Statistical based approaches treated words as atomic symbols:

$\boxed{\text{Animal}}$  $\boxed{\text{Dog}}$  $\boxed{\text{Trick}}$

Or in vector space:

$[0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ \ldots\ 0]$

- Also known as 'one hot' representation.

|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | Sparse representation |
| 0 | 0 | 1 | 0 | capture each word; but very |
| 0 | 0 | 0 | 1 | expensive. |

BoW → term frequency  ⎤ always sparse

Tf-IDF  ⎦ can't understand context

Animal: $[0001000 00..] = V_1$

Dog: $[0000001000] = V_2$

Truck: $[000010000...] = V_3$

* Distributional Representation      you shall know a word by
  diff meaning  in  diff  context.    Company it keeps: John Rupert.
  → One  of  the  most  successful  ideas  of  modern  NLP.
  → Linguistic units with similar distributions  have
    similar meaning.
  → A  bank is a financial institution  that  accepts
    deposits from  the  public & creates credit.
  → The  northern  bank  of  the  river is flooded.

* <u>Co-occurrence matrix</u>

  → A  co-occurrence matrix is a (terms × terms) matrix
    which captures the number of times  a  term  appears
    in the  context of  another term.
  → The  context is defined as  a  window  of <u>K words</u>
    around  the  terms

  Text Corpus   (like n-gram)  capturing  in form of  a matrix
  1. "ChatGPT is a language model."
  2. "Language models like ChatGPT are powerful"
  3. "The language model ChatGPT is based on GPT-3.5
     arch."

       window  size = 1    (Take  1 NN)  Before & After.
  can  capture context  (too many o's)

**\* Distributed Representations**

→ Compact, dense, low-dimensional and real-valued representations

→ Also, known as word embeddings each single component of the vector representation does not have any meaning of its own.

→ The interpretable features are hidden & distributed among uninterpretable vector components.
(Each word as a dense represent$^n$)

$$animals = \begin{bmatrix} 0.286 \\ 0.792 \\ -0.112 \\ -0.143 \\ 0.341 \\ 0.512 \end{bmatrix}$$

each not interpretable

are they good repr of desired word

man ≡ woman (word embeddings)

what exactly some colors mean.

( closeness of two words)

**\* Word Embeddings** (feed with neural network)

Bert, Glove

→ They are numerical representations of words in a continuous vector space, where each word is mapped to a dense vector of real numbers.

→ These vectors capture the semantic & contextual meaning of words based on their usage in a large corpus of a text.

→ Transformer based: GPT

color coding of vectors help us interpret the similarity b/w 2 words.

* <u>Analogies</u>

→ You can add & subtract word embeddings and see the concept of analogies.

king - man + woman ≅ queen.

model. most_similar (positive = ["king", "woman"], negative = ("man"))

<u>T-SNE visualization</u>

word embeddings

Gave all travel documents to a model

50 dim $\xrightarrow[\text{reduction}]{\text{dim}}$ 2 dim              most words related
                                                                                    to

each word has 50 features              (city)                              (travel)

                                                              (feeling)

info captured for a word              (relative)                        (food)

* <u>vector Embedding of words</u>

1. <u>latent Semantic Analysis/ Indexing</u> (1988)
   • Term weighing based model.
   • Consider occurrences of terms at document level.

2. Word2vec (2013)
   • Prediction-based model.
   • Consider occurrences of terms at context level.

**\* Prediction Based Models - word2vec**

→ Word2vec is a neural network-based model created by Google researchers in 2013 to learn word embeddings (ie; vector representations of words)

→ These embeddings capture semantic relationships b/w words, enabling similar words to have similar vector representations in a high-dimensional space.

→ In 2013, proposed word2vec can be trained using two new models.

**\* CBOWM**

→ CBOWM focusses on guessing a word based on its context.

→ Predicting $n^{th}$ word given the context words.

Ex: The quick brown fox jumped over the lazy dog.
    ↓      ↓      ↓
   100...  010...  0010...

(network is learning that wherever there is quick fox it should be followed by fox).

→ the input (first) layer is represented by a one-hot encoded vector and consists of the context words surrounding the target word.

→ Training data: All n-word windows in the corpus.

→ The hidden (second) layer is where the word embeddings are learned and the size depends on the dimensionality.