

Building Large Language Model Applications

Advanced NLP Techniques: N-grams and Word Embeddings

Hamza Farooq
Dr. Saima Hassan





Recap of Lecture 3

Language representation

- **One-Hot Encoding** – Represents each word as a unique binary vector, ignoring word relationships.
- **Bag-of-Words** – Represents text as word frequency counts, disregarding word order.
- **TF-IDF** – Measures word importance by balancing frequency with document uniqueness.

Limitations:

- **TF-IDF** treats words independently, ignoring the sequence or relationships between words.
- Cannot capture context or meaning in phrases.



Learning outcomes

- N-gram
- Word Embedding
 - Word2Vec
 - CBOW
 - Skip-gram
 - Glove
- Conclusion



Next Word Prediction

What is Next Word Prediction?

Next word prediction is a fundamental concept in NLP where **a model predicts the most probable word following a given sequence of words**

How It Works:

- The model analyzes **previous words** in a sentence.
- It predicts the **next word** based on learned patterns from large text corpora.
- Techniques include **N-grams, Word2Vec, Neural Networks**, and **Transformer-based models (GPT, BERT, etc.)**.

Example: The weather is very...

Prediction:

1. "nice" 
2. "cold" 
3. "hot" 

Real-world applications: Search engines, text messaging, AI-powered writing assistants.



Human Word Prediction

How Do Humans Predict Words?

Humans have the ability to imagine the next word in a sentence based on various types of knowledge:

Domain Knowledge – Understanding common phrases in specific fields.

- **Example:** "Red blood ___" → *cells, count, pressure*

Syntactic Knowledge – Recognizing grammatical patterns.

- **Example:** "The ___" → *adj (beautiful), noun (dress/location)*

Lexical Knowledge – Knowing which words frequently appear together.

- **Example:** "Baked ___" → *potato, bread, goods*



Applications

Predicting words based on context is essential for various applications:

Autocorrect & Text Suggestions

*"He will **recieve**/**receive** the email soon."*

Speech Recognition

"I ate a cherry" is a more likely sentence than "Eye eight uh Jerry"

Handwriting Recognition

*"Order 3 more **bottles**/**battles** of water."*

Machine Translation

- **English:** "She has a big heart."
- **French:** "Elle a un grand cœur." (instead of "cardiaque")



Probabilistic Language Models

Assign a probability to a sentence

- **Machine Translation:**

$P(\text{high winds tonight}) > P(\text{large winds tonight})$

- **Spell Correction**

The office is about fifteen minuets from my house

$P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$

- **Speech Recognition**

$P(\text{I saw a van}) > P(\text{eyes awe of an})$

Summarization, question-answering, etc., etc.!!



Probabilistic Language Models

Goal: Compute the probability of a sentence or sequence of words:

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

Probability of an upcoming word:

$$P(w_n \mid w_1, w_2, \dots w_{n-1})$$

What is a Language Model?

A model that estimates:

$P(W) \rightarrow$ Probability of a full sentence

$P(w_n \mid w_1, w_2, \dots w_{n-1}) \rightarrow$ Probability of the next word



How to Compute $P(W)$?

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

To compute the probability of a sequence of words $P(W)$, we use the **Chain Rule of Probability**.

Chain Rule of Probability:

Recall the definition of **conditional probability**:

$$P(A,B) = P(A) P(B|A)$$

More variables:

$$P(A,B,C,D) = P(A) P(B|A) P(C|A,B) P(D|A,B,C)$$

The Chain Rule in General Form:

For a sequence of words $(w_1, w_2, w_3, \dots, w_n)$, we expand using the Chain Rule:

$$P(w_1, w_2, w_3, w_4, w_5 \dots w_n) = P(w_1) P(w_2|w_1) P(w_3|w_1, w_2) \dots P(w_n|w_1, w_2, \dots, w_{n-1})$$



Chain Rule Application

Applying the Chain Rule to compute joint probability of words in sentence:

$$P(w_1, w_2, w_3, w_4, w_5 \dots w_n) = P(w_n | w_1, w_2, \dots, w_{n-1}) =$$

$$= \prod_{k=1}^n P(w_k | w_1, w_2, \dots, w_{k-1})$$

The joint probability of an entire sequence of words can be estimated by multiplying together a number of conditional probabilities.



