

PRECONDITIONING OF OPTIMAL TRANSPORT*

MAX KUANG[†] AND ESTEBAN G. TABAK[†]

Abstract. A preconditioning procedure is developed for the L_2 and more general optimal transport problems. The procedure is based on a family of affine map pairs which transforms the original measures into two new measures that are closer to each other, while preserving the optimality of solutions. It is proved that the preconditioning procedure minimizes the remaining transportation cost among all admissible affine maps. The procedure can be used on both continuous measures and finite sample sets from distributions. In numerical examples, the procedure is applied to multivariate normal distributions and to a two-dimensional shape transform problem.

Key words. preconditioning, optimal transport

AMS subject classifications. 65F08, 15A23, 49K45

1. Introduction. The original optimal transport problem, proposed by Monge in 1781 [14], asks how to move a pile of soil between two locations with minimal cost. Giving the cost $c(x, y)$ of moving a unit mass from point x to point y , one seeks the map $y = T(x)$ that minimizes its integral. After normalizing the two piles so that each has total mass one and can be regarded as a probability measure, the problem adopts the form

$$(1) \quad \min_{T_{\#}\mu=\nu} \int c(x, T(x)) d\mu(x),$$

where μ and ν are the original and target measures, and $T_{\#}\mu$ denotes the push forward measure of μ by the map T .

In the 20th century, Kantorovich [10] relaxed Monge's definition, allowing the movement of soil from one location to multiple destinations and vice versa. Denoting the mass moved from x to y by $\pi(x, y)$, we can rewrite the minimization problem as

$$(2) \quad \min_{\pi} \int c(x, y) \pi(x, y) dx dy$$

among couplings $\pi(x, y)$ satisfying the marginal constraints

$$\int \pi(x, y) dy = \mu(x)$$
$$\int \pi(x, y) dx = \nu(y).$$

Since the second half of the 20th century, mathematical properties of the optimal transport solution have been studied extensively, as well as applications in many different areas (see for instance [16, 12, 3, 7, 8, 4], or [20] for a comprehensive list.). Since closed-form solutions of the multi-dimensional optimal transport problems are relatively rare, a number of numerical algorithms have been proposed. We reference below some recent representatives of the different approaches taken:

PDE methods: Benamou and Brenier [2] introduced a computational fluid approach to solve the problem with continuous distributions $\mu_{1,2}$, exploiting the structure of the interpolant of the optimal map to solve the PDE corresponding to the optimization problem in the dual variables.

*This work was funded by the Office of Naval Research.

[†]Courant Institute, New York University, New York, NY 10012 (kuang@cims.nyu.edu, tabak@cims.nyu.edu).

39 **Adaptive Linear Programming:** Oberman and Ruan [15] discretized the given
 40 continuous distributions and solved the resulting linear programming problem
 41 in an adaptive way that exploits the sparse nature of the solution (the fact
 42 that the optimal plan has support on a map.)

43 **Entropy Regularization:** The discrete version of optimal transport is the earth
 44 mover’s problem in image processing [17], a linear programming problem
 45 widely used to measure the distance between images and in networks. Recent
 46 development on entropy regularization [18] introduced effective algorithms to
 47 solve regularized versions of these problems.

48 **Data-driven Formulations:** Data-driven formulations take as input not the distri-
 49 butions $\mu_{1,2}$ but sample sets from both. Methodologies proposed include a
 50 fluid-flow-like algorithm [19], an adaptive linear programming approach [5],
 51 and a procedure based on approximating the interpolant in a feature-space
 52 [11].

53 In this paper, we introduce a novel procedure to precondition the input probability
 54 measures or samples thereof, so that the resulting measures or sample sets are closer
 55 to each other while preserving the optimality of solutions. The procedure and its
 56 properties are discussed for both L_2 and more general cost functions induced by an
 57 inner product.

58 In theoretical applications, the preconditioning procedure is used to give alterna-
 59 tive derivations of a lower bound for the total transportation cost and of the optimal
 60 map between multivariate normal distributions. For practical applications, we use
 61 the procedure on sample sets to get preconditioned sets, which are then given as
 62 input to optimal transport algorithms to calculate the optimal map. Inverting the
 63 the preconditioning map pairs used, we recover the optimal map between the original
 64 distributions.

65 **2. Optimal Transport.** Let μ and ν be two probability measures on the same
 66 sample space \mathcal{X} . Optimal transport asks how to optimally move the mass from μ to
 67 ν , given a function $c(x, y)$ represents the cost of moving a unit of mass from point
 68 x to point y . Monge’s formulation seeks a map $y = T(x)$ that minimizes the total
 69 transportation cost:

$$70 \quad (3) \quad \min_{T_{\#}\mu=\nu} \mathbb{E}_{\mu} c(X, T(X)),$$

71 where $T_{\#}\mu$ represents the pushforward measure of μ through the map T .

72 A transfer plan $\pi(x, y)$ is the law of a coupling (X, Y) between the two measures
 73 μ and ν . For any measurable set $E \subset \mathcal{X}$,

$$74 \quad \pi(E \times \mathcal{X}) = \mu(E), \quad \pi(\mathcal{X} \times E) = \nu(E).$$

75 Denoting the family of all transfer plans by $\Pi(\mu, \nu)$, Kantorovich’s relaxation of the
 76 optimal transport problem is

$$77 \quad (4) \quad \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{\pi} c(X, Y).$$

78 Since the maps $Y = T(X)$ represent a subset of all couplings between μ and ν , the
 79 feasible domain for (3) lies within the one for (4).

80 While there are many results on the general optimal transport problem, a par-
 81 ticularly well-studied and useful case is the L_2 optimal transport on \mathcal{R}^N , in which μ
 82 and ν are probability measures on \mathcal{R}^N and the cost function $c(x, y)$ is given by the

83 squared Euclidean distance $\|x - y\|^2$. In this case, with moderate requirements, one
 84 can prove that the solution to Kantorovich's relaxation (4) is unique and agrees with
 85 the solution to Monge's problem (3). In other words, the unique optimal coupling
 86 (X, Y) corresponds to a map $Y = T(X)$. Moreover, this optimal map is the gradient
 87 of a convex potential ϕ , so we have the following statement:

88 **THEOREM 1.** *For Kantorovich's relaxation (4) with the L_2 cost function and ab-*
 89 *solute continuous measures μ and ν , the optimal coupling (X, Y) is a map $Y = T(X)$,*
 90 *where $T : \mathcal{R}^N \rightarrow \mathcal{R}$ is defined by*

91 (5)
$$T(x) = \nabla\phi(x)$$

92 where $\phi(x)$ is convex and $T_{\#}\mu = \nu$.

