

Beyond Euclidean: A Conceptual Review of Distance Metrics and Dissimilarity Measures in Contemporary Machine Learning

This paper presents a broad conceptual review of distance metrics and dissimilarity measures fundamental to the field of Machine Learning (ML). While ubiquitous algorithms like k -nearest neighbours, clustering methods, and support vector machines heavily rely on measuring similarity, the default choice of the Euclidean distance often obscures a deeper, application-specific consideration of data geometry. We aim to provide a broad perspective across diverse data types and domains, moving from classical vector norms to covariance-aware, angular, temporal, topological, probabilistic, image-based, and manifold-aware measures. Beyond a simple catalogue, we address several theoretical implications in the ML pipeline. These include the relationship between distance metrics and positive definite kernels, distance metric learning as a subfield dedicated to learning task-specific metrics, latent space fidelity as a way to analyse what distances in low-dimensional latent spaces reveal about semantic separation in the observed space, induced metrics arising from feature maps and learned representations, identifiability of distances as an advanced perspective on latent geometry, and computational complexity aspects. Ultimately, this paper argues that researchers and practitioners should critically evaluate the underlying geometric assumptions in their models, and that the selection of the appropriate distance metric is, in fact, an explicit modelling decision.

Keywords: *Distance Metrics, Dissimilarity Measures, Similarity Learning, Machine Learning Geometry, Metric Learning*

Introduction

Distance measures are central to the conceptual and operational foundations of Machine Learning. Many widely used algorithms, such as k -nearest neighbours, clustering methods (e.g., k -means), and support vector machines, rely explicitly or implicitly on the notion of similarity between data points. Applied clustering studies, including fuzzy cluster analysis of educational data, illustrate how the chosen similarity structure shapes the resulting grouping (Bedalli & Ninka, 2015). In practice, however, the Euclidean distance is frequently used as a default choice, often without critical examination of its suitability for the data at hand.

Formally, a *distance metric* $d: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ satisfies four properties for all $x, y, z \in \mathcal{X}$:

$$d(x, y) \geq 0 \quad (1)$$

$$d(x, y) = 0 \Leftrightarrow x = y \quad (2)$$

$$d(x, y) = d(y, x) \quad (3)$$

$$d(x, z) \leq d(x, y) + d(y, z) \quad (4)$$

These correspond to non-negativity, identity of indiscernibles, symmetry, and the triangle inequality. A *dissimilarity measure*, also called *generalised metric*, may

1 relax one or more of these conditions. A *divergence* is often understood as a non-
 2 negative dissimilarity that vanishes exactly when the two inputs are identical,
 3 although the term is used somewhat differently across fields. A distance that relaxes
 4 (2) is a *pseudometric*. One that relaxes (3) is a *quasimetric*. One that relaxes (4) is a
 5 *semimetric*. Dissimilarity measures allow for greater flexibility in modelling
 6 complex data relationships. This distinction is particularly relevant in modern ML
 7 applications, where data often reside in structured, high-dimensional, or non-
 8 Euclidean spaces.

9 A growing body of research highlights that the choice of distance metric
 10 fundamentally shapes the geometry of the data space and, consequently, the
 11 behaviour of learning algorithms (Kaya & Bilge, 2019). This has led to the
 12 development of Distance Metric Learning (DML) (Suárez et al., 2018), which seeks
 13 to learn task-specific distance functions directly from data.

14 Comprehensive mathematical catalogues of distances, such as the work of
 15 Deza and Deza (Deza & Deza, 2016), provide a broad reference for definitions,
 16 properties, and examples across many areas. Recent guides also provide broad
 17 overviews of similarity measures and their data-science applications (Levy et al.,
 18 2025). The present review has a different focus: it does not aim to catalogue
 19 distances exhaustively, but to explain how metric choice operates as a modelling
 20 decision in contemporary ML pipelines, including kernel methods, learned
 21 embeddings, latent spaces, induced geometries, and practical evaluation settings.

22 A well-known phenomenon in high-dimensional spaces is the *distance*
 23 *concentration effect*, where the relative contrast between nearest and farthest
 24 neighbours diminishes as dimensionality grows (Aggarwal et al., 2001). This
 25 undermines the intuitive meaning of proximity and helps explain why Euclidean
 26 distance can become unreliable for nearest-neighbour search, clustering, and related
 27 distance-based methods. Section 2 revisits this effect when discussing classical
 28 vector norms.

29 Several theoretical considerations arise when selecting or designing a distance
 30 metric:

- 31
- 32 • **Kernel Compatibility:** The formal relationship between distance measures
 33 and positive definite kernels, which determines their compatibility with
 34 kernel methods.
- 35 • **Metric Learning:** The mechanism for learning optimised distance metrics
 36 directly from data distributions.
- 37 • **Latent Space Fidelity:** The fidelity of distances within learned, low-
 38 dimensional latent representations.
- 39 • **Induced Geometries:** The role of explicit feature transformations in
 40 inducing entirely new geometric structures.
- 41 • **Scale-Invariance:** The mathematical handling of distances within
 42 projective spaces where vectors exhibit inherent scale-invariance.
- 43 • **Identifiability:** The identifiability of learned representations, ensuring that
 44 an inferred geometry is uniquely determined rather than an artefact of
 45 optimisation or parametrisation.
- 46

1 The remainder of this paper is structured as follows. Section 2 presents a
 2 taxonomy of distance metrics and dissimilarity measures across domains. Section 3
 3 discusses theoretical implications, including kernel connections, metric learning,
 4 latent space representations, geometric transformations, and identifiability.
 5 Section 4 addresses practical computational constraints and limitations to have in
 6 mind regarding each distance. Section 5 translates these themes into practical
 7 selection criteria, and Section 6 concludes.

8 **Classical and Specialised Distance Metrics**

9 *Classical Vector Norms*

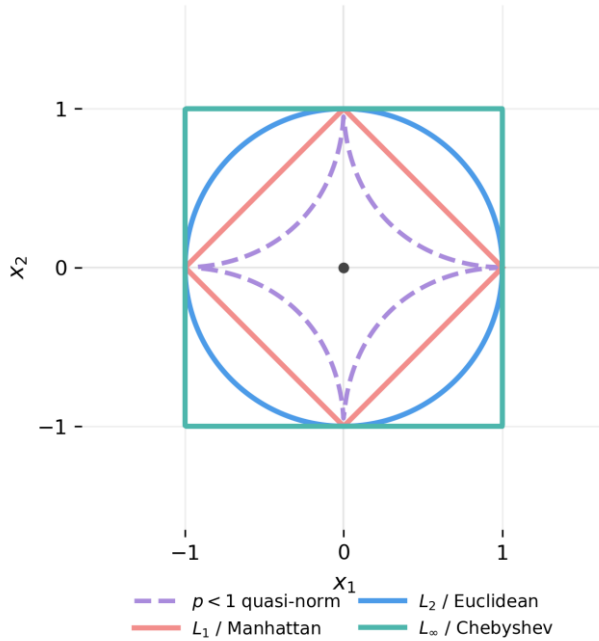
10
 11 The most fundamental class of distance metrics arises from the Minkowski L_p
 12 norm:

$$13 \quad \| \mathbf{x} - \mathbf{y} \|_p = \left(\sum_i |x_i - y_i|^p \right)^{\frac{1}{p}}. \quad (5)$$

14
 15
 16 Special cases of this formulation include the Euclidean (L_2), Manhattan (L_1),
 17 and Chebyshev (L_∞) distances. These metrics differ in how they aggregate
 18 coordinate-wise differences, leading to distinct geometric interpretations. For
 19 example, Euclidean distance assumes isotropy, while Manhattan distance is more
 20 robust for outlier detection or when working in high-dimensional sparse settings
 21 (Zhou et al., 2025). Figure 1 visualises these differences through the boundaries of
 22 the corresponding unit balls.

23
 24 Despite their simplicity, these norms implicitly assume that features are
 25 independent and equally scaled. This assumption is rarely satisfied in real-world
 26 datasets, where correlations and heterogeneous feature scales are common.

1 **Figure 1.** Boundaries of L_p unit balls in \mathbb{R}^2 . The L_1 boundary is diamond-shaped,
 2 reflecting Manhattan distance. The L_2 boundary is circular and isotropic. The L_∞
 3 boundary is square, reflecting dependence on the largest coordinate-wise
 4 difference. The dashed $p < 1$ curve illustrates a non-convex quasi-norm shape, for
 5 which the triangle inequality no longer holds



6

7

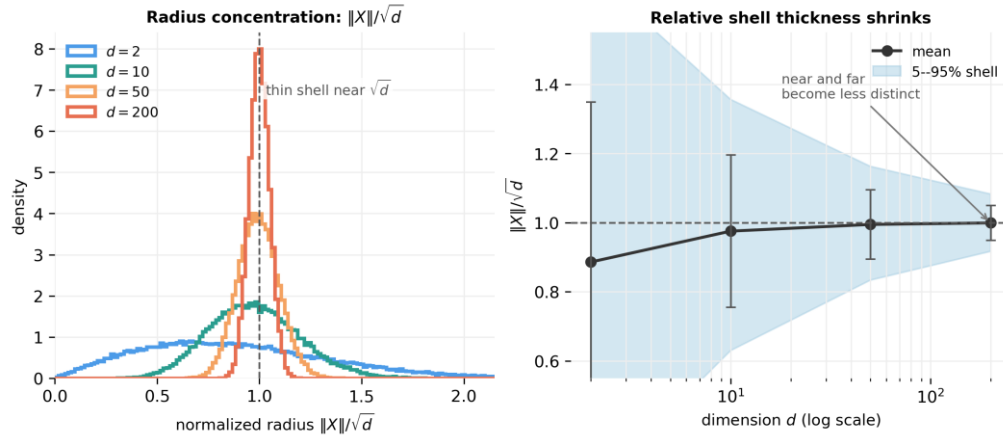
8 The limitations of these classical norms become especially visible in high
 9 dimensions. As dimensionality increases, most pairwise Euclidean distances can fall
 10 within a narrow range, reducing the relative contrast between near and far points.
 11 This phenomenon can be understood through concentration of measure. For
 12 example, in a high-dimensional Gaussian distribution, probability mass concentrates
 13 in a thin shell at radius approximately \sqrt{d} rather than near the origin. In such
 14 settings, almost all sampled points lie at nearly the same distance from the centre
 15 and, by extension, from each other. The result is that proximity becomes less
 16 informative for methods such as nearest-neighbour search and clustering.

17 Figure 2 illustrates this effect empirically. The plot was generated by drawing
 18 45,000 samples from standard Gaussian distributions $X \sim \mathcal{N}(0, I_d)$ for
 19 dimensions $d \in \{2, 10, 50, 200\}$, computing the Euclidean radius $\|X\|$, and
 20 normalising it by \sqrt{d} . As d increases, the distribution of $\|X\|/\sqrt{d}$ becomes
 21 increasingly concentrated around one, making the thin-shell geometry visible.

22

23

1 **Figure 2.** Concentration of measure for high-dimensional Gaussian data. For $X \sim$
 2 $\mathcal{N}(0, I_d)$, the radius $\|X\|$ concentrates near \sqrt{d} . After normalisation by \sqrt{d} , the
 3 empirical distribution of radii collapses into an increasingly thin shell around one.
 4 This is the geometric intuition behind the “soap bubble” analogy: most mass lies
 5 near a shell rather than near the origin, making relative notions of near and far less
 6 informative in high dimensions



7
8

9 Covariance-Aware Metrics

10

11 The Mahalanobis distance generalises the Euclidean distance by explicitly
 12 incorporating the covariance structure of the data:

13

$$14 \quad D_M(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})}. \quad (6)$$

15

16 Unlike the Euclidean distance, which assumes that all features are uncorrelated
 17 and equally scaled, the Mahalanobis distance accounts for both variance and
 18 correlation between features through the covariance matrix Σ . Given observations
 19 $\mathbf{x}_1, \dots, \mathbf{x}_n$ with sample mean $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$, this covariance is commonly estimated
 20 as

21

$$22 \quad \Sigma = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T. \quad (7)$$

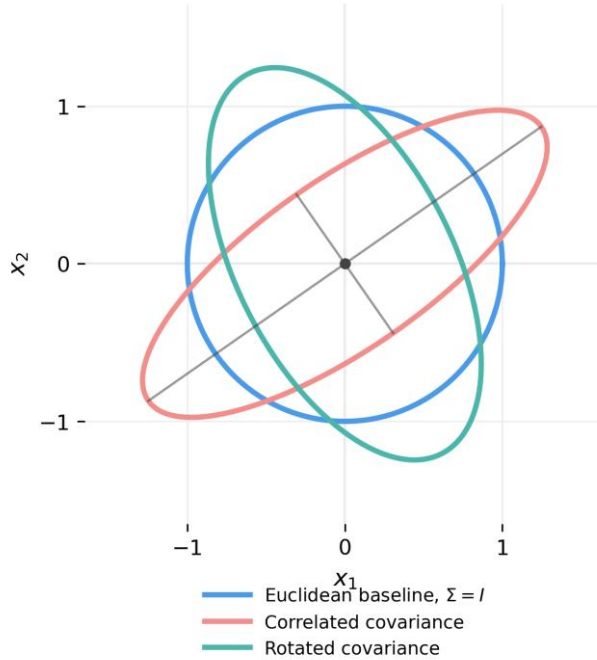
23

24 In practice, this means that differences along directions with high variance are
 25 considered less informative and are therefore down-weighted, while differences
 26 along low-variance directions are amplified. Intuitively, this reflects the idea that
 27 variation commonly observed in the data should contribute less to the notion of
 28 dissimilarity than rare or atypical variation.

29

30

1 **Figure 3.** Mahalanobis distance contours for different covariance structures. While
 2 Euclidean distance treats all directions equally, Mahalanobis distance rescales and
 3 rotates the effective geometry according to the covariance matrix Σ , down-
 4 weighting high-variance directions and amplifying low-variance directions.



