

# Dimensionality Reduction for NLP

## Lecture 13: PCA, t-SNE, and UMAP

PSYC 51.07: Models of Language and Communication - Week 4

Winter 2026

# Today's Lecture



1.  The Curse of Dimensionality
2.  Principal Component Analysis (PCA)
3.  t-SNE: Stochastic Neighbor Embedding
4.  UMAP: Uniform Manifold Approximation
5.  Comparison & Best Practices
6.  Visualization Techniques

*Goal: Learn to visualize and explore high-dimensional embeddings*

# The Curse of Dimensionality



## The Problem:

### Word Embeddings:

- Word2Vec: 300 dimensions
- GloVe: 300 dimensions
- FastText: 300 dimensions
- BERT: 768 dimensions
- GPT-3: 12,288 dimensions!

### Challenges:

- Can't visualize 300D space

# Why Dimensionality Reduction?

## Visualization:

- 2D/3D plots
- Explore semantic structure
- Identify clusters
- Discover patterns
- Quality assurance

## Computation:

- Faster algorithms
- Less memory
- Enable real-time systems

# Principal Component Analysis (PCA)



The classic linear dimensionality reduction method

Intuition:

- Find directions of maximum variance
- Project data onto these directions
- First PC: most variance
- Second PC: second most (orthogonal)
- And so on...

Properties:

- Linear transformation

# PCA: The Algorithm



## Step-by-Step Worked Example:

Original Data (3D  $\rightarrow$  2D):

Word	x1	x2	x3
king	2.0	1.5	0.1
queen	1.8	1.6	0.2
man	1.0	0.5	0.1
woman	0.9	0.6	0.2

**Step 1:** Center data (subtract mean)

**Step 2:** Compute covariance matrix C

# PCA in Practice



```
1from sklearn.decomposition import PCA
2from sklearn.preprocessing import StandardScaler
3import numpy as np
4import matplotlib.pyplot as plt
5
6# Assume we have word embeddings: shape (vocab_size, 300)
7# For example, Word2Vec embeddings
8
9# Standardize the data (optional but recommended)
10scaler = StandardScaler()
11embeddings_scaled = scaler.fit_transform(embeddings)
12
13# Apply PCA
14pca = PCA(n_components=2) # Reduce to 2D for visualization
15embeddings_2d = pca.fit_transform(embeddings_scaled)
16
17# Check variance explained
18print(f"Variance explained: {pca.explained_variance_ratio_}")
```

# PCA in Practice



```
21# Visualize
22plt.figure(figsize=(10, 8))
23plt.scatter(embeddings_2d[:, 0], embeddings_2d[:, 1], alpha=0.5)
24
25# Annotate some words
26words = ['king', 'queen', 'man', 'woman', 'cat', 'dog']
27for word in words:
28    idx = word_to_idx[word]
29    plt.annotate(word, (embeddings_2d[idx, 0], embeddings_2d[idx, 1]))
30
31plt.xlabel('PC1')
32plt.ylabel('PC2')
33plt.title('Word Embeddings – PCA Projection')
34plt.show()
```

...continued

Winter 2026

# PCA Limitations

## Assumes Linearity:

- Only captures linear relationships
- Word embeddings often have non-linear structure
- May miss important patterns

1Non-linear manifold  $\rightarrow$  PCA struggles here!

## Other Issues:

- **Variance  $\neq$  Importance:**
- Preserves variance, not semantic structure

Winter 2026 • Noisy dimensions may have high variance

# t-SNE: t-Distributed Stochastic Neighbor Embedding



**Non-linear dimensionality reduction for visualization**

**Key Idea:**

- Model similarity as probability
- High-D: Gaussian similarity
- Low-D: t-distribution similarity
- Minimize divergence between them
- Preserves

**Advantages:**

# t-SNE: Key Concepts



Why t-distribution in low-D?

The Crowding Problem Illustrated:

1High-D: 10 points can each have  
2                  9 equidistant neighbors

3  
4          \*      \*      \*  
5          \*     \*     \*  
6          \*      \*     \*

7  
8Low-D (2D): Can't fit 9 equidistant  
9                  neighbors around 1 point!

