

What is the distance between pt. a and pt. b?

The N-dimensional real Cartesian space,

denoted \Re^N is the collection of all N-dimensional vectors with real elements. A metric, or distance measure, is a real-valued function with three properties:

 $\forall \bar{x}, \bar{y}, \bar{z} \in \Re^{N}:$ 1. $d(\bar{x}, \bar{y}) \ge 0$. 2. $d(\bar{x}, \bar{y}) = 0$ if and only if $\bar{x} = \bar{y}$

 $2. a(x, y) = 0 \qquad \text{If and only II} \qquad x = 0$

3. $d(\bar{x}, \bar{y}) \leq d(\bar{x}, \bar{z}) + d(\bar{z}, \bar{y})$

The Minkowski metric of order s, or the l_s metric, between x and y is:

$$d_{s}(\bar{x}, \bar{y}) \equiv \sqrt[s]{\sum_{k=1}^{N} |x_{k} - y_{k}|^{s}} = ||\bar{x} - \bar{y}||_{s}$$

(the norm of the difference vector).

Important cases are:

1. l1 or city block metric (sum of absolute values),

$$d_1(\bar{x}, \bar{y}) = \sum_{k=1}^{N} |x_k - y_k|$$

2. 12, or Euclidean metric (mean-squared error),

$$d_2(\bar{x}, \bar{y}) = \sqrt{\sum_{k=1}^{N} |x_k - y_k|^2}$$

3. 1 or Chebyshev metric,

$$d_{\infty}(\bar{x}, \bar{y}) = \max_{k} \left| x_{k} - y_{k} \right|$$



We can similarly define a weighted Euclidean distance metric:

$$d_{2w}(x, y) = \sqrt{|x - y|^T} \underline{W} |x - y|$$

where:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_k \end{bmatrix}, \ \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_k \end{bmatrix}, \text{ and } \underline{W} = \begin{bmatrix} w_{11} \ w_{12} \ \dots \ w_{1k} \\ w_{21} \ w_{22} \ \dots \ w_{2k} \\ \dots \ \dots \ \dots \\ w_{k1} \ w_{k2} \ \dots \ w_{kk} \end{bmatrix}.$$

Why are Euclidean distances so popular?

One reason is efficient computation. Suppose we are given a set of M reference vectors, \bar{x}_m , a measurement, \bar{y} , and we want to find the nearest neighbor:

$$NN = \min_{m} d_2(\bar{x}_m, \bar{y})$$

This can be simplified as follows:

We note the minimum of a square root is the same as the minimum of a square (both are monotonically increasing functions):

$$d_{2}(\bar{x}_{m}, \bar{y})^{2} = \sum_{j=1}^{k} (x_{m_{j}} - y_{j})^{2} = \sum_{j=1}^{k} x_{m_{j}}^{2} - 2x_{m_{j}}y_{j} + y_{j}^{2}$$
$$= \left\| \bar{x}_{m} \right\|^{2} - 2\bar{x}_{m} \cdot \bar{y} + \left\| \bar{y} \right\|^{2}$$
$$= C_{m} + C_{y} - 2\bar{x}_{m} \cdot \bar{y}$$

Therefore,

$$NN = \min_{m} d_2(\bar{x}_m, \bar{y}) = C_m - 2\bar{x}_m \bullet \bar{y}$$

Thus, a Euclidean distance is virtually equivalent to a dot product (which can be computed very quickly on a vector processor). In fact, if all reference vectors have the same magnitude, C_m can be ignored (normalized codebook).

Prewhitening of Features

Consider the problem of comparing features of different scales:

Suppose we represent these points in space in two coordinate systems using the transformation:

$$\bar{z} = \underline{V}\bar{x}$$

System 1:

$$\beta_1 = 1\hat{i} + 0\hat{j}$$
 and $\beta_2 = 0\hat{i} + 1\hat{j}$

$$a = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \qquad b = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$d_2(\bar{a},\bar{b}) = \sqrt{0^2 + 1^2} = 1$$

System 2:

$$\gamma_1 = -2\hat{i} + 0\hat{j}$$
 and $\gamma_1 = -1\hat{i} + 1\hat{j}$



The magnitude of the distance has changed. Though the rank-ordering of distances under such linear transformations won't change, the cumulative effects of such changes in distances can be damaging in pattern recognition. Why?