5
6

7 Geometrically, this induces a transformation of the feature space in which the
 8 data is linearly transformed so that the covariance becomes the identity matrix. In
 9 this transformed space, the Mahalanobis distance reduces to the standard Euclidean
 10 distance. Consequently, distance contours that are spherical under Euclidean
 11 distance become ellipsoidal in the original space, aligning with the principal
 12 directions of the data distribution. This makes the Mahalanobis distance particularly
 13 suitable for anisotropic data, where the assumption of isotropy is violated. Figure 3
 14 illustrates this covariance-aware deformation of distance contours.

15 The Mahalanobis distance plays a central role in statistical modelling, anomaly
 16 detection, and classification. For example, in multivariate Gaussian models, it
 17 naturally arises in the exponent of the likelihood function and is therefore directly
 18 linked to probabilistic interpretations of distance. In anomaly detection, points with
 19 large Mahalanobis distance from the mean are considered outliers relative to the
 20 underlying distribution. Furthermore, in Distance Metric Learning (DML), many
 21 approaches aim to learn a matrix $M = \Sigma^{-1}$ (or a related positive semi-definite
 22 matrix), thereby adapting the geometry of the space to optimise task-specific notions
 23 of similarity (Weinberger & Saul, 2009).

24 Several extensions have been proposed. The local Mahalanobis distance
 25 (Ghosh et al., 2025) adapts the matrix M depending on the region of the space,
 26 allowing for non-linear and heterogeneous data structures. Similarly, kernelised
 27 Mahalanobis distances implicitly define such transformations in high-dimensional
 28 feature spaces via kernel functions, bridging the gap between metric learning and
 29 kernel methods (Schölkopf et al., 2001). These extensions illustrate that the

1 Mahalanobis distance is not a single isolated metric but rather a representative of a
 2 broader class of geometry-aware distances that adapt to the structure of the data,
 3 either through statistical estimation or learning.

4
 5 *Cosine Similarity, Angular Distance and Pearson's Distance*

6
 7 Cosine similarity measures the directional alignment between two non-zero
 8 vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$:

$$10 \quad S_{\cos}(\mathbf{x}, \mathbf{y}) = \cos\theta = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}. \quad (8)$$

11
 12 As a similarity score, S_{\cos} reflects angle rather than magnitude. A common
 13 dissimilarity derived from it is the cosine dissimilarity

$$15 \quad D_{\cos}(\mathbf{x}, \mathbf{y}) = 1 - S_{\cos}(\mathbf{x}, \mathbf{y}), \quad (9)$$

16
 17 but D_{\cos} is not a true metric because it can violate the triangle inequality. When
 18 vectors are normalised to unit length ($\|\mathbf{x}\| = \|\mathbf{y}\| = 1$), Euclidean distance is
 19 proportional to cosine similarity via

$$21 \quad \|\mathbf{x} - \mathbf{y}\| = \sqrt{2(1 - S_{\cos}(\mathbf{x}, \mathbf{y}))}, \quad (10)$$

22
 23 so cosine-based comparisons can be viewed as chordal or angular constraints on a
 24 spherical manifold.

25 The Pearson (correlation) distance is a centred analogue of cosine dissimilarity
 26 and accounts for mean offsets. Using the vector means $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$, the mean-adjusted
 27 vectors are $\mathbf{x} - \bar{\mathbf{x}}$ and $\mathbf{y} - \bar{\mathbf{y}}$. The Pearson (correlation) distance can be written as

$$29 \quad D_{\text{corr}}(\mathbf{x}, \mathbf{y}) = 1 - \frac{(\mathbf{x} - \bar{\mathbf{x}}) \cdot (\mathbf{y} - \bar{\mathbf{y}})}{\|\mathbf{x} - \bar{\mathbf{x}}\| \|\mathbf{y} - \bar{\mathbf{y}}\|}. \quad (11)$$

30
 31 Cosine- and correlation-based measures are particularly useful in high-
 32 dimensional embedding spaces (e.g., word and sentence embeddings (Mikolov et
 33 al., 2013; Pennington et al., 2014; Reimers & Gurevych, 2019)) because they
 34 emphasise angular relationships and are less sensitive to variations in vector norms.
 35 Directional high-dimensional models, such as discriminant analysis with the von
 36 Mises-Fisher distribution, provide a statistical example in which angular structure is
 37 central rather than incidental (Romanazzi, 2014). The choice among raw cosine
 38 similarity, angular distance, or Pearson distance should be guided by whether
 39 magnitude, absolute direction, or mean-centred pattern is most relevant to the task.

40
 41

1 *Point Cloud Dissimilarity*

2
3 Before considering finite point cloud losses, it is useful to recall a classical
4 distance between sets. Given a pointwise Minkowski norm $\|\cdot\|_p$, the distance from
5 a point x to a non-empty set Y is $D_p(x, Y) = \inf_{y \in Y} \|x - y\|_p$. This induces the
6 Hausdorff-type set distance
7

$$8 \quad D_H^{(p)}(X, Y) = \max \left\{ \sup_{x \in X} D_p(x, Y), \sup_{y \in Y} D_p(y, X) \right\}. \quad (12)$$

9
10 For non-empty compact sets this is a genuine metric, and it is not restricted to
11 finite point clouds. For finite point sets, however, it is governed by the worst nearest-
12 neighbour discrepancy, making it highly sensitive to outliers (Deza & Deza, 2016).

13 For comparing unordered 3D point sets, a common alternative is the Chamfer
14 distance. Such point-set comparisons are central in 3D reconstruction settings,
15 including object reconstruction pipelines and broader surveys of deep-learning-
16 based architectural reconstruction (Enesi & Kuqi, 2023; Gorjian et al., 2025). Given
17 two point sets X and Y , it measures bidirectional nearest-neighbour discrepancies:
18

$$19 \quad D_{CD}(X, Y) = \sum_{x \in X} \min_{y \in Y} \|x - y\|^2 + \sum_{y \in Y} \min_{x \in X} \|x - y\|^2. \quad (13)$$

20
21 Intuitively, this measure penalises each point in one structure based on its
22 closest counterpart in the other structure, capturing local geometric consistency
23 without requiring explicit point correspondences. Naively, the nearest-neighbour
24 searches require comparing all pairs of points, although spatial indexes or
25 approximate search can reduce the practical cost. This makes it computationally
26 efficient relative to full optimal transport and well-suited for large, unordered point
27 clouds. Recent variants also attempt to learn or reweight Chamfer-style matching
28 costs for point cloud reconstruction tasks (Huang et al., 2024). In the form given
29 above, the Chamfer quantity is best understood as a dissimilarity measure rather
30 than a metric in the strict axiomatic sense. In particular, the nearest-neighbour
31 assignments are recomputed independently for each pair of point sets, so the
32 resulting quantity does not in general satisfy the triangle inequality.

33 However, because the Chamfer distance relies purely on nearest-neighbour
34 matching, it is primarily sensitive to local geometric deviations and does not
35 explicitly encode global structural relationships. As a result, it may assign low
36 distances to structures that share similar local features but differ significantly in their
37 overall fold or global topology.

38 An alternative to the Chamfer distance is the Earth Mover’s Distance (EMD),
39 a point-set dissimilarity grounded in optimal transport theory (Peyré & Cuturi, 2019;
40 Villani, 2009). To compare two point sets, EMD treats them as discrete distributions
41 of mass rather than as collections of independent nearest-neighbour queries. Let
42 $X = \{x_i\}_{i=1}^m$ and $Y = \{y_j\}_{j=1}^n$. Assign non-negative weights $a_i \geq 0$ to the points of
43 X and $b_j \geq 0$ to the points of Y , with $\sum_i a_i = \sum_j b_j = 1$. These weights specify

1 how much mass is located at each point. Also let $c(x_i, y_j) \geq 0$ be the ground cost
2 of moving one unit of mass from x_i to y_j .

3 The admissible transport plans are collected in

$$4 \quad \Pi(a, b) = \{\pi \in \mathbb{R}_{\geq 0}^{m \times n} : \sum_{j=1}^n \pi_{ij} = a_i, \sum_{i=1}^m \pi_{ij} = b_j\}. \quad (14)$$

6 Thus each coefficient $\pi_{ij} \geq 0$ specifies how much mass is transported from x_i to
7 y_j , while the row and column constraints ensure that all mass from X is sent and all
8 mass required by Y is received. With these definitions, EMD is the minimum total
9 transport cost

$$10 \quad D_{\text{EMD}}(X, Y) = \min_{\pi \in \Pi(a, b)} \sum_{i=1}^m \sum_{j=1}^n \pi_{ij} c(x_i, y_j), \quad (15)$$

11 so the distance reflects a globally optimal redistribution rather than a set of local
12 nearest-neighbour decisions. If we define the ground cost as $c(x, y) = \|x - y\|^p$,
13 we obtain the common formulation of the discrete Wasserstein distance:

$$14 \quad W_p(X, Y) = \left(\min_{\pi \in \Pi(a, b)} \sum_{i=1}^m \sum_{j=1}^n \pi_{ij} \|x_i - y_j\|^p \right)^{1/p}. \quad (16)$$

15 In this sense, EMD is the optimal transport formulation underlying
16 Wasserstein-type distances. In common usage, EMD often refers specifically to the
17 W_1 case with ground cost $\|x - y\|$. Optimal transport ideas have also been
18 extended beyond Euclidean ground spaces, for example to tropical geometric
19 settings where Wasserstein-type distances are defined on tropical projective spaces
20 (Lee et al., 2022).

21 In protein backbone comparison, for example, each structure can be represented
22 as a point set of 3D coordinates, typically corresponding to $C\alpha$ atoms. Chamfer
23 distance can provide a simple geometric baseline, but biologically meaningful
24 similarity often depends on global fold, alignment, and length normalisation rather
25 than local nearest-neighbour proximity (Zhang & Skolnick, 2004).

31 *Information-Theoretic Measures for Probability Distributions*

32 When the objects to be compared are probability distributions rather than
33 individual data points, classical vector-space metrics are generally not applicable.
34 Instead, a family of information-theoretic divergences and distances has been
35 developed to quantify distributional discrepancy.

36 The Kullback–Leibler (KL) divergence (Cover & Thomas, 2006; Kullback &
37 Leibler, 1951) measures the information lost when a distribution Q is used to
38 approximate a reference distribution P :

$$39 \quad D_{\text{KL}}(P \parallel Q) = \int_{\text{supp}(P)} P(x) \log \frac{P(x)}{Q(x)} dx. \quad (17)$$

1 To ensure the integral is well-defined, it is assumed that the support of P is
 2 contained within the support of Q ($\text{supp}(P) \subseteq \text{supp}(Q)$). Otherwise, the
 3 divergence becomes infinite. It answers the question: *how much additional*
 4 *information is required to encode samples from P using a code optimised for Q ?*
 5 Crucially, D_{KL} is not a metric: it is asymmetric ($D_{\text{KL}}(P \parallel Q) \neq D_{\text{KL}}(Q \parallel P)$) and
 6 unbounded, which limits its use in settings that require a proper distance structure.

7 The Jensen–Shannon (JS) divergence (Lin, 1991) remedies the asymmetry of KL
 8 by measuring how far P and Q each deviate from their mixture $M = 1/2 (P + Q)$:
 9

$$10 \quad D_{\text{JS}}(P \parallel Q) = \frac{1}{2} D_{\text{KL}}(P \parallel M) + \frac{1}{2} D_{\text{KL}}(Q \parallel M). \quad (18)$$

11
 12 The JS divergence is symmetric and bounded in $[0,1]$ (when using the base-2
 13 logarithm), and its square root satisfies the triangle inequality, making it a true
 14 metric. It finds widespread application in generative modelling and natural language
 15 processing.

16 The Hellinger distance (Hellinger, 1909) is defined as
 17

$$18 \quad D_{\text{H}}(P, Q) = \frac{1}{\sqrt{2}} \|\sqrt{P} - \sqrt{Q}\|_2. \quad (19)$$

19
 20 Here \sqrt{P} and \sqrt{Q} denote the pointwise square-root representations of the two
 21 distributions. For discrete probability vectors $P = (p_k)$ and $Q = (q_k)$, this is
 22

$$23 \quad D_{\text{H}}(P, Q) = \left(\frac{1}{2} \sum_k (\sqrt{p_k} - \sqrt{q_k})^2 \right)^{1/2}. \quad (20)$$

24
 25 The square-root transformation embeds probability distributions into a
 26 Euclidean geometry in which ordinary ℓ_2 distance becomes a bounded statistical
 27 distance. The Hellinger distance is therefore a genuine metric that is symmetric,
 28 bounded in $[0,1]$, and particularly robust for comparing distributions with differing
 29 supports. It is often preferred over the KL divergence in practical data-science
 30 applications, for instance in the presence of class imbalance, because it does not
 31 diverge when one distribution assigns zero probability to events that the other does
 32 not. Figure 4 contrasts this bounded behaviour with the unbounded growth of KL
 33 divergence in a simple two-outcome case.

