# Vector Similarity Search Algorithms



[0.34, 2.35, 8.34, ...]
300 dimensions

Chicken
Wolf
Dog
Cat
Banana
Apple

Vector: Mathematical representation of data in an n-dimensional space ( each dimension == a feature of the data)

Embeddings: map raw data into a vector space (preserves semantic relationships - meaning)

# Vector Similarity Search: Foundations and Applications

## Introduction

Vector Similarity Search (VSS) represents one of the most transformative technological frameworks in modern data science and artificial intelligence. At its core, VSS enables machines to understand, organize, and retrieve information based on meaning and contextual relationships rather than exact matches or predefined categories. This paradigm shift has fundamentally altered how we interact with large-scale information systems and has become the backbone of numerous applications that we encounter daily.

In a world where data grows exponentially in both volume and complexity, traditional search and retrieval methods have reached their limitations. The increasing sophistication of user expectations— from finding visually similar images to discovering conceptually related documents without matching

keywords—necessitates approaches that can capture the nuanced relationships inherent in our data. Vector Similarity Search addresses this need by representing data items as mathematical vectors in high-dimensional spaces where similarity becomes a measurable distance.

## The Mathematics of Meaning

The power of VSS lies in its ability to translate abstract concepts such as semantic meaning, visual similarity, or musical resemblance into precise mathematical relationships. When we convert data into vectors (a process called embedding), we create a representation where:

- Each dimension corresponds to a feature or attribute of the data
- The position of a vector in this high-dimensional space encodes semantic information
- The geometric distance between vectors correlates with their conceptual similarity
- Complex relationships between concepts emerge as patterns in the vector space

This mathematical framework allows us to perform operations on meaning itself. We can find items similar to a query by identifying its nearest neighbors in vector space. More remarkably, we can perform algebraic operations on these vectors that translate to semantic operations—adding and subtracting concepts, identifying analogies, and discovering latent relationships.

## The Technological Ecosystem

Vector Similarity Search exists within a broader ecosystem of technologies and methods:

1. **Embedding Models**: Neural networks and other algorithms that transform raw data (text, images, audio) into vector representations that preserve semantic relationships.
2. **Distance Metrics**: Mathematical functions that quantify similarity between vectors, with different metrics (Euclidean, Manhattan, Cosine) capturing different aspects of similarity.
3. **Indexing Structures**: Specialized data structures and algorithms that organize vectors for efficient retrieval, overcoming the "curse of dimensionality" that would otherwise make high-dimensional search prohibitively expensive.
4. **Vector Databases**: Purpose-built systems that store, index, and retrieve vectors at scale, often incorporating approximation techniques to balance accuracy and performance.

## Applications Transforming Industries

The practical applications of Vector Similarity Search extend across virtually every domain where information retrieval matters:

- **Information Retrieval**: Modern search engines use VSS to understand user intent beyond keywords, finding conceptually relevant results even when terminology differs.
- **Recommendation Systems**: Streaming services, e-commerce platforms, and content sites leverage VSS to suggest items based on complex patterns of user preferences.
- **Computer Vision**: Image recognition, visual search, and content-based image retrieval all rely on vector representations of visual features.
- **Natural Language Processing**: Language models, translation systems, and sentiment analysis tools use word and sentence embeddings to capture semantic relationships.
- **Anomaly Detection**: Financial systems detect fraudulent transactions by identifying statistical outliers in vector space.
- **Drug Discovery**: Pharmaceutical research uses molecular embeddings to find compounds with similar properties or potential interactions.

$$\text{Cosine Similarity}(A, B) = \frac{A \cdot B}{\|A\|\|B\|}$$

**Best use:**
- Topic modeling
- Document similarity
- Collaborative filtering

$$\text{Euclidean Distance}(A, B) = \sqrt{\sum_{i=1}^{n}(A_i - B_i)^2}$$

**Best use:**
- Clustering analysis
- Anomaly & fraud detection

$$\text{Dot Product}(A, B) = A \cdot B = \sum_{i=1}^{n} A_i B_i$$

**Best use:**
- Image retrieval & matching
- Neural networks & Deep Learning
- Music Recommendation

# Algorithms

# Dot Product

**Overview:**

The Dot Product is a straightforward, yet powerful, similarity measure used extensively in machine learning, data mining, and statistics for finding the similarity between two vectors. It's particularly important in the realm of cosine similarity and is fundamental for algorithms employed in search engines, recommendation systems, and other data-driven applications.

