

A Supplementary Results and Comments

A.1 Proof that eqs. (4) and (5) are equivalent when $N_x = N_y = N$

This is the result of a straightforward calculation, exploiting several elementary facts from linear algebra. First, for any matrix \mathbf{A} we have that $\|\mathbf{A}\|_F^2 = \text{Tr}[\mathbf{A}^\top \mathbf{A}]$. Second, for any matrices $\{\mathbf{A}, \mathbf{B}, \mathbf{C}\}$ with appropriate dimensions such that the product \mathbf{ABC} is defined, we have that $\text{Tr}[\mathbf{ABC}] = \text{Tr}[\mathbf{CAB}] = \text{Tr}[\mathbf{BCA}]$, which is called the *cyclic trace property*. Finally, for any orthogonal matrix $\mathbf{Q} \in \mathcal{O}(N)$ we have $\mathbf{Q}^\top \mathbf{Q} = \mathbf{Q}\mathbf{Q}^\top = \mathbf{I}$.

With these ingredients we can manipulate the squared Procrustes distance as follows:

$$\begin{aligned}
d_{\mathcal{O}}^2(\mathbf{X}, \mathbf{Y}) &= \min_{\mathbf{Q} \in \mathcal{O}(N)} \|\mathbf{X} - \mathbf{Y}\mathbf{Q}\|_F^2 \\
&= \min_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr}[\mathbf{X}^\top \mathbf{X} + \mathbf{Q}^\top \mathbf{Y}^\top \mathbf{Y} \mathbf{Q} - 2\mathbf{X}^\top \mathbf{Y} \mathbf{Q}] && (\|\mathbf{A}\|_F^2 = \text{Tr}[\mathbf{A}^\top \mathbf{A}]) \\
&= \min_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr}[\mathbf{X}^\top \mathbf{X}] + \text{Tr}[\mathbf{Q}^\top \mathbf{Y}^\top \mathbf{Y} \mathbf{Q}] - 2\text{Tr}[\mathbf{X}^\top \mathbf{Y} \mathbf{Q}] && (\text{Tr}[\cdot] \text{ is linear}) \\
&= \min_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr}[\mathbf{X}^\top \mathbf{X}] + \text{Tr}[\mathbf{Y}^\top \mathbf{Y} \mathbf{Q} \mathbf{Q}^\top] - 2\text{Tr}[\mathbf{X}^\top \mathbf{Y} \mathbf{Q}] && (\text{cyclic trace property}) \\
&= \min_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr}[\mathbf{X}^\top \mathbf{X}] + \text{Tr}[\mathbf{Y}^\top \mathbf{Y}] - 2\text{Tr}[\mathbf{X}^\top \mathbf{Y} \mathbf{Q}] && (\text{orthogonality}) \\
&= \text{Tr}[\mathbf{X}^\top \mathbf{X}] + \text{Tr}[\mathbf{Y}^\top \mathbf{Y}] - 2 \max_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr}[\mathbf{X}^\top \mathbf{Y} \mathbf{Q}] && (\text{first two terms are constant}) \\
&= \text{Tr}[\mathbf{X}^\top \mathbf{X}] + \text{Tr}[\mathbf{Y}^\top \mathbf{Y}] - 2\|\mathbf{X}^\top \mathbf{Y}\|_* && (\text{discussed below})
\end{aligned}$$

as claimed in the main text. The final step is the comes from the celebrated closed form solution to the orthogonal Procrustes problem, which is comprehensively reviewed by Gower and Dijksterhuis [20]. Briefly, the result can be understood by considering the singular value decomposition $\mathbf{X}^\top \mathbf{Y} = \mathbf{USV}^\top$. Then, due to the cyclic trace property,

$$\text{Tr}[\mathbf{X}^\top \mathbf{Y} \mathbf{Q}] = \text{Tr}[\mathbf{USV}^\top \mathbf{Q}] = \text{Tr}[\mathbf{SV}^\top \mathbf{Q} \mathbf{U}] \quad (16)$$

This final expression is maximized by setting $\mathbf{Q} = \mathbf{V}\mathbf{U}^\top$, resulting in:

$$\text{Tr}[\mathbf{SV}^\top \mathbf{V}\mathbf{U}^\top \mathbf{U}] = \text{Tr}[\mathbf{S}] = \|\mathbf{X}^\top \mathbf{Y}\|_* \quad (17)$$

Since the sum of the diagonal elements of \mathbf{S} is simply the sum of the singular values of $\mathbf{X}^\top \mathbf{Y}$ (i.e. equal to the nuclear norm of this matrix).