Markov Assumption

The intuition of the n -gram model is that instead of computing the probability of a word given its entire history, we can approximate the history by just the last few words.

$P(\text{delicious} | \text{The cake with chocolate frosting looks absolutely})$

Bigram model

$P(\text{delicious} | \text{absolutely})$

When using a bigram model to predict the conditional probability of the next word, make the following approximation:

$$P(w_n | w_{n-1})$$

Generalize the bigram to the trigram and then to the n -gram



Language Models

A language model is a machine learning model **LM** that predicts upcoming words.

A **LM** assigns a **probability to each possible next word**.

What is an N-gram?

N-gram is the simplest kind of **LM**

Or **A sequence of n words used in language modeling**

Types of N-Grams:

- 📌 **Unigram (1-gram):** "Learning"
- 📌 **Bigram (2-gram):** "Machine learning"
- 📌 **Trigram (3-gram):** "Deep learning models"
- 📌 **4-gram:** "Artificial intelligence is evolving"



N-gram Example

"She enjoys drinking hot coffee."

Bigram Representation:

("She enjoys"), ("enjoys drinking"), ("drinking hot"), ("hot coffee")

Trigram Representation:

("She enjoys drinking"), ("enjoys drinking hot"), ("drinking hot coffee")

N-grams: Examples



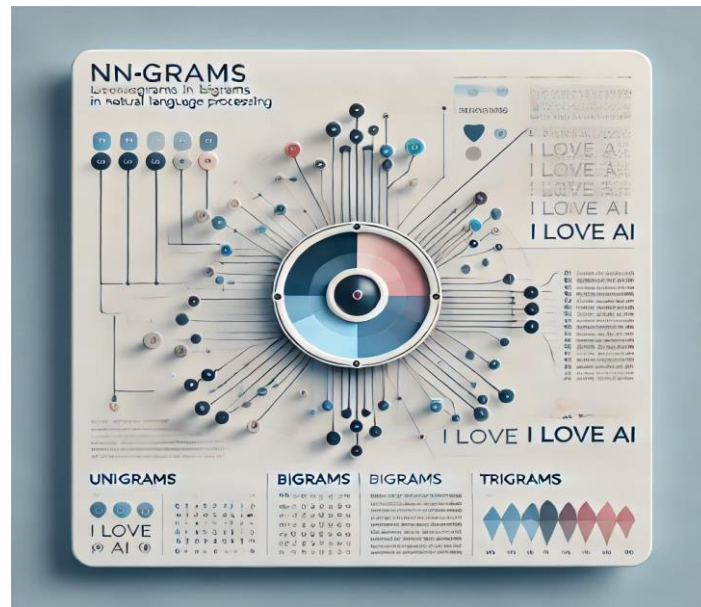
"In language modeling, language modeling is essential."

N-Gram Breakdown:

Unigram (N=1): ['In', 'language', 'modeling', 'language', 'modeling', 'is', 'essential']

Bigram (N=2): ['In language', 'language modeling', 'modeling language', 'language modeling', 'modeling is', 'is essential']

Trigram (N=3): ['In language modeling', 'language modeling language', 'modeling language modeling', 'language modeling is', 'modeling is essential']





Estimating Probabilities

N-gram conditional probabilities can be estimated from raw text based on the ***relative frequency*** of word sequences.

Bigram:

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

N-gram:

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}w_n)}{C(w_{n-N+1}^{n-1})}$$

To have a consistent probabilistic model, append a unique start (<s>) and end (</s>) symbol to every sentence and treat these as additional words.



N-grams: Probability Calculation Example

<S> I am Sam </S>

<S> Sam I am </S>

<S> I do not like green eggs and jam </S>

Some of the bigram probabilities from this corpus can be as:

