

Distance Metrics in Machine Learning

Created by Daniel Zaldaña
<https://x.com/ZaldanaDaniel>

Euclidean Distance

Formula

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

When to use:

- Continuous data in low-dimensional space
- When scale and magnitude matter
- Default choice for clustering algorithms

Properties:

- Symmetric: $d(x, y) = d(y, x)$
- Non-negative: $d(x, y) \geq 0$
- Sensitive to outliers and scale

```
1 from sklearn.metrics.pairwise import (
2     euclidean_distances,
3     manhattan_distances,
4     cosine_distances
5 )
6 import numpy as np
7 from typing import ndarray
8
9 def euclidean(
10     x: ndarray,
11     y: ndarray
12 ) -> float:
13     """Calculate Euclidean distance using
14         sklearn."""
15     # Reshape if needed for single vectors
16     X = x.reshape(1, -1) if x.ndim == 1 else x
17     Y = y.reshape(1, -1) if y.ndim == 1 else y
18     return euclidean_distances(X, Y)[0, 0]
```

Manhattan Distance

Formula

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n |x_i - y_i|$$

When to use:

- Grid-like patterns (e.g., city blocks)
- When diagonal movement costs more
- Robust to outliers

```
1 def manhattan(
2     x: ndarray,
3     y: ndarray
4 ) -> float:
5     """Calculate Manhattan distance using
6         sklearn."""
7     X = x.reshape(1, -1) if x.ndim == 1 else x
8     Y = y.reshape(1, -1) if y.ndim == 1 else y
9     return manhattan_distances(X, Y)[0, 0]
```

Cosine Similarity

Formula

$$\text{similarity}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

When to use:

- Text analysis and document similarity
- High-dimensional sparse data
- When direction matters more than magnitude

```
1 def cosine_similarity(
2     x: ndarray,
3     y: ndarray
4 ) -> float:
5     """Calculate cosine similarity using
6         sklearn."""
7     X = x.reshape(1, -1) if x.ndim == 1 else x
8     Y = y.reshape(1, -1) if y.ndim == 1 else y
9     return 1 - cosine_distances(X, Y)[0, 0]
```

Mahalanobis Distance

Formula

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})}$$

When to use:

- Correlated features
- Anomaly detection
- Scale-invariant clustering

```
1 from sklearn.covariance import
2     EmpiricalCovariance
3
4 def mahalanobis(
5     x: ndarray,
6     y: ndarray,
7     cov: ndarray = None
8 ) -> float:
9     """Calculate Mahalanobis distance using
10         sklearn."""
```

```

9     X = x.reshape(1, -1) if x.ndim == 1 else x
10    Y = y.reshape(1, -1) if y.ndim == 1 else y
11
12    if cov is None:
13        # Estimate covariance from data
14        cov_estimator = EmpiricalCovariance()
15        cov_estimator.fit(np.vstack([X, Y]))
16        cov = cov_estimator.covariance_
17
18    diff = X - Y
19    inv_covmat = np.linalg.inv(cov)
20    return np.sqrt(
21        diff.dot(inv_covmat).dot(diff.T)
22    )[0, 0]

```

Minkowski Distance

Formula

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$$

When to use:

- Generalizing distance metrics
- When you need to tune the influence of large differences
- Experimenting with different p-norms

```

1 from sklearn.metrics import pairwise_distances
2
3 def minkowski(
4     x: ndarray,
5     y: ndarray,
6     p: float = 2
7 ) -> float:
8     """Calculate Minkowski distance using
9         sklearn."""
10    X = x.reshape(1, -1) if x.ndim == 1 else x
11    Y = y.reshape(1, -1) if y.ndim == 1 else y
12    return pairwise_distances(
13        X, Y, metric='minkowski', p=p
14    )[0, 0]

```

Jaccard Distance

Formula

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{|\mathbf{x} \cap \mathbf{y}|}{|\mathbf{x} \cup \mathbf{y}|}$$

When to use:

- Binary or set-based data
- Comparing discrete features
- Document similarity with word sets

```

1 def jaccard(
2     x: ndarray,
3     y: ndarray
4 ) -> float:
5     """Calculate Jaccard distance using
6         sklearn."""
7     X = x.reshape(1, -1) if x.ndim == 1 else x
8     Y = y.reshape(1, -1) if y.ndim == 1 else y
9     return pairwise_distances(
10        X, Y, metric='jaccard'
11    )[0, 0]

```

Hamming Distance

Formula

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \mathbb{1}_{x_i \neq y_i}$$

When to use:

- Categorical data
- Error detection in communication
- Comparing equal-length strings

```

1 def hamming(
2     x: ndarray,
3     y: ndarray
4 ) -> float:
5     """Calculate Hamming distance using
6         sklearn."""
7     X = x.reshape(1, -1) if x.ndim == 1 else x
8     Y = y.reshape(1, -1) if y.ndim == 1 else y
9     return pairwise_distances(
10        X, Y, metric='hamming'
11    )[0, 0]

```