**Definition and Formula:**

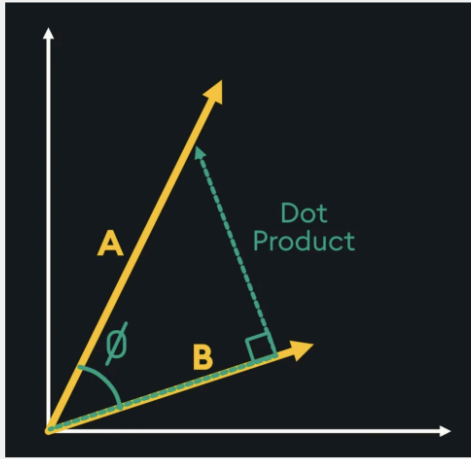Given two vectors, A and B, the Dot Product is computed as follows:

$$A = [a_1, a_2, ..., a_n]$$

Vector A

$$B = [b_1, b_2, ..., b_n]$$

Vector B

$$A \cdot B = a_1 \cdot b_1 + a_2 \cdot b_2 + ... + a_n \cdot b_n = \sum_{i=1}^{n} a_i \cdot b_i$$

**Best use:**
- Image retrieval & matching
- Music recommendation
- ...

$$\text{Dot Product}(A, B) = A \cdot B = \sum_{i=1}^{n} A_i B_i$$

**Applications in Similarity Search:**

- **Cosine Similarity Basis**: The Dot Product is fundamental for calculating cosine similarity, which is a widely used metric in similarity search. When the vectors are normalized, the dot product essentially provides the cosine of the angle between the two vectors, offering an effective measure of similarity.
- **Efficient Retrieval**: It allows for efficient retrieval of similar vectors in large databases, making it vital for search engines and recommendation systems where speed is crucial.
- **Machine Learning Models**: Used in training machine learning models, especially in deep learning, where the dot product operation is employed extensively during the forward and backward propagation stages.

**Advantages:**

- **Speed**: It's computationally efficient, providing fast calculations which are crucial for real-time applications.
- **Simplicity**: Due to its straightforward formula, it's easy to implement and understand.

**Limitations:**

- **Magnitude Sensitivity**: The dot product is sensitive to the magnitude of the vectors, and it may not accurately represent similarity if the magnitudes of the vectors being compared are significantly different.

**Conclusion:**

The Dot Product is a foundational vector similarity search algorithm with widespread applications and relevance in various fields. Its simplicity and efficiency make it a go-to choice for professionals working with large datasets and requiring fast similarity computations.

**Implementation:**

```python
def dot_product(v1, v2):
    """
    Calculate the dot product between two vectors.

    Parameters:
    v1 (list or array): First vector
    v2 (list or array): Second vector

    Returns:
    float: The dot product of the two vectors
```

```python
    """
    dot_prod = 0.0
    for i in range(len(v1)):
        dot_prod += v1[i] * v2[i]
    return dot_prod

# Example with word embeddings (simplified)
happy = [0.8, 0.5, 0.3]    # Embedding for "happy"
joyful = [0.7, 0.6, 0.4]   # Embedding for "joyful"
computer = [0.1, 0.2, 0.9] # Embedding for "computer"

# Calculate dot products
similarity_happy_joyful = dot_product(happy, joyful)
similarity_happy_computer = dot_product(happy, computer)

print(f"Dot product between 'happy' and 'joyful': {similarity_happy_joyful}")
print(f"Dot product between 'happy' and 'computer': {similarity_happy_computer}")

# OR

import numpy as np

def dot_product_numpy(v1, v2):
    return np.dot(v1, v2)
```

## Cosine Similarity

Cosine Similarity is a widely used metric for measuring similarity between two vectors, often employed in the fields of information retrieval, text mining, and machine learning. It gauges the cosine of the angle between two vectors, providing insights into their orientation in a multi-dimensional space.

