# Language Representation

- Context Free Embeddings

- Context Based Embeddings

By: Kanav Bansal
(That AI Guy)

# Language Representation

By: Kanav Bansal
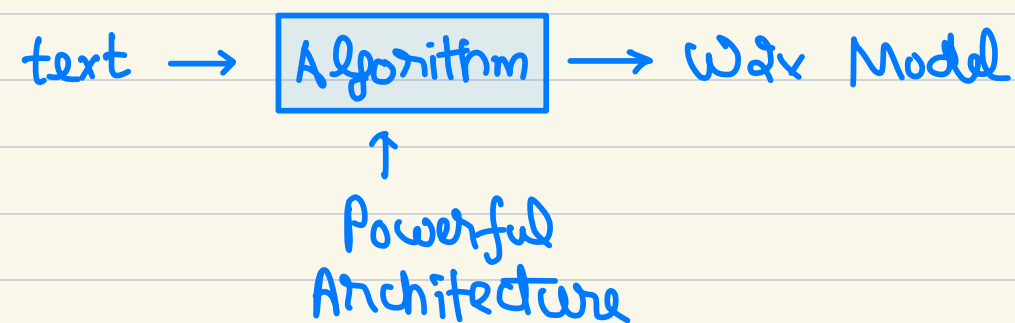(That AI Guy)

Various Text Representation Techniques:

1. Bag of words $\longrightarrow$ Captures frequency

2. TFIDF $\longrightarrow$ Captures Word Importance

3. W2v $\longrightarrow$ Captures Semantic Meaning $\longrightarrow$ This is done in a way that words which are
by learning the relationships similar will have vector representation
b/w different words. closer to each other.

text $\longrightarrow$ | Algorithm | $\longrightarrow$ W2v Model

$\uparrow$
Powerful
Architecture

Ques: What do you mean by Semantic Meaning?
Ans: Semantic meaning is about understanding
a deeper meaning & relationships b/w words.
Bow & TFIDF treat words as independent
entities, whereas W2v creates numerical
representations based upon the context in
which the word occured.

GloVe $\longrightarrow$ Based on Matrix Factorization
of the word Co-occurance matrix

FastText $\longrightarrow$ It extends W2v by representing
words as Character n-gram. for eg:

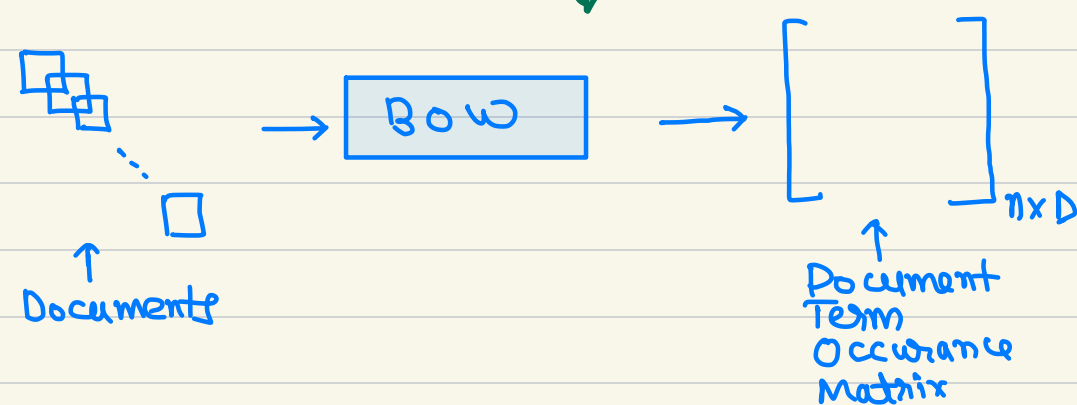'where' $\rightarrow$ '< wh', 'whe', 'her', 'ere', 're >'

This helps capture the Sub-words information
& hence handles OOV words effectively.