10  
11Solution: t-distribution has  
12                  heavier tails → allows  
13                  moderately-distant points  
14                  to spread further apart

# t-SNE in Practice



```
1from sklearn.manifold import TSNE
2import matplotlib.pyplot as plt
3
4# Apply t-SNE (can be slow for large datasets)
5tsne = TSNE(
6    n_components=2,          # 2D visualization
7    perplexity=30,          # try 5-50
8    n_iter=2000,            # iterations
9    random_state=42,        # reproducibility
10   verbose=1               # show progress
11)
12
13embeddings_2d = tsne.fit_transform(embeddings)
14
15# Visualize with labels
16plt.figure(figsize=(12, 10))
17
18# Color by semantic category (if available)
```

# t-SNE in Practice



```
21
22for cat, color in zip(categories, colors):
23    mask = labels == cat
24    plt.scatter(
25        embeddings_2d[mask, 0],
26        embeddings_2d[mask, 1],
27        c=color,
28        label=cat,
29        alpha=0.6
30    )
31
32# Annotate some words
33for word, idx in word_to_idx.items():
34    if word in important_words:
35        plt.annotate(
36            word,
37            (embeddings_2d[idx, 0], embeddings_2d[idx, 1])
38        )
```

## t-SNE in Practice

```
41plt.title('Word Embeddings - t-SNE Projection')
42plt.show()
```

...continued

# t-SNE Limitations and Caveats !

## 1. Slow: $O(n^2)$ complexity

```
1# Timing comparison
2n=1000: ~10 seconds
3n=5000: ~4 minutes
4n=10000: ~15 minutes
```

## 2. Non-deterministic:

```
1tsne1 = TSNE(random_state=42)
2tsne2 = TSNE(random_state=123)
3# Different layouts!
```

## 3. No out-of-sample:

# UMAP: Uniform Manifold Approximation & Projection



**Modern alternative to t-SNE (2018)**

**Key Advantages over t-SNE:**

- $(O(n \log n) \text{ vs } O(n^2))$
- Preserves local and global structure
- Can transform new points
- Theoretically grounded (topology)
- Scales to millions of points
- More robust hyperparameters

**When to Use:**

# UMAP in Practice



```
1import umap
2import matplotlib.pyplot as plt
3
4# Apply UMAP
5reducer = umap.UMAP(
6    n_components=2,          # 2D visualization
7    n_neighbors=15,          # local/global balance (try 5-50)
8    min_dist=0.1,            # cluster tightness (try 0.0-0.99)
9    metric='cosine',         # good for word embeddings
10   random_state=42
11)
12
13embeddings_2d = reducer.fit_transform(embeddings)
14
15# Can also transform new points! (unlike t-SNE)
16new_embeddings_2d = reducer.transform(new_embeddings)
17
18# Visualize
```

# UMAP in Practice



```
21    embeddings_2d[:, 0],  
22    embeddings_2d[:, 1],  
23    c=cluster_labels,      # color by cluster  
24    cmap='Spectral',  
25    s=5,  
26    alpha=0.6  
27)  
28plt.colorbar(scatter)  
29  
30# Annotate  
31for word in important_words:  
32    idx = word_to_idx[word]  
33    plt.annotate(  
34        word,  
35        (embeddings_2d[idx, 0], embeddings_2d[idx, 1]),  
36        fontsize=12  
37    )  
38
```

# PCA vs. t-SNE vs. UMAP



Feature	PCA	t-SNE	UMAP
Speed	Very Fast	Slow	Fast
Scalability	Excellent	Poor (<10k)	Excellent
Global structure	Yes	No	Yes
Local structure	Partial	Yes	Yes
Deterministic	Yes	No	Partial
Out-of-sample	Yes	No	Yes

Concrete Timing Comparison (10,000 word embeddings):

# Visualization Best Practices



## 1. Preprocessing:

- Standardize features (mean=0, std=1)
- Remove outliers (optional)
- For very large datasets: sample or use PCA first

## 2. Try multiple hyperparameters:

- t-SNE perplexity: 5, 10, 30, 50
- UMAP n\_neighbors: 5, 15, 30, 50
- Different views reveal different structure