93 While this characterization of the solution is attractively simple, closed-form solutions
 94 of the L_2 optimal transport on \mathcal{R}^N are rare for $N > 1$. The difficulties of deriving
 95 closed-form solutions boosted research to solve the optimal transport problem numer-
 96 ically. An incomplete list of formulations and methods can be found in section 1.

97 The goal of this paper is not to provide a complete numerical recipe to solve L_2
 98 optimal transport problems, but to introduce a practical preconditioning procedure.
 99 This procedure transforms the original measures μ and ν into two new measures, so
 100 that the optimal transport problems between the new measures is easier to solve,
 101 while the optimality of solutions is preserved by the transformation. The procedure
 102 extends beyond L_2 to any cost function induced by an inner product.

103 **3. Admissible Map Pairs.** The basic framework of the preconditioning proce-
 104 dure is as follows:

105
$$\begin{array}{ccc} X \sim \mu & \xrightarrow{Y=G^{-1}(T(F(X)))} & Y \sim \nu \\ \downarrow \tilde{X}=F(X) & & \downarrow \tilde{Y}=G(Y) \\ \tilde{X} \sim \tilde{\mu} & \xrightarrow{\tilde{Y}=T(\tilde{X})} & \tilde{Y} \sim \tilde{\nu} \end{array}$$

106 Suppose that we transform μ and ν into two new measures $\tilde{\mu}$ and $\tilde{\nu}$ via some
 107 invertible maps F and G and that the solution to the new L_2 optimal transport
 108 problem between $\tilde{\mu}$ and $\tilde{\nu}$ is given by $\tilde{Y} = T(\tilde{X})$. Then the map

109 (6)
$$Y = G^{-1}(T(F(X)))$$

110 pushes forward μ into ν . We call the pair of invertible maps (F, G) an *admissible map*
 111 *pair* if the resulting map (6) is optimal for the original problem between μ and ν .

112 There are several simple *admissible map pairs*.

113 **DEFINITION 2 (Translation Pairs).** *Given two vectors m_1, m_2 in \mathcal{R}^N , a Transla-*
 114 *tion Pair (F, G) is defined by*

115 (7)
$$F(X) = X - m_1, \quad G(Y) = Y - m_2.$$

116 If $\tilde{Y} = T(\tilde{X})$ is an optimal map, then $T = \nabla\phi$ for some convex function ϕ , which
 117 implies that

118 (8)
$$Y = m_2 + T(X - m_1) = \nabla [m_2X + \phi(X - m_1)],$$

119 so $Y = B^{-1}(T(A(X)))$ is indeed the optimal map between μ and ν . Thus *translation*
 120 *pairs are admissible map pairs.*

121 Among all *translation pairs*, we can minimize the total transportation cost in the
 122 new problem:

$$\begin{aligned} 123 \quad \mathbb{E}\|\tilde{X} - \tilde{Y}\|^2 &= \mathbb{E}\|X - m_1 - Y + m_2\|^2 \\ 124 \quad &= \mathbb{E}\|X - \mathbb{E}X - Y + \mathbb{E}Y\|^2 + \|\mathbb{E}X - m_1 - \mathbb{E}Y + m_2\|^2 \\ 125 \quad &\geq \mathbb{E}\|X - \mathbb{E}X - Y + \mathbb{E}Y\|^2 \end{aligned}$$

127 This shows that the transportation cost between \tilde{X} and \tilde{Y} is minimized when $\mathbb{E}X -$
 128 $\mathbb{E}Y = m_1 - m_2$. In particular, we can adopt $m_1 = \mathbb{E}X$ and $m_2 = \mathbb{E}Y$, which gives
 129 both measures a zero mean. We call the corresponding *translation pair* the *mean*
 130 *translation pair.*

131 **DEFINITION 3** (Scaling Pairs). *Given two nonzero numbers α, β in \mathcal{R} , the Scaling*
 132 *Pair (F, G) is defined by:*

$$133 \quad (9) \quad F(X) = \alpha X, \quad G(Y) = \beta Y.$$

134 Clearly if $\tilde{Y} = T(\tilde{X}) = \nabla\phi(\tilde{X})$ is an optimal map,

$$135 \quad (10) \quad Y = \frac{1}{\beta}T(\alpha X)$$

136 is also an optimal map. So all the *scaling pairs* are *admissible map pairs*. In particular,
 137 one can choose

$$138 \quad \alpha = \frac{1}{\sqrt{\mathbb{E}\|X\|^2}}, \quad \beta = \frac{1}{\sqrt{\mathbb{E}\|Y\|^2}},$$

139 so that

$$140 \quad \mathbb{E}\|\tilde{X}\|^2 = \mathbb{E}\|\tilde{Y}\|^2 = 1.$$

141 We call this specific *scaling pair* the *normalizing scaling pair*.

142 Next we discuss general linear *admissible map pairs*. We will think of X as row
 143 vectors, so the matrices representing linear transformations act on X on the right.

144 **THEOREM 4.** *Let $F(X) = XA$ and $G(Y) = YB$, where $A, B \in \mathcal{R}^{N \times N}$ are invert-*
 145 *ible matrices. Denote by $\tilde{Y} = T(\tilde{X})$ the optimal map from $\tilde{\mu}$ to $\tilde{\nu}$. If $B = (A^T)^{-1}$,*
 146 *the induced map between μ and ν is also optimal.*

147 *Proof.* The induced map can be written as

$$148 \quad Y = T(XA)B^{-1} = T(XA)A^T$$

149 Let $T(X) = \nabla\phi(X)$ and $\psi(X) = \phi(XA)$ we have

$$150 \quad (11) \quad Y_i = \sum_{j=1}^N \phi_j(XA)(A^T)_{ij} = \frac{\partial}{\partial X_i} \phi(XA) \Rightarrow Y = \nabla\psi(X).$$

151 Since ψ is also a convex function, the induced map $Y = T(XA)B^{-1}$ is also an optimal
 152 map. \square

153 *Remark 5.* Another way to understand this theorem is to consider map pairs
 154 (F, G) that do not alter the inner product. In fact, the theorem holds if, for any
 155 $x, y \in \mathcal{R}^N$,

$$156 \quad (12) \quad xy^T = F(x)G(y)^T.$$

157 This observation implies that the same result holds for more general cost functions:
 158 as long as the metric $d(x, y)$ is induced by an inner product $\langle x, y \rangle$, we only need the
 159 pair F and G to be adjoint operators to guarantee they form an *admissible map pair*.

