

Autocomplete & Language Models

Word Embeddings

Word2Vec, GloVe & Neural Language Models

Contents

1

Language Models

N-gram models, Smoothing, Autocomplete

4

Word2Vec

CBOW, Skip-gram, Negative Sampling

2

Perplexity & Evaluation

Đánh giá chất lượng language model

5

GloVe

Global Vectors, Co-occurrence matrix

3

Word Embeddings

Distributional semantics, Vector space

6

Lab & Final Assignment

Release Final Assignment (50%)

Review: Autocorrect & HMM

Autocorrect System

Edit Distance (Levenshtein)

$$D[i,j] = \min(D[i-1,j]+1, D[i,j-1]+1, D[i-1,j-1]+cost)$$

HMM & POS Tagging

HMM Parameters: $\lambda = (A, B, \pi)$

A: Transition | B: Emission | π : Initial

Noisy Channel Model

$$\hat{w} = \operatorname{argmax} P(w) \times P(x|w)$$

Language Model Error Model

Viterbi Algorithm

Dynamic Programming cho sequence decoding

$$v[t,j] = \max(v[t-1,i] \times a[i,j]) \times b[j,o[t]]$$

PART 1

Language Models

N-gram Models & Autocomplete

Unigram, Bigram, Trigram

Smoothing Techniques

Autocomplete System

What is a Language Model

Definition

A Language Model (LM) is defined as a probabilistic model that assigns probabilities to sequences of words. It identifies which sentences have a higher 'likelihood of occurrence' within natural language

Core Task

$$P(W) = P(w_1, w_2, \dots, w_n)$$

Prob of a sequence

Insight

$P(\text{"The cat sat"}) >> P(\text{"Sat cat the"})$

LM captures grammar, meaning, world knowledge

Applications

Autocomplete

"I want to eat ____" → pizza, lunch, dinner

Speech Recognition

Choose transcription that has highest prob

Machine Translation

Choose the most fluent transcription

Spelling Correction

Noisy channel: $\operatorname{argmax} P(w) \times P(x|w)$

Text Generation

ChatGPT, GPT-4, Claude, ...

Chain Rule of Probability

CHAIN RULE

$$\begin{aligned} P(w_1, w_2, \dots, w_n) &= P(w_1) \times P(w_2 | w_1) \times P(w_3 | w_1, w_2) \times \dots \times P(w_n | w_1, \dots, w_{n-1}) \\ &= \prod P(w_k | w_1, \dots, w_{k-1}) \end{aligned}$$



Example: "The students opened their books"

$P(\text{The})$

$P(\text{students} | \text{The})$

$P(\text{opened} | \text{The students})$

$P(\text{their} | \text{The students opened})$

$P(\text{books} | \text{The students opened their})$

Problem

More history length → harder to estimate prob

We don't have enough data to count long sequences

Solution: Markov Assumption

Limit history length

$P(w_n | w_1 \dots w_{n-1}) \approx P(w_n | w_{n-k} \dots w_{n-1})$

N-gram Models

Unigram (n=1)

$$P(w_n | w_1 \dots w_{n-1}) \approx P(w_n)$$

Don't have context – depend on word frequency

Bigram (n=2)

$$P(w_n | w_1 \dots w_{n-1}) \approx P(w_n | w_{n-1})$$

Consider 1 previous word

Trigram (n=3)

$$P(w_n | w_1 \dots w_{n-1}) \approx P(w_n | w_{n-2}, w_{n-1})$$

Consider 2 previous words - balance accuracy vs sparsity

N-gram (general)

$$P(w_n | w_1 \dots w_{n-1}) \approx P(w_n | w_{n-n+1} \dots w_{n-1})$$

4-gram, 5-gram: need more data

Maximum Likelihood Estimation

$$P(w_n | w_{n-1}) = C(w_{n-1}, w_n) / C(w_{n-1})$$

$C(\cdot)$ = count in training corpus

Example: Bigram

Corpus: "I want to eat Chinese food"

Trade-off

Larger n → More context → Better accuracy

Larger n → More sparse data → Less reliable

Bigram Example

Training Corpus

"I am Sam"

"Sam I am"

"I do not like green eggs and ham"

Unigram Counts $C(w)$

Word	Count	Word	Count	Word	Count
I	3	like	1	and	1
am	2	green	1	ham	1
Sam	2	eggs	1	SUM	14
do	1	not	1		

