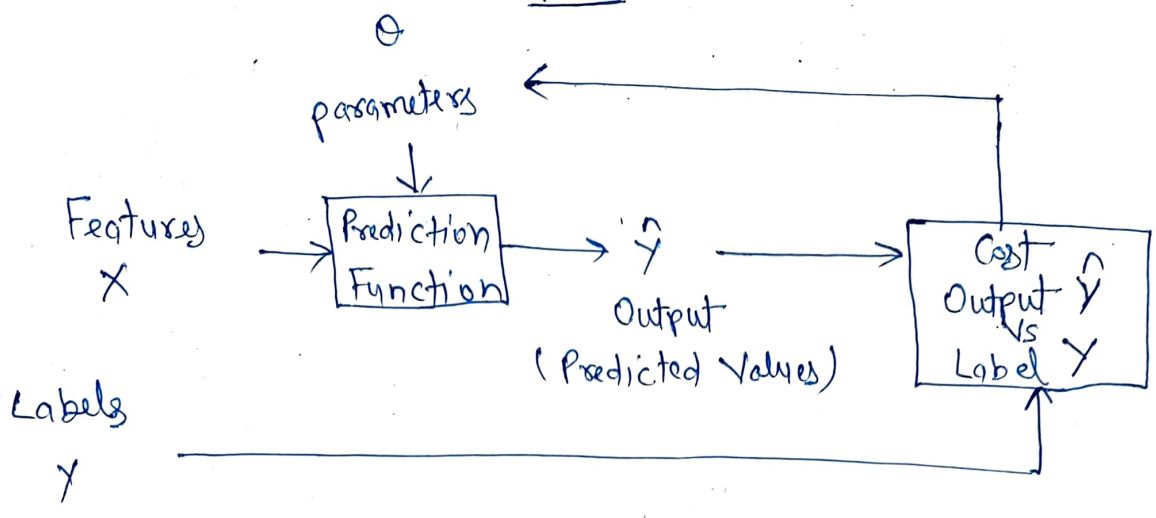


NLP



Tweet: I am happy because I am learning NLP

Logistic Regression
 ↗ Positive 1
 ↘ Negative 0

Vocabulary & Feature Extraction

Tweets
 [tweet1, tweet2 -- tweetm]

because
 I am happy I am learning NLP
 I hated the movie

$V = [I, am, happy, because \dots \dots \dots hated, the movie]$

I am happy because I am learning NLP

[I am happy because MLP -- hated the movie]

[1, 1, 1, 1, 1, \dots, 0, 0, 0]

A lot of zeros, sparse representation.

I am happy I am learning NLP
 All zeros

[1, 1, 1, 1, 1, 1, 1, 0, 0, 0]

| ← ————— |N|

[0, 0, 0, 1, 0, 2, \dots, 0, n]

Large Training Time
 Large Prediction Time
 $n \geq |N|$

Corpus: Positive & negative counts

I am happy because I am learning NLP } +ve class
 I am happy
 I am sad, I am not learning NLP } -ve class
 I am sad

Positive Vocab	Post Freq
I	3
am	3
happy	2
because	1
learning	1
NLP	1
sad	0
not	0

Negative Table Vocab	Post Freq
I	3
am	3
happy	0
because	0
learning	1
NLP	1
sad	2
not	1

Encode the words using frequency of the class

Feature Extraction

freqs: dictionary mapping from (word, class) to frequency

$$X_m = [1, \sum_w \text{freqs}(w, 1), \sum_w \text{freqs}(w, 0)]$$

I am sad, I am not learning NLP.

$$X_m = [1, 8, 11]$$

Now vector size reduced to 3.

Vocabulary

I
am
happy
because
learning
NLP
sad
not

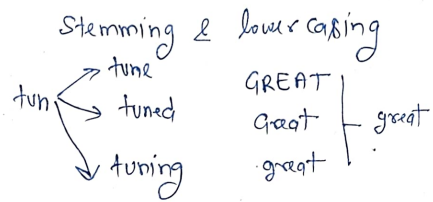
Preprocessing: Stop words & punctuation

@youshuo Bengio and @Geoffrey Hinton are tuning a ~~great~~ GREAT AI model at <https://www.umontreal.ca>!!!