160 The above theorem gives us a family of new *admissible map pairs*.

161 **DEFINITION 6 (Linear Pairs).** *Let A be an invertible matrix in $\mathcal{R}^{N \times N}$, the linear*
 162 *pair (F, G) is defined by:*

$$163 \quad (13) \quad F(X) = XA, \quad G(Y) = Y(A^T)^{-1}$$

164 We first give some examples of common *linear pairs*,

165 **DEFINITION 7 (Orthogonal Pairs).** *For any orthogonal matrix A ,*

$$166 \quad (14) \quad F(X) = XA, \quad G(Y) = YA$$

167 *is called a orthogonal map pair.*

168 For orthogonal pairs, we have $(A^T)^{-1} = A$. This means that performing the
 169 same orthogonal linear transformation on both measures preserves the optimality of
 170 solutions. The interpretation of this result is straightforward, as an orthogonal map
 171 yields a distance-preserving coordinate change which does not alter the cost function.

172 **DEFINITION 8 (Stretching Pairs).** *For any unit vector d and scalar α , we can*
 173 *stretch X by a factor of α along d , and at the same time stretch Y by a factor of $1/\alpha$*
 174 *along the same direction:*

$$175 \quad (15) \quad \begin{aligned} F(X) &= X - (Xd^T)d + \alpha(Xd^T)d = X(I + (\alpha - 1)d^T d) \\ G(Y) &= Y - (Yd^T)d + 1/\alpha(Yd^T)d = X(I + (1/\alpha - 1)d^T d) \end{aligned}$$

176 *We call such map pairs stretching pairs.*

177 It can be verified this is indeed a *linear pair*, and thus an *admissible map pair*.

178 Composing translation and linear pairs, one obtains a more general class of *affine*
 179 *pairs*. Among all *affine pairs*, we seek the optimal one for our preconditioning proce-
 180 *dure*. We first state a linear algebra result:

181 **THEOREM 9.** *For any two positive-definite matrices Σ_1 and Σ_2 in $\mathcal{R}^{N \times N}$, there*
 182 *exists an invertible matrix $A \in \mathcal{R}^{N \times N}$ such that*

$$183 \quad (16) \quad D = A^T \Sigma_1 A = A^{-1} \Sigma_2 (A^T)^{-1}$$

184 *where D is a diagonal matrix with entries satisfying*

$$185 \quad (17) \quad d_1 \geq d_2 \geq \dots \geq d_N > 0.$$

186 *In addition, D is unique.*

187 *Proof.* We first prove the existence of A . Since $\Sigma_1^{1/2}$ is invertible, we can replace
 188 A by a matrix B satisfying

$$189 \quad B = \Sigma_1^{1/2} A$$

190 and

$$191 \quad D = B^T B = B^{-1} \Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} (B^T)^{-1}.$$

192 Because $\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2}$ is positive definite, it admits an eigenvalue decomposition of the
193 form

$$194 \quad (18) \quad \Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} = Q \Lambda Q^T,$$

195 with Q orthogonal and Λ diagonal with sorted, positive diagonal entries. Setting
196 $B = Q \Lambda^{1/4}$, we have

$$197 \quad B^T B = \Lambda^{1/2}$$

198 and

$$199 \quad B^{-1} \Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} (B^T)^{-1} = \Lambda^{-1/4} Q^T Q \Lambda Q^T Q \Lambda^{-1/4} = \Lambda^{1/2}.$$

200 Thus the conditions of the theorem are satisfied with

$$201 \quad (19) \quad D = \Lambda^{1/2}, \quad A = \Sigma_1^{-1/2} Q \Lambda^{1/4}.$$

202 To prove the uniqueness of D , suppose that there are D_1, A_1 and D_2, A_2 such that

$$203 \quad D_1 = A_1^T \Sigma_1 A_1 = A_1^{-1} \Sigma_2 (A_1^T)^{-1}$$

$$204 \quad D_2 = A_2^T \Sigma_1 A_2 = A_2^{-1} \Sigma_2 (A_2^T)^{-1}.$$

206 Then

$$207 \quad D_1^2 = A_1^{-1} \Sigma_2 \Sigma_1 A_1$$

$$208 \quad D_2^2 = A_2^{-1} \Sigma_2 \Sigma_1 A_2,$$

210 implying that D_1^2 , $\Sigma_2 \Sigma_1$ and D_2^2 are similar to each other. Since D_1 and D_2 are
211 positive diagonal matrices with sorted entries, they must be identical, which proves
212 the uniqueness of D . \square

213 Using the theorem above, we can define the following *optimal linear pair*:

214 DEFINITION 10 (*Optimal Linear Pair*). Assume that μ and ν are mean-zero mea-
215 sures with covariance matrices Σ_1 and Σ_2 , and let A be a $N \times N$ matrix that satisfies
216 (16). We define the optimal linear pair (F, G) through:

$$217 \quad (20) \quad F(X) = XA, \quad G(Y) = Y(A^T)^{-1}.$$

218 (Notice that the matrix A can be constructed following (18) and (19) in the proof of
219 Theorem 9.)