$$P(I|<S>) = 2/3 = 0.67$$

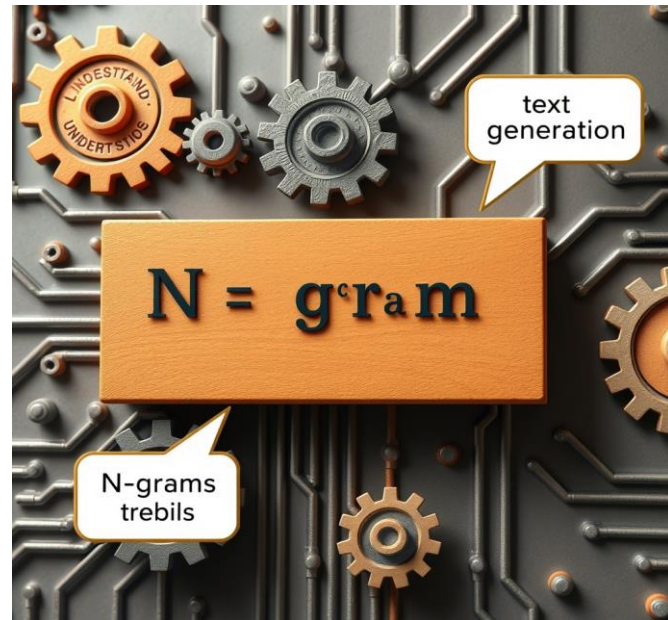
$$P(\text{Sam}|<S>) = 1/3 = 0.33$$

$$P(\text{am}|I) = 2/3 = 0.67$$

$$P(</S>|\text{Sam}) = 1/2 = 0.5$$

$$P(\text{Sam}|\text{am}) = 1/2 = 0.5$$

$$P(\text{do}|I) = 1/3 = 0.33$$





N-grams: Limitations

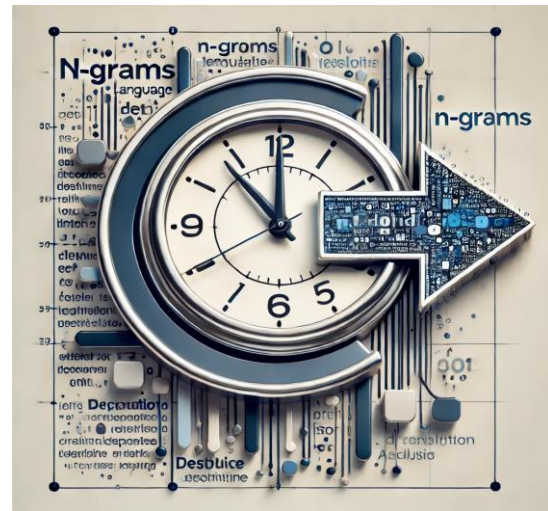
Lacks Long-Range Context: N-grams only consider neighboring words, ignoring distant dependencies.

Data Sparsity: Higher-order N-grams often have insufficient data for accurate probabilities.

Exponential Growth: Vocabulary size increases rapidly with larger N, requiring more storage and computation.

No Semantic Understanding: N-grams rely on word patterns without understanding meaning.

Fixed Window Size: Important context outside the N-gram window is ignored.





Word Embeddings



From One-Hot Encoding to Continuous Representation

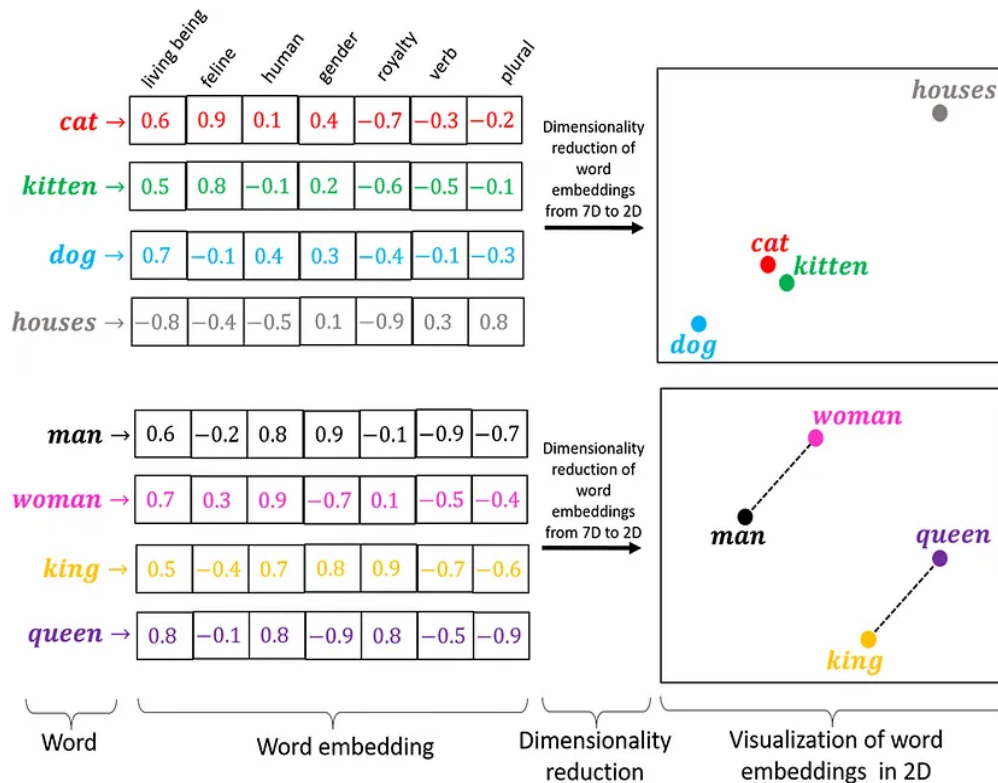
Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$



