

# Neural Language Models

## Week 2 - Word Embeddings and Word2Vec

NLP Course 2025

September 28, 2025

Using Optimal Readability Template

## Week 2: The Semantic Revolution

### From Words as IDs to Words as Meanings

#### The Problem

- Words are just **IDs**
- No semantic similarity
- "cat" and "dog" equally different as "cat" and "democracy"
- Can't generalize

#### The Solution

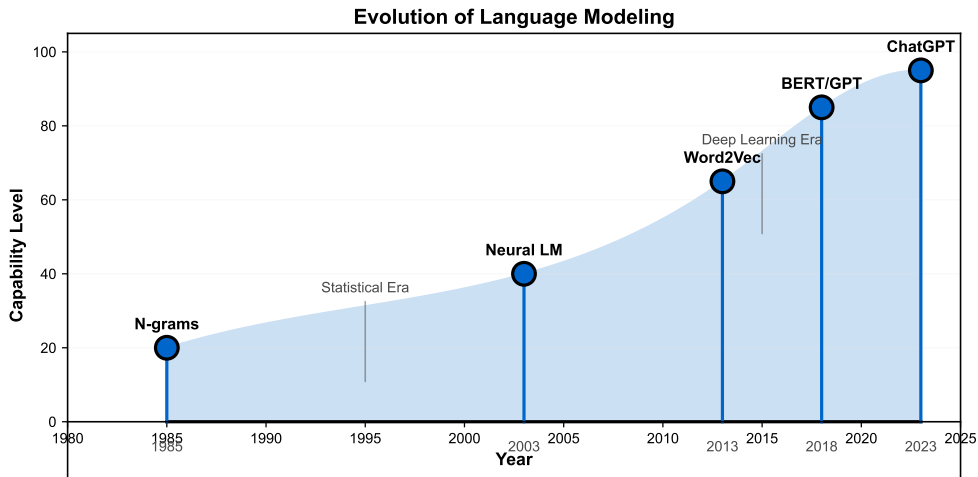
- Words as **vectors**
- Similar words nearby
- Math operations work!
- $\text{King} - \text{Man} + \text{Woman} = \text{Queen}$

#### The Impact

- Powers all modern NLP
- **1M+** developers use
- Semantic search
- Foundation for GPT/BERT

**Core Insight: You shall know a word by the company it keeps**

# The Evolution of Language Modeling



# Interactive: Word Association Game

Fill in the blank:

1. The cat sat on the \_\_\_\_\_
2. I drink my coffee with milk and \_\_\_\_\_
3. The capital of France is \_\_\_\_\_
4. She opened the door with her \_\_\_\_\_

How did you know?

- Context provides meaning
- Similar contexts → similar words
- Pattern recognition

This is Word2Vec's insight:

- Learn from **billions** of contexts
- Words in similar contexts get **similar vectors**
- Mathematics captures semantics

# Where Word Embeddings Power Your Life (2025)

## Search Engines

- Google semantic search
- Bing neural matching
- DuckDuckGo instant answers

## Translation

- Google Translate
- DeepL
- Microsoft Translator

## Business Tools

- Grammarly corrections
- Resume matching
- Customer support bots

## Virtual Assistants

- Siri/Alexa understanding
- Google Assistant
- ChatGPT responses

## Recommendations

- Netflix shows
- Spotify Discover
- YouTube suggestions

## Market Size

- **\$2.7B** by 2025
- **1M+** developers
- **500M+** daily users

# The 2013 Breakthrough: Mathematical Semantics

$$\text{King} - \text{Man} + \text{Woman} = \text{Queen}$$

## The Discovery:

- Vectors encode **relationships**
- Arithmetic operations preserve meaning
- Geometry captures semantics

## Why This Matters:

- Computers understand **analogies**
- Transfer learning possible
- One model, many tasks
- Foundation for all modern NLP

## More Examples:

- Paris - France + Italy = **Rome**
- Bigger - Big + Small = **Smaller**
- Walking - Walk + Swim = **Swimming**