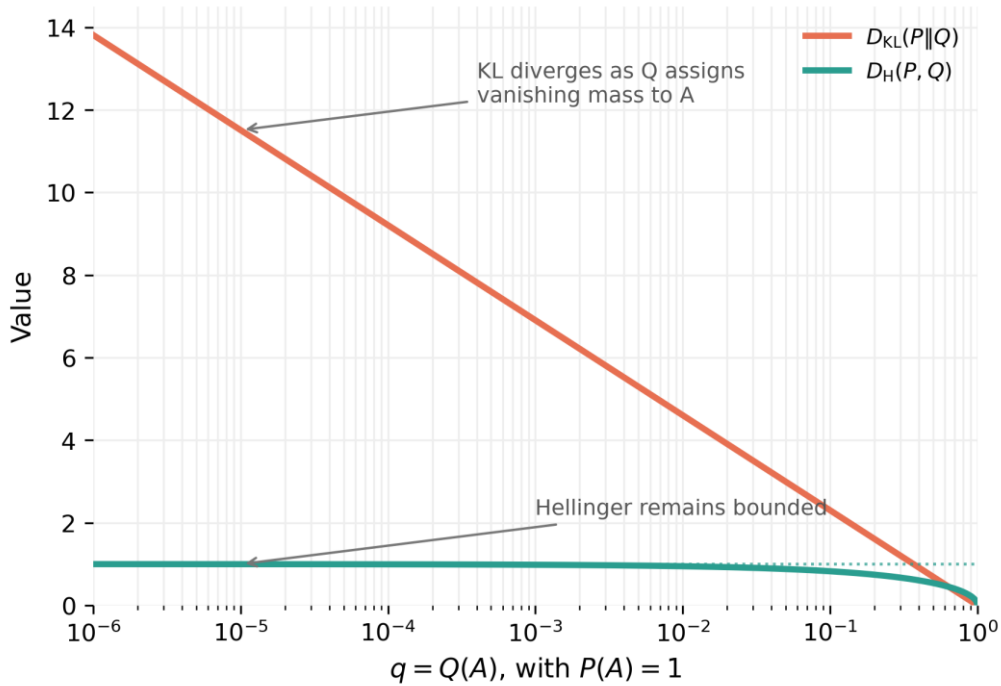
34 The Wasserstein distance extends the optimal-transport perspective introduced
 35 for point sets in Section 2.4 to probability distributions. Unlike KL, JS, and
 36 Hellinger distances, which compare probability mass through density values,
 37 Wasserstein distance also uses a ground metric on the sample space. It therefore
 38 measures how much work is required to move the mass of one distribution into the
 39 other. Formally, the p -Wasserstein distance is
 40

$$41 \quad W_p(P, Q) = \left(\inf_{\gamma \in \Gamma(P, Q)} \int \|x - y\|^p d\gamma(x, y) \right)^{1/p}, \quad (21)$$

42

1 where $\Gamma(P, Q)$ denotes the set of all couplings of P and Q (Peyré & Cuturi, 2019;
 2 Villani, 2009). The W_2 distance, often referred to as the *Fréchet distance* in the
 3 context of feature distributions, is of particular practical importance: it underlies the
 4 Fréchet Inception Distance (FID), a standard evaluation metric for generative image
 5 models (Heusel et al., 2017). A key advantage of the Wasserstein distance over KL
 6 and JS divergences is that it can remain meaningful even when two distributions
 7 have little or no overlap in support, a common situation in high-dimensional
 8 generative modelling. The trade-off is computational: exact optimal transport can
 9 be expensive for large empirical distributions.

10
 11 **Figure 4.** Comparison of KL divergence and Hellinger distance in a two-outcome
 12 example. Let $P(A) = 1$ and $Q(A) = q$. As q approaches zero, $D_{\text{KL}}(P \parallel Q) =$
 13 $\log(1/q)$ diverges, while $D_{\text{H}}(P, Q) = \sqrt{1 - \sqrt{q}}$ remains bounded by one. This
 14 illustrates why bounded distances can be more stable than KL divergence when
 15 distributions assign very different probability mass to the same event



16
 17
 18 Together, these measures span a spectrum of trade-offs between theoretical
 19 rigour, computational tractability, and sensitivity to distributional properties, and the
 20 appropriate choice depends heavily on the application at hand.

21 *Topological Data Analysis*

22
 23
 24 Topological Data Analysis provides a framework for characterising the shape
 25 of data by tracking the birth and death of topological features (connected
 26 components, loops, and voids) across a range of scales (Chazal & Michel, 2021).
 27 This range is usually organised as a filtration, meaning a nested sequence of spaces

1 built by gradually increasing a scale parameter, for example by connecting points
 2 that lie within a growing distance threshold. These features are summarised in a
 3 persistence diagram \mathcal{D} , a multiset of points $(b, d) \in \mathbb{R}^2$ with $b < d$, where b and
 4 d denote the birth and death scales of a topological feature respectively. Comparing
 5 two such diagrams requires metrics that respect this combinatorial structure. The
 6 two most widely used are the bottleneck distance and the diagram Wasserstein
 7 distance. The latter uses the same optimal-matching intuition as the transport
 8 distances discussed in Section 2.4 and Section 2.5, but the objects being matched
 9 are topological features in persistence diagrams rather than mass elements in the
 10 original data space.

11 The bottleneck distance between two persistence diagrams \mathcal{D}_1 and \mathcal{D}_2 is
 12 defined as

$$14 \quad D_\infty(\mathcal{D}_1, \mathcal{D}_2) = \inf_{\eta} \sup_{x \in \mathcal{D}_1} \|x - \eta(x)\|_\infty, \quad (22)$$

15 where the infimum is taken over all bijections $\eta: \mathcal{D}_1 \rightarrow \mathcal{D}_2$, with points on the
 16 diagonal $\{(a, a): a \in \mathbb{R}\}$ included in both diagrams to account for unmatched
 17 features. The bottleneck distance is thus governed by the *worst-case*
 18 correspondence: it records the largest displacement required to match any single
 19 feature across the two diagrams. This gives it strong stability guarantees under small
 20 perturbations of the filtration, but also means that a single large unmatched feature
 21 can dominate the distance.
 22

23 The q -Wasserstein distance on persistence diagrams generalises this by
 24 aggregating the cost over *all* matched pairs:

$$26 \quad D_{W,q}(\mathcal{D}_1, \mathcal{D}_2) = \left(\inf_{\eta} \sum_{x \in \mathcal{D}_1} \|x - \eta(x)\|_\infty^q \right)^{1/q}, \quad (23)$$

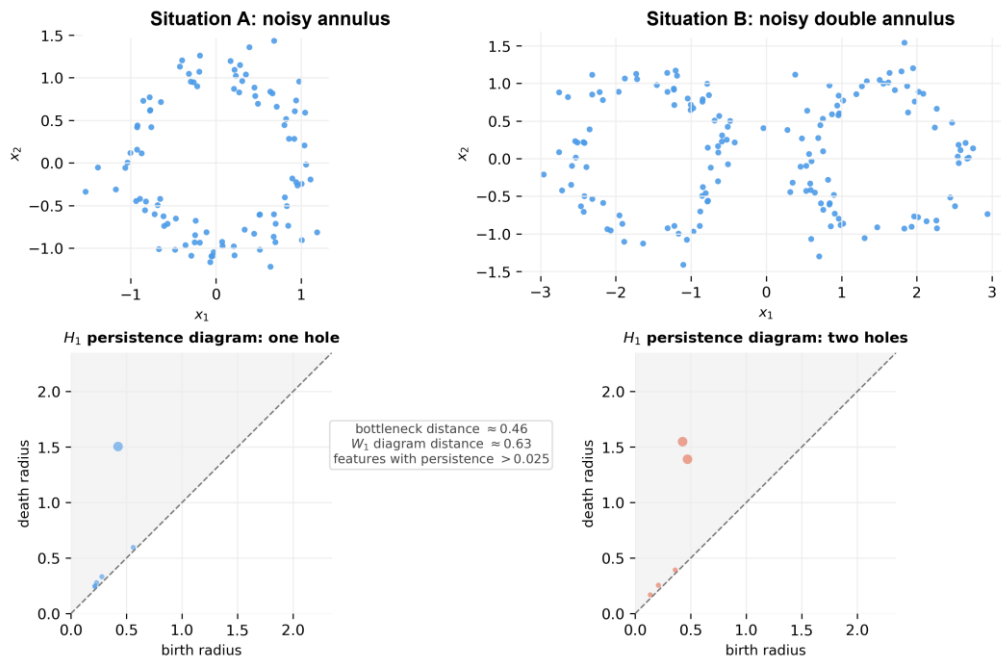
27 where again the infimum ranges over bijections that allow matching to the diagonal.
 28 In the limit $q \rightarrow \infty$ one recovers the bottleneck distance, while finite q penalises the
 29 total transport cost across the entire matching.
 30

31 Crucially, the use of the term ‘‘Wasserstein’’ here refers to a discrete optimal
 32 assignment problem between point sets (extending the perspective of Section 2.4)
 33 rather than the comparison of probability distributions detailed in Section 2.5. While
 34 both formulations share the foundational concept of minimising a global transport
 35 or matching cost, the diagram Wasserstein distance operates directly on the
 36 combinatorial structure of persistence diagrams, preserving its status as a proper
 37 metric without requiring an information-theoretic or probabilistic framework.

38 Figure 5 illustrates this construction on a small synthetic experiment. We
 39 sampled two noisy point clouds in the plane: one annular cloud with a single central
 40 hole and a second cloud formed from two separated annular components. A
 41 Vietoris–Rips filtration was then built by increasing the distance threshold at which
 42 sampled points are connected, and the resulting one-dimensional persistence
 43 diagrams were computed from the birth and death of loops. The plotted points retain
 44 features with persistence greater than 0.025, so short-lived near-diagonal

1 fluctuations are suppressed while the dominant holes remain visible. The bottleneck
 2 and diagram Wasserstein distances are then obtained by optimally matching these
 3 diagram points, allowing unmatched features to be assigned to the diagonal.

4
 5 **Figure 5.** Persistence diagrams for two noisy annulus-like point cloud examples. The
 6 top row shows the sampled point clouds: one annular cloud with a single hole and a
 7 second cloud containing two separated annular components. The bottom row shows
 8 the corresponding H_1 persistence diagrams after filtering to features with persistence
 9 greater than 0.025. The single-annulus example has one dominant off-diagonal
 10 point, while the double-annulus example has two dominant off-diagonal points. The
 11 displayed bottleneck and W_1 diagram distances summarise the cost of optimally
 12 matching these topological features, including possible matches to the diagonal



13
 14
 15 Both distances are stable with respect to small perturbations of the underlying
 16 data: a small change in the input filtration induces a correspondingly small change
 17 in the persistence diagram under either metric. This *stability property* is a
 18 cornerstone result of Topological Data Analysis and underpins the practical utility
 19 of these distances in noisy settings. Furthermore, both metrics are invariant to certain
 20 geometric transformations, such as rigid motions, provided the filtration is defined
 21 in a transformation-invariant manner. These properties make persistence-diagram
 22 distances well suited to applications in computational geometry, material science,
 23 and biological shape analysis, where data exhibit complex global structure that local,
 24 point-wise metrics fail to capture.

25 26 *Temporal Alignment*

27
 28 Standard point-wise distances, such as the Euclidean distance, are ill-suited for
 29 comparing time series that exhibit temporal misalignment: two sequences that are

1 otherwise identical but shifted in phase, or scaled in speed, will incur a large
 2 Euclidean distance despite being semantically similar. Dynamic Time Warping
 3 (DTW) (Sakoe & Chiba, 1978) addresses this limitation by seeking an optimal non-
 4 linear alignment between two sequences before measuring their discrepancy.
 5 Figure 6 shows how DTW avoids the misleading point-wise mismatch produced by
 6 a simple Euclidean comparison.

7 Formally, let $\mathbf{x} = (x_1, \dots, x_m)$ and $\mathbf{y} = (y_1, \dots, y_n)$ be two time series. A
 8 *warping path* $\mathcal{W} = (w_1, \dots, w_K)$ is a sequence of index pairs $w_k = (i_k, j_k) \in$
 9 $\{1, \dots, m\} \times \{1, \dots, n\}$ satisfying:

- 11 • *Boundary conditions*: $w_1 = (1, 1)$ and $w_K = (m, n)$.
- 12 • *Monotonicity*: $i_k \leq i_{k+1}$ and $j_k \leq j_{k+1}$ for all k .
- 13 • *Step continuity*: $i_{k+1} - i_k \leq 1$ and $j_{k+1} - j_k \leq 1$ for all k .

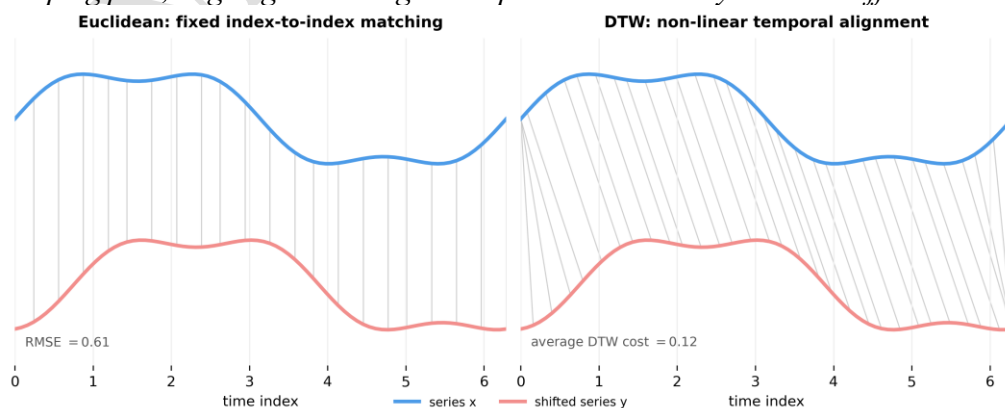
15 The DTW distance is then defined as the cost of the optimal warping path:

$$17 \quad D_{\text{DTW}}(\mathbf{x}, \mathbf{y}) = \min_{\mathcal{W}} \sum_{k=1}^K D(x_{i_k}, y_{j_k}), \quad (24)$$

19 where $D(\cdot, \cdot)$ is a local distance measure, typically the squared or absolute
 20 difference. This optimisation is solved exactly in $\mathcal{O}(mn)$ time via dynamic
 21 programming, making DTW tractable for moderate sequence lengths.

22 The warping path encodes a many-to-one alignment: a single time step in one
 23 sequence may be matched to multiple consecutive steps in the other, allowing DTW
 24 to absorb differences in speed, rhythm, or phase that would otherwise inflate a point-
 25 wise distance. It should be noted, however, that DTW does not satisfy the triangle
 26 inequality in general and is therefore a dissimilarity measure rather than a true
 27 metric.

29 **Figure 6.** *Euclidean matching versus Dynamic Time Warping for two shifted time*
 30 *series. A point-wise Euclidean comparison aligns samples at the same time index, so*
 31 *phase shifts create large apparent discrepancies. DTW instead allows a non-linear*
 32 *warping path, aligning similar signal shapes even when they occur at different times*



33
34

DTW has found widespread application across domains where temporal structure is central: in speech recognition it enables robust comparison of utterances spoken at different rates. It has also been used for pattern discovery in time series databases, approximate matching in vocal-music retrieval and scoring, and analogy-based pandemic forecasting (Berndt & Clifford, 1994; Ye, 2026; Zhang & Ji, 2026). In bioinformatics it aligns gene-expression time courses or electrocardiographic signals. In finance it measures the similarity of price trajectories that may lead or lag one another. These diverse use cases illustrate that DTW is not merely a technical fix for phase misalignment, but a principled modelling choice reflecting the belief that *when* events occur matters less than *whether* they occur and in what order.