Word Embedding

What is Word Embedding?

- Word Embedding is a technique for **representing words and documents in a numerical format**.
- It transforms **words into low-dimensional vectors** that capture their **meaning and relationships**.
- **Words with similar meanings have similar vector representations**, making them useful for AI models.





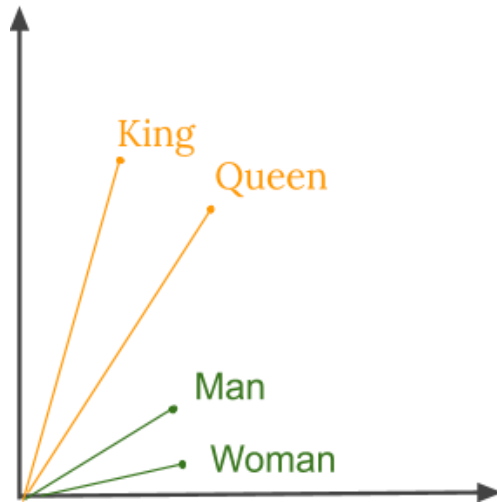
Word Embedding

- **Why Use Word Embeddings?**

- Reduce Dimensionality
- Overcomes limitations of one-hot encoding and TF-IDF.
- Enables understanding of relationships like:
 - "king - man + woman = queen"
 - Synonyms and similar words are closer in vector space.
- Efficient for NLP tasks like classification, translation, and more.

- **Applications:**

- Sentiment Analysis, Question Answering, Machine Translation





Word2Vec

Word2vec is a technique in NLP for obtaining **vector representations of words**. These vectors capture information about the **meaning of the word based on the surrounding words**. The word2vec algorithm estimates these representations by modeling text in a **large corpus**. Once trained, such a model can detect synonymous words or suggest additional words for a partial sentence. Word2vec was developed by **Tomáš Mikolov and colleagues at Google and published in 2013**

Word2vec is a two-layer network where there is input one hidden layer and output



Word2Vec

Word2Vec –learn through training on a large text corpus. The features are learned automatically based on the word's context in the corpus.

1. Each word starts as a one-hot vector

For a vocabulary of 10k words, Each word is represented as a **10k-dimensional one-hot vector**

2. Hidden Layer (Embedding Layer)

The one-hot vector is multiplied by a weight matrix (**size: $V \times N$** where V = vocabulary size, N = embedding size).

The resulting output is a **low-dimensional dense vector** (word embedding).

The values in this vector act as the **features of the word**.

3. Features are learned through training

Word2Vec uses either **CBOW** or **Skip-gram**.

The model adjusts the weight matrix based on how words **co-occur** in a given context.



Word2Vec

Example:

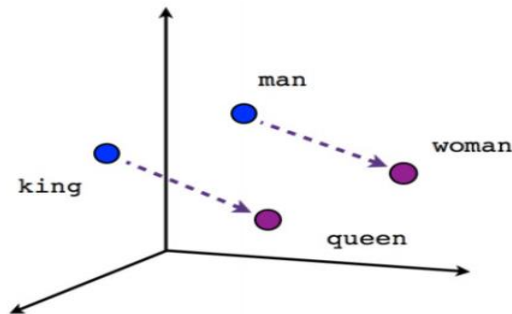
Let's say the word "king" gets converted to a 100-dimensional vector:

king=[0.21,-0.12,0.75,...,0.34]

These values represent different hidden features such as:

Gender, Royalty, Age ...

