

CHAT-BOT

* **CHAT-BOT** :- chatbot are soft. application that use artificial intelligence and NLP technique, in which you ask a question and based on some rules, this system provide you an answer.

* Types of CHAT-BOT :-

→ Retrieval Based (Rule Based) :- with this type of chat-bot, communication work by having a pre-set of question/Answer rules. Based on similarity in input que. and data set que., model provide answers.

User Question:

what is c?

NLP Preprocessing

Similarity matrix (user que == Dataset Que)

Based on similarity

c is manish

Answer from CHAT-BOT

Pre-Define Data-set

	Que	Ans
1.	what is a?	a is Ram
2.	what is b?	b is shyam
3.	what is c?	c is manish

→ Generative Model Based (Self-learning) :- This chat-bot work 2)
by generating new response using ML/DL on a lot of
historical data and previous conversation, without
needing a pre-define dataset.

3>

* Basic NLP Pre-Processing Steps *

* Tokenization * Stemming * Lemmatization * Vectorization
* word Embedding * NER * Text classification

↳ **Tokenization** : Tokenization is the task of breaking or separating the whole sentence into words / apostrophy or punctuation.

Exp:- This is a dog → This is a dog

split the string (text)
using space ("/s")

Exp: We're friends → we 're friends

Appostopy Problem

Exp: what is your name? → what is your name ?

Punctuation
Problem
separate
token
consider

↳ **Stemming** : → This take input as tokenized output words in form of array / list.

→ stemming remove suffix / pre-fix from word and improve the quality of word which is in list.

4) Ex: This dog was coming
↓ ↓ ↓ ↓
Thi dog wa Com

→ **Lemmatization** :- Lemmatization find the root of the words, not just remove the affix from the words.

Ex Birds is am are
↓
Bird be

→ **Named Entity Recognition** :- A task of NLP, whose purpose is to locate and classify named entities in string into pre-defined categories.

Ex: President Abdul Kalam of India on his first state visit to the United State on Tuesday night.
Person Loc ordinal Location Date Time

5)

→ **Vectorization** : vectorization is a process which is used to turning a text document into a numerical vector.

→ Basic approach are:

↳ Bag of words (count vectorizer)

↳ TF-IDF (Term frequency inverse term frequency)

↳ word2vec

ex: "Bag of word" -

1. first define fixed length vocabulary.

["I", "am", "you", "are", "manish"]

2. map each word to an index in voc.

[I=1, am=2, you=3, are=4, manish=5]

3. Based on this index, construct a vector in which the word index is a 1 if the word in document else 0.

ex: "I am manish" → [1, 1, 0, 0, 1]
 ↑ ↑ ↑ ↑ ↑
 I am you are manish