A.2 Proof that eqs. (9) and (10) are equivalent

Recall that we are in the setting where $N_x = N_y = N$. First, for any matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ with columns $\{\mathbf{a}_1, \dots, \mathbf{a}_N\}$ we have $\|\mathbf{A}\|_F^2 = \sum_{i=1}^N \|\mathbf{a}_i\|^2$. Since $\sum_j \mathbf{P}_{ij} \mathbf{y}_j$ gives column i of the matrix product $\mathbf{Y}\mathbf{P}$, we have:

$$\|\mathbf{X} - \mathbf{Y}\mathbf{P}\|_F^2 = \sum_{i=1}^N \|\mathbf{x}_i - \sum_{j=1}^N \mathbf{P}_{ij} \mathbf{y}_j\|^2 \quad (18)$$

Recall that \mathbf{P} is a permutation matrix in the present context. Thus, let $\sigma(i) \in \{1, \dots, N\}$ denote the index of the unique nonzero element of row i in \mathbf{P} . Intuitively, $\sigma(i)$ defines the permutation in which column i of \mathbf{X} is matched to column $\sigma(i)$ of \mathbf{Y} . With this notation, we can re-write the expression above:

$$\sum_{i=1}^N \|\mathbf{x}_i - \sum_{j=1}^N \mathbf{P}_{ij} \mathbf{y}_j\|^2 = \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{y}_{\sigma(i)}\|^2 \quad (19)$$

Now define $\delta[i, j]$ as a function that takes in two integers and evaluates to one if $i = j$ and evaluates to zero if $i \neq j$. (This is often called the Kronecker delta function.) We can write:

$$\sum_{i=1}^N \|\mathbf{x}_i - \mathbf{y}_{\sigma(i)}\|^2 = \sum_{i=1}^N \sum_{j=1}^N \delta[\sigma(i), j] \cdot \|\mathbf{x}_i - \mathbf{y}_j\|^2 \quad (20)$$

since the inner sum will evaluate to zero whenever $j \neq \sigma(i)$. Finally, we argue that $\mathbf{P}_{ij} = \delta[\sigma(i), j]$. Indeed, $\sum_j \delta[\sigma(i), j] \mathbf{y}_j = \mathbf{y}_{\sigma(i)}$ which agrees with $\sum_j \mathbf{P}_{ij} \mathbf{y}_j = \mathbf{y}_{\sigma(i)}$. Thus, we have shown that:

$$d_{\mathcal{P}}^2(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{P} \in \mathcal{P}(N)} \|\mathbf{X} - \mathbf{Y}\mathbf{P}\|_F^2 = \min_{\mathbf{P} \in \mathcal{P}(N)} \sum_{i,j} \mathbf{P}_{ij} \|\mathbf{x}_i - \mathbf{y}_j\|_2^2 \quad (21)$$

To arrive at eq. (10), we need to show that we can relax the constraint of the minimization over the permutation group to over the Birkhoff polytope—i.e. to show that:

$$\min_{\mathbf{P} \in \mathcal{P}(N)} \sum_{i,j} \mathbf{P}_{ij} \|\mathbf{x}_i - \mathbf{y}_j\|_2^2 = \min_{\mathbf{P} \in \mathcal{B}(N)} \sum_{i,j} \mathbf{P}_{ij} \|\mathbf{x}_i - \mathbf{y}_j\|_2^2 \quad (22)$$

Here we evoke two well-known results. First, the celebrated Birkhoff–von Neumann theorem states that the vertices of $\mathcal{B}(N)$ are one-to-one matched with the permutation matrices $\mathcal{P}(N)$. Second, the final expression is a linear program since the objective function is linear in \mathbf{P} and the constraints are linear (as can be verified from eq. 2). Thus, we evoke a basic fact from the theory of linear programming (see e.g. [5]), which states that, assuming that a finite solution exists, at least one vertex of the feasible set is a solution. Any such vertex is called a “basic feasible solution” and the fact that such solutions exist motivates the well known simplex algorithm for linear programming. Thus, we conclude that relaxing the constraints from $\mathbf{P} \in \mathcal{P}(N)$ to $\mathbf{P} \in \mathcal{B}(N)$ does not allow us to further minimize the objective, and so eq. (22) is valid. Taking square roots on both sides of eq. (22) proves the result claimed in the main text.