**Best Use:**

- When vector magnitude need to be considered
- Ideal for clustering tasks (k-means clustering)
- Topic modeling
- Document similarity
- Collaborative filtering

**Definition and Formula:**

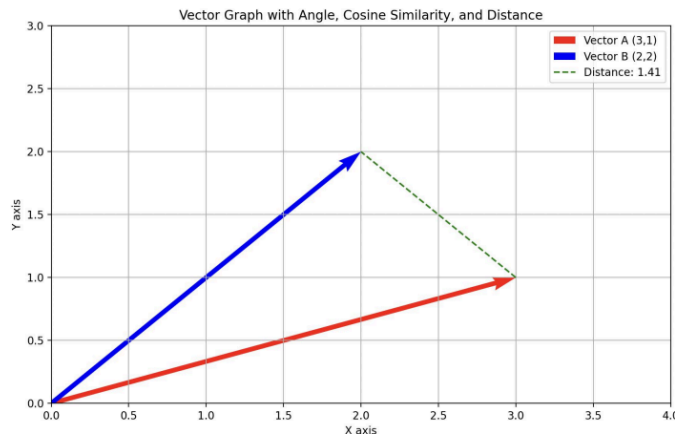Given two vectors, A and B, the Cosine Similarity between these vectors is computed as:

$$A = [a_1, a_2, ..., a_n]$$

Vector A

$$B = [b_1, b_2, ..., b_n]$$

Vector B

$$S(A, B) = \frac{A \cdot B}{||A|| \; ||B||} = \frac{\sum\limits_{i=1}^{n} a_i \cdot b_i}{\sqrt{\sum\limits_{i=1}^{n} a_i^2} \cdot \sqrt{\sum\limits_{i=1}^{n} b_i^2}}$$

Vector Graph with Angle, Cosine Similarity, and Distance

- Vector A (3,1)
- Vector B (2,2)
- Distance: 1.41

$$\text{Cosine Similarity}(A, B) = \frac{A \cdot B}{||A||||B||}$$

- ○ Dot prod For Arrow A and Arrow B, that's (3 * 2) + (1 * 2) = 6 + 2 = 8.
- ○ Magnitude For Arrow A, For Arrow A, it's the square root of $(3^2 + 1^2) = \sqrt{(9 + 1)} = \sqrt{10}$
- ○ Magnitude For Arrow B, it's $\sqrt{(2^2 + 2^2)} = \sqrt{(4 + 4)} = \sqrt{8}$.

- Cosine Similarity $= \frac{8}{\sqrt{10} * \sqrt{8}}$

Result: ~0.894

**Applications in Similarity Search:**

- **Document Similarity:** Often used in natural language processing to measure similarity between texts, aiding in document retrieval and clustering.
- **Recommendation Systems:** Cosine Similarity is employed in collaborative filtering to generate recommendations by comparing user or item profiles.
- **Image Comparison:** It's applied in computer vision for comparing feature vectors of images.

**Advantages:**

- **Angle Measurement:** It measures the cosine of the angle between vectors, making it effective for comparing documents of different lengths.
- **Normalization:** The metric inherently normalizes vector lengths, making it sensitive to the direction of the data, not the magnitude.

**Limitations:**

- **Zero Vector Issue:** It does not handle zero vectors well, as it becomes undefined.
- **Not a Metric:** Cosine Similarity doesn't satisfy the triangle inequality and therefore is not a true metric.

**Conclusion:**

Cosine Similarity is a robust and versatile vector similarity search algorithm widely applied across various domains. Though it offers several advantages, like insensitivity to magnitude, understanding its limitations is crucial for effective application. With careful implementation, it remains a valuable tool for professionals working on similarity search, document retrieval, and recommendation systems in the data science and AI fields.