Word2Vec paper: 16,000+ citations

Semantic relationships become vector arithmetic

# The Distributional Hypothesis

## Linguistic Foundation (1954):

"You shall know a word by the company it keeps" - J.R. Firth

## What it means:

- Words with similar **contexts** have similar **meanings**
- Context = surrounding words
- Meaning emerges from usage

## Example contexts for "bank":

- "deposit money in the bank"
- "sitting by the river bank"
- Different contexts → different meanings

## How Word2Vec uses this:

1. Scan billions of sentences
2. Track which words appear together
3. Words in similar contexts get similar vectors
4. Geometry encodes semantics

## The Magic:

- No human labeling needed
- Learns from raw text
- Scales to millions of words
- Works for any language

# From Sparse to Dense: The Representation Revolution

## One-Hot Encoding (Old Way):

- Vocabulary size: 50,000
- cat = [0,0,1,0,0,...,0]
- dog = [0,0,0,1,0,...,0]

## Problems:

- 50,000 dimensions!
- No similarity:  $\text{cat} \cdot \text{dog} = 0$
- Can't generalize
- Massive memory usage

## Dense Embeddings (Word2Vec):

- Embedding size: 300
- cat = [0.2, -0.4, 0.7, ...]
- dog = [0.3, -0.3, 0.6, ...]

## Benefits:

- 99.4% smaller!
- Similarity:  $\text{cat} \cdot \text{dog} = 0.8$
- Generalizes to new contexts
- Efficient computation

From 50,000 sparse dimensions to 300 dense dimensions

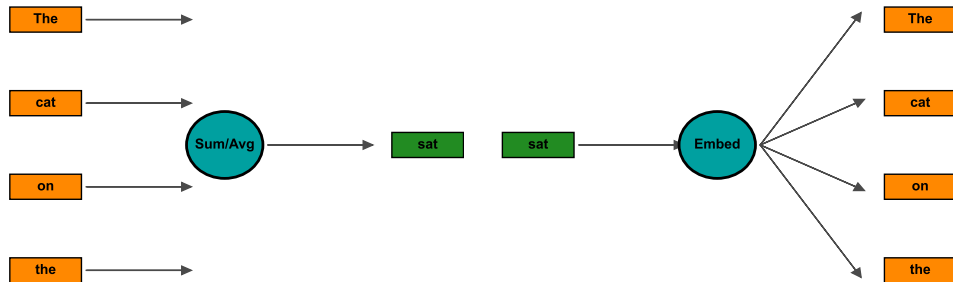


# Word2Vec: Two Architectures

## Word2Vec Architecture Comparison

CBOW: Context  $\rightarrow$  Center

Skip-gram: Center  $\rightarrow$  Context



# Skip-gram: The Architecture That Won

**Training Objective:** Predict context from center word

Given: "The cat sat on the mat"

**Input:** "cat" (center word)

**Outputs to predict:**

- "The" (position -1)
- "sat" (position +1)
- Sometimes: "on" (+2), "the" (-2)

**Window size = 2:**

- Look 2 words left/right
- 4 predictions per center word
- More context = better vectors

**Why Skip-gram wins:**

- Better on **rare words**
- More training examples
- Superior semantic quality
- Used by Google, Facebook

**Training data from one sentence:**

- (cat, The)
- (cat, sat)
- (sat, cat)
- (sat, on)
- ... many more pairs

# Building Word2Vec in PyTorch

```
1 import torch
2 import torch.nn as nn
3
4 class Word2Vec(nn.Module):
5     def __init__(self, vocab_size, embed_dim):
6         super().__init__()
7         # Two embedding matrices
8         self.center_embeddings = nn.Embedding(
9             vocab_size, embed_dim
10        )
11        self.context_embeddings = nn.Embedding(
12            vocab_size, embed_dim
13        )
14
15    def forward(self, center, context):
16        # Get embeddings
17        center_embeds = self.center_embeddings(center)
18        context_embeds =
19            self.context_embeddings(context)
20
21        # Dot product = similarity
22        scores = torch.sum(
23            center_embeds * context_embeds, dim=1
24        )
25
26    return scores
```