## 3. Color intelligently:

- By semantic category

- By cluster assignment

# Interactive Visualization



```
1import plotly.express as px
2import plotly.graph_objects as go
3import pandas as pd
4
5# Prepare data
6df = pd.DataFrame({
7    'x': embeddings_2d[:, 0],
8    'y': embeddings_2d[:, 1],
9    'word': words,
10   'category': categories,
11   'frequency': frequencies
12})
13
14# Create interactive plot with plotly
15fig = px.scatter(
16    df,
17    x='x',
18    y='y',
```

# Interactive Visualization



```
21  hover_data=['word', 'frequency'],
22  title='Word Embeddings (UMAP)',
23  width=1000,
24  height=800
25)
26
27# Add text annotations for important words
28for word in important_words:
29    row = df[df['word'] == word].iloc[0]
30    fig.add_annotation(
31        x=row['x'],
32        y=row['y'],
33        text=word,
34        showarrow=False,
35        font=dict(size=12)
36    )
37
38# Show
```

# Interactive Visualization



```
41# Can also save as HTML  
42fig.write_html('embeddings_viz.html')
```

...continued

# Applications in NLP

## 1. Embedding Quality Check:

```
1# Quick sanity check
2words = ['dog', 'cat', 'fish', # animals
3          'car', 'bus', 'train'] # vehicles
4
5# If good embeddings: two clusters!
6# If bad: random scatter
```

## 2. Finding Polysemous Words:

1 "apple" appears in TWO clusters:  
2- Near: orange, banana, fruit  
3- Near: microsoft, google, tech  
4  
5→ Word has multiple senses!

# Case Study: Visualizing BERT Layers

**Question:** What do different BERT layers capture?

**Approach:**

1. Extract embeddings from each layer
2. Apply UMAP to each layer separately
3. Visualize and compare

**Typical Findings:**

- **Layer 0-2:** Syntactic (POS tags cluster)
- **Layer 3-8:** Semantic (meaning clusters)
- **Layer 9-12:** Task-specific

# Case Study: Visualizing BERT Layers

```
21     viz_2d = reducer.fit_transform(avg_embeddings)
22
23     # Plot
24     plt.figure()
25     plt.scatter(viz_2d[:, 0], viz_2d[:, 1], c=pos_tags)
26     plt.title(f'Layer {layer_idx}')
27     plt.show()
```

...continued

Reference: Jawahar et al. (2019). "What Does BERT Learn about the Structure of Language?"

# Discussion Question



**What does a 2D visualization actually tell us about 300D space?**

**Consider:**

- We compress 300 dimensions into 2
- Massive information loss (>99%)
- Different methods show different views
- Hyperparameters change the story

**What we CAN conclude:**

- Rough semantic groupings
- Relative proximities

# Practical Tips



## 1. Choose the right tool:

- Default: UMAP (fast, balanced)
- Need speed: PCA
- Small dataset, only local: t-SNE

## 2. Preprocessing matters:

- StandardScaler for PCA
- Consider normalization for UMAP/t-SNE
- Remove extreme outliers

## 3. Try multiple settings:

- Run with different hyperparameters

# Summary

## What we learned today:

1. **Curse of Dimensionality:** High-D space is weird, need reduction

2. **PCA:**

- Linear, fast, preserves global structure
- Good for preprocessing and quick exploration

3. **t-SNE:**

- Non-linear, preserves local structure
- Beautiful visualizations but slow
- Sensitive to hyperparameters

4. **UMAP:**

- Fast, scalable, balanced local+global

# Key References



## Foundational Papers:

- Pearson, K. (1901). "On Lines and Planes of Closest Fit to Systems of Points in Space" (PCA)
- van der Maaten & Hinton (2008). "Visualizing Data using t-SNE"
- McInnes et al. (2018). "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction"

## Applications:

- Jawahar et al. (2019). "What Does BERT Learn about the Structure of Language?"
- Coenen et al. (2019). "Visualizing and Measuring the Geometry of BERT"

## Guides & Tutorials:

Winter 2026

# Questions?

**Next Lecture:**

Cognitive Models of Semantic Representation

*How do humans represent meaning?*