GREAT

Stop words	Punctuation
and	,
is	.
are	:
at	!
has	"
for	,
a	

tuning GREAT AI model



Preprocessed tweet:

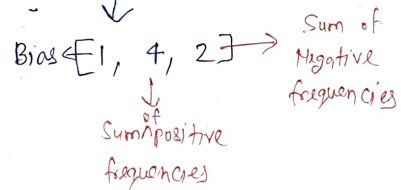
[tun, great, ai, model]

General Overview:

I am Happy Because i am learning NLP @ AI

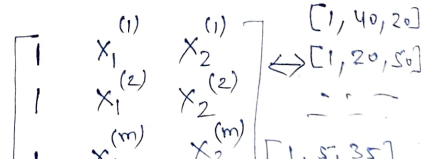
↓ Preprocessing

[happy, learn, nlp]



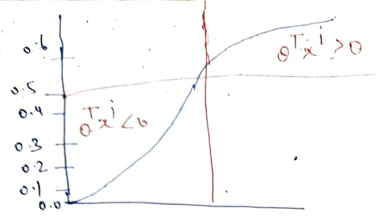
X': Set of tweets

Our X will look like



Logistic Regression Overview

$$h(x^{(i)}, \theta) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$$



@Vishva Bengio @Geoffrey Hinton are training a GREAT model.

[+un, ai, great, model]

$$\theta = \begin{bmatrix} -0.0003 \\ -0.00150 \\ -0.00120 \end{bmatrix}$$

$$x^{(i)} = \begin{bmatrix} 1 \\ 3476 \\ 245 \end{bmatrix}$$

Naive Bayes classifier:

Corpus of tweets

+	+	+	+	+
+	+	+	+	+
+	+	+	-	-
-	-	-	-	-

Tweets containing the words happy

+	+	+	+	+
+	+	+	+	+
+	+	+	ha	ppx
-	-	-	-	-

A → Positive tweet

$$P(A) = 13/20 = 0.65$$

$$P(\text{Negative}) = 1 - P(\text{positive}) = 1 - 0.65 = 0.35$$

Tweets containing the words happy

h	a	pp	x	

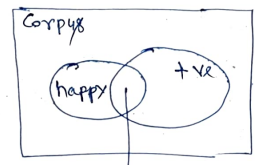
B → Tweets contains happy

$$P(\text{happy} | x) = M_{\text{happy}} / N$$

$$P(B) = 4/20 = 0.20$$

Probability of intersection

$$P(A \cap B) = P(A, B) = 3/20 = 0.15$$



Conditional Probabilities:

+	+	+	ha	ppx
h	a			

$$P(A|B) = P(\text{Positive} | \text{happy})$$

$$P(A|B) = 3/4 = 0.75$$

$$P(B|A) = P(\text{happy} | \text{Positive}) = 3/13 = 0.231$$

$$P(\text{Positive} | \text{happy}) = \frac{P(\text{Positive} \cap \text{happy})}{P(\text{happy})}$$

Bayes's Rule:

$$P(\text{Positive} | \text{happy}) = \frac{P(\text{Positive} \cap \text{happy})}{P(\text{happy})}$$

$$P(\text{happy} | \text{Positive}) = \frac{P(\text{happy} \cap \text{Positive})}{P(\text{Positive})}$$

Summary:

Conditional Probability \rightarrow Bayes' Rule

$$P(X/Y) = \frac{P(Y/X) \times P(X)}{P(Y)}$$

Naive Bayes for Sentiment Analysis

Corpus:

+ve tweets

I am happy because I am learning NLP.
I am happy, not sad

-ve tweets

I am sad, I am not learning NLP.
I am sad, not happy

Words	P(wi/class)	
	Pos	Neg
I	3	3
am	3	3
happy	2	1
→ because	1	0
learning	1	1
NLP	1	1
sad	1	2
not	1	2
	13	12

