

1-4 Intro. to main concepts

5-8 basics of dataset and tokenizers.

9/12 Transformer models in CV and speech recognition

Geekradio -

NLP - field of linguistic and ML, understanding everything related to human language.

NLP tasks

1. Classifying whole sentences:- SA, spam, grammatically correct

2. " each word :- gram. components, NER

3. Generating text content:-

4.

Why challenging?

Pipeline () fn:- Selects a model pretrained for a particular task.

3 main steps

- ① text preprocessed in a way that model can understand
- ② preprocessed ip → to the model.
- ③ predictions are processed so you can make sense.

Available pipelines:-

feature-extraction

fill-mask

NER

qa

SA

Summarization

text-generation

translation

text-generation
translation
zero-shot classification.

Mask-filling:- is used to train BERT models

Zero-shot learning:- need for fine-tuning model on data to use.

Text-generation:- provide a prompt, model will auto-complete it.

Grouped-entities = true, regroup together parts of the sentence

Hugging Face into one

How Transformer Works?

June 2018

GPT-family:- Auto-regressive

BERT-family:- Auto-encoding

BART/T-5 ... Seq-to-seq

Trained as lang. models

large amount of raw texts

Self-supervised fashion.

Objectives automatically computed

humans to label the data.

Statistical understanding of the language.

Transfer learning

fine-tuning in a supervised way.

Caused language modeling

O/P depends on present and past i/p

Masked language modeling:-

To achieve better performance by ↑ size and amount of data in which models are pretrained.

Tgt

CO₂ footprint of LLMs:-

^{Q1}
CO₂ footprint of LLMs: —

Type of energy: — (non)-renewable

Training Time

H/W

I/O

Data

CO₂ emission per kWh wise vary लग्जी ।

Mumbai → 920 gm

Montreal → 20 gm

almost 90X

Other elements: —

pretraining → recycling

fine-tuning

Starting with smaller exp. & Debugging

Lit. reviews for selecting hyperparameter range

Random search Vs Grid search

ML emissions calculator

Code carbon

Emissions trackers

large corpuses

Pretaining: — Training from scratch.

Weights initialized randomly

training w/o any prior knowledge.

Large dataset

Training time Weeks.

Fine-tuning: — Training after pretaining

" on dataset specific task

" " " - fine-tune (fin-tuning)

" on dataset specific task

why we directly train on specific tasks (fine-tuning):—

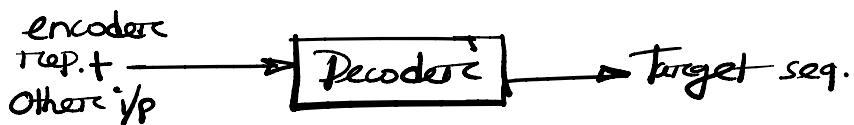
advantage of taking knowledge acquired during pretraining.

e.g. less data for decent result.

Time & resources ↴

Encoder:— Self-Attention, Bi-dir.

Decoder:— Uni-dir., Auto-regressive, Masked Self-Attention



Encoder only models:— req. understanding of i/p.
Sentence classification
NER

Decoder " " :— generative tasks
Text generation.

En-Decoder model / seq-seq:— generative tasks req. on i/p
translation, summarization.

Attention layer:—

Tell the model to give more attention to certain words.

Word \Rightarrow meaning একটি বা অন্যের word \Rightarrow depend করে।
Context

Original architecture of Transformer:—

Designed for translation.

— সূত্রান্তি lang. → **Encoder**

Same Target lang. → **Decoder**

Same target lang. 

Encoder \rightarrow word (here আয়োজন)

Decoder seq. করতে পারতে। Future words Training এর সময় দ্রব্যে
আয়োজন।

1st Attn. layer in decoder pays attn. to all past i/p to decoder.
2nd u u u u can access o/p of encoder. Thus whole
seq. to predict next word.

Architecture: — skeleton of the model.

def. of layers
op. happen w/i : BERT

checkpoints: — weights loaded in a archi., bert-base-cased

Model: — Architecture / checkpoint both.

When to use encoders: —

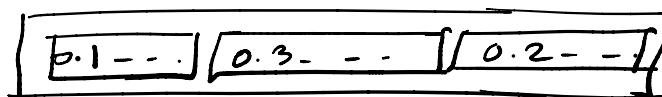
Bi-directional

Good at extracting meaningful info.