**Reference Implementation:**

```python
import numpy as np

def cosine_similarity(v1, v2):
    """
    Calculate the cosine similarity between two vectors.
    Parameters:
    v1 (list or array): First vector
    v2 (list or array): Second vector
    Returns:
    float: Cosine similarity value between -1 and 1
    """

    dot_prod = sum(v1[i] * v2[i] for i in range(len(v1)))

    v1_sqr_sum = sum(v1[i] ** 2 for i in range(len(v1)))
    v2_sqr_sum = sum(v2[i] ** 2 for i in range(len(v1)))

    return dot_prod / (np.sqrt(v1_sqr_sum) * np.sqrt(v2_sqr_sum))

# Example usage
embedding1 = [0.2, 0.5, 0.8, 0.1]
embedding2 = [0.3, 0.4, 0.7, 0.2]
embedding3 = [0.9, 0.1, 0.2, 0.5]

# Calculate similarities
similarity_1_2 = cosine_similarity(embedding1, embedding2)
similarity_1_3 = cosine_similarity(embedding1, embedding3)

print(f"Similarity between embedding1 and embedding2: {similarity_1_2:.4f}")
print(f"Similarity between embedding1 and embedding3: {similarity_1_3:.4f}")
```

# Distance Metrics in Vector Spaces

## Introduction to Distance Metrics

Distance metrics are mathematical functions that measure how far apart two points or vectors are in a geometric space. These metrics form the foundation of various computational techniques in machine learning, information retrieval, and data analysis.

## Manhattan Distance (L1 Norm)

Manhattan distance, also known as L1 distance, taxicab distance, or city block distance, measures the sum of the absolute differences between the coordinates of two points.

### Definition and Formula

For two points A and B in n-dimensional space:

$$d_{Manhattan}(A, B) = \sum_{i=1}^{n} |a_i - b_i|$$

In expanded form:

$$d_{Manhattan}(A, B) = |a_1 - b_1| + |a_2 - b_2| + \ldots + |a_n - b_n|$$

## Python Implementation

```python
import numpy as np

def manhattan_distance(vector_a, vector_b):
    a = np.array(vector_a)
    b = np.array(vector_b)
    return np.sum(np.abs(a - b))

# Alternative implementation using numpy's built-in function
def manhattan_distance_numpy(vector_a, vector_b):
    return np.linalg.norm(np.array(vector_a) - np.array(vector_b), ord=1)
```

## Geometric Interpretation

The Manhattan distance gets its name from the grid layout of Manhattan streets. In a 2D grid, you can only move horizontally or vertically (not diagonally), so the shortest path between two points follows the grid lines. The distance represents the total number of blocks you would need to travel.

## Applications in Embeddings and Vector Search

Manhattan distance is particularly useful in high-dimensional spaces and vector embeddings for several reasons:

1. **Computational Efficiency**: Manhattan distance requires only addition and subtraction operations, making it computationally less expensive than Euclidean distance, which requires squaring and square root operations.
2. **Robustness to Outliers**: The L1 norm is less sensitive to outliers compared to the L2 norm, as it doesn't square the differences (which would amplify large deviations).
3. **Feature Selection**: In sparse high-dimensional spaces (like text embeddings), Manhattan distance can better capture semantic differences because it considers all dimensions more equally.
4. **Interpretability**: The Manhattan distance provides a more interpretable measure in certain domains. Each unit represents a direct step along a single dimension.
5. **Curse of Dimensionality**: In high-dimensional spaces, the Manhattan distance can sometimes be more effective than Euclidean distance due to the "curse of dimensionality" phenomenon.

## Advantages in Vector Search

When working with vector embeddings in search applications:

- **Semantic Relationships**: Manhattan distance can effectively capture semantic relationships in word embeddings and document vectors.
- **Efficiency in Approximate Nearest Neighbor (ANN) Algorithms**: Many fast ANN algorithms work well with Manhattan distance metrics.
- **Distinctiveness**: In high-dimensional spaces, Manhattan distance often provides better discrimination between vectors than Euclidean distance.

## Limitations

- **Non-Euclidean Geometry**: Manhattan distance doesn't correspond to our intuitive understanding of distance in physical space.
- **Rotation Invariance**: Unlike Euclidean distance, Manhattan distance is not invariant to rotations in the vector space.