Bigram Counts $C(w_{n-1}, w_n)$

Bigram	Count	Bigram	Count
(I, am)	2	(not, like)	1
(I, do)	1	(like, green)	1
(am, Sam)	1	(green, eggs)	1
(Sam, I)	1	(eggs, and)	1
(do, not)	1	(and, ham)	1

Bigram Probabilities

$$P(\text{am} | \text{I}) = C(\text{I}, \text{am}) / C(\text{I})$$

$$= 2 / 3 = 0.67$$

$$P(\text{do} | \text{I}) = C(\text{I}, \text{do}) / C(\text{I})$$

$$= 1 / 3 = 0.33$$

$$P(\text{Sam} | \text{am}) = C(\text{am}, \text{Sam}) / C(\text{am})$$

$$= 1 / 2 = 0.50$$

Sentence Probability

$$P(\text{"I am Sam"}) = P(\text{I}) \times P(\text{am}|\text{I}) \times P(\text{Sam}|\text{am})$$

$$= (3/14) \times (2/3) \times (1/2) = 0.071$$

Start & End Tokens

Why Special Tokens?

To help the model recognize sentence boundaries

Allows to calculate $P(\text{first word})$ and knows to stop generating

Special Tokens

<**s**>

Start of Sentence

</**s**>

End of Sentence

We can use: BOS/EOS, [START]/[END], [CLS]/[SEP]

Note for N-gram

Trigram need 2 <**s**> tokens at the beginning to calculate $P(w_1 | \langle s \rangle, \langle s \rangle)$

Example with Bigram

Original sentence:

"I love NLP"

With special tokens:

<**s**> I love NLP </**s**>

Bigram probability:

$$P(\text{sentence}) = P(I|\langle s \rangle) \times P(\text{love}|I) \times P(\text{NLP}|\text{love}) \times P(\langle s \rangle|\text{NLP})$$

Benefits

- Can calculate $P(\text{first word})$
- Model know when a sentence ends
- Necessary for text generation

Sparsity Problem

Core Problem

Many N-gram combinations do not appear in training data

→ Count = 0 → Probability = 0

Scale of the Problem

$V = 50,000$ words

$V^2 = 2.5$ billion possible pairs

$V^3 = 125$ trillion combinations!

Catastrophic Consequence

$P(\text{"The cat jumped"}) = \dots \times P(\text{jumped}|\text{cat}) \times \dots$

If "cat jumped" have not seen → $P = 0 \rightarrow \text{all sentences} = 0!$

Solutions: Smoothing

1. Add-k Smoothing (Laplace)

Thêm k vào mọi count

2. Backoff

Fallback to lower-order N-gram

3. Interpolation

Weighted average of N-grams

4. Kneser-Ney

State-of-the-art cho N-gram LMs

Idea

"Steal" probability mass from seen events

Redistribute to unseen events

Add-k Smoothing (Laplace)

Formula

$$P(w_n | w_{n-1}) = (C(w_{n-1}, w_n) + k) / (C(w_{n-1}) + k \times V)$$

k = smoothing parameter (Usually k=1 or k<1)

V = vocabulary size

Example: k = 1

Without smoothing:

$$P(\text{jumped}|\text{cat}) = 0/100 = 0$$

With Add-1 (V = 10,000):

$$P(\text{jumped}|\text{cat}) = (0+1)/(100+10000) = 0.0001$$

Advantages

- Simple to implement
- No zero probabilities
- Easy to understand

Disadvantages

- Add-1 is too aggressive with large V
- Moves too much probability mass
- Not good for text applications

Add-k & Best Practice

Use $k < 1$: $k=0.5$, $k=0.1$, $k=0.01$

Add-k → classification | LM → Interpolation

Backoff & Interpolation

Backoff

If N-gram does not exist → use (N-1)-gram

Algorithm:

```
if C(wn-2, wn-1, wn) > 0:  
    return P_trigram  
elif C(wn-1, wn) > 0:  
    return α × P_bigram  
else:  
    return β × P_unigram
```

Example

$P(\text{food} \mid \text{Chinese}, \text{eat})?$

If "Chinese eat food" have not seen → fallback to $P(\text{food} \mid \text{eat})$

Interpolation

Always mix all N-gram levels

Formula:

$$\hat{P}(w_n | w_{n-2}, w_{n-1}) = \\ \lambda_1 \times P(w_n | w_{n-2}, w_{n-1}) \\ + \lambda_2 \times P(w_n | w_{n-1}) \\ + \lambda_3 \times P(w_n)$$

Constraint: $\lambda_1 + \lambda_2 + \lambda_3 = 1$

Why Better?