Discrete and Set-Based Metrics

Continuous vector-space distances are not meaningful when data are inherently discrete (binary strings, categorical feature vectors, sets of tokens, etc.) and dedicated metrics are required. Three measures are particularly prominent in this setting.

The Hamming distance between two sequences $x, y \in \Sigma^n$ over an alphabet Σ counts the number of positions at which the two sequences disagree:

$$D_H(x, y) = \sum_{i=1}^n \mathbf{1}[x_i \neq y_i]. \quad (25)$$

For binary strings this reduces to the number of bit-flips required to transform one string into the other, which coincides with the L_1 norm on $\{0,1\}^n$. The Hamming distance is a true metric and forms the foundation of classical error-correcting codes, where the minimum Hamming distance between codewords determines the error-detection and correction capacity of a code. In machine learning it is widely used for comparing categorical feature vectors and hashed representations. Generalised Hamming-type distances have also been proposed when near misses or structured symbol relabellings should receive partial credit rather than being counted as exact mismatches (Bookstein et al., 2002; Liu & Na, 2025). In perceptual hashing, recent work has similarly argued that plain Hamming distance can miss spatial structure in image hashes (McKeown, 2025).

The Hamming distance can be applied to text when strings are represented as equal-length sequences of discrete characters. More general string dissimilarities are provided by *edit distances* (Deza & Deza, 2016), which measure the minimum number or cost of operations required to transform one string into another. A widely used example is the Damerau–Levenshtein distance, defined recursively as:

$$D_{a,b}(i, j) = \min \begin{cases} 0 & \text{if } i = j = 0, \\ D_{a,b}(i-1, j) + 1 & \text{if } i > 0, \\ D_{a,b}(i, j-1) + 1 & \text{if } j > 0, \\ D_{a,b}(i-1, j-1) + \mathbf{1}_{(a_i \neq b_j)} & \text{if } i, j > 0, \\ D_{a,b}(i-2, j-2) + \mathbf{1}_{(a_i \neq b_j)} & \text{if } i, j > 1 \text{ and } a_i = b_{j-1} \text{ and } a_{i-1} = b_j, \end{cases} \quad (26)$$

with $\mathbf{1}_{(a_i \neq b_j)}$ being the indicator function equal to 0 when $(a_i = b_j)$ and equal to 1 otherwise. Each recursive call corresponds to some text alteration e.g., deletion, insertion or mismatch of character. Other examples of popular edit distances are the Longest Common Subsequence (LCS) and the Jaro distance. Although they lost relevancy regarding NLP with the development of text tokenizers, edit distances are still used in practical ML applications such as spelling correctors.

When the objects of interest are *sets* rather than ordered sequences, cardinality-based overlap measures are more natural. The Jaccard index quantifies the similarity between two finite sets A and B as the ratio of their intersection to their union:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}, \quad J(A, B) \in [0, 1]. \quad (27)$$

A value of 1 indicates identical sets and 0 indicates disjoint sets. The corresponding *Jaccard distance* $1 - J(A, B)$ is a true metric. The Jaccard index is particularly prevalent in information retrieval, document deduplication via MinHash, and ecological studies comparing species compositions. Related overlap measures have also been analysed in recommendation and diversity settings, where the Jaccard index appears alongside Sørensen, cosine, and entropy-derived measures (Jost, 2006; Verma & Aggarwal, 2020).

The Dice coefficient (also known as the Sørensen–Dice index) is a closely related overlap measure:

$$\text{DSC}(A, B) = \frac{2|A \cap B|}{|A| + |B|}, \quad \text{DSC}(A, B) \in [0, 1]. \quad (28)$$

The factor of two in the numerator compensates for the fact that the shared elements $A \cap B$ are counted once in each of $|A|$ and $|B|$, rendering the denominator equivalent to $|A \cup B| + |A \cap B|$. Consequently the Dice coefficient up-weights the contribution of the overlap relative to the Jaccard index, and the two are related by

$$\text{DSC}(A, B) = \frac{2J(A, B)}{1 + J(A, B)}. \quad (29)$$

Unlike the Jaccard distance, $1 - \text{DSC}$ does not satisfy the triangle inequality and is therefore not a metric. Nevertheless, the Dice coefficient is the dominant evaluation criterion in medical image segmentation, where it directly measures the spatial overlap between a predicted and a ground-truth region, and it also serves as a differentiable training objective in segmentation networks when applied to soft, probabilistic predictions. Modified Dice coefficients continue to be proposed for medical-image settings where boundary location and spatial error matter in addition to raw overlap (Rainio & Klén, 2025).

Together, these three measures illustrate a recurring theme in metric design: the appropriate notion of dissimilarity is inseparable from the structure of the data. Imposing a continuous geometry on inherently discrete or combinatorial objects risks obscuring the very relationships one wishes to quantify.

1 *Similarity Measures Between Images and Videos*

2
3 Images and videos highlight the gap between signal-level and perceptual
4 similarity: two images can be close under a pixel norm yet visually different, or
5 visually similar despite small translations or intensity changes. This motivates a
6 hierarchy of measures, from pixel-wise errors to structural, information-theoretic,
7 and distribution-level criteria.

8 Pixel-wise measures such as Mean Squared Error (MSE) and Peak Signal-to-
9 Noise Ratio (PSNR) treat an image as a vector in $\mathbb{R}^{H \times W \times C}$. They are simple and
10 useful in compression or restoration benchmarks, but they correlate poorly with
11 human judgement when small spatial shifts or perceptually minor changes produce
12 large pointwise errors.

13 Structural and statistical measures address part of this limitation. Edge-
14 detection and facial-image-processing applications illustrate why the chosen image
15 similarity measure or feature representation can strongly affect downstream
16 recognition and matching pipelines (Arriagada et al., 2019; Rodrigues et al., 2016).
17 The Structural Similarity Index (SSIM) compares local luminance, contrast, and
18 structure, and was introduced precisely to better align image-quality assessment
19 with human perception (Wang et al., 2004). Normalised cross-correlation (NCC)
20 compares local intensity patterns after centring and scaling, making it useful for
21 template matching and optical flow when brightness changes are approximately
22 affine (Lewis, 1995). Mutual Information (MI) compares joint intensity statistics
23 and is therefore common in multi-modal image registration (Maes et al., 1997; Viola
24 & Wells, 1997).

25 For generative models, evaluation usually compares distributions of images
26 rather than paired examples. The Fréchet Inception Distance (FID) fits Gaussians to
27 real and generated Inception-v3 activations and computes the corresponding
28 Wasserstein-2 distance (Heusel et al., 2017). For video, Fréchet Video Distance
29 (FVD) applies the same principle to spatio-temporal features, making temporal
30 coherence part of the comparison (Unterthiner et al., 2018). These distribution-level
31 metrics are reference-free in the sense that they compare sets of generated and real
32 samples rather than paired examples, but they depend strongly on the chosen feature
33 encoder and its domain biases. Recent empirical studies of synthetic-image
34 evaluation continue to examine how FID and Inception Score behave for generative
35 models such as VAEs (Chan & Sithungu, 2025).

36 *Geodesic Distances and Manifold-Aware Metrics*

37
38
39 Many real-world datasets do not occupy a flat Euclidean space but instead lie
40 on or near a low-dimensional *manifold* embedded in a high-dimensional ambient
41 space. In such settings, the straight-line Euclidean distance between two points may
42 be geometrically misleading: it cuts through empty regions of the ambient space that
43 contain no data, whereas the true dissimilarity should follow the curvature of the
44 manifold. The geodesic distance addresses this by measuring the length of the
45 shortest path between two points *along the manifold*, rather than through the
46 surrounding space.

1 Formally, given a Riemannian manifold (\mathcal{M}, g) with metric tensor g , the
2 geodesic distance between two points $p, q \in \mathcal{M}$ is:

$$3 \quad 4 \quad D_{\text{geo}}(p, q) = \inf_{\gamma} \int_0^1 \sqrt{g_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))} dt, \quad (30)$$

5 where the infimum is taken over all smooth curves $\gamma: [0,1] \rightarrow \mathcal{M}$ with $\gamma(0) = p$
6 and $\gamma(1) = q$, and $\dot{\gamma}(t)$ denotes the tangent or velocity vector of the curve at time
7 t . When \mathcal{M} is a sphere S^{n-1} , this reduces to the great-circle distance. When $\mathcal{M} =$
8 \mathbb{R}^n with the standard metric, it coincides with the Euclidean distance.

10 Graph-based approximation and Isomap

11 In practice, the manifold \mathcal{M} is unknown and must be inferred from a finite
12 sample of data points. The standard approach is to construct a neighbourhood graph
13 G in which each point x_i is connected to its k nearest Euclidean neighbours (or all
14 neighbours within radius ε), with edge weights equal to the Euclidean distances
15 between connected points. The geodesic distance between any two points is then
16 approximated by the shortest path in G , computed efficiently via Dijkstra’s or
17 Floyd–Warshall’s algorithm. Recent graph-algorithm work continues to refine all-
18 pairs shortest-path computation, including Floyd–Warshall variants for sparse or
19 disconnected graphs (Zugan et al., 2025). This approximation converges to the true
20 geodesic distance as the sample density increases, under mild smoothness
21 conditions on \mathcal{M} .

22 Isomap (Tenenbaum et al., 2000) builds directly on this idea: it replaces the
23 Euclidean pairwise distance matrix used in classical Multidimensional Scaling
24 (MDS) with the matrix of graph-based geodesic distances, then applies MDS to find
25 a low-dimensional embedding that preserves these geodesic distances as faithfully
26 as possible. This yields a nonlinear dimensionality reduction method that correctly
27 “unfolds” curved manifolds where PCA and Euclidean MDS fail entirely. Figure 7
28 illustrates why the distinction between ambient and intrinsic distance matters in such
29 settings on the famous Swiss roll example.

31 Heat kernels and diffusion distances

32 A complementary family of manifold-aware metrics is based on the heat
33 equation on the data graph. The *diffusion distance* at time t between points x_i and
34 x_j is defined as:

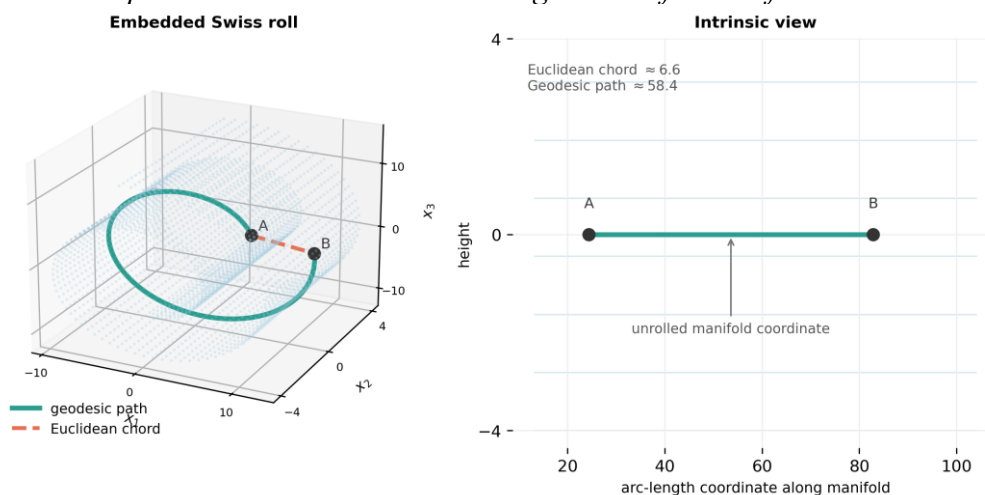
$$35 \quad 36 \quad D_t(x_i, x_j)^2 = \sum_k \lambda_k^{2t} (\phi_k(x_i) - \phi_k(x_j))^2, \quad (31)$$

37 where λ_k and ϕ_k are the eigenvalues and eigenvectors of the normalised graph
38 Laplacian. To construct this operator, one first forms a weighted adjacency matrix
39 W , where W_{ij} records the affinity between neighbouring data points, and a diagonal
40 degree matrix D with $D_{ii} = \sum_j W_{ij}$. The unnormalised graph Laplacian is $L = D -$
41 W , while common normalised variants include $L_{\text{sym}} = I - D^{-1/2}WD^{-1/2}$ and
42 $L_{\text{rw}} = I - D^{-1}W$. These matrices are discrete analogues of the Laplace–Beltrami

1 operator on a manifold: they compare the value of a function at a point with its
 2 values on neighbouring points, and therefore encode how signals, heat, or random
 3 walks propagate across the data graph.

4

5 **Figure 7.** *Ambient Euclidean distance versus intrinsic geodesic distance on a Swiss-
 6 roll manifold. In the embedded view, the Euclidean chord between two points cuts
 7 through the surrounding space and ignores the data manifold. The geodesic path
 8 instead follows the rolled surface. In the intrinsic unrolled coordinate system, this
 9 path corresponds to distance measured along the manifold itself*



10

11

12 Intuitively, the diffusion distance measures how differently heat (or a random
 13 walk) spreads from x_i versus x_j over t steps. This makes it robust to small
 14 perturbations and noise: two points connected by many short paths (i.e., well within
 15 the same region of the manifold) will have a small diffusion distance even if a single
 16 shortest path is corrupted. Diffusion Maps (Coifman & Lafon, 2006) use this
 17 distance to derive a low-dimensional embedding, analogously to Isomap but with
 18 stronger noise robustness and a probabilistic interpretation via random walks.

19

20 Geodesic distances on meshes and point clouds

21

22 When the manifold is explicitly represented as a triangulated surface mesh,
 23 exact geodesic distances can be computed via the Fast Marching Method (Sethian,
 24 1996), which propagates a wavefront across triangle faces in $O(N \log N)$ time, or
 25 via the Heat Method (Crane et al., 2013), which solves two sparse linear systems on
 26 the mesh and is particularly efficient on GPU hardware. These methods are standard
 27 in 3D shape analysis and computational geometry, where geodesic distances
 28 underpin shape descriptors, surface parameterisation, and shape retrieval.