220 This pair has the following useful properties:

221 PROPERTY 11. The resulting random variables \tilde{X}, \tilde{Y} derived from the optimal
222 linear pair have the same diagonal covariance matrix D :

$$223 \quad (21) \quad \mathbb{E} \tilde{X}^T \tilde{X} = A^T \Sigma_1 A = D$$

$$224 \quad (22) \quad \mathbb{E} \tilde{Y}^T \tilde{Y} = A^{-1} \Sigma_2 (A^T)^{-1} = D.$$

226 PROPERTY 12. Among all possible linear pairs $X' = XC, Y' = Y(C^T)^{-1}$ given
 227 by an invertible matrix C , the optimal linear pair minimizes $\mathbb{E}\|X' - Y'\|^2$. In other
 228 words, for any invertible matrix C :

$$229 \quad (23) \quad \mathbb{E}\|X' - Y'\|^2 \geq \mathbb{E}\|\tilde{X} - \tilde{Y}\|^2.$$

230 *Proof.* For any matrix C , we have:

$$\begin{aligned} 231 \quad \mathbb{E}\|X' - Y'\|^2 &= \mathbb{E}X'X'^T + \mathbb{E}Y'Y'^T - 2\mathbb{E}X'Y'^T \\ 232 &= \mathbb{E}XC^T C^T X^T + \mathbb{E}Y(C^T)^{-1}C^{-1}Y^T - 2\mathbb{E}XY^T \\ 233 &= \mathbb{E}\operatorname{tr}(C^T X^T X C) + \mathbb{E}\operatorname{tr}(C^{-1}Y^T Y (C^T)^{-1}) - 2\mathbb{E}XY^T \\ 234 &= \operatorname{tr}(C^T \Sigma_1 C) + \operatorname{tr}(C^{-1} \Sigma_2 (C^T)^{-1}) - 2\mathbb{E}XY^T. \end{aligned}$$

236 On the other hand, (16) is equivalent to

$$237 \quad \Sigma_1 = (A^T)^{-1}DA^{-1}, \quad \Sigma_2 = ADA^T.$$

238 In terms of $S = A^{-1}C$,

$$\begin{aligned} 239 \quad \mathbb{E}\|X' - Y'\|^2 &= \operatorname{tr}(S^T D S) + \operatorname{tr}(S^{-1}D(S^T)^{-1}) - 2\mathbb{E}XY^T \\ 240 &= \operatorname{tr}(S S^T D) + \operatorname{tr}((S S^T)^{-1}D) - 2\mathbb{E}XY^T. \end{aligned}$$

242 Writing $S = (s_1, s_2, \dots, s_N)^T$ and $(S^T)^{-1} = (z_1, z_2, \dots, z_N)^T$, we have

$$\begin{aligned} 243 \quad \mathbb{E}\|X' - Y'\|^2 &= \sum_{i=1}^N d_i s_i^T s_i + \sum_{i=1}^N d_i z_i^T z_i - 2\mathbb{E}XY^T \\ 244 &= \sum_{i=1}^N d_i (s_i^T s_i + z_i^T z_i) - 2\mathbb{E}XY^T \\ 245 &\geq \sum_{i=1}^N d_i (2s_i^T z_i) - 2\mathbb{E}XY^T \\ 246 &= 2 \sum_{i=1}^N d_i - 2\mathbb{E}XY^T \\ 247 &= \mathbb{E}\|\tilde{X} - \tilde{Y}\|^2. \end{aligned}$$

249 Notice that we have the equal sign when $S = I$, which means that $C = A$. Thus □

$$250 \quad \mathbb{E}\|X' - Y'\|^2 \geq 2 \sum_{i=1}^N d_i - 2\mathbb{E}XY^T = \mathbb{E}\|\tilde{X} - \tilde{Y}\|^2.$$

251 Composing the *mean translation pair* and the *optimal linear pair* one obtains the
 252 *optimal affine pair*. It follows from the properties above that the *optimal affine pair*
 253 not only gives the two distributions zero means and transforms the covariance matrices
 254 into diagonal matrices, but also minimizes the distance between $\tilde{\mu}$ and $\tilde{\nu}$ among all
 255 *affine pairs*.

256 **4. Admissible Map Pairs For General Cost Functions.** In [Theorem 4](#), we
 257 introduced a class of affine maps that preserves the optimality of solutions for L_2 cost.
 258 As mentioned in the remark, similar results hold for more general cost functions. For
 259 cost functions induced by an inner product, we have the following generalization of
 260 [Theorem 4](#):

261 **THEOREM 13.** *Let $\langle \cdot, \cdot \rangle$ be an inner product in \mathcal{R}^N . For the optimal transport*
 262 *problem with cost*

$$263 \quad (24) \quad c(x, y) = \langle x - y, x - y \rangle,$$

264 *we have (F, G) is an admissible map pair if F and G are adjoint operators with respect*
 265 *to inner product $\langle \cdot, \cdot \rangle$.*

Proof. It follows from the fact that $c(x, y) = \|x\|^2 + \|y\|^2 - 2\langle x, y \rangle$, where only the last term depends on the actual coupling between X and Y , that

$$\operatorname{argmin} E [c(X, Y)] = \operatorname{argmax} E [\langle X, Y \rangle].$$

266 Since this applies to both the original and the preconditioned problems, their optimal
 267 solutions satisfy

$$268 \quad (X^*, Y^*) = \operatorname{argmax} [E \langle X, Y \rangle] \quad \text{and} \quad (\tilde{X}^*, \tilde{Y}^*) = \operatorname{argmax} [E \langle \tilde{X}, \tilde{Y} \rangle].$$

269 But if F and G are adjoint,

$$270 \quad \langle \tilde{X}, \tilde{Y} \rangle = \langle F(X), G(Y) \rangle = \langle X, Y \rangle,$$

271 so

$$272 \quad (\tilde{X}^*, \tilde{Y}^*) = (F(X^*), G(Y^*)),$$

273 proving the conclusion. □

274 Any inner product on \mathcal{R}^N can be written in terms of the standard vector multi-
 275 plication, through the introduction of a positive definite kernel matrix K :

$$276 \quad (25) \quad \langle x, y \rangle = xKy^T,$$

277 so stating that the linear operators $F(X) = XA, G(Y) = YB$ are adjoint is equivalent
 278 to

$$279 \quad (26) \quad AKB^T = K.$$

280 We can also derive the *optimal linear pair* for general cost functions. Here we only
 281 state without proof the core linear algebra theorem.