P(wi/class)

Pos	Neg
$\frac{3}{13} = 0.24$	0.25
0.24	0.25
0.15	0.08
0.08	0.0
0.08	0.08
0.08	0.08
0.08	0.17
0.08	0.17
1	1

← Sum

Naive Bayes:

Tweet: I am happy, today, I am learning

$$\prod_{i=1}^m \frac{P(w_i/Pos)}{P(w_i/Neg)} = \frac{0.14}{0.10} = 1.4 > 1$$

$$\frac{P(I/Pos)}{P(I/Neg)}$$

$$P(+ve/w)$$

$$= \frac{P(w/Positive) P(+ve)}{P(w)}$$

$$\frac{0.26}{0.20} * \frac{0.20}{0.20} * \frac{0.14}{0.10} * \frac{0.26}{0.20} * \frac{0.26}{0.20} * \frac{0.10}{0.10}$$

⊗

NB: Based on Bayes Theorem widely used for text classification:
Spam detection, sentiment analysis, topic classification, document filtering.

Naive \rightarrow assume that feature (words) are conditionally independent given the class.

$$P(c/x) = \frac{P(x/c) P(c)}{P(x)}$$

c: class, x = document

P(c): Prior probability of class

P(x/c) = likelihood

P(c/x) = posterior probability

P(x) is constant for all the class

$$P(c/x) \propto P(c) \prod_{i=1}^n P(w_i/c)$$

Naive Independence Assumption

for a document $x = \{w_1, w_2, \dots, w_n\}$

$$P(x|c) = \prod_{i=1}^n P(w_i|c)$$

Interpretation

The presence of one word does not affect the probability of another word given the class.

- Steps:
- 1) Text preprocessing (tokenization, stopword removal)
 - 2) Feature Extraction (Bag of words / word counts)
 - 3) Compute!
 - class prior $P(c)$
 - word likelihood $P(w|c)$
 4. Apply Bayes Rule
 5. choose class with maximum posterior probability.

Example:

Document	Text	Class
D1	free money offer	Spam
D2	free prize money	Spam
D3	meeting schedule today	Ham
D4	Project meeting today	Ham

Vocab = { free, money, offer, prize, meeting, schedule, today, Project }

Vocab Size = 8

class prior $\Rightarrow P(\text{Spam}) = \frac{2}{4} = 0.5$

$P(\text{Ham}) = \frac{2}{4} = 0.5$

Word Counts per class:

Word	Count	Class
free	2	Spam Documents
money	2	
offer	1	
prize	1	

Words	Count	Class
meeting	2	Ham Documents
today	2	
schedule	1	
project	1	

Total words in Spam = 6

Total words in Ham = 6

Likelihood with Laplace Smoothing:

$$P(w|c) = \frac{\text{Count}(w,c) + 1}{\text{total words in } c + |V|}$$

Ex: $P(\text{free}|\text{Spam}) = \frac{2+1}{6+8} = \frac{3}{14}$

$P(\text{free}|\text{Ham}) = \frac{0+1}{8+8} = \frac{1}{14}$

Test Document: free money

Spam Probability:

$$P(\text{Spam}|x) < P(\text{Spam}) \times P(\text{free}|\text{Spam}) \times P(\text{money}|\text{Spam})$$

$$= 0.5 \times \frac{2}{14} \times \frac{2}{14} = 0.0229$$

$$P(\text{Ham}|x) < 0.5 \times \frac{1}{14} \times \frac{1}{14} = 0.00255$$

Final Prediction:
Spam (✓)

$$P(\text{Spam}|x) > P(\text{Ham}|x)$$

Advantage:

- Simple & easy to implement
- Very fast
- works well with high dimensional data
- performs well with small datasets
- strong baseline for NLP tasks.

Limitation!

- Independence assumption is unrealistic
- Can't capture word order semantics
- poor performance when features are correlated
- lower accuracy compared DL models.