4. ELMO $\longrightarrow$ Captures contextual representation & word sequences

5. BERT $\longrightarrow$ Captures contextual representation & word sequences

By: Kanav Bansal
(That AI Guy)

# BoW Approach



↑
Documents

BoW

$n \times D$

↑
Document
Term
Occurance
Matrix

Problems:
1. High Dimensional & Sparsity
2. Sequence is not captured
3. Can't handle OOV words
4. Semantics is not captured

---

Distributed Representation:
Earlier approaches were high-dim & sparse.
In distributed representation we try to compress the
dimensionality which results in a compact &
dense vectors.
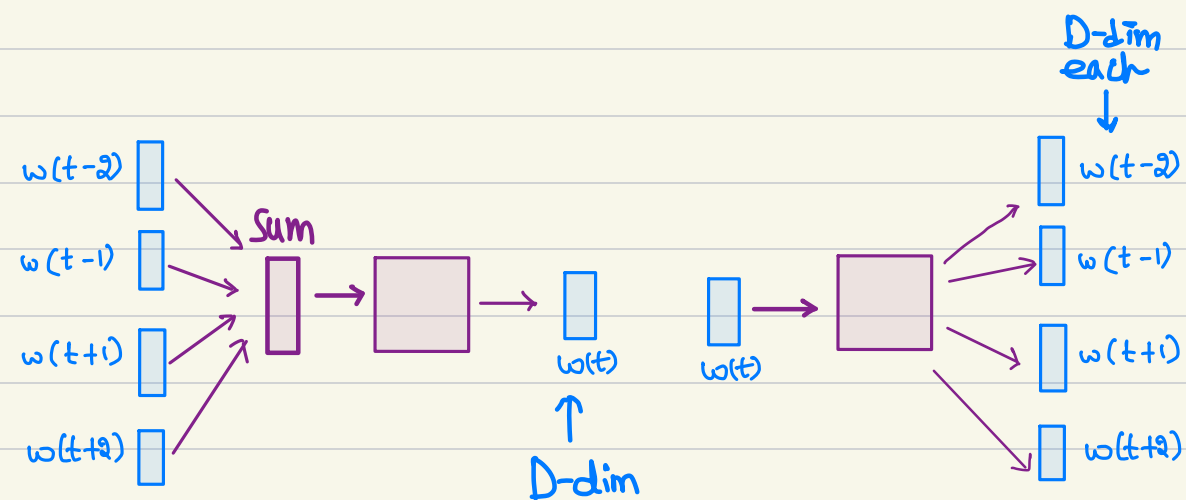
W2v is a distributed representation.

# W2V ↓

By: Kanav Bansal
(That AI Guy)

How does Word2Vec learns Embeddings?

W2V utilizes either of the following Architectures

1. CBOW

2. Skip Gram

3. Skip Gram + Negative Sampling

$doc_i \rightarrow \dots \quad w_5 \quad w_6 \quad w_7 \quad w_8 \quad w_9 \quad \text{-----}$

→ Training Data pairs for CBOW

| X | y |
|---|---|
| $w_5$ | $w_7$ |
| $w_6$ | $w_7$ |
| $w_8$ | $w_7$ |
| $w_9$ | $w_7$ |
| ⋮ | ⋮ |
| ↑ Context | ↑ Target |

D-dim each ↓

| w(t-2) | | | w(t-2) |
|---|---|---|---|
| w(t-1) | Sum | | w(t-1) |
| w(t+1) | → w(t) | w(t) → | w(t+1) |
| w(t+2) | ↑ D-dim | | w(t+2) |

### 1. CBOW
→ Not used in real time
→ Bad Representations for rare words
→ Overfits on frequent words
→ D-dim Classifier is not trivial