We can simplify the distance calculation in the transformed space:

$$d_{2}(\underline{V}\bar{x}, \underline{V}\bar{y}) = \sqrt{\left[\underline{V}\bar{x} - \underline{V}\bar{y}\right]^{T}\left[\underline{V}\bar{x} - \underline{V}\bar{y}\right]}$$
$$= \sqrt{\left[\bar{x} - \bar{y}\right]^{T}\underline{V}^{T}\underline{V}\left[\bar{x} - \bar{y}\right]}$$
$$= d_{2}w(\bar{x}, \bar{y})$$

This is just a weighted Euclidean distance.

Suppose all dimensions of the vector are not equal in importance. For example, suppose one dimension has virtually no variation, while another is very reliable. Suppose two dimensions are statistically correlated. What is a statistically optimal transformation?

Consider a decomposition of the covariance matrix (which is symmetric):

$$\underline{C} = \underline{\Phi} \underline{\Lambda} \underline{\Phi}^T$$

where Φ denotes a matrix of eigenvectors of \underline{C} and $\underline{\Lambda}$ denotes a diagonal matrix whose elements are the eigenvalues of \underline{C} . Consider:

$$z = \underline{\Lambda}^{-1/2} \underline{\Phi} \bar{x}$$

The covariance of z, \underline{C}_z is easily shown to be an identity matrix (prove this!) We can also show that:

$$d_2(\bar{z}_1, \bar{z}_2) = \sqrt{[\bar{x}_1 - \bar{x}_2]^T \underline{C}_{\bar{x}}^{-1} [\bar{x}_1 - \bar{x}_2]}$$

Again, just a weighted Euclidean distance.

- If the covariance matrix of the transformed vector is a diagonal matrix, the transformation is said to be an orthogonal transform.
- If the covariance matrix is an identity matrix, the transform is said to be an orthonormal transform.
- A common approximation to this procedure is to assume the dimensions of x
 are uncorrelated but of unequal variances, and to approximate <u>C</u> by a diagonal matrix, <u>Λ</u>. Why? This is known as variance-weighting.

"Noise-Reduction"

The prewhitening transform, $\bar{z} = \Delta^{-1/2} \Phi \bar{x}$, is normally created as a $k \times k$ matrix in which the eigenvalues are ordered from largest to smallest:

$$\begin{bmatrix} z_1 \\ z_2 \\ \dots \\ z_k \end{bmatrix} = \begin{bmatrix} \lambda_1^{-1/2} & & \\ & \lambda_2^{-1/2} & & \\ & & \ddots & \\ & & & \lambda_k^{-1/2} \end{bmatrix} \begin{bmatrix} v_{11} & v_{12} & \dots & v_{13} \\ v_{21} & v_{22} & \dots & v_{2k} \\ \dots & & & \ddots & \dots \\ v_{k1} & v_{k2} & \dots & v_{kk} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_k \end{bmatrix}$$

where

$$\lambda_1 > \lambda_2 > \dots > \lambda_k$$
.

In this case, a new feature vector can be formed by truncating the transformation matrix to l < k rows. This is essentially discarding the least important features.

A measure of the amount of discriminatory power contained in a feature, or a set of features, can be defined as follows:

$$\% \text{ var } = \frac{\sum_{j=1}^{l} \lambda_j}{\sum_{j=1}^{k} \lambda_j}$$

This is the percent of the variance accounted for by the first *l* features.

Similarly, the coefficients of the eigenvectors tell us which dimensions of the input feature vector contribute most heavily to a dimension of the output feature vector. This is useful in determining the "meaning" of a particular feature (for example, the first decorrelated feature often is correlated with the overall spectral slope in a speech recognition system — this is sometimes an indication of the type of microphone).

Maximum Likelihood Classification

Consider the problem of assigning a measurement to one of two sets:



What is the best criterion for making a decision?