Hidden Markov Model (HMM)

HMM is a probabilistic sequence model where

- System is assumed to be a Markov process (current state depends only on the previous state)
- The states are hidden, but
- Each hidden state emits an observable symbol.

In NLP:

- Hidden States: linguistic labels (POS tags, NER tags, phonemes)
- Observations \rightarrow words (or characters)

Why HMM in NLP?

Natural language is

- Sequential
- Context dependent
- Uncertain

HMMs model:

- Sequential dependency: via transition probabilities.
- Word-label uncertainty: via emission probabilities.

Core Assumptions of HMM \leftarrow

1. Markov Assumption

$$P(s_t | s_{t-1}, s_{t-2}, \dots) = P(s_t | s_{t-1})$$

2. Output Independence Assumption

$$P(o_t | s_t, s_{t-1}, \dots) = P(o_t | s_t)$$

Component of HMM:

An HMM is defined by $\lambda \Rightarrow (A, B, \pi)$

(a) State (S)

Hidden labels:

Example (POS tagging): $S = \{ \text{Noun (N)}, \text{Verb (V)} \}$

(b) Observation (O)

observed words $O = \{ \text{dog, barks} \}$

(c) Transition Probabilities (A): $A_{ij} = P(S_j / S_i)$

(d) Emission Probabilities (B):

$$B_j(o) = P(o / S_j)$$

(e) Initial Probabilities (π)

$$\pi_i \Rightarrow P(S_i = 1)$$

Application of HMM in NLP

Task	Hidden State	Observation
POS Tagging	POS Tags	words
NER	Entity Tag	words
Speech Recognition	Phonemes	Acoustic signal
Text Segmentation	Segment label	characters
Language modeling	latent state	words

Time/NN Flies/VB

Fruit/NN flies/NNS

Time/NN like/IN arrow/~~IN~~ NN.

state = {NN, VB, NNS, IN}

observation (words) = {Time, flies, Fruit, like, arrow}

Initial Probability $P(NN) = 1.0$

Transition Probability:

NN \rightarrow VB: 1 $P(VB/NN) = \frac{1}{3}$

NN \rightarrow NNS: 1

$P(NNS/NN) = \frac{1}{3}$

NN \rightarrow IN: 1

$P(IN/NN) = \frac{1}{3}$

Emission Probability

Tag	word	probability
NN	Time	2/4
NN	Fruit	1/4
NN	arrow	1/4
VB	Flies	1
NNS	flies	1
IN	like	1

word 1: Time

$$V_1(NN) = P(NN) \times P(\text{Time}/NN) = 1.0 \times 0.5$$

$$= 0.5$$

ex 2: fli's

NN → VB

$$0.5 \times \frac{1}{3} \times 1 = 0.1667$$

NN → NNS

$$0.5 \times \frac{1}{3} \times 1 = 0.1667$$

Both path tie

Time / NN fli's / VB

Time / NN fli's / NNS

Q.2 Given Training corpus!

I / PRP like / VB dogs / NNS

You / PRP like / VB cats / NNS

- 1) Find the state & observation
 - 2) Compute Initial & transition probability
 - 3) Compute the probability of sentence
- "I like dogs"

Solⁿ: States & observations (PRP)

states {PRP, VB, NNS}

observations {I, You, like, dogs, cats}

Probability Estimation

$$P(\text{PRP}) = 1$$

Transition Probabilities

Transition	Prob
PRP → VB	1
NB → NNS	1

Emission Probabilities

Tag	word	Probability
PRP	I	0.5
PRP	You	0.5
VB	like	1
NNS	dogs	0.5
NNS	Cats	0.5

Sentence Prob Calculation

⇒: I like dogs

$$P = P(\text{PRP}) \times P(\text{I} / \text{PRP}) \times P(\text{VB} / \text{PRP}) \times P(\text{like} / \text{VB}) \\ \times P(\text{NNS} / \text{VB}) \times P(\text{dogs} / \text{NNS})$$

$$= 1 \times 0.5 \times 1 \times 1 \times 1 \times 0.5$$

$$= 0.25$$

Limitation of HMM

Strong independence assumptions (output depends only on current state).

