

Lecture 10: X-Hour Embeddings Workshop

Week 3: Hands-On Dimensionality Reduction and Word Vectors

PSYC 51.07: Models of Language and Communication

Learning Objectives

By the end of this session, you will:

1. Implement classic dimensionality reduction (LSA, LDA)
2. Train and analyze Word2Vec embeddings
3. Visualize high-dimensional embeddings using UMAP
4. Compare different embedding methods on real data
5. Understand semantic relationships captured by embeddings

Workshop format: Hands-on coding with the 20 Newsgroups dataset

Workshop Overview

Today's Agenda:

1. **Part 1:** Why embeddings? From sparse to dense representations
2. **Part 2:** LSA - Latent Semantic Analysis with SVD
3. **Part 3:** LDA - Latent Dirichlet Allocation for topic modeling
4. **Part 4:** Word2Vec - Neural word embeddings
5. **Part 5:** Visualizing embeddings with UMAP
6. **Part 6:** Comparing methods and document classification

Companion notebook: `xhour_embeddings_demo.ipynb`

Part 1: Why Embeddings?

The problem with sparse representations:

Last week (BoW, TF-IDF):

- High dimensional (vocab size)
- Sparse (mostly zeros)
- No semantic similarity
- "dog" and "puppy" are orthogonal

Embeddings:

- Low dimensional (50-300 dims)
- Dense (all non-zero)
- Similar words cluster together

Sparse vs Dense: Concrete Comparison

```
1# Sparse representation (one-hot / BoW)
2# Vocabulary: [cat, dog, puppy, car, truck, vehicle]
3
4cat_sparse = [1, 0, 0, 0, 0, 0] # 6 dimensions, 5 zeros
5dog_sparse = [0, 1, 0, 0, 0, 0]
6puppy_sparse = [0, 0, 1, 0, 0, 0]
7
8# Cosine similarity: cat-dog = 0, dog-puppy = 0 (orthogonal!)
9
10# Dense embedding (learned from data)
11cat_dense = [0.8, -0.2, 0.5] # 3 dimensions, all non-zero
12dog_dense = [0.7, -0.1, 0.6] # Similar to cat!
13puppy_dense = [0.75, -0.15, 0.55] # Very similar to dog!
14car_dense = [-0.3, 0.9, -0.4] # Different cluster
15
16# Cosine similarity: cat-dog = 0.98, dog-puppy = 0.99
```

Key insight: Dense embeddings capture that "dog" and "puppy" are semantically related

The Magic of Word Vectors

Famous example: king - man + woman = queen

Vector Arithmetic

Word embeddings capture semantic relationships as directions in space:

- Gender direction: woman - man
- Royalty direction: king - queen
- Pluralization: words - word

Key insight: Meaning encoded as geometry!

Vector Arithmetic: Step-by-Step Example

```
1import numpy as np
2
3# Pretend embeddings (simplified to 3D for illustration)
4embeddings = {
5    'king': np.array([0.9, 0.8, 0.2]),
6    'queen': np.array([0.85, 0.75, 0.7]),
7    'man': np.array([0.7, 0.6, 0.1]),
8    'woman': np.array([0.65, 0.55, 0.6]),
9}
10
11# The analogy: king - man + woman = ?
12result = embeddings['king'] - embeddings['man'] + embeddings['woman']
13# result = [0.9-0.7+0.65, 0.8-0.6+0.55, 0.2-0.1+0.6]
14#           = [0.85, 0.75, 0.7] ← Very close to 'queen'!
15
16# Why does this work?
17# king - man = "royalty" direction = [0.2, 0.2, 0.1]
18# woman + royalty = queen
```

Part 2: Latent Semantic Analysis (LSA)

Using SVD to find latent topics:

$$X \approx U_k \Sigma_k V_k^T$$

Algorithm:

1. Build TF-IDF matrix X
2. Apply Singular Value Decomposition
3. Keep top k dimensions
4. Use U_k as word embeddings

Interpretation:

- U : word-topic associations
- Σ : topic strengths

LSA in Code

```
1from sklearn.decomposition import TruncatedSVD
2from sklearn.feature_extraction.text import TfidfVectorizer
3
4# Build TF-IDF matrix
5tfidf = TfidfVectorizer(max_features=5000, stop_words='english')
6tfidf_matrix = tfidf.fit_transform(documents)
7
8# Apply LSA
9lsa = TruncatedSVD(n_components=100, random_state=42)
10doc_embeddings = lsa.fit_transform(tfidf_matrix)
11word_embeddings = lsa.components_.T
12
13print(f"Explained variance: {lsa.explained_variance_ratio_.sum():.2%}")
```

Try it: Find similar words using cosine similarity!