### 2. Skip Gram
→ Even more challenging than CBOW because with a single word now we have to predict the whole context.

→ Training Data pairs for Skip Gram

| X | y |
|---|---|
| $w_7$ | $w_5$ |
| $w_7$ | $w_6$ |
| $w_7$ | $w_8$ |
| $w_7$ | $w_9$ |
| ⋮ | ⋮ |
| ↑ Target | ↑ Context |

There is no negative sampling here.
We will always pass words within a context

By: Kanav Bansal
(That AI Guy)

# 3. Skip-Gram with Negative Sampling

## Using NN to compress the dimensionality ↓

Let's learn how to take each word's D-dim
one hot vector & compress it to d-dim



One-hot D-dim word vector

Skip-Gram with Negative Sampling
(Best till 2015)

if $w_i$ & $w_J$ co-occur

if $w_i$ & $w_J$ donot co-occur

A fundamental assumption here is:

If a Machine Learning architecture can start
predicting if two words co-occur or not,

then the d-dim representation we get must
have captured the semantic meaning of the
word.

* Negative Sampling:

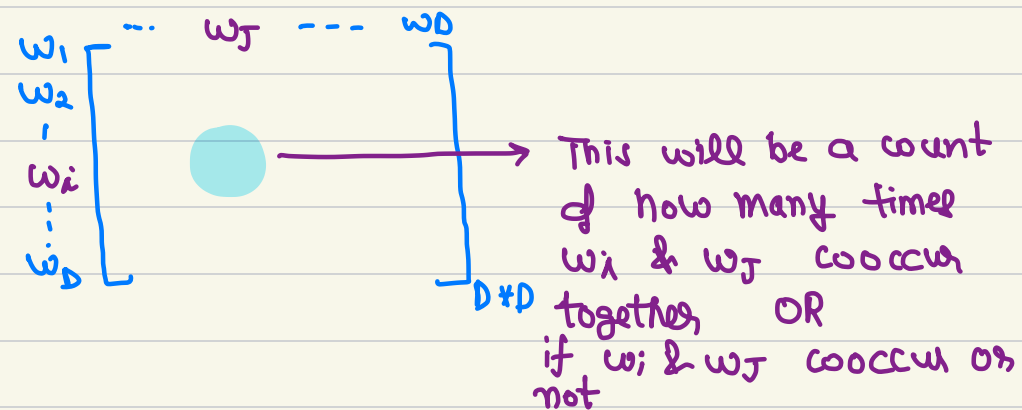$w_i$ & $w_J$ can be the words which are not
even in the context.

It's like teaching by showing both what's
right & wrong.

# GloVe 7

By: Kanav Bansal
(That AI Guy)

* Let's now have a look at Co-occurance Matrix

$$doc_i \rightarrow \quad \dots \quad w_5 \quad w_6 \quad w_7 \quad w_8 \quad w_9 \quad \text{----}$$

- Observe that $w_7$ cooccurs with $w_5, w_6,$ $w_8$ & $w_9$. (consider a window of size 5)

- Now we can create a cooccurance matrix instead of document term matrix.

- This won't preserve the sequence information but some context of each word will be preserved.

$$\begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_i \\ \vdots \\ w_D \end{array} \begin{array}{c} \dots \quad w_J \quad \text{---} \quad w_D \end{array}$$

$D * D$

This will be a count of how many times $w_i$ & $w_J$ cooccur together, OR if $w_i$ & $w_J$ cooccur or not

Observe that the above co-occurance matrix is very large.

# Distributed VS Contextual Representation

By: Kanav Bansal
(ThatAIGuy)

## Distributed Representation ⟶

Dense & Low Dimensional Representation of text

## Contextual Representation ⟶

Each word can have multiple meanings based upon the context in which it is used.

Ex:

I went to a river bank for a bath. Later I went to a bank to discuss regarding my loan application.

# Example of Contextual Representation & Long term Dependencies

Anaconda ▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓

▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓ 3.

▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓

4. ▓▓▓▓▓▓▓▓ Amazon.

Anaconda ▓▓▓▓▓▓▓▓▓▓▓▓

▓▓▓▓▓▓▓▓▓▓ Python. 2.