28

29 Limitations

30

31 Geodesic approximations via neighbourhood graphs are sensitive to the choice
 32 of k or ϵ . If it is too small, the graph becomes disconnected. If it is too large, short-
 33 circuit edges bridge distinct regions of the manifold, corrupting the distance.
 34 Furthermore, geodesic distances are expensive to compute for large datasets as all-
 pairs shortest paths scale as $O(N^2 \log N)$ and are not straightforwardly

1 differentiable, complicating their use as training objectives in deep learning
2 pipelines.

5 **Theoretical Implications**

7 *Positive Definite Kernels and Their Relationship to Distance Metrics*

9 Kernel methods occupy a central place in machine learning, underpinning
10 support vector machines, Gaussian processes, kernel PCA, and a broad class of non-
11 parametric estimators, with kernel-trick formulations also appearing in regression
12 and classification schemes beyond standard SVMs (Huh, 2015). Their power
13 derives from the *kernel trick*: by implicitly mapping data into a high-dimensional
14 (possibly infinite-dimensional) feature space \mathcal{H} , linear operations in \mathcal{H} capture
15 nonlinear structure in the original input space without ever computing the map
16 explicitly. The central object is therefore a positive definite kernel, and the key
17 question for this review is when a distance can be used to construct one.

19 Positive definite kernels and RKHS embeddings

20 A symmetric function $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a positive definite kernel if, for any
21 finite collection of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{X}$ and coefficients $c_1, \dots, c_n \in \mathbb{R}$,

$$23 \quad \sum_{i=1}^n \sum_{j=1}^n c_i c_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0. \quad (32)$$

24
25 Equivalently, the Gram matrix G with entries $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ is positive semi-
26 definite for every finite sample. A canonical example is the Gaussian radial basis
27 function (RBF) kernel

$$29 \quad K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{y}\|_2^2}{2\sigma^2}\right), \quad (33)$$

30
31 where nearby vectors receive large similarity values and distant vectors receive
32 values close to zero. Under standard conditions, such kernels admit a feature map
33 $\phi: \mathcal{X} \rightarrow \mathcal{H}$ into a reproducing kernel Hilbert space (RKHS) such that

$$35 \quad K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_{\mathcal{H}}. \quad (34)$$

36
37 This factorisation is what makes the kernel trick possible: inner products in \mathcal{H}
38 can be evaluated directly through K without explicit construction of ϕ (Berg et al.,
39 1984; Berlinet & Thomas-Agnan, 2004).

41 From distance metrics to kernels: conditionally negative definite functions

42 Not every distance metric induces a valid PD kernel, and understanding which
43 distances do is essential for the correct application of kernel methods. The key
44 concept is that of a *conditionally negative definite* (CND) distance: a symmetric
45 function $d: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ with $d(\mathbf{x}, \mathbf{x}) = 0$ is CND if

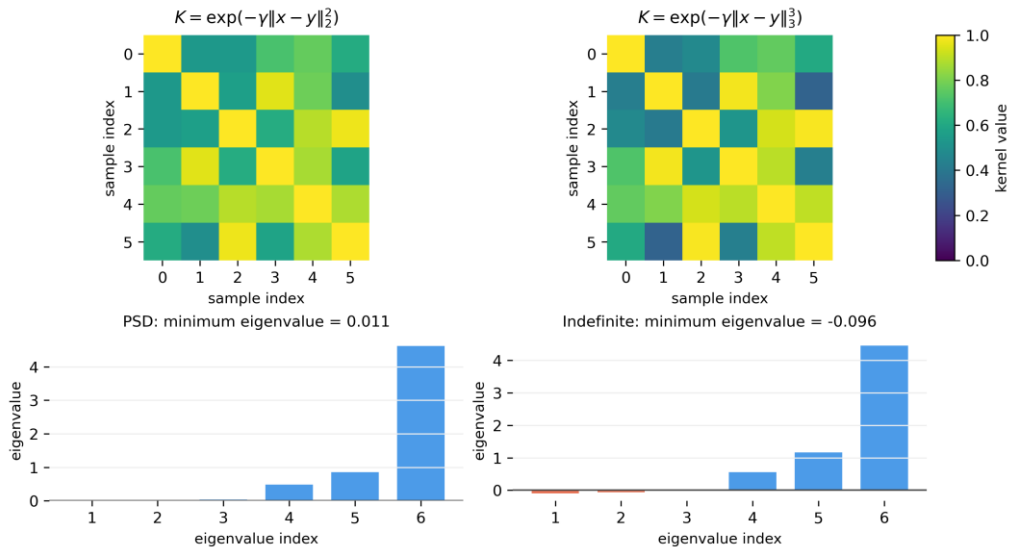
1
$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j d(\mathbf{x}_i, \mathbf{x}_j) \leq 0 \quad \text{whenever} \quad \sum_{i=1}^n c_i = 0. \quad (33)$$

2
3 Schoenberg’s theorem states that d is *CND* exactly when

4
5
$$K_t(\mathbf{x}, \mathbf{y}) = \exp(-t d(\mathbf{x}, \mathbf{y})) \quad (34)$$

6
7 is *PD* for every $t > 0$ (Berg et al., 1984; Ressel, 1976; Schoenberg, 1938). This
8 result has profound practical consequences: it precisely characterises which distance
9 functions can be “kernelized” via the ubiquitous Gaussian (RBF) kernel, and
10 explains why the RBF kernel constructed from a non-*CND* distance may yield an
11 indefinite Gram matrix, rendering the kernel method ill-posed. Figure 8 illustrates
12 this failure mode for an L_3 -based radial kernel.

13
14 **Figure 8.** Kernel validity depends on the distance used to construct the Gram
15 matrix. For the same finite point set, a radial kernel built from squared Euclidean
16 distance remains positive semi-definite, as expected from the *CND* property of that
17 distance. Replacing it with an analogous construction based on an L_3 distance can
18 produce a negative eigenvalue, showing that the resulting Gram matrix is indefinite
19 and therefore not a valid positive definite kernel



20
21
22 Which distances are conditionally negative definite?

23 Several canonical distances are known to be *CND*, and hence admit valid RBF
24 kernelizations:

- 25
26 • The squared Euclidean distance $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ is *CND*.
27 Consequently, the standard Gaussian kernel $K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|_2^2 /$
28 $2\sigma^2)$ is *PD*.

- 1 • For $0 < p \leq 2$, the L_p distance $\| \mathbf{x} - \mathbf{y} \|_p^p$ is *CND*. For $p > 2$ it is not, so
- 2 the RBF kernel built on $\| \mathbf{x} - \mathbf{y} \|_p$ with $p > 2$ may not be *PD*.
- 3 • The geodesic distance on a sphere is *CND*, enabling valid kernel methods on
- 4 spherical data. Related work extends Schoenberg-type characterisations of
- 5 positive definite kernels on spherical and bundled domains (Kuryatnikova &
- 6 Vera, 2019).
- 7 • The Hamming distance on binary strings is *CND*, and the corresponding RBF
- 8 (diffusion) kernel is widely used in bioinformatics.

9
10 Conversely, distances that are *not* *CND*, such as the L_p norm for $p > 2$, or

11 certain graph-theoretic distances, do not directly yield valid *PD* kernels via the RBF

12 construction. In such cases, practitioners must either modify the distance (e.g., by

13 taking a fractional power d^θ with $\theta \in (0,1)$, which can restore the *CND* property),

14 use indefinite kernel machinery, or seek an alternative explicit feature map.

15 The practical implication for metric selection is clear. Choosing a distance

16 metric is not merely a modelling decision about dissimilarity. It also determines

17 whether the metric is compatible with the kernel machinery one intends to deploy

18 downstream. A distance that faithfully captures domain geometry but fails the *CND*

19 condition will corrupt the Gram matrix and invalidate theoretical guarantees for

20 kernel SVMs, Gaussian processes, and related methods. Verifying or enforcing the

21 *CND* property should therefore be considered an integral part of metric design in

22 kernel-based pipelines (Berg et al., 1984; Schoenberg, 1938).

23 *Distance Metric Learning*

24
25
26 The distance metrics surveyed in Section 2 are largely hand-crafted, relying on

27 domain knowledge or geometric intuition to define dissimilarity. Distance Metric

28 Learning (DML) takes a data-driven perspective: rather than specifying a fixed

29 metric, it *learns* a distance function from labelled or weakly labelled data such that

30 the induced geometry reflects task-specific notions of similarity (Suárez et al.,

31 2021). The core premise is that performance in downstream similarity-based tasks

32 (nearest neighbour classification, clustering, etc.) is directly determined by the

33 geometry of the distance, and that this geometry can and should be optimised jointly

34 with or prior to the task itself.

35 Classical Mahalanobis metric learning

36
37 The most natural parametric family for DML is the class of Mahalanobis

38 distances

$$39 \quad D_M(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^\top \mathbf{M} (\mathbf{x} - \mathbf{y})}, \quad \mathbf{M} \succeq 0, \quad (37)$$

40 where the positive semi-definite matrix \mathbf{M} plays the role of Σ^{-1} in Equation 6.

41 Learning \mathbf{M} amounts to finding a linear transformation \mathbf{L} (with $\mathbf{M} = \mathbf{L}^\top \mathbf{L}$) that

42 maps the input space into a new Euclidean space in which the standard distance is

43 task-appropriate (Goldberger et al., 2004; Weinberger & Saul, 2009).

44
45

1 Deep metric learning

2 Classical Mahalanobis methods are limited to linear transformations of the
3 input space. Deep Metric Learning (DeepML) replaces the linear map \mathbf{L} with a deep
4 neural network $f_\theta: \mathcal{X} \rightarrow \mathbb{R}^d$, learning a non-linear embedding in which distances
5 reflect semantic similarity (Hoffer & Ailon, 2015; Kaya & Bilge, 2019).

7 Angular margin losses

8 A further line of development replaces Euclidean distance in the embedding
9 space with angular (cosine) distance and incorporates the margin directly into the
10 classification objective. Methods such as ArcFace (Deng et al., 2022) and CosFace
11 normalise both embeddings and weight vectors to lie on a hypersphere, then
12 introduce an additive angular margin m into the softmax logit:

$$14 \quad \mathcal{L}_{\text{Arc}} = -\log \frac{e^{s \cos(\theta_{y_i+m})}}{e^{s \cos(\theta_{y_i+m})} + \sum_{j \neq y_i} e^{s \cos \theta_j}}, \quad (38)$$

15
16 where θ_j is the angle between the embedding and the j -th class centre, and s is a
17 scale factor. Because the angular margin corresponds directly to a geodesic distance
18 on the hypersphere, ArcFace enforces tighter intra-class compactness and larger
19 inter-class separation than purely Euclidean objectives, and has become the de facto
20 standard loss in large-scale face recognition (Deng et al., 2022).

22 Supervision regimes and sampling

23 DML methods differ not only in their loss function but also in the form of
24 supervision they require. Pair- and triplet-based methods rely on relative constraints
25 (“ x_i is more similar to x_j than to x_k ”), which are often cheaper to obtain than
26 absolute class labels. This has motivated *self-supervised* and *contrastive*
27 *representation learning* variants (such as SimCLR and MoCo) that generate
28 positive pairs through data augmentation without any manual annotation, effectively
29 learning a metric purely from the structure of the data (Chen et al., 2020; He et al.,
30 2020).

32 Applications and broader impact

33 DML has achieved state-of-the-art results in face verification and recognition,
34 person re-identification, image retrieval, few-shot learning, and recommendation
35 systems (Kaya & Bilge, 2019). In contemporary ML, related metric choices also
36 appear in transformer attention, where scaled dot products define token-level
37 similarity, and in multimodal contrastive systems such as CLIP, where image and
38 text embeddings are aligned by cosine-like objectives (Radford et al., 2021;
39 Vaswani et al., 2017). Diffusion models likewise rely on feature-space and
40 distributional distances for training diagnostics and evaluation, even when their
41 generative objective is not itself a distance metric (Ho et al., 2020). Beyond
42 performance, DML represents a conceptual shift in how distance is treated in the
43 ML pipeline: rather than a fixed geometric assumption about the data, the metric
44 becomes a *learned inductive bias*, shaped by the task, the data distribution, and the
45 available supervision signal. This view reinforces the central argument of the

1 present paper: metric selection and design are modelling decisions with
2 consequences for learning.

3 The same issue reappears in representation learning more broadly. DML
4 explicitly trains an embedding so that distances serve a task objective. Latent-
5 variable and generative models may also produce embeddings, but without a metric-
6 learning objective there is no guarantee that Euclidean latent distance captures
7 semantic similarity. This motivates the latent-space fidelity question considered
8 next.

9 10 *Latent Space Fidelity*

11
12 Representation learning methods such as autoencoders and variational
13 autoencoders (VAE) compress high-dimensional observations $x \in \mathcal{X}$ into latent
14 codes $z \in \mathcal{Z}$. A central question is whether distances in \mathcal{Z} reflect meaningful
15 semantic or geometric relationships in \mathcal{X} . This issue is directly relevant to
16 downstream tasks such as interpolation, nearest-neighbour search, clustering, and
17 anomaly detection.