If we compare this with "queen", their vectors would be similar but differ in specific features (e.g., gender).





Types of Word Embeddings

1. **Word2Vec** – A neural network-based model that learns word representations by analyzing context.
 - a. **Skip-gram**: Predicts surrounding words given a target word.
 - b. **CBOW (Continuous Bag of Words)**: Predicts a target word based on surrounding words.

2. GloVe (Global Vectors for Word Representation)

Based on word co-occurrence statistics.

Captures both global (document-wide) and local (sentence-level) relationships.

Generates word vectors by analyzing word frequency patterns in large corpora.



Word2Vec: Continuous Bag of Words (CBOW)

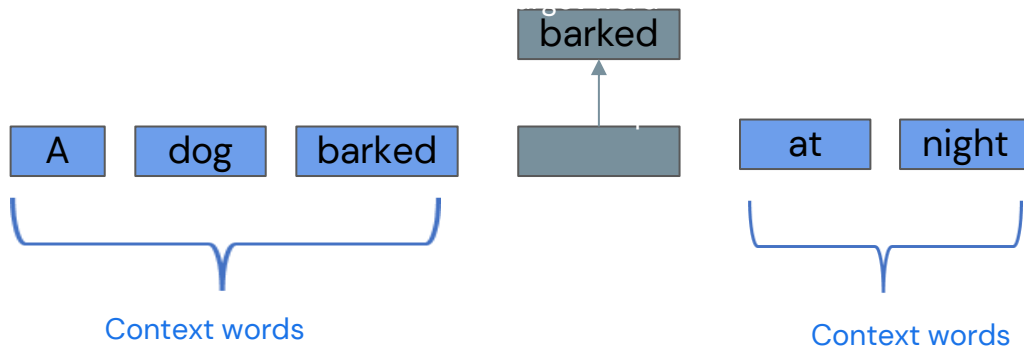
Predict a **target word** w_t from its surrounding **context words**.

Efficient for frequent words.

Example:

- Sentence: A dog barked loudly at night.
 - **Context words:** "dog," "barked," "at," "night."
 - **Predict:** "loudly."

Mathematical Objective: Maximize: $P(w_t | w_{t-n}, \dots, w_{t+n})$





Word2Vec: Skip-gram

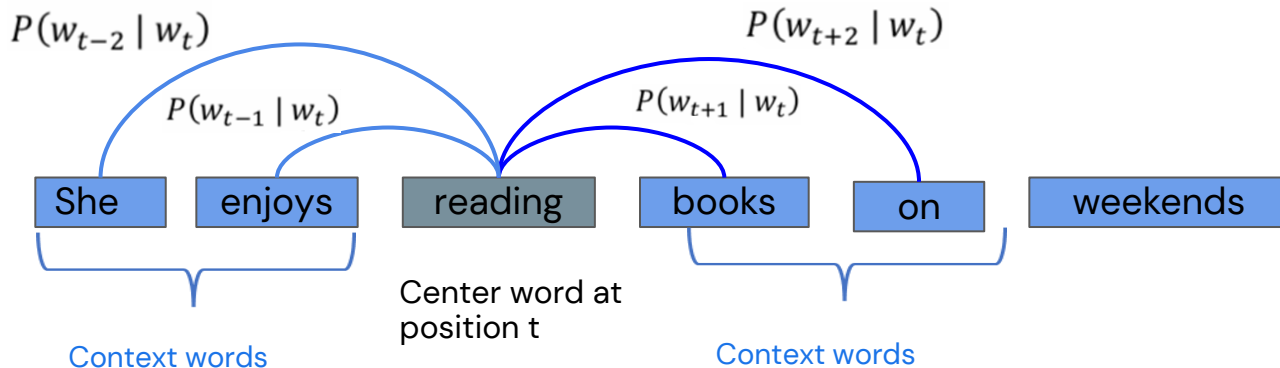
Predicts the surrounding **context** given a **target** word.