1.

Can you tell the meaning of the words from above paragraph?

Example of Contextual Representation &
Long term Dependencies

Anaconda can grow to be more than
30 feet & weigh over 500 pounds. 3.

They live in swamps in the rain—
4. —forests of the Amazon.

Anaconda is also a popular open

Source distribution of Python. 2.

1.

Can you tell the meaning
of the words from above
paragraph?

w2v / Glove / fastText → Numerical Representation
captures Semantics
ie Relationships b/w
words

Note: Every word will have one numerical
representation.

* **IDEAS:**

1. One word can have multiple
*Contextual Representation* meanings. This meaning can
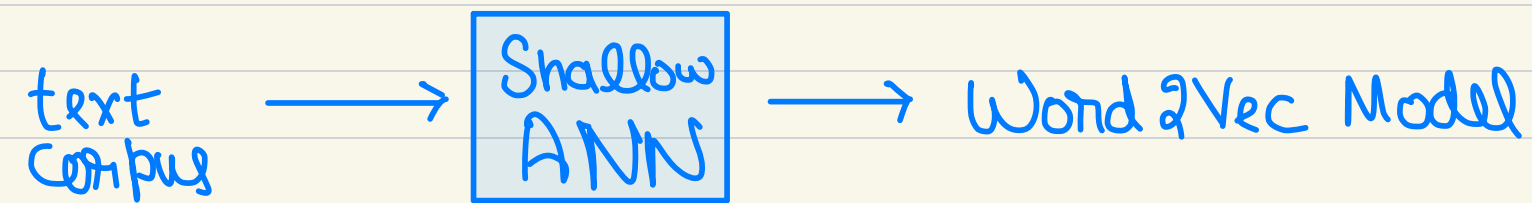be represented by looking at
the near-by words of context.

2. In order to understand the meaning
*Attention* we don't need to look at all the
surrounding words. Just by focusing
on few important words, we can
understand more about each word.

* Attention mechanism plays a vital
role by learning context-Based
Language Representations.

* Attention mechanism also help with
long term dependencies.

# Context Free Embeddings Summary 1

By: Kanav Bansal
(That AI Guy)