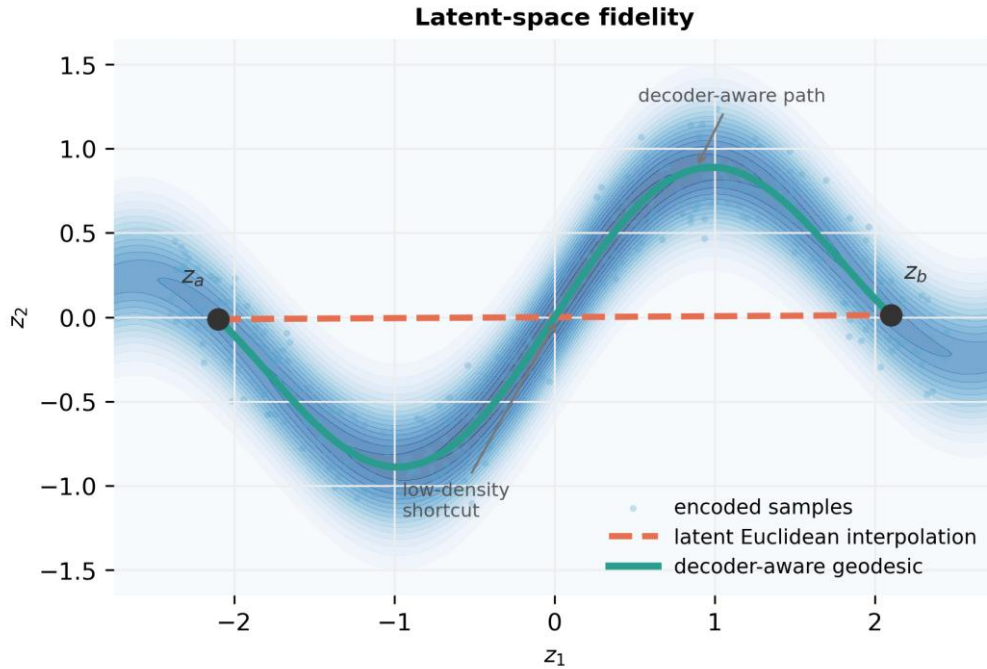
18 The simplest assumption is that Euclidean distance in latent space is
19 meaningful. In practice this often fails because the decoder is nonlinear, so straight
20 lines in \mathcal{Z} need not correspond to natural paths on the data manifold in observation
21 space. As a result, latent Euclidean distance may reflect artefacts of the
22 parametrisation rather than genuine similarity.

23 A principled alternative is to endow the latent space with the pullback
24 Riemannian metric induced by the generator. Arvanitidis, Hansen, and Hauberg
25 (Arvanitidis et al., 2018) show that this leads to geodesics that better respect the
26 learned manifold geometry than naive Euclidean interpolation. Related work further
27 argues that uncertainty is crucial for recovering meaningful latent geometry and that
28 purely deterministic embeddings may distort or obscure the manifold structure
29 (Arvanitidis et al., 2021; Hart et al., 2009). Figure 9 illustrates why a straight latent
30 interpolation can differ from a decoder-aware geodesic path. The figure is synthetic:
31 latent samples were drawn around a one-dimensional curved manifold $z_2 =$
32 $0.72\sin(1.75z_1) + 0.18z_1$ with additive Gaussian noise, and the background
33 density was computed from distance to that same curve.

34 A separate failure mode is posterior collapse in VAEs, where the latent
35 variables become uninformative because the decoder can model the data without
36 using them. In that regime, latent distances lose semantic content altogether
37 (Bowman et al., 2016; Ichikawa & Hukushima, 2024; Lucas et al., 2019). The
38 broader lesson is that the geometry of a learned representation should not be
39 assumed. It must be justified by the model and, when necessary, replaced by a
40 geometry induced by the decoder or by explicit metric-preservation objectives.

41
42

1 **Figure 9.** *Latent-space fidelity in a learned representation. A straight Euclidean*
 2 *interpolation between latent codes z_a and z_b can pass through low-density or*
 3 *semantically invalid regions of the latent space. A decoder- or pullback-aware*
 4 *geodesic instead follows the learned high-density geometry, better preserving*
 5 *meaningful variation in the observation space*



6
7
8
9

Geometric Transformations and Induced Metrics

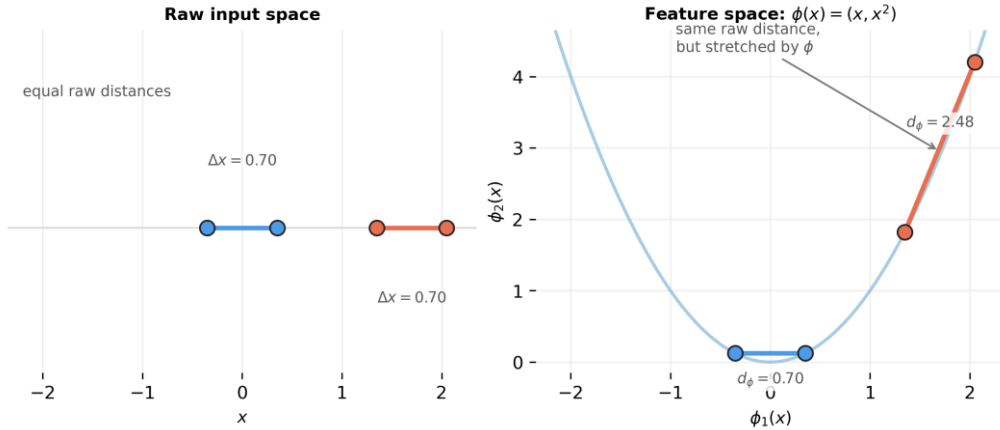
10 A central but often overlooked insight in machine learning is that every feature
 11 transformation implicitly defines a new distance. If data are mapped from an input
 12 space \mathcal{X} into a feature space \mathcal{F} by $\phi: \mathcal{X} \rightarrow \mathcal{F}$, then many algorithms compare $\phi(x)$
 13 and $\phi(y)$ rather than x and y :

$$14 \quad 15 \quad D_\phi(x, y) = \|\phi(x) - \phi(y)\|_2 \quad (39).$$

16 Thus the modelling question is not only which distance is used, but in which
 17 representation it is used.

18 As a compact example, a polynomial feature map sends an input vector to
 19 monomials of its coordinates. Geometrically, this changes the notion of proximity
 20 by comparing points after nonlinear expansion rather than in the original
 21 coordinates. Two inputs that are moderately separated in the raw coordinates may
 22 become much farther apart after quadratic or higher-order terms are introduced,
 23 while other directions may be compressed. This is why a linear separator after a
 24 polynomial expansion corresponds to a nonlinear decision boundary in the original
 25 space: the algorithm is operating in the geometry of the transformed representation.
 26 Figure 10 illustrates this induced-distance effect for a simple quadratic feature map.
 27

1 **Figure 10.** A nonlinear feature map induces a new distance geometry. In the raw
 2 input space, two point pairs have the same Euclidean separation. After the
 3 quadratic feature map $\phi(x) = (x, x^2)$, Euclidean distance in feature space assigns
 4 different separations to the two pairs. The operative distance is therefore
 5 $D_\phi(x, y) = \|\phi(x) - \phi(y)\|_2$, not the raw input distance



6
7

8 Kernel methods make the same idea implicit. A kernel $K(x, y) =$
 9 $\langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$ defines the RKHS distance

10

$$11 \quad D_K(x, y) = \sqrt{K(x, x) - 2K(x, y) + K(y, y)}, \quad (40)$$

12

13 which may differ substantially from Euclidean distance in the input space
 14 (Schölkopf et al., 2001). For the Gaussian RBF kernel, for example, $K(x, y) =$
 15 $\exp(-\|x - y\|^2 / 2\sigma^2)$ maps Euclidean separation into a bounded similarity scale
 16 controlled by σ : small changes inside the bandwidth are treated as highly similar,
 17 while points beyond a few bandwidths become nearly indistinguishable from one
 18 another in terms of kernel similarity. The chosen bandwidth is therefore a geometric
 19 modelling choice, not merely a numerical hyperparameter.

20

21 Deep representations follow the same principle: contrastive, classification, and
 22 reconstruction losses can induce different notions of proximity. Even standard
 23 architectural operations such as batch normalisation and layer normalisation alter
 24 the scale and orientation of intermediate representations, and therefore change
 25 which directions in activation space are emphasised by subsequent dot products or
 26 Euclidean comparisons (Ba et al., 2016; Ioffe & Szegedy, 2015). The practical
 27 implication is straightforward: when a model transforms the data, the operative
 28 distance is the one induced by that transformation, not necessarily the Euclidean
 29 distance in the raw input coordinates.

29

30 *Distances in Projective Geometry*

31

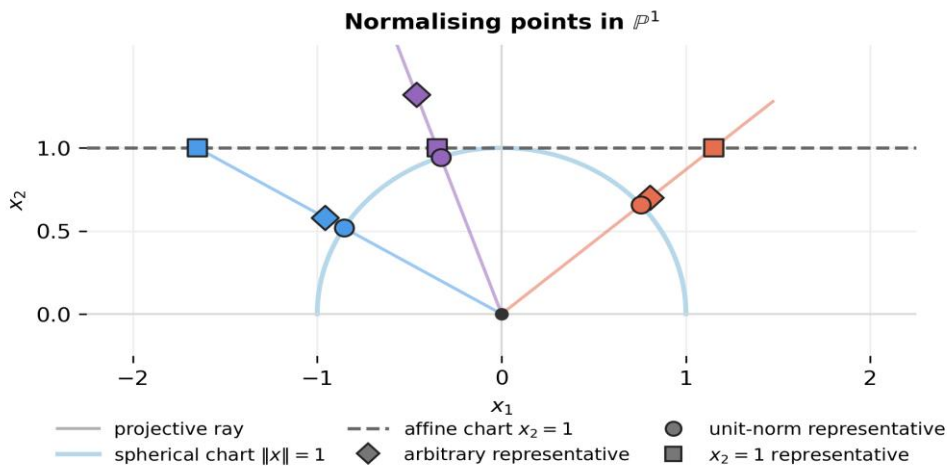
32 Projective representations illustrate a recurring theme of this review: the
 33 coordinates used by an algorithm are not always coordinates in which Euclidean
 34 distance is meaningful. Straight lines in 3-space, for example, can be represented by
 35 Plucker coordinates in \mathbb{P}^5 (De Boi et al., 2023; Pottmann et al., 1999), while planar

1 homographies are represented by 3×3 matrices defined only up to scale (Hartley
2 & Zisserman, 2004). In both cases, multiplying the coordinate vector or matrix by a
3 non-zero scalar leaves the represented geometric object unchanged.

4 This scale invariance makes naive Euclidean distances coordinate dependent.
5 Normalising homogeneous coordinates by unit norm, by a chosen component such
6 as $h_{33} = 1$, or by dehomogenisation leads to different numerical distances and may
7 therefore change the outcome of matching, estimation, or clustering. Figure 11
8 illustrates this dependence on the chosen representative in the simple case of \mathbb{P}^1 .
9 Alternatives, such as angular distances (see also Section 2.3), geometry-aware line
10 distances, cross-ratios (computed for 4 points rather than 2) or weighted projective
11 geometry (in case the scale has a relevant interpretation) are often preferable when
12 the task at hand is to compare the represented objects rather than their arbitrary
13 coordinate representatives (Pottmann & Wallner, 2009).

14 In computer vision, homography estimation, line matching, and multi-view
15 reconstruction compare objects that are naturally defined only up to projective scale.
16 A distance on coordinate arrays can therefore disagree with the geometric error
17 relevant to the task. Related line-geometric calibration and reconstruction settings
18 illustrate the same point: the metric should compare the represented geometric
19 object, not merely its coordinates (De Boi et al., 2021; De Boi et al., 2022). This
20 section is included as a reminder that some data representations carry invariances
21 before any learning algorithm is applied, and that those invariances should be
22 reflected in the chosen dissimilarity. Similar representation-dependent issues arise
23 in camera modelling and line-calculus-based optical calibration, where the
24 operational distance is tied to the chosen geometric parametrisation rather than to
25 raw image or actuator coordinates (De Boi et al., 2024; Penne et al., 2023).

26
27 **Figure 11.** *Different representatives of points in \mathbb{P}^1 . Each coloured ray represents*
28 *one projective point. Diamonds show arbitrary homogeneous representatives,*
29 *circles spherical normalisation, and squares affine normalisation. Distances*
30 *depend on the chosen representation: after affine normalisation, the distance*
31 *between the blue and purple representatives is close to the distance between the*
32 *purple and orange representatives, whereas after spherical normalisation these two*
33 *distances become substantially different*



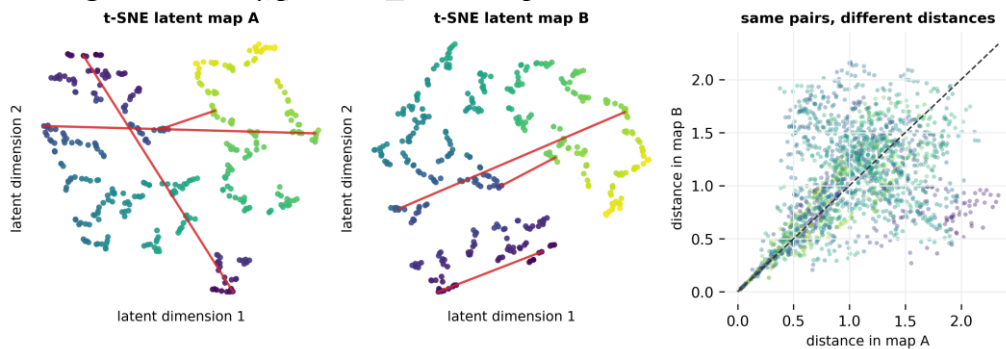
34

1 *Identifiability of Distance Metrics*

2
3 Identifiability asks whether geometric quantities inferred from data are
4 uniquely determined, at least up to an accepted equivalence. In latent variable
5 models, different latent coordinate systems may represent the same observed data
6 distribution, so Euclidean distances between latent codes are not automatically
7 interpretable. Related issues arise in nonlinear independent component analysis and
8 identifiable variational autoencoders, where additional structure or auxiliary
9 variables are needed to recover meaningful latent factors (Hyvärinen & Morioka,
10 2016; Khemakhem et al., 2020).

11 A principled response is to focus on geometry induced by the model, such as
12 pullback metrics and geodesic distances, rather than on raw coordinates. Recent
13 work shows that, under suitable regularity conditions, such metric structure can be
14 identifiable even when the latent parametrisation is not (Syrota et al., 2025).
15 Figure 12 illustrates this issue with two t-SNE embeddings of the same high-
16 dimensional Swiss-roll data. Although both embeddings are generated from the
17 same observations, different random initialisations lead to different two-
18 dimensional coordinates and different Euclidean distances between matched point
19 pairs (Maaten & Hinton, 2008). Related uncertainty-aware extensions of
20 multidimensional scaling reinforce the same caution for unstable or uncertain
21 embeddings (Hagele et al., 2023).