Limited context modelling due to first order Markov assumption.

Poor handling of long range dependencies.

Data sparsity & zero probability problem.

Generative nature, model $P(\text{words} | \text{tags})$ instead of $P(\text{tags} | \text{words})$.

Modern Solutions:

LSTM/GRU \rightarrow Capture long term dependencies

Transformer-based models (BERT, RoBERTa)

\Rightarrow Context representations

Hybrid HMM-Neural models - useful for low-resource languages.

Training Corpus:

- | | DET | NOUN | VERB |
|----|-----|------|-------|
| 1. | the | dog | barks |
| 2. | a | dog | runs |
| 3. | the | cat | runs |

Test:

the dog runs

Tag Set: DET, NOUN, VERB

Vocab: the, a, dog, cat, barks, runs

Frequency Counts:

Tag Counts: DET = 3, NOUN = 3, VERB = 3

Transition Count:

START \rightarrow DET = 3, DET \rightarrow NOUN = 3, NOUN \rightarrow VERB = 3

Emission:

DET: the = 2, a = 1

NOUN: dog = 2, cat = 1

VERB: barks = 1, runs = 2

Laplace's Smoothing:

Number of tags = 3, Vocab size = 6

$$P(\text{NOUN} | \text{DET}) = \frac{3+1}{3+3} = 4/6$$

$$P(\text{DET} | \text{DET}) = 0+1/3+3 = 1/6$$

Emission Prob!

$$P(\text{the} / \text{DET}) = \frac{2+1}{3+6} = 3/9$$

$$P(\text{dog} / \text{Noun}) = \frac{2+1}{3+6} = 3/9$$

$$P(\text{runs} / \text{VERB}) = \frac{2+1}{3+6} = 3/9$$

Viterbi (Normal Probability)

Sentence: the dog runs.

Initialization: $\text{DET} = 4/6 \times 3/9 = 0.222$

Recursion:

$$\text{Noun} = 0.2222 \times \frac{4}{6} \times \frac{3}{9} = 0.0494$$

$$\text{VERB} = 0.0494 \times \frac{4}{6} \times \frac{3}{9} = 0.0110$$

Best path = DET \rightarrow Noun \rightarrow VERB

Viterbi (log probabilities)

$$\ln(4/6) = -0.405$$

$$\ln(3/9) = -1.099$$

$$\text{DET} = -1.504$$

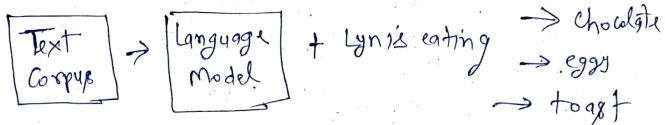
$$\text{Noun} = -3.008$$

$$\text{VERB} = -4.512$$

Language Model

N-gram based language model:

- Create a language model from text corpus to
- estimate the probability of word sequences
- estimate the prob of a word following a sequence of words.
- Apply this concept to autocomplete a sentence



other applications:

- Speech Recognition
- Spelling Correction
- Argumentative Communication

N-gram:

An N-gram is a sequence of N words.

Corpus: I am happy because I am learning

Unigram: { I, am, happy, because, learning }

Bigram: { I am, am happy, happy because }

Trigram: { I am happy, am happy because }

Sequence Notation:

Corpus: This is great
 $w_1 w_2$

teacher drinks tea.
 w_{500}

$m=500$

$$w_1^m = w_1 w_2 \dots w_m$$

Unigram Probability:

Corpus: I am happy because I am learning.