text Corpus → [ Shallow ANN ] → Word2Vec Model

How to train?
- CBOW
- SkipGram
- SkipGram with Negative Sampling

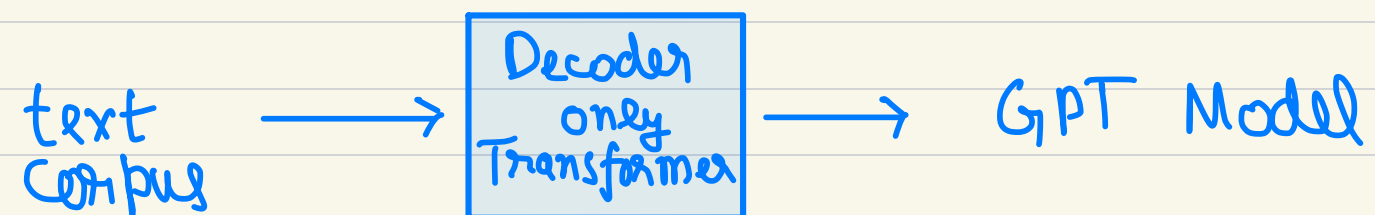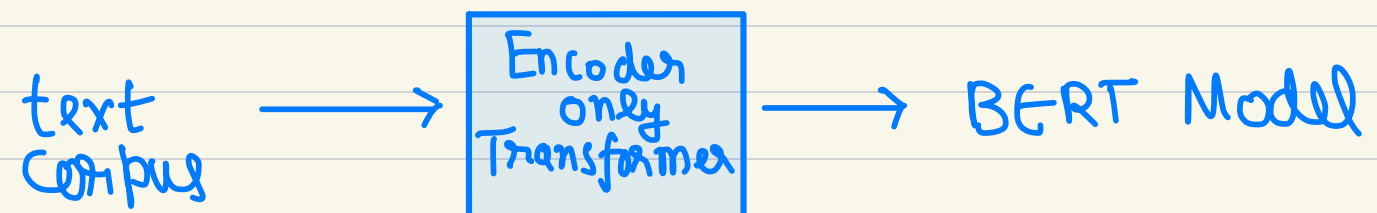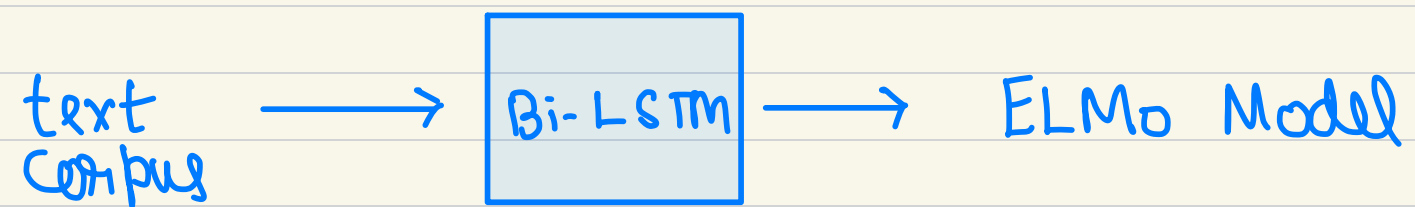→ Helps in mapping the problem statement to classification task.

Pro:
- Dense & Low Dim Vectors
- Captures word Semantics
- FastText can capture sub-level word info & hence handling OOV words

Cons:
- Sequence info is not captured
- Contextual Representations are not captured (ie Context Free Embeddings OR Static Embeddings)

# Context Based Embeddings Summary

By: Kanav Bansal
(ThatAI Guy)

text Corpus $\longrightarrow$ | Bi-LSTM | $\longrightarrow$ ELMo Model

text Corpus $\longrightarrow$ | Encoder only Transformer | $\longrightarrow$ BERT Model

text Corpus $\longrightarrow$ | Decoder only Transformer | $\longrightarrow$ GPT Model

By: Kanav Bansal
(That AI Guy)

**Basic Vectorization**
BoW    TFIDF

**Distributed Representation**
Word2Vec/GloVe/FastText

**Universal Lang. Rep.**
ELMo

BoW → Word Occurence

TFIDF → word importance

Word2Vec/GloVe/FastText → word's semantic

ELMo → Word's sequence + Contextual Embedding

**BERT**

**Universal Language Representation**
+ Bi-directional Contextual embeddings
+ Parallel Training & Scalable
+ Takes care of Long-Term dependencies
↑ Using Attention Mechanism

| Vectorization Techniques | BoW | TfIDF | W2V | ELMo | BERT |
|---|---|---|---|---|---|
| 1. Dimensionality | ↑ | ↑ | ↓ | ↓ | ↓ |
| 2. Sparsity | ↑ | ↑ | ↓ | ↓ | ↓ |
| 3. Semantic | ✗ | ✗ | ✓ | ✓ | ✓ |
| 4. Sequence info | ✗ | ✗ | ✗ | ✓ | ✓ |
| 5. OOV | ✗ | ✗ | ✗ | ✓ | ✓ |
| 6. Contextual Embeddings | ✗ | ✗ | ✗ | ✓ | ✓ |

Lots of Text Data
+
Approach: CBow/SkipGram
+
Architecture: Shallow ANN
• No seq info
• No contextual info
• OOV can be handled by FastText

Lots of Text Data
+
Approach: Language Model
+
Architecture: Bi-LSTM
• Slow to train
• Long-Term Dependencies

Lots of Text Data
+
Approach: Masked LM & NSP
+
Architecture: Encoders Only (Transformers)