- Always uses all available information
- No zero probabilities
- λ weights learned from held-out data

Typical λ values: $\lambda_1=0.6$, $\lambda_2=0.3$, $\lambda_3=0.1$

Autocomplete System

Autocomplete Pipeline

Input Processing

Get context (previous N-1 words)

Candidate Generation

Find all words following context

Probability Scoring

Calculate $P(\text{word}|\text{context})$ for each

Ranking & Return

Return top-k suggestions

Example: Trigram Autocomplete

User types:

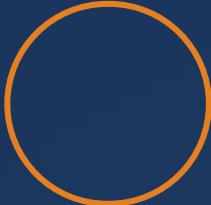
"I want to |"

Context = ("want", "to")

eat

go

buy



PART 2

Perplexity

Evaluating Language Models



Definition & Intuition



Mathematical Formula



Interpretation



What is Perplexity?

Definition

Perplexity (PP) is a metric for evaluating the quality of a language model. It measures how 'surprised' the model is when it encounters test data.

Formula

$$\begin{aligned} \text{PP}(W) &= P(w_1, w_2, \dots, w_n)^{-1/N} \\ &= \sqrt[N]{1 / P(W)} \end{aligned}$$

N = total number of words in test set

Log Form (numerical stability)

$$\text{PP} = \exp(-1/N \times \sum \log P(w_i | \text{context}))$$

Intuition: "Weighted Average Branching Factor"

Perplexity = The average number of words the model must 'choose' from at each position.

Model "confused" between 50 words

Model "confused" between 500 words

Lower is Better!

Low PP = Model predict better

PP = 1 → Perfect prediction (impossible)

Perplexity: Example

Test Sentence (Bigram)

"I want to eat"

N = 4 words (excluding <s>)

Bigram Probabilities für LM:

$$P(I | \text{<s>}) = 0.20$$

$$P(\text{want} | I) = 0.30$$

$$P(\text{to} | \text{want}) = 0.65$$

$$P(\text{eat} | \text{to}) = 0.10$$

Step-by-Step Calculation

Step 1: Joint probability

$$\begin{aligned} P(W) &= 0.20 \times 0.30 \times 0.65 \times 0.10 \\ &= 0.0039 \end{aligned}$$

Step 2: Inverse probability

$$1/P(W) = 1/0.0039 = 256.41$$

Step 3: N-th root (N=4)

$$\begin{aligned} PP &= \sqrt[4]{256.41} \\ &= 4.0 \checkmark \end{aligned}$$

Interpretation

PP = 4 means that the model "confused" around 4 words at each position.
→ It is good for a simple bigram model!

Log form: $\log PP = -1/N \times \sum \log P$ = more stable

PART 3

Word Embeddings

Distributional Semantics & Dense Representations

From Sparse to Dense

Semantic Similarity

Word Analogies

Distributional Semantics

Core Idea

"You shall know a word by the company it keeps"

— J.R. Firth, 1957

Distributional Hypothesis

Words with similar meanings will appear in similar contexts.

Context of "dog":

"The **dog** barks", "I pet my **dog**", "The **dog** ran"

Context of "cat":

"The **cat** meows", "I pet my **cat**", "The **cat** ran"

→ Similar contexts → Similar meanings!

Evolution of Word Representations

1. One-Hot Encoding

[0, 0, 1, 0, 0, ..., 0] - Sparse, no similarity

2. TF-IDF / Co-occurrence

Count-based, still high-dimensional

3. Dense Embeddings (Word2Vec, GloVe)

[0.2, -0.5, 0.8, ...] - Low-dim, captures semantics

4. Contextualized (BERT, GPT)

Different vector for each context!

Today: Word2Vec & GloVe (Static Embeddings)

Word Embedding Properties

Semantic Similarity

Similar words → Nearby vectors in space

Cosine Similarity:

$\text{sim}(\text{king}, \text{queen})$

= 0.85

$\text{sim}(\text{king}, \text{apple})$

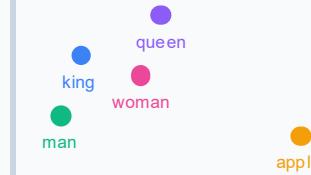
= 0.12

Word Analogies

Semantic relationships are encoded as vector operations!