Size of corpus $\Rightarrow m=7$

$$P(I) = 2/7$$

$$P(\text{happy}) = 1/7$$

Bigram Prob

$$P(\text{am} | I) = \frac{C(I \text{ am})}{C(I)} = 2/2 = 1$$

$$P(\text{happy} | I) = \frac{C(I \text{ happy})}{C(I)} = 0/2 = 0$$

$$P(y|x) = \frac{C(x,y)}{\sum_w C(x,w)} = \frac{C(x,y)}{C(x)}$$

Trigram Prob

$$P(\text{happy} | I \text{ am}) = \frac{C(I \text{ am happy})}{C(I \text{ am})}$$

$$P(w_3 | w_1^2) = \frac{C(w_1^2 w_3)}{C(w_1^2)}$$

$P(w_i | w_1^{i-1})$ $N=i$

$\langle s \rangle$ I am Sam $\langle /s \rangle$

$\langle s \rangle$ Sam I am $\langle /s \rangle$

$\langle s \rangle$ I do not like Sam $\langle /s \rangle$

Calculate the prob of $P(\text{am}/I)$ using the corpus.

$$P(w_n/w_{n-1}) = \frac{\text{Count}(w_{n-1}, w_n)}{\text{Count}(w_{n-1})}$$

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

Prob of sentence $\langle s \rangle$ I am Sam $\langle /s \rangle$ according to bigram model.

Decomposition

$$P(I/\langle s \rangle) \times P(\text{am}/I) \times P(\text{Sam}/\text{am}) \times P(\langle /s \rangle/\text{Sam})$$

$$P(\text{Sentence}) = \frac{2}{3} \times \frac{2}{3} \times \frac{1}{2} \times \frac{1}{3} = \frac{4}{54} = 0.074$$

③ Calculate $P(\text{Sam}/\text{do})$ using Add-1 Laplacian Smoothing.

$$P = \frac{\text{Count}(w_{n-1}, w_n) + 1}{\text{Count}(w_{n-1}) + V}$$

$$P(\text{Sam}/\text{do}) = \frac{0+1}{1+6} = \frac{1}{7}$$

perplexity: a standard metric used to evaluate how well a language model predicts a sample.

Ex: $S = I$ like Sam

$$P(I/\langle s \rangle) = 0.5$$

$$P(\text{like}/I) = 0$$

$$P(\text{Sam}/\text{like}) = 0.2$$

$$P(\langle /s \rangle/\text{Sam}) = 0.4$$

Perplexity is the inverse probability of test set , normalized by the no of words, for a sentence W with N words

$$PP(W) = P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}}$$

$$PP = \sqrt[N]{\frac{1}{P(w_1/w_0) \times P(w_2/w_1) \dots P(w_N/w_{N-1})}}$$

Note Usually N include end of sentence token $\langle /s \rangle$ not $\langle s \rangle$ (start to end)

For a sentence $w = w_1 w_2 \dots w_N$, the perplexity of a language model is defined as:

$$PP(w) = [P(w)]^{-\frac{1}{N}}$$

or equivalent in log form

$$PP(w) = \exp\left(-\frac{1}{N} \sum_{i=1}^N \log P(w_i | w_{1-T(i)})\right)$$

Lower Perplexity \Rightarrow better language model.

Ex: Sentence: <S> I love MLP </S>

No of words (including end token)
 $N = 4$

Given Bigram	Probability
Bigram	Prob
$P(I <S>)$	0.5
$P(\text{love} I)$	0.4
$P(\text{MLP} \text{love})$	0.1
$P(</S> \text{MLP})$	0.2

$$P(\text{Sentence}) = P(w) = 0.5 \times 0.4 \times 0.1 \times 0.2 = 0.004$$

Perplexity Calculation

$$PP(w) = (0.004)^{-\frac{1}{4}} = \sqrt[4]{250} = 3.97$$

log-prob version (not in practice)

$$\begin{aligned} \log P(w) &= \log 0.50 + \log 0.40 \\ &\quad + \log 0.10 + \log 0.20 \\ &= -(0.6931) + (-0.9163) + (-2.3026) \\ &\quad + (-1.6094) \\ &= -5.5214 \end{aligned}$$

$$PP(w) = \exp\left(-\frac{-5.5214}{4}\right) = \exp(1.38035) = 3.97$$

$Perplexity = 3.97$

Sentences

$\langle S \rangle \langle S \rangle$ | love NLP $\langle /S \rangle$

$\langle S \rangle \langle S \rangle$ | love AI $\langle /S \rangle$

$$N = 4 + 4 = 8$$

Given trigram prob

Trigram

Prob

$$P(I | \langle S \rangle \langle S \rangle)$$

Basic word Representations!

