# Word Embeddings: A Guided Discovery
## Pre-Class Worksheet

## 1  Character-Based Similarity

### Exploration 1: Comparing Words by Letters

Let's explore whether comparing words by their shared letters captures semantic similarity.

---

**Exercise:** Calculate the character overlap between word pairs.

For each pair, identify:

- Common letters (ignore duplicates)

- Total unique letters in the longer word

- Overlap percentage

| Word 1 | Word 2 | Common | Total | Overlap % |
|--------|--------|--------|-------|-----------|
| cat | car | _____ | 3 | _____ |
| cat | kitten | _____ | 6 | _____ |
| bank | tank | _____ | 4 | _____ |
| dog | puppy | _____ | 5 | _____ |

**Observation:** Which word pair has the highest character overlap? _____

**Question:** Does this match semantic similarity? _____

---

## 2  Understanding Dot Product as Similarity

### Mathematical Foundation

The dot product is fundamental to measuring similarity in vector spaces.

---

**Definition:** For vectors $\vec{a} = [a_1, a_2, ..., a_n]$ and $\vec{b} = [b_1, b_2, ..., b_n]$:

$$\vec{a} \cdot \vec{b} = \sum_{i=1}^{n} a_i \times b_i = a_1 b_1 + a_2 b_2 + ... + a_n b_n$$

**Geometric Interpretation:**

$$\vec{a} \cdot \vec{b} = |\vec{a}| \times |\vec{b}| \times \cos(\theta)$$
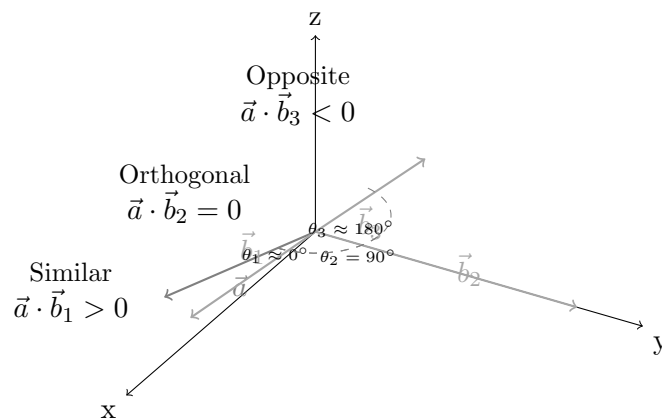
where $\theta$ is the angle between vectors.

**Key Insights:**

- When $\theta = 0°$ (parallel): $\cos(0°) = 1 \rightarrow$ Maximum similarity

---

- When $\theta = 90°$ (orthogonal): $\cos(90°) = 0 \rightarrow$ No similarity

- When $\theta = 180°$ (opposite): $\cos(180°) = -1 \rightarrow$ Maximum dissimilarity

**Practice:** Calculate the dot product:

- $[1, 0, 1] \cdot [1, 1, 0] = $ _____

- $[2, 3] \cdot [1, 2] = $ _____

- $[1, 0, 0] \cdot [0, 1, 0] = $ _____



# 3  One-Hot Encoding

**Exploration 2: Vector Representation Attempt**

One-hot encoding assigns each word a unique position in a high-dimensional space.

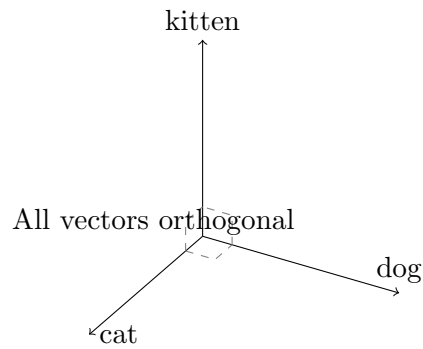**Vocabulary:** {cat, dog, kitten, car, truck}

**One-Hot Vectors:**

| Word | Vector Components | | | | |
|------|---|---|---|---|---|
| cat | 1 | 0 | 0 | 0 | 0 |
| dog | 0 | 1 | 0 | 0 | 0 |
| kitten | 0 | 0 | 1 | 0 | 0 |
| car | 0 | 0 | 0 | 1 | 0 |
| truck | 0 | 0 | 0 | 0 | 1 |

**Calculate Similarities:** Using dot product

$$\text{cat} \cdot \text{dog} = (1 \times 0) + (0 \times 1) + (0 \times 0) + (0 \times 0) + (0 \times 0) = \underline{\quad}$$
$$\text{cat} \cdot \text{kitten} = (1 \times 0) + (0 \times 0) + (0 \times 1) + (0 \times 0) + (0 \times 0) = \underline{\quad}$$
$$\text{cat} \cdot \text{car} = \underline{\hspace{4cm}} = \underline{\quad}$$
$$\text{dog} \cdot \text{truck} = \underline{\hspace{4cm}} = \underline{\quad}$$

**Discovery:** What pattern do you notice? _____
**Problem:** What angle exists between all word pairs? _____ degrees

kitten

All vectors orthogonal

dog

cat

# 4 Dense Vector Representations

## Solution: Distributed Representations

Instead of one-hot vectors, we use dense vectors where each dimension captures semantic properties.