$\text{king} - \text{man} + \text{woman} \approx \text{queen}$

Vector Space Visualization



Gender direction: $\text{man} \rightarrow \text{woman} \approx \text{king} \rightarrow \text{queen}$

Benefits of Word Embeddings ??

PART 4

Word2Vec

Neural Word Embeddings

Word2Vec: Introduction

What is Word2Vec?

Developed by **Mikolov et al. (2013)** at Google

A family of neural network models that learn dense word vectors from large unlabeled text corpora.

Core Idea

"Predict a word from its context, or predict context from a word"

→ The hidden layer weights become the word embeddings!

💡 Key insight: Self-supervised learning

No labeled data needed - the text itself provides supervision!

Two Architectures

1 CBOW

Continuous Bag Of Words

Context → Target word

"The cat sat on the ___" → "mat"

2 Skip-gram

Inverse of CBOW

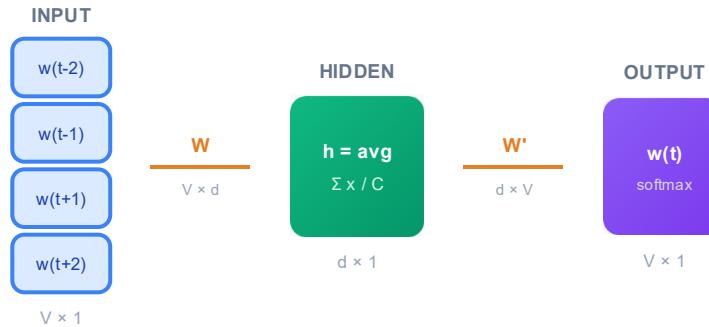
Target word → Context

"mat" → "cat", "sat", "on", "the"

Skip-gram: works better for rare words | CBOW: faster, better for frequent words

CBOW Architecture

Neural Network Architecture



V = vocabulary size (10K-1M) | d = embedding dimension (50-300) | C = context window size

Training Process

- 1 One-hot encode context words
- 2 Lookup embeddings from W
- 3 Average to get hidden state h
- 4 Softmax over $W' \times h$
- 5 Cross-entropy loss + backprop

Example

Context: "the cat ___ on mat"

Target: "sat"

Window = 2: [the, cat, on, mat] → sat

After training

Matrix W is word embeddings!

CBOW: Mathematical Formulation

Forward Pass

1. Hidden layer (average of context embeddings)

$$h = (1/C) \times \sum W T_{xc}$$

2. Output scores (dot product with all word vectors)

$$u = W^T \times h$$

3. Softmax probability

$$P(wt|context) = \exp(u_t) / \sum \exp(u_j)$$

Objective: Maximize Log-Likelihood

$$J = (1/T) \times \sum \log P(wt|wt-c, \dots, wt+c)$$

T = total training examples, c = context window size

Intuition: What does training optimize?

- Words appearing in similar contexts → **similar embeddings**
- High dot product with true target → increase probability
- Push apart vectors of unrelated words

Softmax Bottleneck

Tính softmax over V words rất expensive!

V = 1,000,000 → 1M operations per word!

→ Solution: Negative Sampling hoặc Hierarchical Softmax

Key Parameters

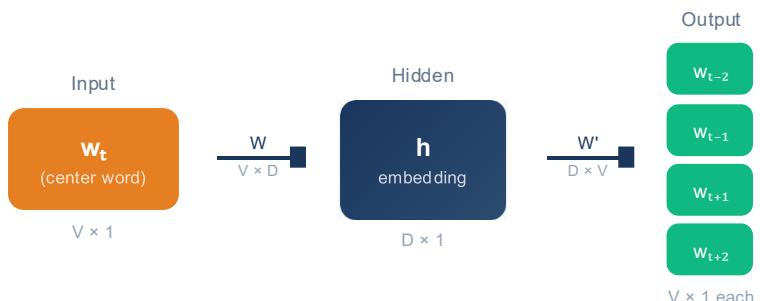
d: 50-300 dims

c: 2-10 window

lr: 0.025

Skip-gram Architecture

Architecture: Center → Context



Insight: Predict MULTIPLE context words from ONE center word. Better for rare words!

