→ NLP where Embeddings

  → Prereq: Processing the text data.

  [ _____ ] ~ Non-numeric Data

  → Math model ~ Numeric Data.

→ Tokenization : [        CORPUS ~ 1000 lines ]

  [ -ˌ-ˌˌˌˌˌ-ˌˌˌ ] ~ Parameters

→ Preprocessing :    STOPWORDS

[ → (The) propensity (of) action (has) consequences . ]

         Semantic      ~    propensity action
                                         consequences
Connecting sentence grammatically.

→ TF-IDF does not have correlation index of words

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 1.2 \\ 22.3 & 0 & 0 & 4.5 \\ . & . & -- & = \end{bmatrix}$$    ~ Sparsity ( Sparse )

                              [ Very big & has multiple ]
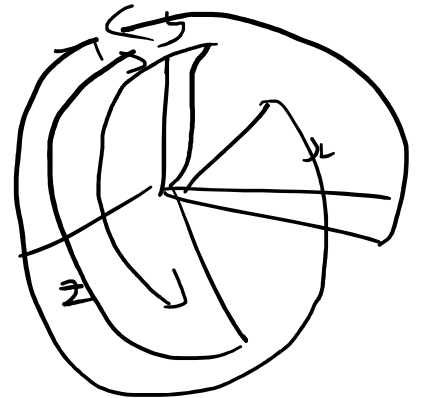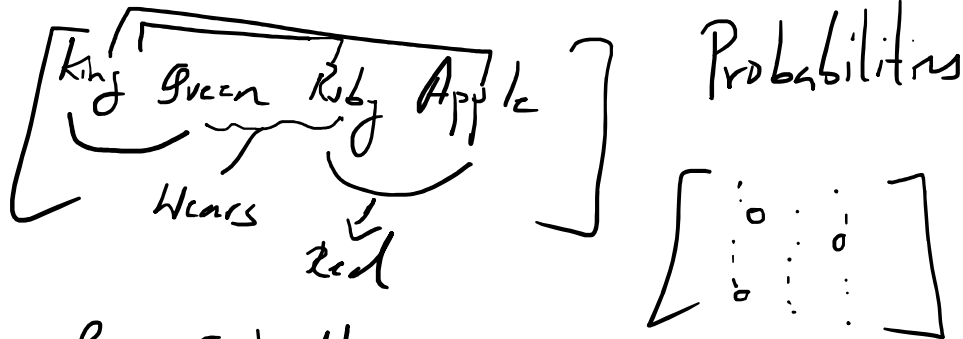                                         zeroes
( 5572 , 8600 )         16 to 32 dimensions
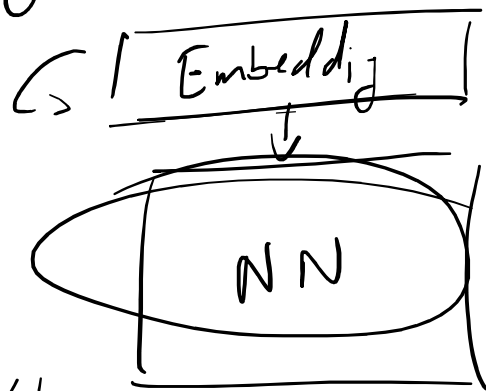
                                    Density of words

→ Word Embeddings (Deep Learning)

[ King Queen Ruby Apple ]
   Wears    Red

Probabilities

[ . ° . ° ]
  ° ° ° .

Density 4 words



→ Train Embedding

∨ Transfer

Learning

⟶ | Embedding |
       ↓
   ( N N )

∨ Traing overhead
   (System)

Word 2 Vec

[ GloVe ]

tokens
500 words vocabs ⟋ 32 dimension

6 billion tokens

500
∿ [ 6 0 0 0 0 0 0 0 0 ]
   100 / 200 / 300 dimensions

→ LSTM

Bidirectional

LSTM

| Embedding |
      ⇓

o/p

→ Creating Embedding on words



[unkomplete (incomplete)] ∨ Character Level Semantics

↝ Character Level Embedding — 1-D Convolutional Layer

 2-D Conv   1-D [▯▯▯▯▯▯...▯] 1-D Conv

→ Parallel Processing (Siamese style n/w)

Token Embd
character Emb

 ⊕ → [] → o/p