- integer
- One hot vectors
- word embeddings

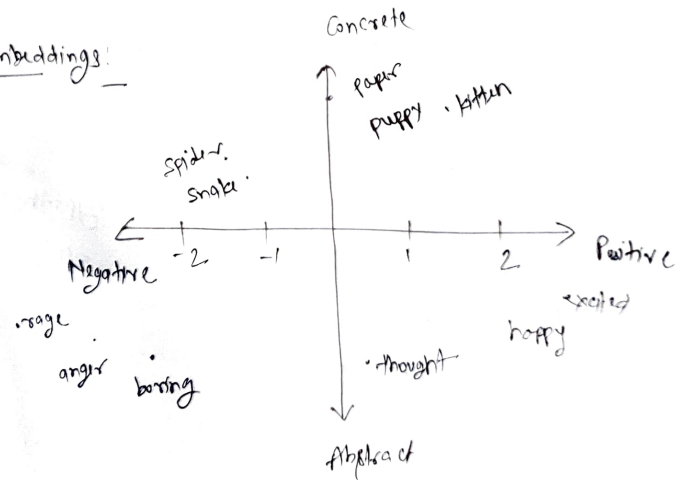
word	Number
a	1
able	2
about	3
⋮	⋮
hand	615
happy	621
zebra	1000

happy

1	0	9
2	0	able
3	0	about
615	0	hand
621	1	happy
1000		zebra

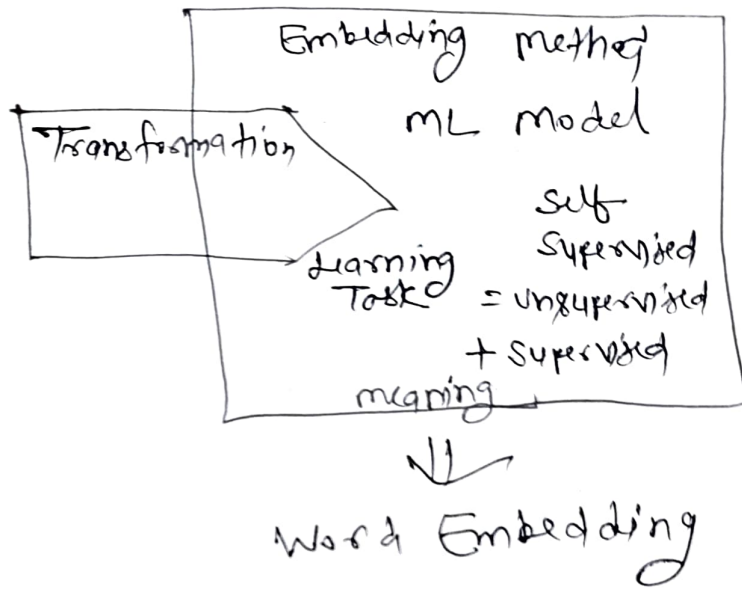
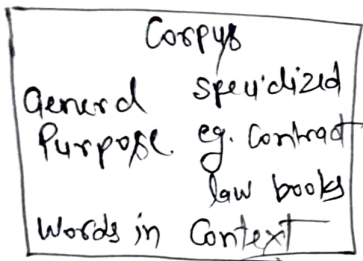
⇒ One hot encoding
 Pros: Simple & require no implied ordering.
 Cons: one hot vectors: huge & encode no meaning.

Embeddings!



Word Embedding Process

Hyperparameters
word Embedding size



classical method

• Word2Vec (Google 2013)

Continuous bag of word (CBOW) the model learn to predict the center word given some context words.

• Word 2 Vec operates on Distributional hypothesis words that appear in similar context show similar meaning

The cat sit on the mat

The dog sits on the mat

CBOW

input [The _ sits on]

• Target: cat