A.3 Relation between soft matching distance and correlation score

Here we show that the optimal soft permutation matrix $\mathbf{P} \in \mathcal{T}(N_x, N_y)$ that minimizes the expression in eq. 11 equals the one which maximizes the expression in eq. 12. First, beginning with the minimization problem in eq. 11, we can break the expression into three terms:

$$\operatorname{argmin}_{\mathbf{P} \in \mathcal{T}(N_x, N_y)} \sum_{i,j} \mathbf{P}_{ij} \|\mathbf{x}_i - \mathbf{y}_j\|^2 \quad (23)$$

$$= \operatorname{argmin}_{\mathbf{P} \in \mathcal{T}(N_x, N_y)} \sum_{i,j} \mathbf{P}_{ij} (\mathbf{x}_i^\top \mathbf{x}_i + \mathbf{y}_j^\top \mathbf{y}_j - 2\mathbf{x}_i^\top \mathbf{y}_j) \quad (24)$$

$$= \operatorname{argmin}_{\mathbf{P} \in \mathcal{T}(N_x, N_y)} \underbrace{\sum_{i,j} \mathbf{P}_{ij} \mathbf{x}_i^\top \mathbf{x}_i}_{(A)} + \underbrace{\sum_{i,j} \mathbf{P}_{ij} \mathbf{y}_j^\top \mathbf{y}_j}_{(B)} - 2 \underbrace{\sum_{i,j} \mathbf{P}_{ij} \mathbf{x}_i^\top \mathbf{y}_j}_{(C)} \quad (25)$$

Considering term (A) first, we argue that this term is constant with respect to any feasible \mathbf{P} since:

$$\sum_{i,j} \mathbf{P}_{ij} \mathbf{x}_i^\top \mathbf{x}_i = \sum_{i=1}^{N_x} \left(\sum_{j=1}^{N_y} \mathbf{P}_{ij} \right) \mathbf{x}_i^\top \mathbf{x}_i = \sum_{i=1}^{N_x} \left(\frac{1}{N_x} \right) \mathbf{x}_i^\top \mathbf{x}_i \quad (26)$$

where the final equality follows from definition of the transportation polytope in eq. (3)—namely, the rows of \mathbf{P} each sum to $1/N_x$. We then can make a similar argument for term (B). In particular, since the columns of \mathbf{P} each sum to $1/N_y$, we have:

$$\sum_{i,j} \mathbf{P}_{ij} \mathbf{y}_j^\top \mathbf{y}_j = \sum_{j=1}^{N_y} \left(\sum_{i=1}^{N_x} \mathbf{P}_{ij} \right) \mathbf{y}_j^\top \mathbf{y}_j = \sum_{j=1}^{N_y} \left(\frac{1}{N_y} \right) \mathbf{y}_j^\top \mathbf{y}_j \quad (27)$$

In summary, we see that only term (C) is non-constant. That is, we have

$$\operatorname{argmin}_{\mathbf{P} \in \mathcal{T}(N_x, N_y)} \sum_{i,j} \mathbf{P}_{ij} \|\mathbf{x}_i - \mathbf{y}_j\|^2 = \operatorname{argmin}_{\mathbf{P} \in \mathcal{T}(N_x, N_y)} -2 \sum_{i,j} \mathbf{P}_{ij} \mathbf{x}_i^\top \mathbf{y}_j + \text{const.} \quad (28)$$

$$= \operatorname{argmax}_{\mathbf{P} \in \mathcal{T}(N_x, N_y)} \sum_{i,j} \mathbf{P}_{ij} \mathbf{x}_i^\top \mathbf{y}_j \quad (29)$$

as we have claimed.

A.4 Computational complexity

Computing the soft-matching distance requires solving an optimal transport problem in the discrete setting. The solution to the transportation problem, which is a linear program, can be derived using the network simplex algorithm. With efficient implementations of the simplex algorithm as in the Python Optimal Transport Library, the complexity of solving the linear program is $O(n^3 \log n)$, assuming that the two representations being compared have n units. Such efficient implementations enable broad application of optimal transport-based solutions in real-world settings.

A.5 Relevance of the Soft-Matching Metric to Disentangled Representation Learning Metrics

It is worth noting that the proposed soft-matching metric can also serve as a valuable tool in the field of disentangled representation learning (DRL) [4] due to its sensitivity to the representational basis. In DRL, the objective is to learn a model that effectively disentangles and makes the underlying generative factors of the data explicit in representational form (i.e. aligned with representational units). Within the DRL literature, various measures have been developed to quantify the alignment between learned representations and ground truth generative factors. Typically, the desiderata involve a combination of different criteria, encompassing the similarity in information content (explicitness) and the degree of one-to-one correspondence between the representational units and generative factors (modularity and compactness) [15]. The soft-matching distance metric offers a unique advantage by simultaneously capturing sensitivity to both of these critical properties.