## Key Design Choices:

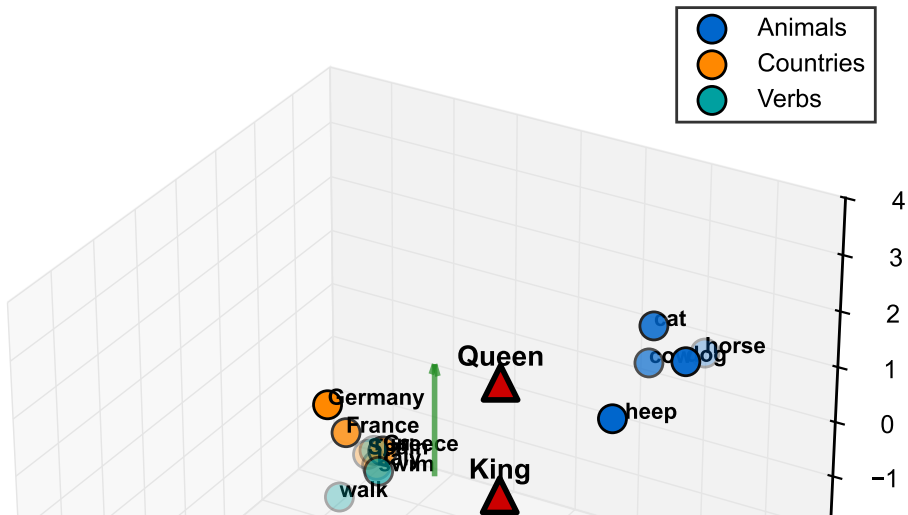
- **Two matrices:** center and context
- Embedding dim: typically 300
- Dot product for similarity
- Simple = fast training

## Training Process:

1. Sample (center, context) pairs
2. Compute similarity scores
3. Maximize correct pairs
4. Minimize random pairs

Full implementation: 50 lines of code!

## Word Embeddings in 3D Space



# The Softmax Challenge

Converting scores to probabilities:

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^V e^{z_j}}$$

## The Problem:

- Vocabulary size  $V = 50,000$
- Must compute all 50,000 scores
- Denominator sums **50,000 exponentials**
- Every training step!

## Computational cost:

- Per sample:  $O(V \cdot d)$
- 1B training samples
- = **15 trillion operations**

## The Solution: Negative Sampling

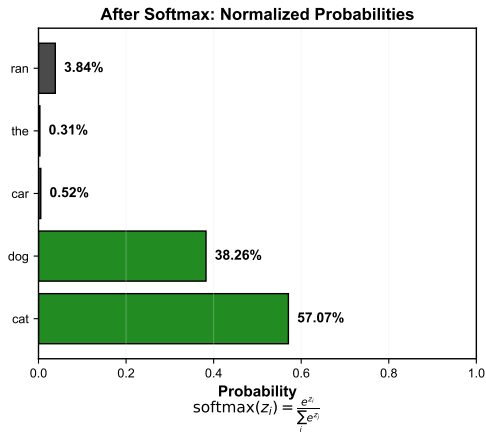
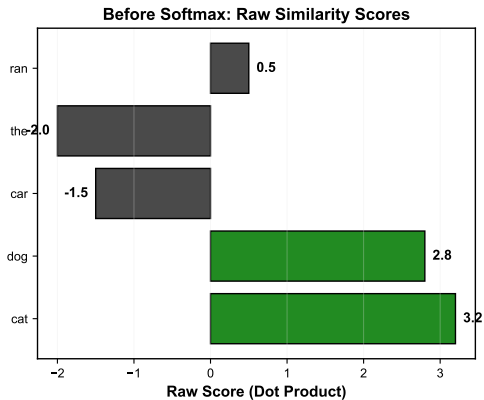
- Don't compute all 50,000
- Just sample **5-20 negatives**
- 99.96% speedup!
- Quality stays the same