22
23 **Figure 12.** *Two latent representations of the same observed high-dimensional points.*
24 *The data are sampled from a Swiss-roll manifold embedded in 20 dimensions, then*
25 *reduced to two dimensions by two t-SNE runs with different random initialisations.*
26 *The colour encodes the same underlying manifold parameter in both maps. The right*
27 *panel compares distances for the same pairs of observations in the two latent maps,*
28 *showing that not every pairwise distance is preserved*



29 30 31 32 **Computational Aspects of Distance Functions**

33
34 Beyond their distinct use cases, the distances covered in this paper exhibit
35 significant variation in computational cost. This is especially relevant for ML
36 algorithms and models that are based on similarity measures, such as k-means or
37 Gaussian processes, as they require computing distances not only between two
38 specific points but between every possible pair of points in a potentially large

1 dataset. As a result, even inexpensive distances with pairwise complexity $\mathcal{O}(d)$
 2 ultimately incur a total cost of $\mathcal{O}(N^2d)$ when used in such models, where N is the
 3 size of the dataset or batch and d is the dimensionality of a sample. It is therefore
 4 important to understand both the cost of a single comparison and the number of
 5 times such comparisons are performed within a model. This section focuses on
 6 pairwise comparison complexities.

7 Providing a uniform frame for analysing the computational complexity of these
 8 distances is nearly impossible, as computational costs depend on multiple factors.
 9 The time required to compute a distance depends on the specific algorithm used, the
 10 structure and shape of the data on which it is computed, and the hardware used to
 11 perform the computations. For computationally expensive metrics specifically, the
 12 literature is constantly evolving, with new approximation methods, parallel
 13 algorithms optimised to run on GPUs, and memory optimisations. Therefore, this
 14 section only aims to provide a broad overview of computational considerations to
 15 keep in mind when choosing a metric for a specific application.

16 Minkowski L_p norms, cosine similarity, angular distances and Pearson distance
 17 are among the simplest distances in terms of computation. They require a single pass
 18 over the vector components, leading to an $\mathcal{O}(d)$ complexity, where d is the
 19 dimensionality of the vectors. GPU-optimised algorithms are available (Eslami et al.,
 20 2017).

21 Set-based metrics can be computed efficiently in $\mathcal{O}(|A| + |B|)$, where $|A|$ and
 22 $|B|$ are the cardinalities of the two sets A and B . For very large sets, this complexity
 23 can be drastically reduced using hashing algorithms such as MinHash (Broder, 1997).

24 Image distances that operate pixel-wise naturally incur an $\mathcal{O}(WHC)$
 25 computational complexity, where W and H denote the width and height of the
 26 images, respectively, and C represents the number of colour channels. Because the
 27 Structural Similarity Index Measure (SSIM) relies on a sliding window, its
 28 complexity escalates to $\mathcal{O}(WHk^2)$, where k is the window’s side length (Levy et
 29 al., 2025). To alleviate these computational overheads, domain-specific
 30 optimisations can be deployed, such as *Fast SSIM* (Chen & Bovik, 2011).
 31 Conversely, the Fréchet Inception Distance (FID) represents one of the most
 32 computationally demanding metrics, as its evaluation requires computing the matrix
 33 square root of a covariance product. This matrix operation scales cubically, resulting
 34 in an $\mathcal{O}(d^3)$ complexity, where d is the dimensionality of the underlying Inception-
 35 v3 latent embeddings.

36 Metrics that involve a minimisation problem often require more complex
 37 computations. Edit distances can be computed using a dynamic programming
 38 algorithm in $\mathcal{O}(l^2)$ time, where l is the length of the strings. DTW computation also
 39 has quadratic complexity, although constraints such as the Sakoe–Chiba
 40 band (Sakoe & Chiba, 1978) can provide cheaper lower-bound approximations in
 41 practice. Chamfer distance, as a sum of nearest-neighbour distances, can also be
 42 computed in $\mathcal{O}(N^2)$, where N is the number of samples.

43 For EMD, the optimisation problem involves a computation of complexity
 44 $\mathcal{O}(n^3 \log n)$, where n is the number of bins. This complexity can be reduced when
 45 working with one-dimensional data, or by using one of many approximation

1 methods (Shirdhonkar & Jacobs, 2008), such as the Sinkhorn algorithm (Cuturi,
2 2013).

3 As seen with the FID, metrics that involve covariance matrices often incur
4 higher computational costs. The Mahalanobis distance requires a Cholesky
5 decomposition of the covariance matrix, implying an $\mathcal{O}(d^3)$ preprocessing cost,
6 where d is the dimensionality of the vectors. However, this decomposition only
7 needs to be computed once and can then be reused for subsequent comparisons,
8 reducing the complexity of each comparison to $\mathcal{O}(d^2)$ due to the matrix–vector
9 multiplication.

10 Measures for probability distributions provide good examples of the difficulty
11 of summarising the complexity of a metric with a single figure. In the case of
12 discrete distributions, these measures can often be computed in linear time.
13 Computing them on continuous distributions typically involves numerical
14 integration, making them intractable in many cases (Pérez-Cruz, 2008). Recent
15 literature has explored estimation methods based on k-nearest-neighbours to
16 approximate the KL divergence (Zhao & Lai, 2020). On the other hand, restricting
17 the problem to specific distributions (e.g., Gaussian distributions) often leads to
18 analytical solutions that can be computed efficiently. For example, the KL
19 divergence between two d -dimensional Gaussian distributions can be computed in
20 $\mathcal{O}(d^3)$ time for full covariance matrices, but this reduces to $\mathcal{O}(d)$ for isotropic
21 Gaussians.
22
23

24 **Practical Synthesis: Choosing a Metric**

25
26 The central practical question is not which distance is best in the abstract, but
27 which assumptions about the data and task should be made explicit. Table 1
28 summarises choices that are mainly driven by the data representation, while Table 2
29 summarises choices that are mainly driven by modelling or theoretical constraints.

30 A hand-crafted metric is often appropriate when the data type already suggests
31 a clear invariance: coordinate-wise comparison for vectors, covariance-aware
32 comparison for correlated features, angular comparison for embeddings, projective
33 normalisation for homogeneous coordinates, Chamfer or EMD for point clouds,
34 divergence or transport distances for distributions, diagram distances for topological
35 summaries, DTW for temporally misaligned sequences, Jaccard or Hamming
36 distance for discrete objects, perceptual or feature-based criteria for images and
37 videos, and geodesic or diffusion distances for manifold-like data. These choices are
38 attractive because their assumptions are transparent, but they can fail when the task-
39 specific notion of similarity differs from the generic geometry of the data type.

40 Metric learning becomes preferable when supervision or reliable weak
41 supervision is available and the goal is predictive performance in a specific task.
42 Kernel-compatible distances become important when the chosen dissimilarity will
43 be used to form a positive definite kernel. Pullback and induced metrics become
44 preferable when the central assumption is geometric rather than label-based:
45 similarity should follow a feature map, a learned representation, or a decoder-
46 induced surface. Common mistakes are to ignore feature scaling, apply Euclidean

1 distance after a nonlinear transformation without asking what geometry was
 2 induced, or use a computationally expensive metric without checking whether its
 3 additional structure changes the downstream decision.

4
 5
 6

Table 1. Compact guide for selecting representation-driven metrics or dissimilarity measures

Data or task	Useful choice	Metric status	Main caution
Scaled vectors	Euclidean, Manhattan, Chebyshev, Minkowski	Yes for $p \geq 1$, no for $p < 1$ because triangle inequality fails	Sensitive to scaling
Correlated features	Mahalanobis / learned PSD metric	Yes if matrix is positive definite, pseudo-metric if only semi-definite	Requires stable estimation
Sparse text or embeddings	Cosine similarity / angular distance	Cosine similarity is not a distance, angular distance is a metric	Scale information is removed
Projective data	Normalisation-aware projective distance	Conditional, coordinate distances depend on the chosen normalisation	Coordinates depend on gauge
Point clouds	Chamfer / EMD	Chamfer is not a metric because triangle inequality can fail, EMD is a metric under standard transport assumptions	Local and global costs differ
Distributions	KL, JS, Hellinger, Wasserstein	KL is not, JS is not unless square-rooted, Hellinger and Wasserstein are metrics	Support and symmetry differ
Persistence diagrams	Bottleneck / diagram Wasserstein	Yes on persistence diagrams with diagonal matching	Sensitive to filtration
Time series	Dynamic Time Warping	No, triangle inequality can fail	Can over-warp sequences
Sets or binary data	Jaccard / Hamming	Jaccard distance and Hamming distance are metrics	Ignores feature dependence
Strings or discrete ordered sequences	Damerau–Levenshtein / LCS / Jaro	Levenshtein is a metric. Damerau–Levenshtein and LCS-based distances depend on the precise variant, while Jaro relaxes the triangle inequality	Purely character-based, no grammar or semantic encoded
Images or videos	MSE, PSNR, NCC, SSIM, MI, FID, FVD	Root-MSE is a metric, most others are not because symmetry, identity, or triangle inequality fail	Encoder bias can dominate
Manifold-like data	Geodesic / diffusion distance	Usually yes when the graph or manifold construction is valid	Graph construction is fragile

1 **Table 2.** *Compact guide for model-dependent and theoretical metric choices.*

Data or task	Useful choice	Metric status	Main caution
Kernel methods	CND-compatible distances	Conditional, CND guarantees kernel validity but not all metric axioms	Indefinite Gram matrices
Supervised embeddings	Learned metric	Conditional, depends on whether the learned form enforces metric axioms	May overfit supervision
Latent generative models	Pullback / decoder metric	Conditional, valid when the induced geometry is regular	Euclidean latent distance can mislead
Feature-transformed data	Induced feature-space metric	Metric if the feature map is injective, otherwise pseudo-metric	Raw-space intuition may fail
Identifiability questions	Geometry-level comparison	Not a single metric, but a check on whether the distance is uniquely determined	Identifiability must be checked

2
3 The tables should be read as decision aids rather than rankings. The most
4 reliable choice is usually the simplest metric whose invariances match the scientific
5 or engineering question. More complex learned or geometric distances are justified
6 when they encode information that the simpler metric cannot.

9 Conclusions

10
11 This review has highlighted the diversity and importance of distance metrics
12 and dissimilarity measures in modern Machine Learning. Moving beyond Euclidean
13 distance is not merely a matter of replacing one formula by another. It changes the
14 geometry in which learning takes place, the invariances encoded by the model, and
15 the meaning of similarity used by downstream algorithms.

16 Across the examples considered here, a common pattern emerges: each metric
17 becomes appropriate only under particular assumptions about the data. When those
18 assumptions are violated, the resulting notion of distance may be computationally
19 convenient but semantically misleading.

20 The theoretical perspectives discussed in this paper reinforce the same point. In
21 modern ML pipelines, the operative metric is often induced by a feature map, a
22 learned embedding, a kernel, or a decoder-defined geometry rather than chosen
23 directly in the raw input space. Open problems remain in making metric choice more
24 auditable, especially for large foundation models, multimodal embeddings, and
25 generative models whose learned spaces are difficult to interpret. For both
26 researchers and practitioners, the main message is straightforward: distance is a
27 modelling choice, and it should be justified with the same care as architecture, loss,
28 or prior.

29
30

1 **References**

- 2
- 3 Aggarwal CC, Hinneburg A, Keim DA (2001) On the surprising behavior of distance
4 metrics in high dimensional space. Lecture notes in computer science (including
5 subseries lecture notes in artificial intelligence and lecture notes in bioinformatics)
6 1973
- 7 Arriagada RC, Apablaza RC, Pinninghoff MA (2019) Efficient Edge Detection by Adapting
8 Artificial Bee Colony Algorithm. Athens Journal of Technology and Engineering 6(3)
9 : 179–194
- 10 Arvanitidis G, Hansen LK, Hauberg S (2018) Latent Space Oddity: On the Curvature of
11 Deep Generative Models. 6th international conference on learning representations
- 12 Arvanitidis G, Hauberg S, Schölkopf B (2021) Geometrically Enriched Latent Spaces.
13 Proceedings of the 24th international conference on artificial intelligence and statistics
14 130 : 631–639
- 15 Ba JL, Kiros JR, Hinton GE (2016) Layer Normalization. arXiv preprint arXiv:1607.06450
- 16 Bedalli E, Ninka I (2015) Exploring an Educational System's Data through Fuzzy Cluster
17 Analysis. Athens Journal of Sciences 2(1) : 33–44
- 18 Berg C, Christensen JPR, Ressel P (1984) Harmonic Analysis on Semigroups: Theory of
19 Positive Definite and Related Functions. Springer
- 20 Berline A, Thomas-Agnan C (2004) Reproducing Kernel Hilbert Spaces in Probability and
21 Statistics. Kluwer Academic Publishers
- 22 Berndt D, Clifford J (1994) Using Dynamic Time Warping to Find Patterns in Time Series.
23 Workshop on knowledge discovery in databases 398
- 24 Bookstein A, Kulyukin VA, Raita T (2002) Generalized Hamming Distance. Information
25 Retrieval 5(4)
- 26 Bowman SR, Vilnis L, Vinyals O, Dai AM, Jozefowicz R, Bengio S (2016) Generating
27 Sentences from a Continuous Space. Proceedings of the 20th SIGNLL conference on
28 computational natural language learning : 10–21
- 29 Broder AZ (1997) On the resemblance and containment of documents. Proceedings.
30 Compression and complexity of SEQUENCES 1997 (cat. no.97TB100171) : 21–29
- 31 Chan DA, Sithungu SP (2025) Evaluating the Suitability of Inception Score and Fréchet
32 Inception Distance as Metrics for Quality and Diversity in Image Generation.
33 Proceedings of the 7th international conference on computational intelligence and
34 intelligent systems
- 35 Chazal F, Michel B (2021) An Introduction to Topological Data Analysis: Fundamental and
36 Practical Aspects for Data Scientists. Frontiers in Artificial Intelligence 4
- 37 Chen M-J, Bovik AC (2011) Fast structural similarity index algorithm. Journal of Real-
38 Time Image Processing 6(4) : 281–287
- 39 Chen T, Kornblith S, Norouzi M, Hinton G (2020) A Simple Framework for Contrastive
40 Learning of Visual Representations. Proceedings of the 37th international conference
41 on machine learning 119 : 1597–1607
- 42 Coifman RR, Lafon S (2006) Diffusion Maps. Applied and Computational Harmonic
43 Analysis 21(1) : 5–30
- 44 Cover TM, Thomas JA (2006) Elements of Information Theory. Wiley-Interscience
- 45 Crane K, Weischedel C, Wardetzky M (2013) Geodesics in heat: A new approach to
46 computing distance based on heat flow. ACM Transactions on Graphics 32(5)
- 47 Cuturi M (2013) Sinkhorn distances: Lightspeed computation of optimal transport.
48 Advances in neural information processing systems : 2292–2300
- 49 De Boi I, Ek CH, Penne R (2023) Surface Approximation by Means of Gaussian Process
50 Latent Variable Models and Line Element Geometry. Mathematics 11(2)

