

02 Embeddings

What are Embeddings?

Embeddings are dense vector representations of data in a continuous vector space. They capture semantic meaning and contextual relationships by mapping discrete entities (like words, sentences, images, or user behaviors) to points in a multi-dimensional space where similar items are positioned closer together.

Unlike traditional sparse representations (like one-hot encoding), embeddings are:

- **Dense:** They use a fixed number of dimensions to represent information efficiently
- **Learned:** They're derived from data rather than manually engineered
- **Continuous:** They exist in a smooth vector space allowing for mathematical operations
- **Semantic:** They capture meaningful relationships between entities

Real-World Example

Consider how we might represent colors. A traditional representation might use discrete RGB values (255, 0, 0) for red. An embedding approach would learn that "scarlet" and "crimson" should be positioned close to "red" in vector space, while "navy" and "azure" would cluster near "blue." This proximity in the embedding space reflects their semantic similarity.

How Embeddings Work

Embeddings transform discrete data into continuous vector representations through a process of learning from context and relationships within data.

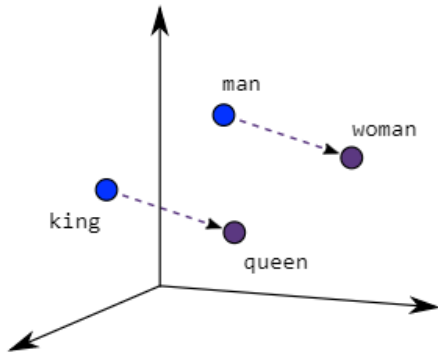
The Training Process

1. **Initialize:** Start with random vectors for each entity (word, image, etc.)
2. **Define an objective:** Create a task where predicting context is required
 - For words: Predict surrounding words (Word2Vec)
 - For sentences: Predict whether one sentence follows another (Sentence-BERT)
 - For images: Predict whether images are similar (CLIP)
3. **Train:** Update embeddings through gradient descent to minimize prediction error
4. **Extract:** The resulting vectors become your embeddings

Vector Operations

Once trained, embeddings enable powerful vector operations:

- **Similarity:** Calculate cosine similarity between vectors to find related items
- **Analogy:** Famous example: "king" - "man" + "woman" \approx "queen"
- **Clustering:** Group similar entities based on embedding proximity



Dimensionality

Most embedding models use vectors with dimensions ranging from 100-1024. Higher dimensions can capture more nuanced relationships but require more data to train effectively and more computational resources to use.

Real-World Example

Consider a product recommendation system. When a user views a red cotton t-shirt, the system doesn't just recommend identical items. Instead, it:

1. Converts the product into its embedding vector
2. Finds other products with similar embeddings
3. Recommends products that may differ in color or exact style but share underlying characteristics the user might appreciate

The system learns these relationships from user behavior patterns rather than explicit rules.

Use Cases of Embeddings

Natural Language Processing

- **Semantic search:** Find documents based on meaning rather than exact keywords
- **Text classification:** Categorize documents into topics or sentiment
- **Machine translation:** Convert words/sentences between languages
- **Question answering:** Understand the intent behind questions

Computer Vision

- **Image retrieval:** Find visually similar images
- **Face recognition:** Match faces across different images
- **Video understanding:** Classify actions or scenes

Recommender Systems

- **Content-based filtering:** Recommend items similar to what users have liked
- **Collaborative filtering:** Find users with similar preferences
- **Hybrid approaches:** Combine multiple embedding types

Graph Analysis

- **Node embeddings:** Represent entities in a network
- **Link prediction:** Identify potential new connections
- **Community detection:** Find clusters of related entities

Multimodal Applications

- **Text-to-image generation:** Transform text descriptions into visual elements
- **Cross-modal retrieval:** Find images matching text queries or vice-versa

Real-World Example

A healthcare organization might use embeddings to:

1. Convert patient medical records to embeddings
2. Identify clusters of patients with similar conditions
3. For a new patient, quickly find similar historical cases
4. Use those similar cases to recommend treatment plans based on what worked previously

This approach can find subtle patterns that rule-based systems might miss.