282 **THEOREM 14.** *Let Σ_1, Σ_2 and K be positive-definite matrices in $R^{N \times N}$. There*
 283 *exist invertible matrices $A, B \in R^{N \times N}$ such that*

$$284 \quad (27) \quad AKB^T = K$$

285 and

$$286 \quad (28) \quad D = K^{1/2}A^T\Sigma_1AK^{1/2} = K^{1/2}B^T\Sigma_2BK^{1/2}$$

287 where D is a unique diagonal matrix with entries satisfying

$$288 \quad (29) \quad d_1 \geq d_2 \geq \dots \geq d_N > 0.$$

289 Matrices constructed so as to satisfy the above theorem give the *optimal linear pairs*
 290 with respect to the corresponding cost. Notice that in this case the resulting measures
 291 no longer have diagonal covariance matrices:

$$292 \quad (30) \quad \mathbb{E}\tilde{X}^T\tilde{X} = \mathbb{E}\tilde{Y}^T\tilde{Y} = K^{-1/2}DK^{-1/2}.$$

293 **5. Preconditioning Procedure and Its Applications.** We go back to the
 294 L_2 cost case and introduce the full preconditioning procedure using all the *admissible*
 295 *map pairs* discussed in [section 3](#).

296 **DEFINITION 15 (Preconditioning Procedure).** *For two random variables X and*
 297 *Y with probability measures μ and ν , let*

$$298 \quad (31) \quad m_1 = \mathbb{E}X, \quad m_2 = \mathbb{E}Y,$$

$$299 \quad (32) \quad \Sigma_1 = \mathbb{E}[(X - m_1)^T(X - m_1)], \quad \Sigma_2 = \mathbb{E}[(Y - m_2)^T(Y - m_2)].$$

300 *We construct two matrices A and D that satisfy (16), and apply the preconditioning*
 301 *procedure:*

$$302 \quad (33) \quad \tilde{X} = (X - m_1)A, \quad \tilde{Y} = (Y - m_2)(A^T)^{-1}.$$

303 *If the optimal map between $\tilde{\mu}$ and $\tilde{\nu}$ is $\tilde{Y} = T(\tilde{X})$, the optimal map between $X \sim \mu$*
 304 *and $Y \sim \nu$ is*

$$305 \quad (34) \quad Y = [m_2 + T((X - m_1)A)A^T].$$

306 This preconditioning procedure moves the two given measures into new measures
 307 with zero mean and the same diagonal covariance matrix. An extra step that one
 308 can add to the preconditioning procedure uses the scaling pairs to normalize both
 309 measures so that they have total variance one. In the numerical experiments for this
 310 article we do not perform this extra step.

311 One straightforward theoretical application of the procedure is a simple derivation
 312 of the optimal map between multivariate normal distributions. If $X \sim N(m_1, \Sigma_1)$ and
 313 $Y \sim N(m_2, \Sigma_2)$, the \tilde{X} and \tilde{Y} resulting from the application of the preconditioning
 314 procedure have the same distribution $N(0, D)$. Since the optimal coupling between
 315 identical measures is the identity map, the optimal map between $N(m_1, \Sigma_1)$ and
 316 $N(m_2, \Sigma_2)$ is

$$317 \quad (35) \quad Y = m_2 + (X - m_1)AA^T = m_2 + (X - m_1)\Sigma_1^{-1/2}(\Sigma_1^{1/2}\Sigma_2\Sigma_1^{1/2})^{1/2}\Sigma_1^{-1/2},$$

318 a result that agrees with the one found in [\[9\]](#) through different means.

319 This procedure also gives an alternative proof to the following lower bound intro-
 320 duced in [\[6\]](#):

321 **THEOREM 16.** *Suppose (X, Y) is the optimal coupling between μ and ν . Let $m_1 =$
 322 $\mathbb{E}X$ and $m_2 = \mathbb{E}Y$ and $\Sigma_{1,2}$ be their respective covariance matrices. Denoting the
 323 nuclear norm of a matrix M by $\|M\|_*$, we have the following lower bound for the total
 324 transportation cost:*

$$325 \quad (36) \quad \mathbb{E}\|X - Y\|^2 \geq \|m_1 - m_2\|^2 + \|\Sigma_1\|_* + \|\Sigma_2\|_* - 2\|\Sigma_1^{1/2}\Sigma_2^{1/2}\|_*.$$

326 *Proof.* This bound follows directly from the estimation in the proof of [Property 12](#).
 327 Since

$$328 \quad \|\Sigma_1\|_* = \text{tr}(\Sigma_1), \quad \|\Sigma_2\|_* = \text{tr}(\Sigma_2), \quad \|\Sigma_1^{1/2}\Sigma_2^{1/2}\|_* = \sum_{i=1}^N d_i,$$

329 applying the *optimal affine pair* to general random variables X and Y , we have

$$330 \quad \mathbb{E}\|X - Y\|^2 = \|m_1 - m_2\|^2 + \|\Sigma_1\|_* + \|\Sigma_2\|_* - 2\|\Sigma_1^{1/2}\Sigma_2^{1/2}\|_* + \mathbb{E}\|\tilde{X} - \tilde{Y}\|^2.$$

331 Since clearly $\mathbb{E}\|\tilde{X} - \tilde{Y}\|^2$ is non-negative, we derive the lower bound (36) along with
332 the condition for the bound to be sharp. \square

333 A more general application of this procedure is to precondition measures and
334 datasets before applying any numerical optimal transport algorithm. The new prob-
335 lem is generally easier to solve, as it has a smaller transportation cost than the original
336 one.

337 In practice, instead of continuous probability measures in closed form, one of-
338 ten has only sample points drawn from otherwise unknown distributions. Applying
339 the procedure of this article to precondition a problem posed in terms of samples
340 is straightforward, since the preconditioning maps act on the random variables, and
341 hence on the sample points. The only difference is that, instead of the true mean val-
342 ues and covariance matrices, one uses estimates, such as their empirical counterparts,
343 to define the preconditioning maps.

344 **6. Numerical Experiments.** Our first example concerns optimal transport
345 problems between two-dimensional normal distributions. Consider μ and ν defined
346 by

$$347 \quad (37) \quad \mu = N \left([1, 1], \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix} \right), \quad \nu = N \left([-1, 0], \begin{pmatrix} 1 & -1 \\ -1 & 2 \end{pmatrix} \right).$$

348 We generate $N = 200$ data points $\{x_i\}_{i=1}^{200}$ and $\{y_i\}_{i=1}^{200}$ from each distribution. The
349 distributions and sample sets are shown in figures Figure 1a and Figure 1b.