- 1 De Boi I, Pathak S, Oliveira M, Penne R (2024) How to Turn Your Camera into a Perfect
2 Pinhole Model. *Progress in pattern recognition, image analysis, computer vision, and
3 applications* : 90–107
- 4 De Boi I, Sels S, De Moor O, Vanlanduit S, Penne R (2022) Input and Output Manifold
5 Constrained Gaussian Process Regression for Galvanometric Setup Calibration. *IEEE
6 Transactions on Instrumentation and Measurement* 71 : 1–8
- 7 De Boi I, Sels S, Penne R (2021) Semidata-Driven Calibration of Galvanometric Setups
8 Using Gaussian Processes. *IEEE Transactions on Instrumentation and Measurement*
9 71 : 1–8
- 10 Deng J, Guo J, Yang J, Xue N, Kotsia I, Zafeiriou S (2022) ArcFace: Additive Angular
11 Margin Loss for Deep Face Recognition. *IEEE Transactions on Pattern Analysis and
12 Machine Intelligence* 44(10)
- 13 Deza MM, Deza E (2016) *Encyclopedia of Distances*. Springer Berlin Heidelberg
- 14 Enesi I, Kuqi A (2023) Evaluation of the 3D Reconstruction Performance of Objects in
15 Meshroom: A Case Study. *Athens Journal of Technology and Engineering* 10(1) : 49–
16 70
- 17 Eslami T, Awan MG, Saeed F (2017) GPU-PCC: A GPU based technique to compute
18 pairwise pearson’s correlation coefficients for big fMRI data. *Proceedings of the 8th
19 ACM international conference on bioinformatics, computational biology, and health
20 informatics* : 723–728
- 21 Ghosh A, Ghosh AK, SahaRay R, Sarkar S (2025) Classification using global and local
22 Mahalanobis distances. *Journal of Multivariate Analysis* 207
- 23 Goldberger J, Hinton GE, Roweis S, Salakhutdinov RR (2004) Neighbourhood
24 Components Analysis. *Advances in neural information processing systems* 17
- 25 Gorjian M, Caffey SM, Luhan GA (2025) Exploring Architectural Design 3D
26 Reconstruction Approaches through Deep Learning Methods: A Comprehensive
27 Survey. *Athens Journal of Sciences* 12(3) : 205–234
- 28 Hagele D, Krake T, Weiskopf D (2023) Uncertainty-Aware Multidimensional Scaling.
29 *IEEE Transactions on Visualization and Computer Graphics* 29(1)
- 30 Hart GL, Zach C, Niethammer M (2009) Only Bayes Should Learn a Manifold. 2009 IEEE
31 conference on computer vision and pattern recognition : 1206–1213
- 32 Hartley RI, Zisserman A (2004) *Multiple View Geometry in Computer Vision*. Cambridge
33 University Press, ISBN: 0521540518
- 34 He K, Fan H, Wu Y, Xie S, Girshick R (2020) Momentum Contrast for Unsupervised Visual
35 Representation Learning. *Proceedings of the IEEE/CVF conference on computer
36 vision and pattern recognition* : 9729–9738
- 37 Hellinger E (1909) Neue Begründung der Theorie quadratischer Formen von
38 unendlichvielen Veränderlichen. *Journal für die reine und angewandte Mathematik*
39 136 : 210–271
- 40 Heusel M, Ramsauer H, Unterthiner T, Nessler B, Hochreiter S (2017) GANs Trained by a
41 Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. *Advances in
42 neural information processing systems* 30
- 43 Ho J, Jain A, Abbeel P (2020) Denoising Diffusion Probabilistic Models. *Advances in
44 neural information processing systems* 33 : 6840–6851
- 45 Hoffer E, Ailon N (2015) Deep metric learning using triplet network. *Lecture notes in
46 computer science (including subseries lecture notes in artificial intelligence and lecture
47 notes in bioinformatics)* 9370
- 48 Huang T, Liu Q, Zhao X, Chen J, Liu Y (2024) Learnable Chamfer Distance for Point Cloud
49 Reconstruction. *Pattern Recognition Letters* 178
- 50 Huh M-H (2015) Kernel-Trick Regression and Classification. *Communications for
51 Statistical Applications and Methods* 22(2)

- 1 Hyvärinen A, Morioka H (2016) Unsupervised Feature Extraction by Time-Contrastive
 2 Learning and Nonlinear ICA. *Advances in neural information processing systems* 29
 3 Ichikawa Y, Hukushima K (2024) Learning Dynamics in Linear VAE: Posterior Collapse
 4 Threshold, Superfluous Latent Space Pitfalls, and Speedup with KL Annealing.
 5 *Proceedings of machine learning research* 238
 6 Ioffe S, Szegedy C (2015) Batch Normalization: Accelerating Deep Network Training by
 7 Reducing Internal Covariate Shift. *Proceedings of the 32nd international conference*
 8 *on machine learning* 37 : 448–456
 9 Jost L (2006) Entropy and Diversity. *Oikos* 113(2)
 10 Kaya M, Bilge HŞ (2019) Deep Metric Learning: A Survey. *Symmetry* 11(9)
 11 Khemakhem I, Kingma DP, Monti RP, Hyvärinen A (2020) Variational Autoencoders and
 12 Nonlinear ICA: A Unifying Framework. *Proceedings of the 23rd international*
 13 *conference on artificial intelligence and statistics* 108 : 2207–2217
 14 Kullback S, Leibler RA (1951) On Information and Sufficiency. *The Annals of*
 15 *Mathematical Statistics* 22(1) : 79–86
 16 Kuryatnikova O, Vera JC (2019) Generalizations of Schoenberg’s Theorem on Positive
 17 Definite Kernels. *arXiv*
 18 Lee W, Li W, Lin B, Monod A (2022) Tropical Optimal Transport and Wasserstein
 19 Distances. *Information Geometry* 5(1)
 20 Levy A, Shalom BR, Chalamish M (2025) A guide to similarity measures and their data
 21 science applications. *Journal of Big Data* 12(188)
 22 Lewis JP (1995) Fast Normalized Cross-Correlation. *Vision interface* : 120–123
 23 Lin J (1991) Divergence Measures Based on the Shannon Entropy. *IEEE Transactions on*
 24 *Information Theory* 37(1) : 145–151
 25 Liu P, Na J (2025) Word Motifs and a Generalized Hamming Distance. *Contemporary*
 26 *Mathematics* 6(1)
 27 Lucas J, Tucker G, Grosse R, Norouzi M (2019) Understanding Posterior Collapse in
 28 Generative Latent Variable Models. *Deep generative models for highly structured data,*
 29 *ICLR workshop*
 30 Maaten L van der, Hinton G (2008) Visualizing Data Using t-SNE. *Journal of Machine*
 31 *Learning Research* 9 : 2579–2605
 32 Maes F, Collignon A, Vandermeulen D, Marchal G, Suetens P (1997) Multimodality Image
 33 Registration by Maximization of Mutual Information. *IEEE Transactions on Medical*
 34 *Imaging* 16(2) : 187–198
 35 McKeown S (2025) Beyond Hamming Distance: Exploring Spatial Encoding in Perceptual
 36 Hashes. *Forensic Science International: Digital Investigation* 52
 37 Mikolov T, Chen K, Corrado G, Dean J (2013) Efficient Estimation of Word
 38 Representations in Vector Space. *arXiv preprint arXiv:1301.3781*
 39 Penne R, De Boi I, Vanlanduit S (2023) On the Angular Control of Rotating Lasers by
 40 Means of Line Calculus on Hyperboloids. *Sensors* 23(13)
 41 Pennington J, Socher R, Manning CD (2014) GloVe: Global vectors for word
 42 representation. *EMNLP 2014 - 2014 conference on empirical methods in natural*
 43 *language processing, proceedings of the conference*
 44 Pérez-Cruz F (2008) Kullback-leibler divergence estimation of continuous distributions.
 45 *2008 IEEE international symposium on information theory* : 1666–1670
 46 Peyré G, Cuturi M (2019) Computational Optimal Transport. *Foundations and Trends in*
 47 *Machine Learning* 11(5–6) : 355–607
 48 Pottmann H, Peternell M, Ravani B (1999) An introduction to line geometry with
 49 applications. *Comput. Aided Des.* 31 : 3–16
 50 Pottmann H, Wallner J (2009) *Computational line geometry.* Springer Science & Business
 51 *Media*

- 1 Radford A, Kim JW, Hallacy C, Ramesh A, Goh G, Agarwal S, Sastry G, Askell A, Mishkin
2 P, Clark J, Krueger G, Sutskever I (2021) Learning Transferable Visual Models From
3 Natural Language Supervision. Proceedings of the 38th international conference on
4 machine learning 139 : 8748–8763
- 5 Rainio O, Klén R (2025) Modified Dice Coefficients for Evaluation of Tumor Segmentation
6 from PET Images: A Proof-of-Concept Study. Journal of Imaging Informatics in
7 Medicine
- 8 Reimers N, Gurevych I (2019) Sentence-BERT: Sentence Embeddings Using Siamese
9 BERT-Networks. Proceedings of the 2019 conference on empirical methods in natural
10 language processing : 3982–3992
- 11 Ressel P (1976) A Short Proof of Schoenberg’s Theorem. Proceedings of the American
12 Mathematical Society 57(1)
- 13 Rodrigues M, Kormann M, Al Dulaimi M (2016) Data Protection and Privacy Issues
14 Concerning Facial Image Processing in Public Spaces. Athens Journal of Technology
15 and Engineering 3(1) : 39–52
- 16 Romanazzi M (2014) Discriminant Analysis with High Dimensional von Mises–Fisher
17 Distribution. Athens Journal of Sciences 1(4) : 225–240
- 18 Sakoe H, Chiba S (1978) Dynamic Programming Algorithm Optimization for Spoken Word
19 Recognition. IEEE Transactions on Acoustics, Speech, and Signal Processing 26(1) :
20 43–49
- 21 Schoenberg IJ (1938) Metric Spaces and Positive Definite Functions. Transactions of the
22 American Mathematical Society 44(3) : 522–536
- 23 Schölkopf B, Platt JC, Shawe-Taylor J, Smola AJ, Williamson RC (2001) Estimating the
24 Support of a High-Dimensional Distribution. Neural Computation 13(7) : 1443–1471
- 25 Sethian JA (1996) A fast marching level set method for monotonically advancing fronts.
26 Proceedings of the National Academy of Sciences of the United States of America
27 93(4)
- 28 Shirdhonkar S, Jacobs DW (2008) Approximate earth mover’s distance in linear time. 2008
29 IEEE conference on computer vision and pattern recognition : 1–8
- 30 Suárez JL, García S, Herrera F (2018) A Tutorial on Distance Metric Learning:
31 Mathematical Foundations, Algorithms and Software. arXiv:1812.05944 [cs.LG]
- 32 Suárez JL, García S, Herrera F (2021) A tutorial on distance metric learning: Mathematical
33 foundations, algorithms, experimental analysis, prospects and challenges.
34 Neurocomputing 425
- 35 Syrota S, Zainchkovskyy Y, Xi J, Bloem-Reddy B, Hauberg S (2025) Identifying Metric
36 Structures of Deep Latent Variable Models. Proceedings of machine learning research
37 267
- 38 Tenenbaum JB, De Silva V, Langford JC (2000) A global geometric framework for
39 nonlinear dimensionality reduction. Science 290(5500)
- 40 Unterthiner T, Steenkiste S van, Kurach K, Marinier R, Michalski M, Gelly S (2018)
41 Towards Accurate Generative Models of Video: A New Metric and Challenges.
- 42 Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser L, Polosukhin
43 I (2017) Attention Is All You Need. Advances in neural information processing
44 systems 30
- 45 Verma V, Aggarwal RK (2020) A Comparative Analysis of Similarity Measures Akin to
46 the Jaccard Index in Collaborative Recommendations: Empirical and Theoretical
47 Perspective. Social Network Analysis and Mining 10(1)
- 48 Villani C (2009) Optimal Transport: Old and New. Springer 338
- 49 Viola P, Wells WM (1997) Alignment by Maximization of Mutual Information.
50 International Journal of Computer Vision 24(2) : 137–154

- 1 Wang Z, Bovik AC, Sheikh HR, Simoncelli EP (2004) Image Quality Assessment: From
2 Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing* 13(4)
3 : 600–612
- 4 Weinberger KQ, Saul LK (2009) Distance metric learning for large margin nearest neighbor
5 classification. *Journal of Machine Learning Research* 10
- 6 Ye Q (2026) An Automatic Scoring Method for Vocal Singing Based on Approximate
7 Matching Algorithm. *International Journal of High Speed Electronics and Systems*
8 35(3)
- 9 Zhang C, Ji D (2026) HAL-Net: A Historical Analogy Learning Network for Adaptive and
10 Interpretable Pandemic Forecasting. *Expert Systems with Applications* 299
- 11 Zhang Y, Skolnick J (2004) Scoring Function for Automated Assessment of Protein
12 Structure Template Quality. *Proteins: Structure, Function, and Bioinformatics* 57(4) :
13 702–710
- 14 Zhao P, Lai L (2020) Minimax optimal estimation of KL divergence for continuous
15 distributions. *IEEE Transactions on Information Theory* 66(12) : 7787–7811
- 16 Zhou J, Hodge K, Dong W, Tamakloe E (2025) Distribution-Aware Outlier Detection in
17 High Dimensions: A Scalable Parametric Approach. *Mathematics* 14(1)
- 18 Zupan D, Požar R, Brodnik A (2025) Floyd-Warshall Algorithm for Sparse Graphs.
19 *Algorithms* 18(12)