Implementing Embeddings with Hugging Face (Free Models)

Text Embeddings

Sentence Transformers

PYTHON

```
from sentence_transformers import SentenceTransformer

# https://milvus.io/ai-quick-reference/what-are-some-popular-pretrained-sentence-
# transformer-models-and-how-do-they-differ-for-example-allminilm6v2-vs-
# allmpnetbasev2

# Load a pre-trained model
model = SentenceTransformer('all-MiniLM-L6-v2') # Small, efficient model (384
dimensions)

# Create embeddings for some text
sentences = [
    "This patient shows signs of diabetes",
    "The patient exhibits symptoms consistent with diabetes mellitus",
    "This user frequently purchases outdoor equipment"
]

embeddings = model.encode(sentences)
print(f"Shape of embeddings: {embeddings.shape}") # [3, 384]

# Calculate similarity between first two sentences (medical) vs third (retail)
from sklearn.metrics.pairwise import cosine_similarity
similarity_matrix = cosine_similarity(embeddings)
print(f"Similarity between sentences 1 and 2: {similarity_matrix[0][1]:.4f}") #
High similarity
print(f"Similarity between sentences 1 and 3: {similarity_matrix[0][2]:.4f}") #
Low similarity
```

Max Seq Length

<https://developers.googleblog.com/en/gemini-embedding-available-gemini-api/#:~:text=An%20incredibly%20versatile%20model%2C%20Gemini,down%20from%20the%20default%203072.>

```

model

SentenceTransformer(
  (0): Transformer({'max_seq_length': 256, 'do_lower_case': False, 'architecture': 'BertModel'})
  (1): Pooling({'word_embedding_dimension': 384, 'pooling_mode_cls_token': False, 'pooling_mode_mean_tokens': True,
'pooling_mode_max_tokens': False, 'pooling_mode_mean_sqrt_len_tokens': False, 'pooling_mode_weightedmean_tokens': False,
'pooling_mode_lasttoken': False, 'include_prompt': True})
  (2): Normalize()
)

model.max_seq_length
256

```

Using Transformers Directly

PYTHON

```

from transformers import AutoTokenizer, AutoModel
import torch

# Load BERT model and tokenizer
tokenizer = AutoTokenizer.from_pretrained('bert-base-uncased')
model = AutoModel.from_pretrained('bert-base-uncased')

# Prepare text
text = "Machine learning is transforming healthcare"
inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True)

# Generate embeddings
with torch.no_grad():
    outputs = model(**inputs)

# Use CLS token as sentence embedding
embeddings = outputs.last_hidden_state[:, 0, :]
print(f"Embedding shape: {embeddings.shape}") # [1, 768]

```

Image Embeddings

PYTHON



```
from transformers import CLIPProcessor, CLIPModel
import requests
from PIL import Image

# Load model and processor
model = CLIPModel.from_pretrained("openai/clip-vit-base-patch32")
processor = CLIPProcessor.from_pretrained("openai/clip-vit-base-patch32")

# Load and process an image
url = "http://images.cocodataset.org/val2017/000000039769.jpg"
image = Image.open(requests.get(url, stream=True).raw)
inputs = processor(images=image, return_tensors="pt")

# Generate image embeddings
with torch.no_grad():
    outputs = model.get_image_features(**inputs)

image_embedding = outputs
print(f"Image embedding shape: {image_embedding.shape}") # [1, 512]
```

Multimodal Embeddings (Text and Image)

PYTHON

```
from transformers import CLIPProcessor, CLIPModel

# Load CLIP model and processor
model = CLIPModel.from_pretrained("openai/clip-vit-base-patch32")
processor = CLIPProcessor.from_pretrained("openai/clip-vit-base-patch32")

# Prepare image and text
url = "http://images.cocodataset.org/val2017/000000039769.jpg"
image = Image.open(requests.get(url, stream=True).raw)
texts = ["a photo of a cat", "a photo of a dog"]

# Process inputs
inputs = processor(text=texts, images=image, return_tensors="pt", padding=True)

# Get similarity scores
with torch.no_grad():
    outputs = model(**inputs)

logits_per_text = outputs.logits_per_text # Text-image similarity
print(f"Similarity scores: {logits_per_text}") # Higher score for correct match
```