# Euclidean Distance (L2 Distance)

## Core Concept

Euclidean distance represents the straight-line distance between two points in a multidimensional space. Named after the ancient Greek mathematician Euclid, this fundamental metric extends our intuitive understanding of physical distance to any n-dimensional space. When working with embeddings, this distance quantifies how dissimilar two entities are based on their vector representations.

## Understanding Euclidean Distance

The Euclidean distance measures the straight-line distance between two points in a space. When we have two points:

- Point A = $(a_1, a_2, ..., a_n)$
- Point B = $(b_1, b_2, ..., b_n)$

The distance between them is calculated as:

$$d(A, B) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \ldots + (a_n - b_n)^2}$$

This can indeed be expressed more compactly using summation notation:

$$d(A, B) = \sqrt{\sum_{i=1}^{n} (a_i - b_i)^2}$$

## Geometric Interpretation

This formula is a direct generalization of the Pythagorean theorem. In a 2-dimensional space, the distance between points $(a_1, a_2)$ and $(b_1, b_2)$ is:

$$d(A, B) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

This corresponds to the length of the hypotenuse of a right triangle formed by the horizontal and vertical components of the displacement between the points.

As we extend to higher dimensions, each additional coordinate contributes an additional squared term under the radical.

## Properties of Euclidean Distance

The Euclidean distance has several important properties:

1. Non-negativity: $d(A,B) \geq 0$ for all points A and B
2. Identity of indiscernibles: $d(A,B) = 0$ if and only if A = B
3. Symmetry: $d(A,B) = d(B,A)$
4. Triangle inequality: $d(A,C) \leq d(A,B) + d(B,C)$

These properties make the Euclidean distance a proper metric, which is why it's often called the Euclidean metric.

Where:

- Each term `(aᵢ-bᵢ)²` represents the squared difference between corresponding components

- The sum `∑` aggregates these squared differences across all dimensions
- The square root `√` converts this sum to a true distance measurement

As we move to higher dimensions, the formula continues to capture the shortest path between points, maintaining this geometric interpretation even when we can no longer visualize the space directly.

# Applications in Machine Learning and Data Science

## Vector Similarity Search

Euclidean distance serves as the backbone for many nearest neighbor search algorithms:

- **K-nearest neighbors (KNN)**: Classifies points based on the majority class among k closest neighbors
- **Radius-based search**: Identifies all points within a specified distance threshold
- **Approximate nearest neighbor search**: Techniques like locality-sensitive hashing (LSH) accelerate searches by approximating Euclidean distances

## Clustering Algorithms

Many clustering algorithms rely on Euclidean distance as their primary similarity metric:

- **K-means clustering**: Iteratively assigns points to the nearest cluster center and recalculates centers
- **Hierarchical clustering**: Builds nested clusters by merging or splitting based on distance criteria
- **DBSCAN**: Density-based clustering using distance thresholds to define neighborhoods

## Recommendation Systems

When item or user profiles are represented as feature vectors, Euclidean distance helps identify similarities:

- **Content-based filtering**: Recommends items with feature vectors close to those a user has liked
- **Collaborative filtering**: Finds users with similar preference vectors to generate recommendations

## Computer Vision

Image similarity often relies on Euclidean distance between feature vectors:

- **Facial recognition**: Compares facial embedding vectors to identify or verify individuals
- **Image retrieval**: Finds visually similar images by comparing extracted feature vectors
- **Object detection**: Matches detected objects against known templates based on feature distances

## Anomaly Detection

Points that have large Euclidean distances from clusters or expected patterns can indicate anomalies:

- **Outlier detection**: Identifies data points that deviate significantly from the norm
- **Fraud detection**: Flags transactions with unusual feature vectors compared to typical patterns
- **Predictive maintenance**: Detects when sensor readings diverge from normal operating conditions

# Advantages and Strengths

## Mathematical Properties

Euclidean distance has several desirable mathematical properties:

- **Non-negativity**: Always greater than or equal to zero
- **Identity of indiscernibles**: Distance is zero if and only if the points are identical
- **Symmetry**: Distance from A to B equals distance from B to A
- **Triangle inequality**: Direct distance between points is never greater than going through a third point