## New objective:

- Maximize:  $P(\text{correct context})$
- Minimize:  $P(\text{random words})$
- Binary classification  $\times 6$
- Much faster training

# Softmax Computation Explained

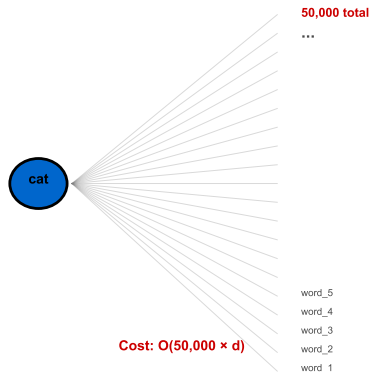
## Softmax: Converting Scores to Probabilities



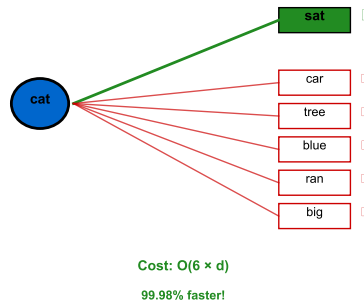
# Negative Sampling: Before and After

## Negative Sampling: The Optimization That Made Word2Vec Practical

Full Softmax: Compute All 50,000 Words



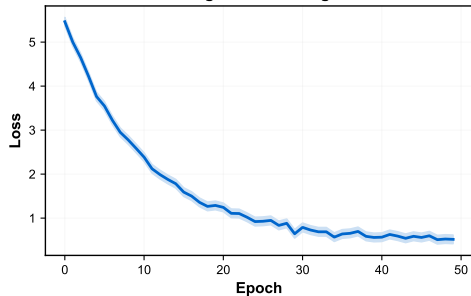
Negative Sampling: Only 5-20 Words



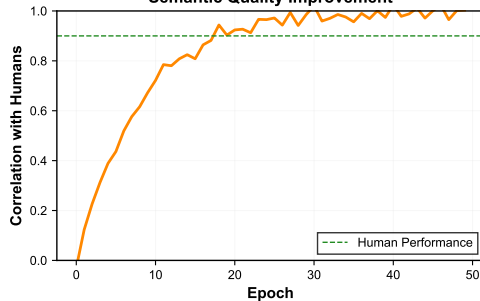
# Training Dynamics and Convergence

## Word2Vec Training Dynamics

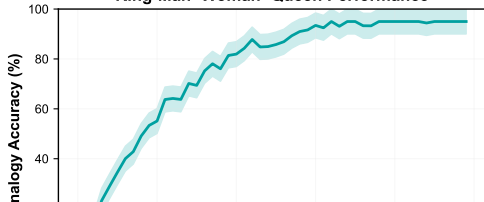
### Training Loss Convergence



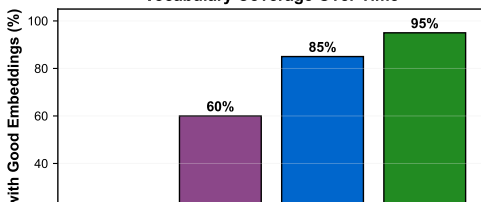
### Semantic Quality Improvement



### King-Man+Woman=Queen Performance



### Vocabulary Coverage Over Time

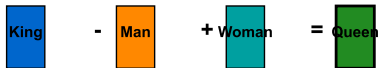




# Semantic Arithmetic in Action

## Semantic Arithmetic: Mathematical Operations on Meaning

### Gender Relationship



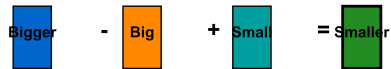
### Capital Cities



### Verb Conjugation



### Comparative Forms



# How Do We Know It Works?

## Intrinsic Evaluation:

### Word Similarity Tasks:

- Human ratings: cat-dog = 7.5/10
- Model similarity: cosine(cat, dog)
- Correlation with humans
- WordSim-353 dataset

### Analogy Tasks:

- a:b :: c:?
- Berlin:Germany :: Paris:?
- Google analogy dataset
- 90%+ accuracy

## Extrinsic Evaluation:

### Downstream Tasks:

- Sentiment analysis
- Named entity recognition
- Machine translation
- Question answering

### Real-world metrics:

- Search relevance ↑15%
- Translation BLEU ↑3.2
- Classification F1 ↑8%
- All from better embeddings!