Example:

- Sentence: She enjoys reading books on weekends.
 - **Center word:** "reading."
 - **Context words for a window size 2:** "She", "enjoy", "books", "on"

Mathematical Objective:

$$\text{Maximize: } \prod_{t=1}^T \prod_{-n \leq j \leq n, j \neq 0} P(w_{t+j} | w_t)$$





CBOW vs Skip-gram

Feature	Skip-gram	CBOW
Input	Target word	Context words
Output	Context words	Target word
Speed	Slower	Faster
Strength	Rare words	Frequent words



Pros of Word2Vec

Captures Word Relationships – Words with similar meanings have similar vector representations.

Low-Dimensional Representation – Uses dense vectors instead of large, sparse one-hot vectors, making it memory-efficient.

Self-Supervised Learning – Does not require labeled data; learns patterns from raw text, making data collection easy.



Cons of Word2Vec

Does Not Preserve Global Information – Word2Vec focuses on local word relationships but does not capture overall document structure.

Solved by GloVe

Limited for Morphologically Rich Languages – Struggles with languages where words have multiple inflected forms (e.g., Arabic, Turkish, Finnish).

Solved by FastText

Lacks Broad Context Awareness – Embeddings are static, meaning the same word has the same representation regardless of context (e.g., "bank" in "river bank" vs. "money bank").

Solved by GPT, BERT, LSTM

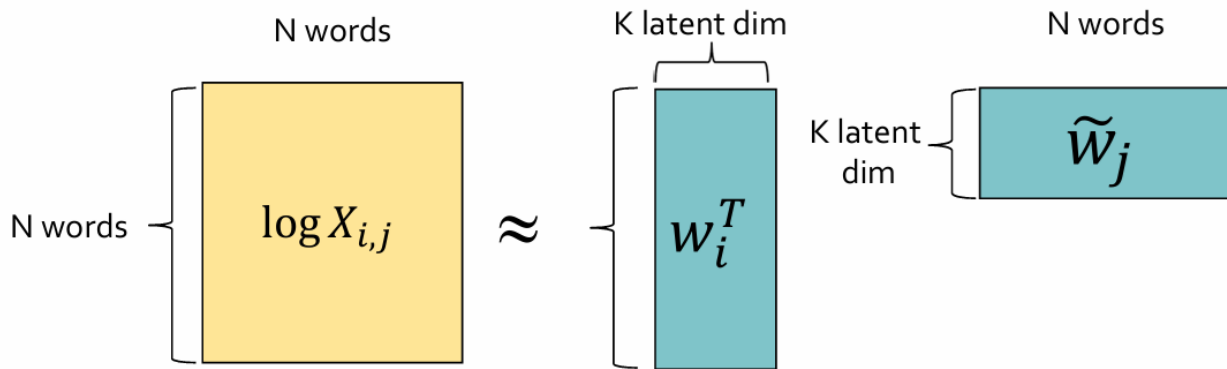


GloVe: Global Vectors for Word Representation

While **word2Vec** is a predictive model — learning vectors to improve the predictive ability, **GloVe** is a **count-based model**, Pennington et al., 2014

Count-based models learn vectors by doing dimensionality reduction on a co-occurrence counts matrix

- Factorize this matrix to yield a lower-dimensional matrix of words and features, where each row yields a vector representation for each word





GloVe: Example

GloVe training starts by forming a co-occurrence matrix X where the ij -th entry is the number of times words on i -th row and j -th column appeared together,

Co-occurrence matrix for a window size of 1 when corpus consists of only the following three sentences

- I like deep learning.
- I like NLP.
- I enjoy flying.

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



The need for Context-Based Models

- To capture **contextual meanings** of words dynamically.
- To understand **sentence structure** and **long-term dependencies**.
- To handle **domain-specific aspects** and evolving meanings effectively.

The Solution

The limitations of Word2Vec led to the rise of advanced models like **BERT** and **GPT**, which offer **contextual embeddings** and excel at understanding language in depth.



Conclusion

Advances in Word Embeddings

- **Word2Vec:** Pioneered contextual word representation with its efficient models:
 - **CBOW:** Predicts a word from its context.
 - **Skip-gram:** Predicts context from a given word.
- **GloVe:** Enhanced embeddings by incorporating global statistical information.



Thank You