## Practical Benefits

- **Intuitive interpretation**: Matches human understanding of physical distance
- **Computational efficiency**: Simple arithmetic operations make it fast to compute
- **Differentiable**: Important for gradient-based optimization in machine learning
- **Preserves relative distances**: Works well when relative magnitudes matter

# Limitations and Challenges

## The Curse of Dimensionality

As dimensionality increases, Euclidean distance faces significant challenges:

- **Distance concentration**: In high dimensions, distances between points tend to become more uniform
- **Sparsity**: Data points become increasingly sparse, requiring more samples for meaningful analysis
- **Computational complexity**: Processing high-dimensional vectors becomes resource-intensive

A practical demonstration of this problem is that as dimensions increase, the ratio of the distance between the closest and farthest points approaches 1, making nearest neighbor searches less meaningful.

## Scale Sensitivity

Euclidean distance is highly sensitive to the scale of features:

- Features with larger scales disproportionately influence the distance calculation
- This necessitates normalization techniques like:
  - Min-max scaling
  - Standardization (z-score normalization)
  - Feature weighting

```
# Example: Impact of feature scaling
# Original features with different scales
feature1 = [100, 200, 300]  # Large scale
feature2 = [0.1, 0.2, 0.3]  # Small scale
```
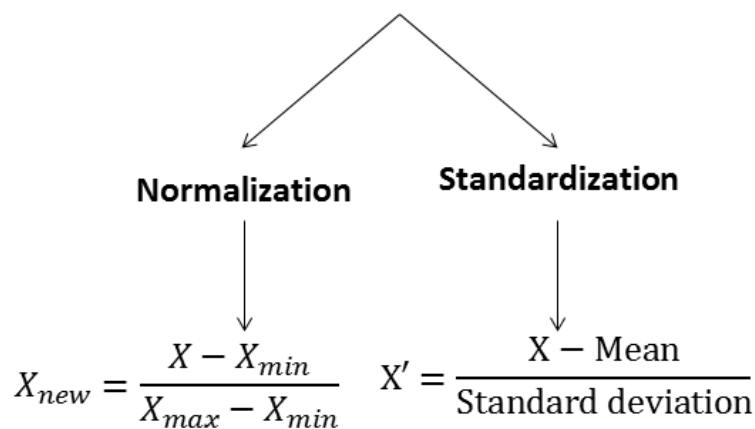
## Outlier Sensitivity

Squaring differences means outliers can disproportionately affect distance calculations:

- A single dimension with a large difference can dominate the overall distance
- This makes Euclidean distance less robust to noisy data compared to alternatives like Manhattan distance

**Side Note:**

## Feature scaling

Normalization          Standardization

$$X_{new} = \frac{X - X_{min}}{X_{max} - X_{min}} \qquad X' = \frac{X - Mean}{Standard\ deviation}$$

**Implementation:**

```python
import math

def euclidean_distance(v1, v2):
    sum_sqr_diff = 0.0
    for i in range(len(v1)):
        sqr_diff = (v1[i] - v2[i]) ** 2
        sum_sqr_diff += sqr_diff
    return math.sqrt(sum_sqr_diff)

# Example with word embeddings (simplified to 3D for clarity)
word1 = [0.2, 0.5, 0.9]  # First word embedding
word2 = [0.4, 0.1, 0.7]  # Second word embedding

# Calculate distance
distance = euclidean_distance(word1, word2)
print(f"Euclidean distance between word1 and word2: {distance:.4f}")
```

## Conclusion

Vector Similarity Search represents a fundamental shift in how machines understand and organize information. By translating meaning into mathematics, it enables a new generation of information systems that can retrieve, categorize, and analyze data based on semantic relationships rather than rigid structures or exact matches. As embedding models become more sophisticated and search algorithms more efficient, VSS will continue to transform how we interact with the growing universe of digital information.

The lectures that follow will explore the mathematical foundations, algorithmic approaches, and practical implementations of Vector Similarity Search in detail, providing you with both theoretical understanding and practical skills to apply these powerful techniques to your own data challenges.