Building a Semantic Search System

PYTHON

```
import numpy as np
from sentence_transformers import SentenceTransformer
from sklearn.metrics.pairwise import cosine_similarity

# Sample database of documents
documents = [
    "Artificial intelligence is revolutionizing healthcare",
    "Machine learning models require significant data for training",
    "Neural networks have multiple layers of computational nodes",
    "Healthcare costs continue to rise in many countries",
    "Data privacy is a major concern in medical applications"
]

# Create model and embeddings
model = SentenceTransformer('all-MiniLM-L6-v2')
document_embeddings = model.encode(documents)

# Function to search documents
def semantic_search(query, top_k=2):
    # Create embedding for the query
    query_embedding = model.encode([query])

    # Calculate similarity scores
    similarity_scores = cosine_similarity(query_embedding, document_embeddings)[0]

    # Get top results
    top_indices = np.argsort(similarity_scores)[::-1][:top_k]

    # Return results
    results = []
    for idx in top_indices:
        results.append({
            "document": documents[idx],
            "score": similarity_scores[idx]
        })
    return results

# Example search
results = semantic_search("AI in medicine", top_k=2)
for i, result in enumerate(results):
    print(f"Result {i+1}: {result['document']} (Score: {result['score']:.4f})")
```

Advanced Concepts in Embeddings

Embedding Spaces and Transfer Learning

Embeddings trained on one task often transfer well to related tasks. This is because the learned vector space captures fundamental patterns that generalize across domains. A model trained on general English text can often produce useful embeddings for specialized domains like legal or medical text with minimal fine-tuning.

Contextual vs. Static Embeddings

Early embedding models like Word2Vec produced static embeddings where each word has a single vector regardless of context. Modern transformer-based models generate contextual embeddings where the same word receives different vectors based on surrounding context:

```
"The bank denied my loan application." → financial meaning
"I sat by the bank of the river."       → geographical meaning
```

Embedding Visualization

Tools like TensorBoard or UMAP allow visualization of high-dimensional embeddings in 2D or 3D space, revealing clusters and relationships:

```
from sklearn.manifold import TSNE
import matplotlib.pyplot as plt

# Create t-SNE model
tsne = TSNE(n_components=2, random_state=42)

# Fit and transform embeddings to 2D
embeddings_2d = tsne.fit_transform(document_embeddings)

# Plot
plt.figure(figsize=(10, 8))
plt.scatter(embeddings_2d[:, 0], embeddings_2d[:, 1])

# Add labels
for i, doc in enumerate(documents):
    plt.annotate(doc[:20] + "...", (embeddings_2d[i, 0], embeddings_2d[i, 1]))

plt.title("t-SNE visualization of document embeddings")
plt.tight_layout()
plt.show()
```

PYTHON

Fine-tuning Embeddings

General-purpose embeddings can be fine-tuned for specific domains:

```
PYTHON
from sentence_transformers import SentenceTransformer, losses, InputExample
from torch.utils.data import DataLoader

# Create a training dataset of similar and dissimilar pairs
train_examples = [
    InputExample(texts=['diabetes symptoms', 'signs of high blood sugar'],
label=0.9),
    InputExample(texts=['heart disease', 'cardiac condition'], label=0.8),
    InputExample(texts=['diabetes symptoms', 'car maintenance'], label=0.1)
]

# Create dataloader
train_dataloader = DataLoader(train_examples, shuffle=True, batch_size=16)

# Load model
model = SentenceTransformer('all-MiniLM-L6-v2')

# Train with cosine similarity loss
train_loss = losses.CosineSimilarityLoss(model)
model.fit(
    train_objectives=[(train_dataloader, train_loss)],
    epochs=1,
    warmup_steps=100
)
```

Ethical Considerations

Embeddings can inherit and amplify biases present in training data. For example:

- Word embeddings trained on news articles have shown gender biases (associating "doctor" with male terms and "nurse" with female terms)
- Face recognition embeddings have demonstrated racial biases in accuracy rates

Mitigation strategies include:

- Careful dataset curation and balancing
- Debiasing techniques during or after training
- Regular auditing of embedding-based systems
- Transparency about limitations

Future Directions

- **Multimodal embeddings:** Unified representations across text, images, audio, and video

- **Sparse embeddings:** Combining benefits of sparse and dense representations for better interpretability
- **Hierarchical embeddings:** Representing information at multiple levels of abstraction
- **Domain-specific architectures:** Models optimized for particular fields like chemistry or genomics
- **Embedding compression:** Creating smaller, more efficient embeddings for resource-constrained environments
- **Foundation model embeddings:** Leveraging increasingly powerful large language models to generate rich, context-aware embeddings

Embeddings continue to evolve as a fundamental building block of modern machine learning systems, enabling machines to understand and process information in ways that more closely resemble human conceptual thinking.