LSA: Finding Similar Words

```
1from sklearn.metrics.pairwise import cosine_similarity
2import numpy as np
3
4# Get vocabulary mapping
5vocab = tfidf.get_feature_names_out()
6word_to_idx = {word: i for i, word in enumerate(vocab)}
7
8def find_similar_words(word, top_n=5):
9    """Find words with similar LSA embeddings."""
10    if word not in word_to_idx:
11        return f"'{word}' not in vocabulary"
12
13    idx = word_to_idx[word]
14    word_vec = word_embeddings[idx].reshape(1, -1)
15
16    # Compute similarities to all words
17    sims = cosine_similarity(word_vec, word_embeddings)[0]
18
```

LSA: Finding Similar Words

```
21     return [(vocab[i], f"{sims[i]:.3f}") for i in top_indices]
22
23print(find_similar_words("computer"))
24# Output: [('software', 0.82), ('program', 0.79),
25#             ('system', 0.71), ('hardware', 0.68), ('disk', 0.65)]
```

...continued

Part 3: LDA for Topic Modeling

A probabilistic approach:

Generative Story

LDA imagines documents are created by:

1. Choosing a mixture of topics
2. For each word, picking a topic
3. Sampling a word from that topic

Key difference from LSA:

- Probabilistic interpretation
- Non-negative weights
- More interpretable topics

LDA Example Output

1Topic 0: hockey, game, team, player, season, nhl, play
2Topic 1: space, nasa, launch, orbit, shuttle, satellite
3Topic 2: computer, software, program, file, windows, system
4Topic 3: medical, doctor, patient, disease, health, treatment
5Topic 4: government, president, congress, law, political

Each document is a **mixture** of topics:

Document #42: 60% Space + 25% Computer + 15% Other

LDA: Complete Working Example

```
1from sklearn.decomposition import LatentDirichletAllocation
2from sklearn.feature_extraction.text import CountVectorizer
3
4# 20 Newsgroups sample documents
5documents = [
6    "The hockey team scored three goals in the game",
7    "NASA launched a new satellite into orbit",
8    "Install the software program on your computer",
9    "The doctor prescribed medicine for the patient",
10   # ... more documents
11]
12
13# Step 1: Create bag-of-words matrix
14vectorizer = CountVectorizer(max_features=1000, stop_words='english')
15bow_matrix = vectorizer.fit_transform(documents)
16
17# Step 2: Fit LDA
18lda = LatentDirichletAllocation(n_components=5, random_state=42)
```

LDA: Complete Working Example

```
21# Step 3: Print topics
22vocab = vectorizer.get_feature_names_out()
23for topic_idx, topic in enumerate(lda.components_):
24    top_words = [vocab[i] for i in topic.argsort()[-7:]]
25    print(f"Topic {topic_idx}: {', '.join(top_words)})
```

...continued

Part 4: Word2Vec

Learning embeddings from context:

Skip-gram:

Given target word, predict context

"The **cat** sat on mat"

- cat → the, sat, on

CBOW:

Given context, predict target

the, sat, on → **cat**

```
1from gensim.models import Word2Vec  
2  
3# 1.1. Model /
```

Word2Vec: Semantic Similarity

```
1# Find similar words
2model.wv.most_similar('computer', topn=5)
3# [('software', 0.82), ('program', 0.79), ('system', 0.75), ...]
4
5# Word analogies
6model.wv.most_similar(
7    positive=['woman', 'king'],
8    negative=['man']
9)
10# [( 'queen', 0.71), ...]
```

Hands-on Exercise

Try creating your own word analogies! What works? What fails?

Word2Vec: Exploring Analogies

```
1# Analogies that typically WORK well:  
2model.wv.most_similar(positive=['paris', 'germany'], negative=['france'])  
3# → 'berlin' (capital cities)  
4  
5model.wv.most_similar(positive=['walking', 'swam'], negative=['swimming'])  
6# → 'walked' (verb tenses)  
7  
8model.wv.most_similar(positive=['bigger', 'cold'], negative=['big'])  
9# → 'colder' (comparatives)  
10  
11# Analogies that often FAIL:  
12model.wv.most_similar(positive=['doctor', 'woman'], negative=['man'])  
13# → might return 'nurse' instead of 'doctor' (reflects bias!)  
14  
15model.wv.most_similar(positive=['sushi', 'italy'], negative=['japan'])  
16# → uncertain results (cultural associations are noisy)
```

Part 5: Visualizing with UMAP

Projecting 100D → 2D:

```
1import umap
2
3reducer = umap.UMAP(
4    n_neighbors=15,
5    min_dist=0.1,
6    metric='cosine'
7)
8
9embeddings_2d = reducer.fit_transform(word_vectors)
```

UMAP advantages:

- Faster than t-SNE
- Preserves global structure

What You Should See

When you visualize embeddings:

Sports cluster:

- hockey, baseball, player, team, game

Space cluster:

- nasa, shuttle, orbit, launch, space

Tech cluster:

- computer, software, program, windows

Medical cluster:

- doctor, patient, hospital, treatment

UMAP Visualization: Complete Code

```
1import umap
2import matplotlib.pyplot as plt
3
4# Get word vectors for a subset of interesting words
5words_to_plot = ['hockey', 'baseball', 'player', 'team', 'game',
6                  'nasa', 'shuttle', 'orbit', 'space', 'satellite',
7                  'computer', 'software', 'program', 'windows', 'disk',
8                  'doctor', 'patient', 'hospital', 'disease', 'treatment']
9
10word_vectors = np.array([model.wv[w] for w in words_to_plot])
11
12# Reduce to 2D with UMAP
13reducer = umap.UMAP(n_neighbors=5, min_dist=0.3, metric='cosine')
14embeddings_2d = reducer.fit_transform(word_vectors)
15
16# Plot
17plt.figure(figsize=(12, 8))
18plt.scatter(embeddings_2d[:, 0], embeddings_2d[:, 1], alpha=0.7)
```

UMAP Visualization: Complete Code

```
21plt.title("Word Embeddings Visualized with UMAP")
22plt.savefig("word_clusters.png")
```

...continued

Part 6: Comparing Methods

Method	Speed	Interpretability	Quality	Data Needed
LSA	Fast	Medium	Medium	Small-Medium
LDA	Medium	High	Medium	Medium
Word2Vec	Medium	Low	High	Large

Recommendations:

- **Quick exploration:** LSA
- **Interpretable topics:** LDA
- **Best semantic quality:** Word2Vec

Document Classification with Embeddings

Using embeddings as features:

```
1def document_vector(doc, model):  
2    """Average word vectors for document."""  
3    tokens = preprocess(doc)  
4    vectors = [model.wv[w] for w in tokens if w in model.wv]  
5    return np.mean(vectors, axis=0) if vectors else np.zeros(100)  
6  
7# Train classifier  
8X_train = [document_vector(doc, w2v) for doc in train_docs]  
9clf = LogisticRegression()  
10clf.fit(X_train, y_train)
```

Compare to TF-IDF baseline!

Classification Comparison: Full Example

```
1from sklearn.linear_model import LogisticRegression
2from sklearn.model_selection import cross_val_score
3
4# Method 1: TF-IDF baseline
5tfidf = TfidfVectorizer(max_features=5000)
6X_tfidf = tfidf.fit_transform(train_docs)
7clf_tfidf = LogisticRegression(max_iter=1000)
8tfidf_scores = cross_val_score(clf_tfidf, X_tfidf, y_train, cv=5)
9
10# Method 2: LSA embeddings
11lsa = TruncatedSVD(n_components=100)
12X_lsa = lsa.fit_transform(X_tfidf)
13clf_lsa = LogisticRegression(max_iter=1000)
14lsa_scores = cross_val_score(clf_lsa, X_lsa, y_train, cv=5)
15
16# Method 3: Word2Vec embeddings
17X_w2v = np.array([document_vector(doc, model) for doc in train_docs])
18clf_w2v = LogisticRegression(max_iter=1000)
```

Classification Comparison: Full Example

```
21print(f"TF-IDF: {tfidf_scores.mean():.3f} (+/- {tfidf_scores.std():.3f})")
22print(f"LSA: {lsa_scores.mean():.3f} (+/- {lsa_scores.std():.3f})")
23print(f"Word2Vec: {w2v_scores.mean():.3f} (+/- {w2v_scores.std():.3f})")
```

...continued

Key Takeaways

- 1. Embeddings capture semantic meaning** - similar words have similar vectors
- 2. Different methods, different strengths:**
 - LSA: Fast, linear, interpretable
 - LDA: Probabilistic, topic-focused
 - Word2Vec: Neural, best for similarity
- 3. Visualization reveals structure** - UMAP shows semantic clusters
- 4. Limitations:**
 - Static (one vector per word, no context)
 - Requires substantial data
 - Can encode biases

Discussion Questions

- 1. Why does vector arithmetic work?** What does "king - man + woman" really mean geometrically?
- 2. Bias in embeddings:** If Word2Vec learns from news articles, what biases might it capture?
- 3. Window size matters:** What happens with window=2 vs window=10?
- 4. Out-of-vocabulary problem:** How do you handle words not in your vocabulary?
- 5. When to use what:** For a sentiment analysis task, would you choose LSA, LDA, or Word2Vec?

Next Steps

For Assignment 2:

- Use embeddings to improve your classifier
- Compare at least 2 embedding methods
- Visualize your embeddings

Coming up in Lecture 11:

- Modern neural word embeddings
- GloVe and FastText
- Subword tokenization

Office hours: Available if you need help with the notebook!