350 We then perform the preconditioning procedure on both the distributions and the
351 sample sets. Notice that the two versions should give slightly different results, because
352 in the sample-based version empirical statistics are used instead of the true ones.
353 The results are shown in Figure 1c and Figure 1d. The preconditioning procedure
354 for continuous measures by definition makes $\tilde{\mu} = \tilde{\nu}$. On the other hand, the two
355 preconditioned sample sets are consistent with the preconditioned measures.

356 In the second example, we test the preconditioning procedure on more complicated
357 distributions. We define both μ and ν to be Gaussian mixtures:

$$358 \quad (38) \quad \begin{aligned} \mu &= \frac{1}{2}N \left([2, -1], \begin{pmatrix} 1/4 & 0 \\ 0 & 1/4 \end{pmatrix} \right) + \frac{1}{2}N \left([2, -3], \begin{pmatrix} 1/2 & 1/4 \\ 1/4 & 1/4 \end{pmatrix} \right) \\ \nu &= \frac{2}{3}N \left([2, 1], \begin{pmatrix} 3 & -1 \\ -1 & 2 \end{pmatrix} \right) + \frac{1}{3}N \left([-2, 1], \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \right) \end{aligned}$$

359 In Figure 2c and Figure 2d the preconditioned datasets have the same diagonal co-
360 variance matrix and are closer to each other than in the original datasets. As in the
361 first example, the preconditioned sample sets are consistent with the corresponding
362 preconditioned measures. This shows numerically that the preconditioning procedure
363 on sample sets is consistent with the procedure on continuous measures.

364 In the third example, we apply the preconditioning procedure along with the
365 sample-based numerical optimal transport algorithm introduced in [11], which takes
366 sample sets as input and compares and transfers them through feature functions.
367 This iterative algorithm approaches the optimal map by gradually approximating
368 the McCann interpolant [13] and updating the local transfer maps. We apply the

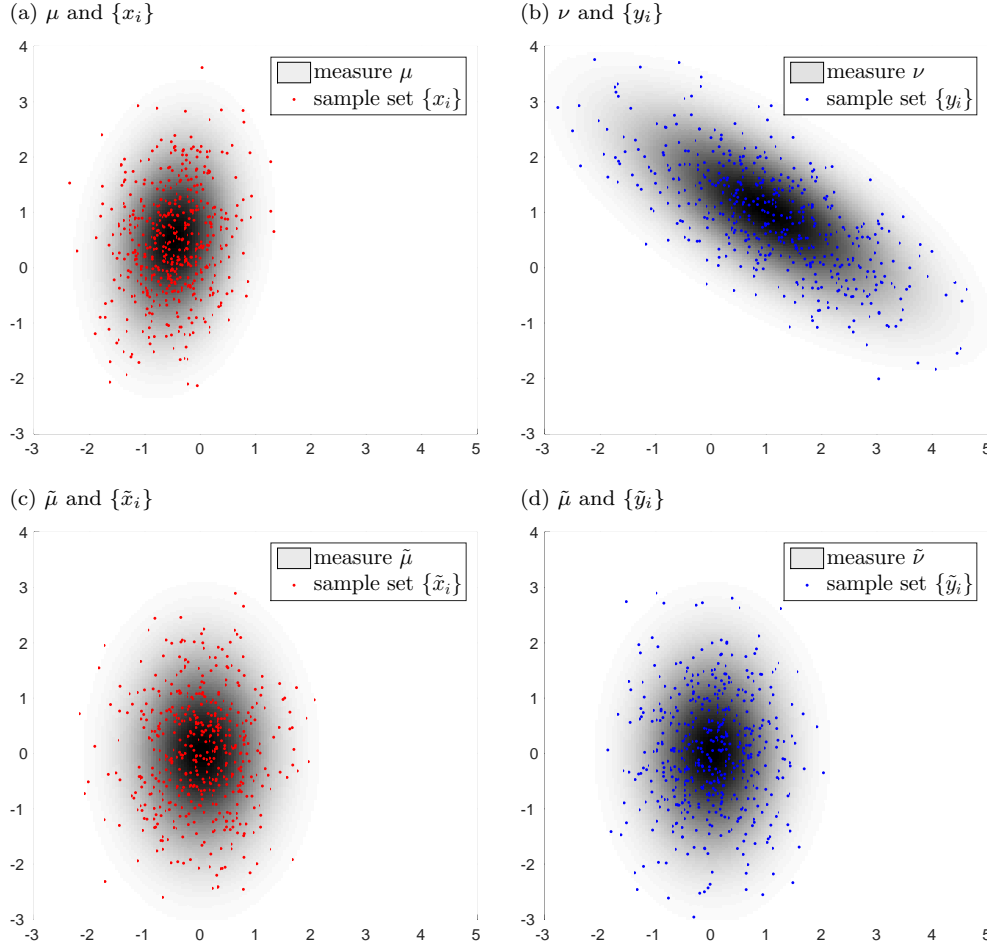


Fig. 1: Preconditioning on the two Gaussian distributions μ and ν defined in (37). Sample sets $\{x_i\}$ and $\{y_i\}$ are sampled from μ and ν respectively, each with sample size 200. In (c)(d), the preconditioned measures $\tilde{\mu}$ and $\tilde{\nu}$ are derived from μ and ν by the preconditioning procedure. $\{\tilde{x}_i\}$ and $\{\tilde{y}_i\}$ are transferred from the original sample sets with maps defined by their empirical mean values and covariance matrices.

369 preconditioning procedure and give the preconditioned sample sets to the algorithm.
 370 Then we take the optimal map from the algorithm's output and transform it to solve
 371 the original problem. The preconditioning procedure is crucial on two grounds: not
 372 only does the algorithm perform better on the preconditioned sample sets, which are
 373 closer to each other than the original ones, but feature selection becomes easier, as
 374 the same features describe the two distributions at similar levels of precision.

