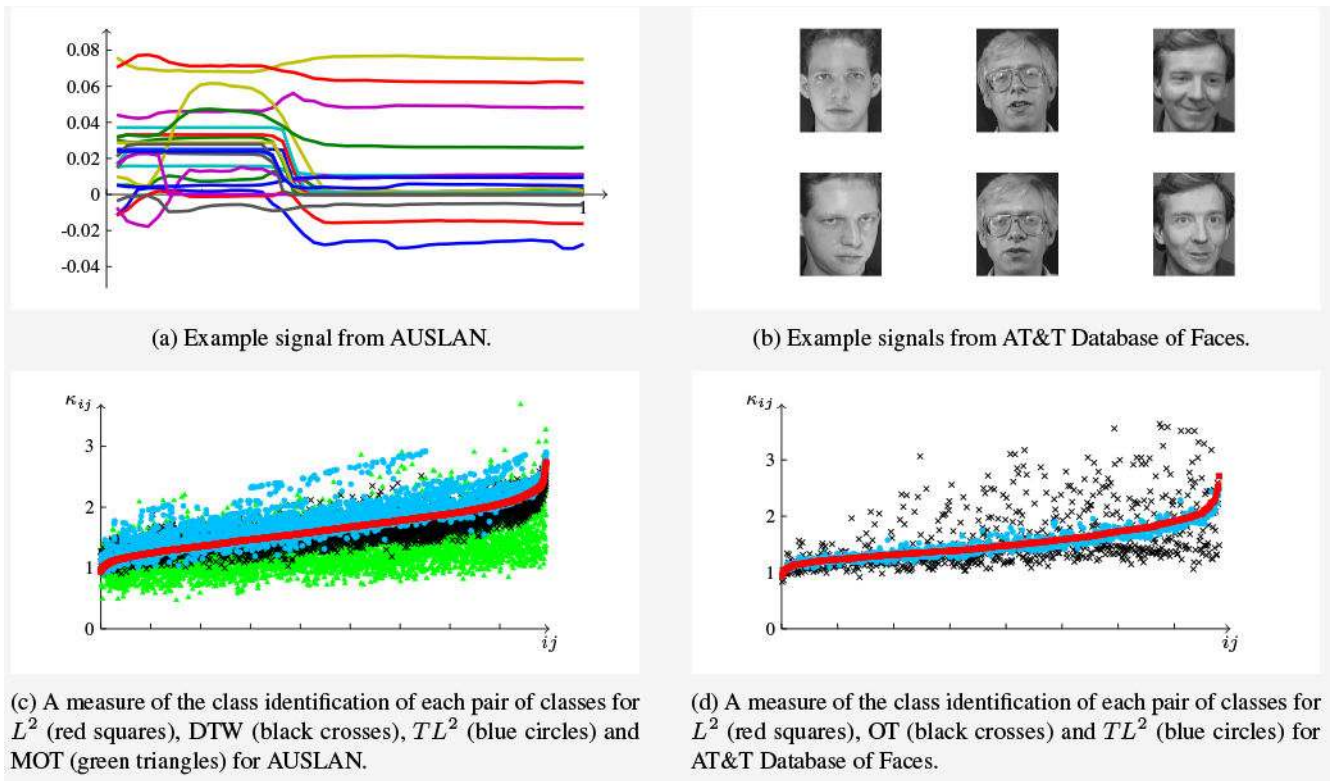


**Figure 4.** For fixed  $\alpha, \beta, \gamma \in (0, 1)$  where  $\beta > \alpha \gg \gamma$  the definition of the classes  $\mathcal{C}_i$ .