Example (window = 2)

Sentence: "The quick brown fox jumps"

Input: "brown" → Predict: The, quick, fox, jumps

Training Objective

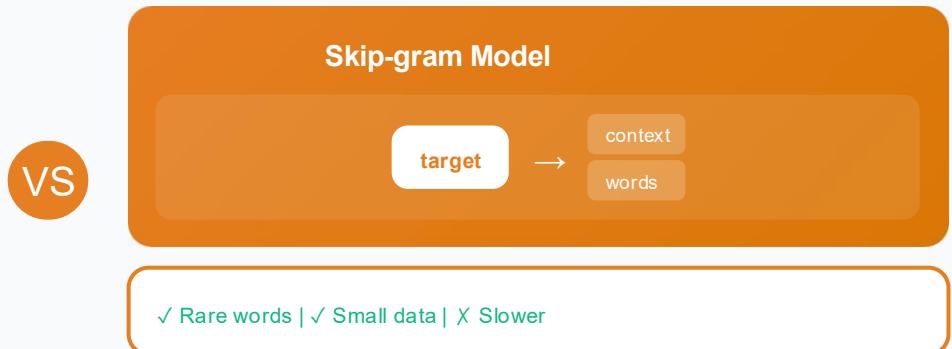
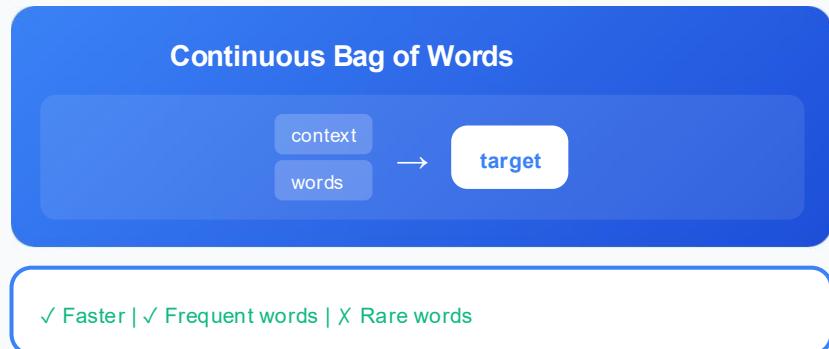
$$\max \sum_t \sum_{-c \leq j \leq c, j \neq 0} \log P(w_{t+j}|w_t)$$

Maximize probability of context words given center word over all positions

Skip-gram Strengths:

- Works well with small data
- Better for rare words
- More training examples per word

CBOW vs Skip-gram



CBOW: large corpus | Skip-gram: small corpus + rare words

The Softmax Challenge

Standard Softmax

$$P(w|w) = \exp(u \cdot v)$$

$$\sum \exp(u \cdot v)$$

u_o
output vector

v_c
center vector

The Bottleneck

Denominator requires summing over **ALL vocabulary words**:

Solutions to Speed Up

Negative Sampling

Instead of all V words, sample k negative words ($k = 5-20$). Binary classification instead of V -way softmax.

Hierarchical Softmax

Build binary tree with words at leaves. Path length = $O(\log V)$ instead of $O(V)$.

Negative Sampling

Idea

Instead of predicting probability over entire vocabulary, convert to **binary classification**:



Positive

Real context pairs



Negative

Random "fake" pairs

Example: "The cat sat on the mat"

Objective Function

$$J = \log \sigma(u_o^T v_c) + \sum_{k=1}^K \log \sigma(-u_k^T v_c)$$

Negative Sampling Distribution

Sample negatives based on word frequency (raised to 3/4 power):

$$P(w) = f(w) / \sum f(w)$$

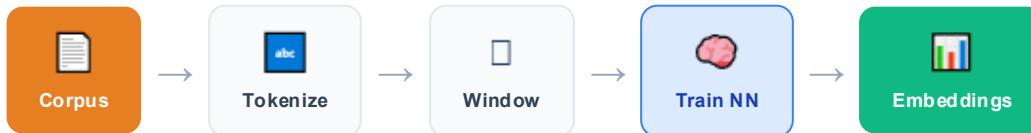
The 3/4 power smooths the distribution — less bias towards very frequent words



Hyperparameters

Benefit: Train on ~15 binary classifiers instead of 100,000-way softmax!