375 We choose a two-dimensional shape transform problem to test the algorithm. The
 376 problem involves finding the optimal transport between two geometrical objects, which
 377 can be described in probabilistic terms by introducing a uniform distribution within
 378 the support of each. For demonstration, consider the specific task of transforming an

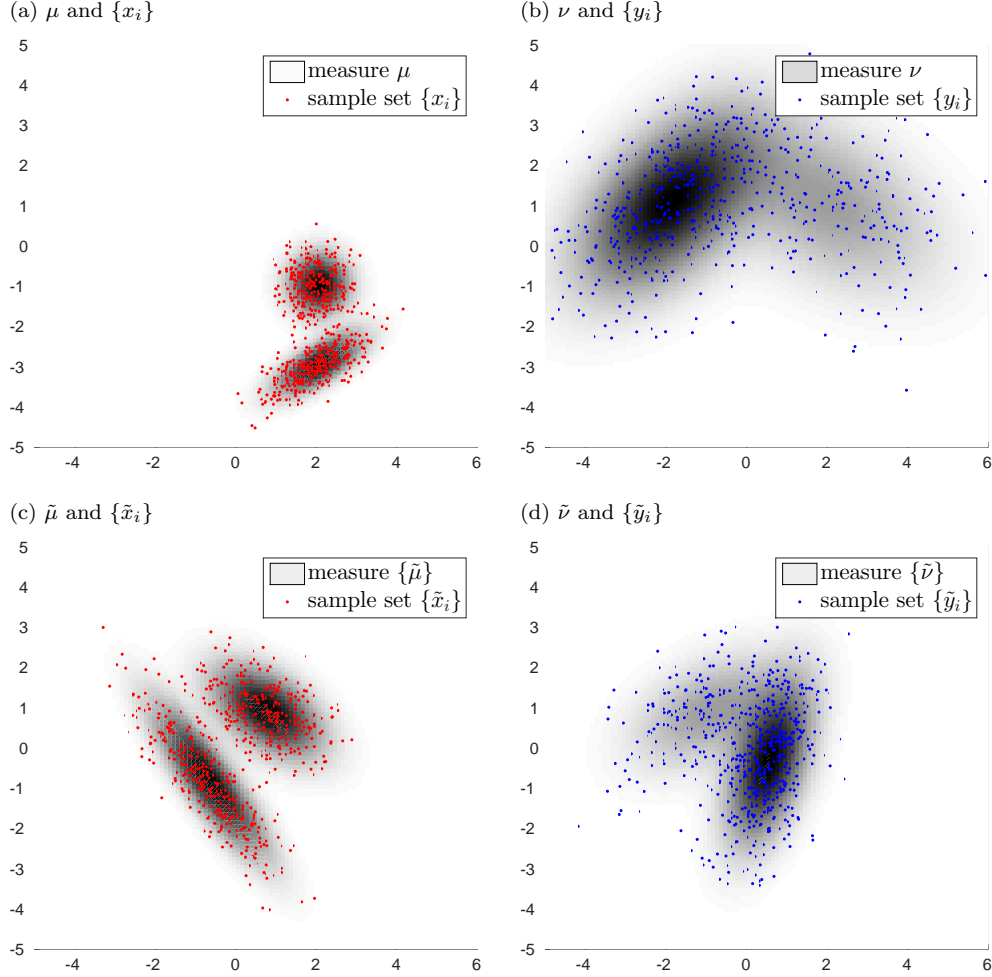


Fig. 2: The distributions μ and ν are the Gaussian mixtures defined in (38). $\{x_i\}$ and $\{y_i\}$ are derived in the same way as in figure Figure 1. In (c)(d), the two preconditioned sample sets $\{\tilde{x}_i\}$ and $\{\tilde{y}_i\}$ are transferred from the original datasets through maps defined in terms of their empirical mean values and covariance matrices.

379 ellipse into a ring (Figure 3a), described by:

$$\begin{aligned}
 \Omega_2 &= \{(x, y) \mid 1 \leq 3(x-5)^2 + 2(y+1)^2 - (x-5)(y+1) \leq 9\} \\
 \Omega_1 &= \{(x, y) \mid (x-1)^2 + 10y^2 \leq 1\}
 \end{aligned}
 \tag{39}$$

381 Both sample sets are drawn from uniform distributions within each region, with
 382 the sample size set to 1000 points per sample set.

383 This is a challenging optimal transport problem, since a) the locations and sizes
 384 of the two regions are different; b) the topological structure of the two regions are
 385 different, as one is simply connected and the other is not; c) both regions have sharp
 386 boundaries, which makes the solution singular; and d) since both shapes are eccentric,
 387 the optimal map between them is not essentially one dimensional as in the transfor-

388 mation between a circle and a circular ring.

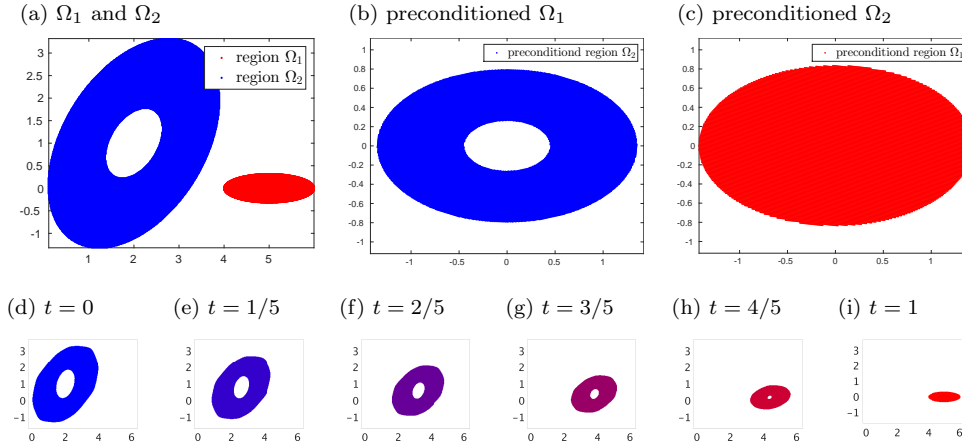


Fig. 3: Shape transformation problem. The two regions Ω_1 and Ω_2 are shown in (a), and their preconditioned images in (b)(c). (d)-(i) illustrate the McCann Interpolation of the optimal map, at times shown in the titles. All computation are carried out on sample sets drawn from the corresponding region. For the plots, we estimate the density function $p(x)$ for each sample set and display the area with $p(x) > \epsilon$, where ϵ is a small constant. The density functions are estimated by kernel density estimator with optimal kernel parameters.

389 The preconditioned regions are shown in Figure 3b and Figure 3c, they share the
 390 same mean and diagonal covariance matrix. The two preconditioned regions are much
 391 closer to each other, the blue one distinguished by its hole and a slightly smaller radius.
 392 Using the sample-based algorithm on the preconditioned sample sets, we find the
 393 optimal map T between the two preconditioned regions. Reversing the preconditioning
 394 step, the map can then be transformed back to the optimal map between Ω_1 and Ω_2 .
 395 The map and its McCann interpolation are shown in the second row of Figure 3.
 396 Without the preconditioning step, the procedure would have produced much poorer
 397 results and at a much higher computational expense.