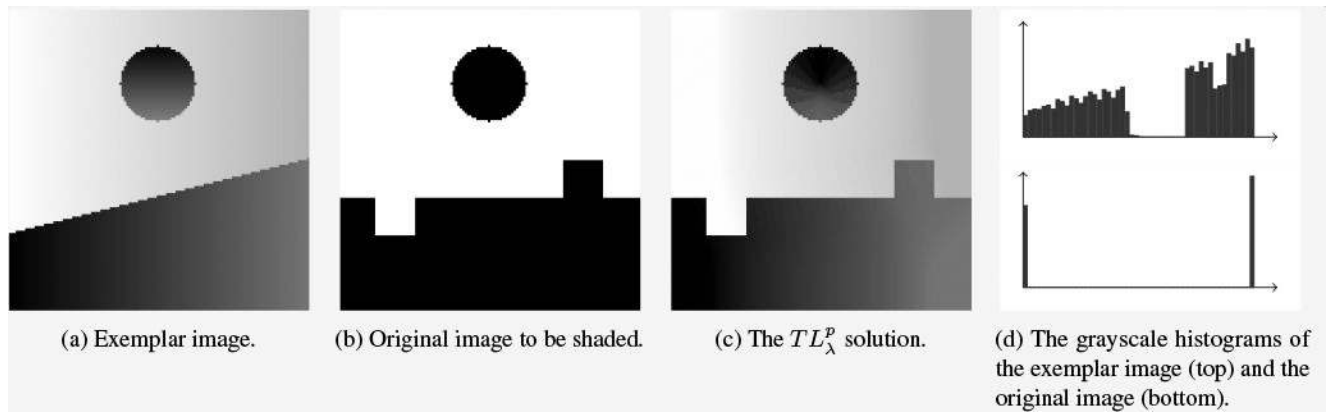


**Figure 5.** Example signals and results for the synthetic data in Section 4.2.



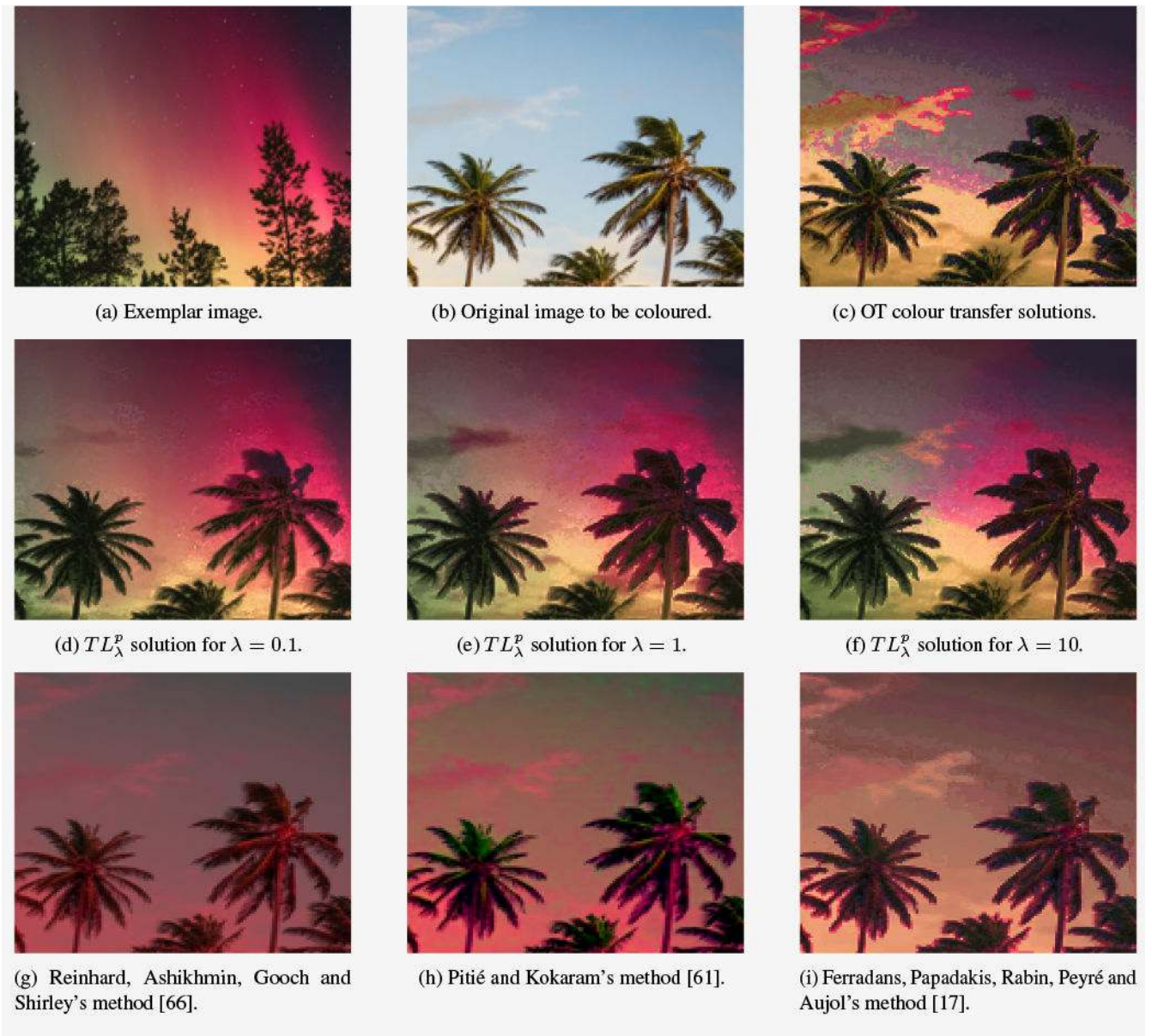


**Figure 6.**  
Example signals and results for the data sets described in Section 4.3.

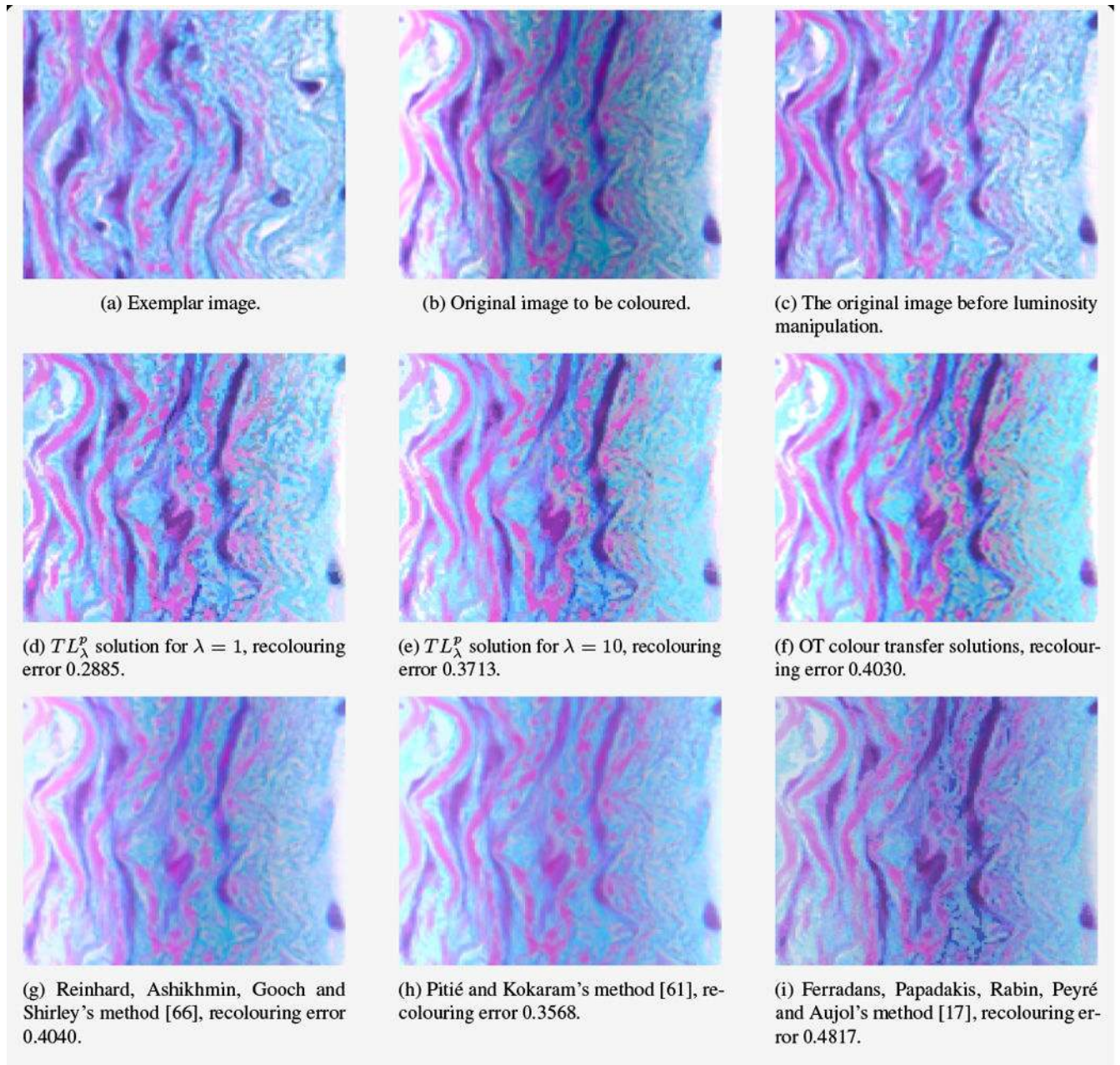


**Figure 7.**

Spatially correlated histogram specification of synthetic grayscale images. An OT induced solution does not exist since there are no transport maps between the image to be shaded and the exemplar image.



**Figure 8.** Spatially correlated histogram specification of real colour images between palm trees and the northern lights.



**Figure 9.** Spatially correlated histogram specification of real colour images between two images of Masson's trichrome staining procedure, one of which the luminosity has been manipulated.

**Table 1**

Error rates for 1NN classification in AUSLAN. The top row corresponds to applying each method to the signal, the second row is each method applied to the derivative of the signal, and the third row is the method applied to a weighted average of the signal and the derivative of the signal.

	$L^2$	DTW	$TL\frac{2}{\lambda}$	MOT
Signal	15.39%	11.45%	12.12%	61.71%
Derivative	22.15%	19.77%	12.63%	10.41%
Weighted Average	8.06%	7.33%	6.70%	10.41%