**Good embeddings improve everything downstream**

# Challenges: Not Everything Is Perfect

## 1. Polysemy Problem:

- "bank" (financial) = "bank" (river)
- One vector for all meanings
- Averages different senses
- Solution: Contextual embeddings (BERT)

## 2. Rare Words:

- Need many examples
- "serendipity" appears rarely
- Poor vectors for rare words
- Solution: Subword embeddings

## 3. Bias Amplification:

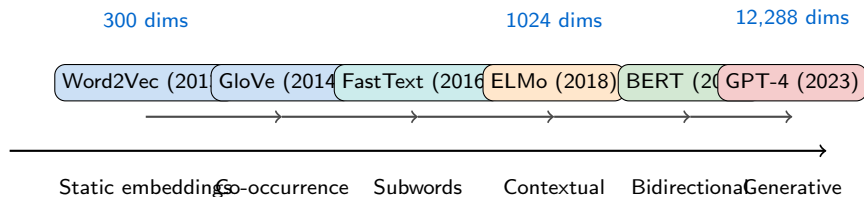
- **Learns societal biases**
- Doctor:Male :: Nurse:Female
- Amplifies stereotypes
- Active research area

## 4. Static Embeddings:

- Fixed after training
- Can't adapt to new contexts
- No fine-tuning possible
- Solution: Transformer models

These limitations led to BERT and GPT development

# From Word2Vec to ChatGPT: The Journey



## Word2Vec's Legacy:

- Proved embeddings work
- Inspired contextual models
- Still used in production
- Foundation for all modern NLP

## What Changed:

- Static → Contextual
- 300 dims → 12,000+ dims
- Word-level → Subword
- Millions → Billions of parameters

# Build It: Semantic Search Engine

```
1 import numpy as np
2 from gensim.models import Word2Vec
3
4 def semantic_search(query, documents, model):
5     """Find semantically similar documents"""
6
7     # Vectorize query
8     query_vec = document_vector(query, model)
9
10    # Vectorize all documents
11    doc_vectors = [
12        document_vector(doc, model)
13        for doc in documents
14    ]
15
16    # Compute similarities
17    similarities = [
18        cosine_similarity(query_vec, doc_vec)
19        for doc_vec in doc_vectors
20    ]
21
22    # Return ranked results
23    ranked = sorted(
24        zip(documents, similarities),
25        key=lambda x: x[1],
26        reverse=True
```

## How it works:

1. Convert query to vector
2. Convert documents to vectors
3. Find nearest neighbors
4. Return ranked results

## Real Examples:

Query: "animal pets"

### Results:

- "dog training tips"
- "cat care guide"
- "hamster habitats"

No keyword matching needed!

## Week 2 Summary: Words Have Meaning!

- Words as **IDs** → Words as **vectors**
- Distributional hypothesis: **Context defines meaning**
- Word2Vec: **Skip-gram** + **Negative sampling**
- Mathematical semantics:  $\text{King} - \text{Man} + \text{Woman} = \text{Queen}$
- From 50,000 sparse → **300 dense** dimensions
- Powers modern NLP: Search, translation, chatbots
- Limitations led to BERT/GPT development

### Key Technical Insights:

- Dot product captures similarity
- Negative sampling avoids softmax bottleneck
- Embeddings are the foundation of all modern NLP

**Next Week:** Recurrent Neural Networks  
How do we process sequences using embeddings?