398 **7. Conclusions and Future Works.** This paper describes a family of *affine*
 399 *map pairs* that preserves the optimality of transport solutions, and finds an optimal
 400 one among them that minimizes the remaining transportation cost. The procedure
 401 extends from the L_2 -cost to more general cost functions induced by an inner product.
 402 Based on these map pairs, we propose a preconditioning procedure which maps input
 403 measures or datasets to preconditioned ones while preserving the optimality of the
 404 solutions.

405 The procedure is efficient, easy to implement and it can significantly reduce the
 406 difficulty of the problem in many scenarios. Using this procedure one can directly solve
 407 the optimal transport problem between multivariate normal distributions. We tested
 408 the procedure both as a stand-alone method and along with a sample-based optimal
 409 transport algorithm. The procedure in all cases successfully preconditioned the input
 410 measures and datasets, making them more regular and closer to their counterparts.

411 For future works, one natural extension is to consider non-linear *admissible map*

412 *pairs*, which can potentially reduce further the total transportation cost and solve
 413 directly a wider class of optimal transport problems. If the family of *admissible*
 414 *map pairs* is rich enough, one can potentially construct a practical optimal transport
 415 algorithm from these map pairs alone.

416 Another possible extension is to the barycenter problem [1]:

$$417 \quad (40) \quad \min_{\pi_k \in \Pi(\mu_k, \nu)} \sum_{k=1}^K w_k \int c(x, y) d\pi_k(x, y),$$

418 where $\mu_1, \mu_2, \dots, \mu_K$ are K different measures with positive weights w_1, w_2, \dots, w_K .
 419 Instead of the two measures of the regular optimal transport problem, we would like
 420 to map K measures simultaneously while preserving the optimality of the solution.
 421 The simplest of such maps is the set of translations that give all measures the same
 422 zero mean.

423 **Acknowledgments.** This work was partially supported by a grant from the
 424 Office of Naval Research and from the NYU-AIG Partnership on Global Resilience.

425

REFERENCES

- 426 [1] M. AGUEH AND G. CARLIER, *Barycenters in the wasserstein space*, SIAM Journal on Mathe-
 427 matical Analysis, 43 (2011), pp. 904–924.
- 428 [2] J.-D. BENAMOU AND Y. BRENIER, *A computational fluid mechanics solution to the monge-*
 429 *kantorovich mass transfer problem*, Numerische Mathematik, 84 (2000), pp. 375–393.
- 430 [3] Y. BRENIER, *Polar factorization and monotone rearrangement of vector-valued functions*, Com-
 431 munications on pure and applied mathematics, 44 (1991), pp. 375–417.
- 432 [4] G. BUTTAZZO, A. PRATELLI, AND E. STEPANOV, *Optimal pricing policies for public transporta-*
 433 *tion networks*, SIAM Journal on Optimization, 16 (2006), pp. 826–853.
- 434 [5] W. CHEN AND E. G. TABAK, *An adaptive linear programming methodology for data driven*
 435 *optimal transport*, Numerische Mathematik, (submitted).
- 436 [6] J. CUESTA-ALBERTOS, C. MATRÁN-BEA, AND A. TUERO-DIAZ, *On lower bounds for the 2-*
 437 *wasserstein metric in a hilbert space*, Journal of Theoretical Probability, 9 (1996), pp. 263–
 438 283.
- 439 [7] M. CULLEN AND R. PURSER, *Properties of the lagrangian semigeostrophic equations*, Journal
 440 of the Atmospheric Sciences, 46 (1989), pp. 2684–2697.
- 441 [8] W. GANGBO, R. J. MCCANN, ET AL., *Shape recognition via wasserstein distance*, (1999).
- 442 [9] C. R. GIVENS, R. M. SHORTT, ET AL., *A class of wasserstein metrics for probability distribu-*
 443 *tions.*, The Michigan Mathematical Journal, 31 (1984), pp. 231–240.
- 444 [10] L. V. KANTOROVICH, *On a problem of monge*, Journal of Mathematical Sciences, 133 (2006),
 445 pp. 1383–1383.
- 446 [11] M. KUANG AND E. G. TABAK, *Sample-based optimal transport and barycenter problems*, (In
 447 preparation).
- 448 [12] J. N. MATHER, *Action minimizing invariant measures for positive definite lagrangian systems*,
 449 Mathematische Zeitschrift, 207 (1991), pp. 169–207.
- 450 [13] R. J. MCCANN, *A convexity principle for interacting gases*, Advances in mathematics, 128
 451 (1997), pp. 153–179.
- 452 [14] G. MONGE, *Mémoire sur la théorie des déblais et des remblais*, De l’Imprimerie Royale, 1781.
- 453 [15] A. M. OBERMAN AND Y. RUAN, *An efficient linear programming method for optimal trans-*
 454 *portation*, arXiv preprint arXiv:1509.03668, (2015).
- 455 [16] S. T. RACHEV AND L. RÜSCHENDORF, *Mass Transportation Problems: Volume I: Theory*, vol. 1,
 456 Springer Science & Business Media, 1998.
- 457 [17] Y. RUBNER, C. TOMASI, AND L. J. GUIBAS, *A metric for distributions with applications to*
 458 *image databases*, in Computer Vision, 1998. Sixth International Conference on, IEEE,
 459 1998, pp. 59–66.
- 460 [18] J. SOLOMON, F. DE GOES, G. PEYRÉ, M. CUTURI, A. BUTSCHER, A. NGUYEN, T. DU, AND
 461 L. GUIBAS, *Convolutional wasserstein distances: Efficient optimal transportation on geo-*
 462 *metric domains*, ACM Transactions on Graphics (TOG), 34 (2015), p. 66.
- 463 [19] E. G. TABAK AND G. TRIGILA, *Data-driven optimal transport*, Commun. Pure. Appl. Math.
 464 doi, 10 (2014), p. 1002.

- 465 [20] C. VILLANI, *Optimal transport: old and new*, vol. 338, Springer Science & Business Media,
466 2008.