---

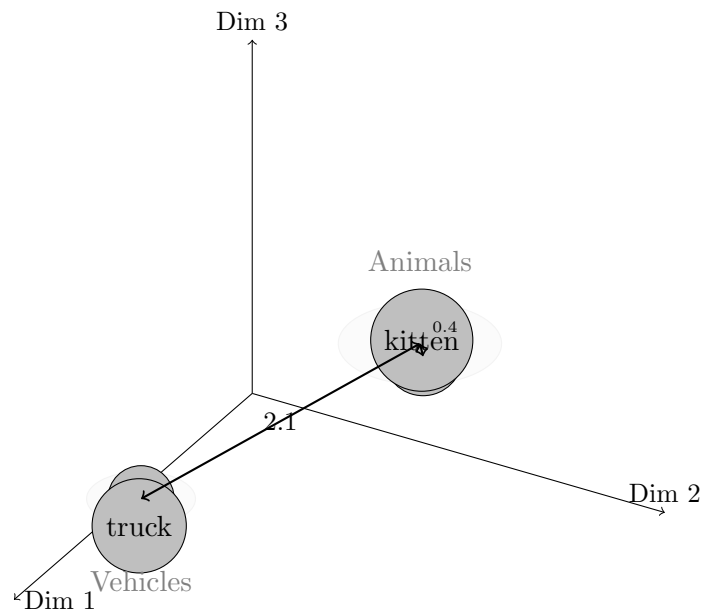**Example Dense Vectors:** (3-dimensional for visualization)

| Word | Dim 1 | Dim 2 | Dim 3 |
|---|---|---|---|
| cat | 0.2 | 0.8 | 0.5 |
| dog | 0.3 | 0.7 | 0.6 |
| kitten | 0.15 | 0.85 | 0.45 |
| car | 0.9 | 0.1 | 0.2 |
| truck | 0.85 | 0.15 | 0.25 |

**Calculate Similarities:**

$$\text{cat} \cdot \text{dog} = (0.2 \times 0.3) + (0.8 \times 0.7) + (0.5 \times 0.6) = \underline{\hspace{2cm}}$$
$$\text{cat} \cdot \text{kitten} = (0.2 \times 0.15) + (0.8 \times 0.85) + (0.5 \times 0.45) = \underline{\hspace{2cm}}$$
$$\text{cat} \cdot \text{car} = (0.2 \times 0.9) + (0.8 \times 0.1) + (0.5 \times 0.2) = \underline{\hspace{2cm}}$$
$$\text{car} \cdot \text{truck} = \underline{\hspace{4cm}} = \underline{\hspace{2cm}}$$
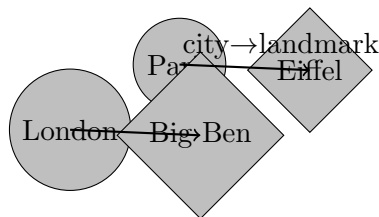
**Rank by Similarity to "cat":**

1. _____ (score: _____)

2. _____ (score: _____)

3. _____ (score: _____)

4. _____ (score: _____)

---

## 5 Vector Relationships

### Discovery: Consistent Relationships

Relationships between concepts form parallel vectors in embedding space.



---

**Vector Arithmetic:**

The relationship "city to its landmark" is consistent across examples:

$$\text{Eiffel Tower} - \text{Paris} \approx \text{Big Ben} - \text{London}$$

Therefore:

$$\text{Paris} - \text{France} + \text{UK} \approx \text{London}$$
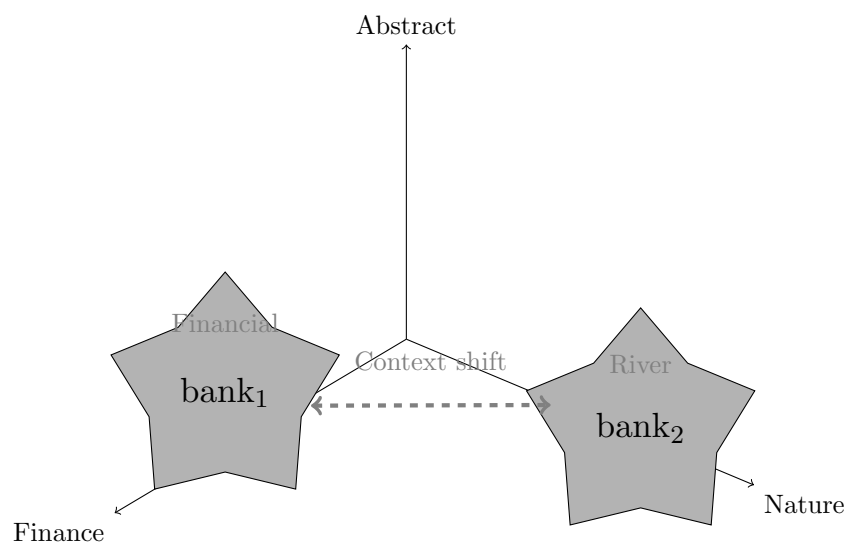
**Complete These Analogies:**

- king - man + woman = _____

- Tokyo - Japan + France = _____

- Einstein - physics + music = _____

**Reflection:** Why might relationships form parallel vectors?

_____

_____

---

# 6 Context-Dependent Representations

**Advanced: Dynamic Embeddings**

Modern systems adjust word positions based on surrounding context.

Abstract

Financial

Context shift

River

bank$_1$

bank$_2$

Finance

Nature

**Examples of Ambiguous Words:**

| Word | Context 1 | Context 2 |
|------|-----------|-----------|
| Apple | fruit, juice, tree | iPhone, Mac, company |
| Java | coffee, beans, island | programming, code, software |
| Python | snake, reptile, zoo | programming, AI, library |
| Spring | season, flowers, warm | coil, bounce, mechanism |

**Question:** How might a system determine which meaning to use?

_____

_____

# Summary

**Key Discoveries:**

1. Character overlap fails to capture semantic similarity

2. Dot product measures vector alignment and similarity

3. One-hot vectors are orthogonal, showing no relationships

4. Dense vectors place similar words nearby in space

5. Vector arithmetic captures analogical relationships

6. Context determines word position in modern systems