Training Process



Prepare Corpus

Tokenize, lowercase, remove rare words (freq < 5)

Generate Training Pairs

Slide window over text, create (center, context) pairs

Initialize Embeddings

Random init $W (V \times D)$ and $W' (D \times V)$

Forward Pass

Compute predictions using current weights

Compute Loss & Gradients

Neg sampling loss, backprop gradients

Update & Repeat

SGD update, repeat for millions of pairs

Final Output

Word Embedding Matrix W :

```
[[0.12, -0.34, ...], # "the"  
 [0.45, 0.23, ...], # "cat"  
 ...]
```

Shape: $(\text{vocab_size} \times \text{embed_dim})$

Use W or $\text{average}(W, W')$

PART 5

GloVe

Global Vectors for Word Representation

Co-occurrence Matrix

Weighted Least Squares

Global Statistics

GloVe: Introduction

What is GloVe?

Global Vectors — an unsupervised learning algorithm that combines:

Global Statistics

Matrix factorization
(like LSA)



Local Context

Window-based learning
(like Word2Vec)

Insight

Word meanings are encoded in **co-occurrence statistics**. The ratio of co-occurrence probabilities captures semantic relationships:

$$P(\text{solid} \mid \text{ice}) / P(\text{solid} \mid \text{steam}) >> 1$$

This ratio reveals "ice" is related to "solid"

Word2Vec vs GloVe Approach

Word2Vec (Predictive)

Learns embeddings by predicting context words. Slides window over corpus → many training pairs.

GloVe (Count-based)

First collects global co-occurrence counts, then learns embeddings to fit this matrix.

GloVe Advantages

- Faster training (process matrix once)
- Uses global corpus statistics
- Better on word analogy tasks
- Parallelizable training

Co-occurrence Matrix

Definition

X = number of times word j appears in context of word i

Corpus: "I like deep learning. I like NLP."

X matrix (symmetric):

	I	like	deep	learning	NLP
I	0	2	0	0	0
like	2	0	1	0	1
deep	0	1	0	1	0
learning	0	0	1	0	0
NLP	0	1	0	0	0

$$X[I, \text{like}] = 2$$

$$X[\text{like}, \text{deep}] = 1$$

Matrix is symmetric: $X_{ij} = X_{ji}$ | Window size affects counts

Key Formulas

Total co-occurrence for word i :

$$X_i = \sum_k X_{ik}$$

Co-occurrence probability:

$$P_{ij} = P(j|i) = X_{ij} / X_i$$

Context Window

Words within k positions of target word count as context (typically $k = 10$). Closer words often weighted more heavily.

Storage

Dense

$$\mathbf{V} \times \mathbf{V}$$

Sparse

~1% non-zero

GloVe Objective Function

GloVe Loss Function

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

w_i
word vec

~w_j
context vec

b_i, ~b_j
biases

f(X_{ij})
weighting

Intuition

We want the dot product of word vectors to approximate the log of co-occurrence:

$$w_i^T \tilde{w}_j \approx \log(X_{ij})$$

Similar words → similar contexts → similar embeddings

Weighting Function f(x)

Give less weight to very frequent co-occurrences (avoid domination by common words):

$$f(x) = \{ (x/x_{\max})^\alpha \text{ if } x < x_{\max}, \text{ else } 1 \}$$

Final embedding: $w + \tilde{w}$ (sum of both vectors)

PART 6

Lab & Assignments

Practice & Final Project Release

Lab: N-gram Language Model

Lab Tasks

Build N-gram Counts

Tokenize corpus, count unigrams, bigrams, trigrams

Implement Smoothing

Add-k smoothing, interpolation

Calculate Perplexity

Evaluate model on test set

Build Autocomplete

Predict next word, rank by probability

Key Takeaways

Language Models

N-gram models for word sequences with smoothing techniques (Add-k, Interpolation)

Perplexity

LM evaluation metric. Lower = better model

Word Embeddings

Dense vectors capturing semantic similarity

Word2Vec

CBOW (context→word) and Skip-gram (word→context) with Negative Sampling

GloVe

Count-based approach using co-occurrence matrix

Final Assignment (50%)

Choose: Chatbot, MT, Summarization, or QA. Present in Session 10.