Ideally, we would select the class for which the conditional probability is highest:

$$c^* = \underset{c}{\operatorname{argmax}} P((c = \hat{c}) | (\bar{x} = \bar{x}))$$

However, we can't estimate this probability directly from the training data. Hence, we consider:

$$c^* = \underset{c}{\operatorname{argmax}} P((\bar{x} = \hat{\bar{x}}) | (c = \hat{c}))$$

By definition

$$P((c = \hat{c}) | (\bar{x} = \hat{\bar{x}})) = \frac{P((c = \hat{c}), (\bar{x} = \bar{x}))}{P(\bar{x} = \hat{\bar{x}})}$$

and

$$P((\bar{x} = \hat{x}) | (c = \hat{c})) = \frac{P((c = \hat{c}), (\hat{x} = \hat{x}))}{P(c = \hat{c})}$$

from which we have

$$P((c = \hat{c})|(\bar{x} = \hat{x})) = \frac{P((\bar{x} = \hat{x})|(c = \hat{c}))P(c = \hat{c})}{P(\bar{x} = \hat{x})}$$

Clearly, the choice of c that maximizes the right side also maximizes the left side. Therefore,

$$c^* = \underset{c}{\operatorname{argmax}} \left[P((\bar{x} = \hat{\bar{x}}) | (c = \hat{c})) \right]$$
$$= \underset{c}{\operatorname{argmx}} \left[P((\bar{x} = \hat{\bar{x}}) | (c = \hat{c})) P(c = \hat{c}) \right]$$

if the class probabilities are equal,

$$c^* = \operatorname{argmx}_{c}[P((\bar{x} = \hat{\bar{x}}) | (c = \hat{c}))]$$

A quantity related to the probability of an event which is used to make a decision about the occurrence of that event is often called a *likelihood measure*.

A decision rule that maximizes a likelihood is called a maximum likelihood decision.

In a case where the number of outcomes is not finite, we can use an analogous continuous distribution. It is common to assume a multivariate Gaussian distribution:

$$\begin{split} f_{\bar{x}|c}(x_1, \dots, x_N \big| c) &= f_{\bar{x}|c}(\hat{x} \big| \hat{c}) \\ &= \frac{1}{\sqrt{2\pi |C_{\bar{x}}|c|}} \exp\left\{ -\frac{1}{2} (\hat{x} - \bar{\mu}_{\bar{x}|c})^T \underline{C}_{\bar{x}|c}^{-1} (\hat{x} - \bar{\mu}_{\hat{x}|c}) \right\} \end{split}$$

We can elect to maximize the log, $\ln[f_{\bar{x}|c}(\bar{x}|c)]$ rather than the likelihood (we refer to this as the log likelihood). This gives the decision rule:

$$c^{*} = \underset{c}{\operatorname{argmin}} \left[(\hat{x} - \bar{\mu}_{\bar{x}|c})^{T} \underline{C}_{\bar{x}|c}^{-1} (\hat{x} - \bar{\mu}_{\hat{x}|c}) + \ln \left\{ \left| \underline{C}_{\bar{x}|c}^{-1} \right| \right\} \right]$$

(Note that the maximization became a minimization.)

We can define a distance measure based on this as:

$$d_{ml}(\bar{x},\bar{\mu}_{\bar{x}|c}) = (\hat{\bar{x}} - \bar{\mu}_{\bar{x}|c})^T \underline{C}_{\bar{x}|c}^{-1} (\hat{\bar{x}} - \bar{\mu}_{\hat{\bar{x}}|c}) + \ln\left\{ \left| \underline{C}_{\bar{x}|c}^{-1} \right| \right\}$$

Note that the distance is conditioned on each class mean and covariance. This is why "generic" distance comparisons are a joke.

If the mean and covariance are the same across all classes, this expression simplifies to:

$$d_{\underline{M}}(\bar{x}, \overline{\mu}_{\bar{x}|c}) = (\hat{\bar{x}} - \overline{\mu}_{\bar{x}|c})^T \underline{C}_{\bar{x}|c}^{-1} (\hat{\bar{x}} - \overline{\mu}_{\hat{\bar{x}}|c})$$

This is frequently called the *Mahalanobis distance*. But this is nothing more than a weighted Euclidean distance.

This result has a relatively simple geometric interpretation for the case of a single random variable with classes of equal variances:



The decision rule involves setting a threshold:

$$a = \left(\frac{\mu_1 + \mu_2}{2}\right) + \frac{\sigma^2}{\mu_1 - \mu_2} \ln\left(\frac{P(c=2)}{P(c=1)}\right)$$

and,

if
$$x < a$$
 $x \in (c = 1)$
else $x \in (c = 2)$

If the variances are not equal, the threshold shifts towards the distribution with the smaller variance.

What is an example of an application where the classes are not equiprobable?