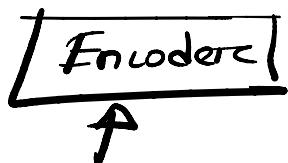
Seq. classification, qa, MLM

NLU

examples — BERT, RoBERTa, ALBERT.



768 Dim Vectors



Welcome to NYC

Use encoders of Transformer

Attn. layers access all the words in initial sentence.

.....

Attn. layers access all the words in initial sentence.

Bi-dir. attn. / auto-encoding

Pretraining by MLM.

Best suited for understanding full sentence.

sentence classification, NER, QA

examples:- (AL) BERT,

Distil ..

ELECTRA

RoBERTa.

Decoder - archi:- GPT-2

Same tasks can be performed as encoder-only models but performance decreased a little.

Masked self-attention mechanism:-

Only see words on the left
right side words are hidden.

When to use?

Uni dir.

Causal tasks; generating seq.

NLG

GPT-2 / Neo.

1024

Only access words before it. Auto-regressive models.

Text generation.

CTRL

GPT-1/2

Transformer-XL

3. En-decoders / seq-toseq models:-

encoder - all the words

decoder - access words positioned before a 'given word.' } attn. layer access.

~~context~~ - access words positioned before a 'given word' to access.

Pre-training more complex

T5 masks random spans of texts. & predicts to pretrain
Best suited for gen. new sentences on a given i/p
Summarization, Translation, Generative QA.

m(BART), Marian, T5.

Bias and Limitations:-

Dataset एके काम के लिये internet मा आकू तोहि दिये pre-train
कर्या रखा। परन्तु कामे अनेक समस्या bias भए जाते।

This man/woman works as a _____.

As a result gender-wise very vary करते।

pretrain model use कर्या समस्या माने बाधते होते ही model
मुला sexist, racist/ homophobic होते होते।

Fine-tuning कराये तो Bias को remove करी possible हो।

Chapter 2:— Using Transformers:-

Intra:— Transformers very large.

- Billions of parameters.

New models released regularly.

Transformer library to load, train & save models.

Features:-

1. Ease of use:— Two line code, यहाँ तक की भी नया model लगाया

2. flexibility:— ~~जैसे~~ other NN model के लिये framework यहाँ तक की लगाया

3. Simplicity:— Understandable & usable by using a single file.

Chapter 2:— Contents

model + Tokenizer = replicate pipeline fn.

model & Config. classes.

Load a model

..... condition.

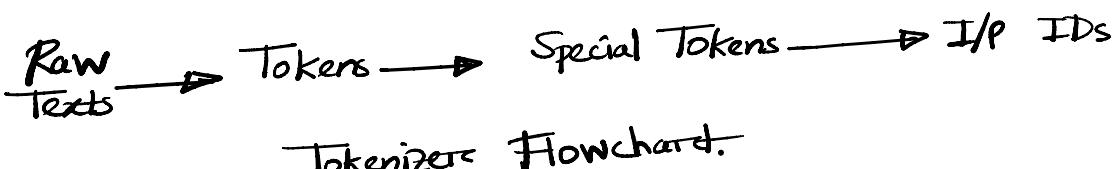
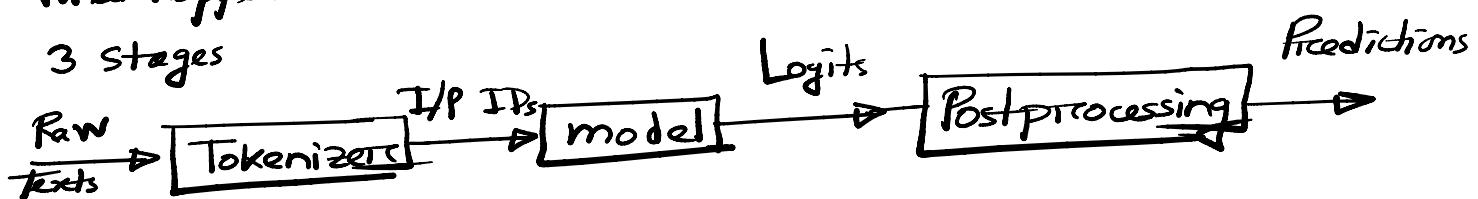
Load a model

Process numerical I/P to o/p prediction.

Tokenizer API:-

What happens inside pipeline fn:-

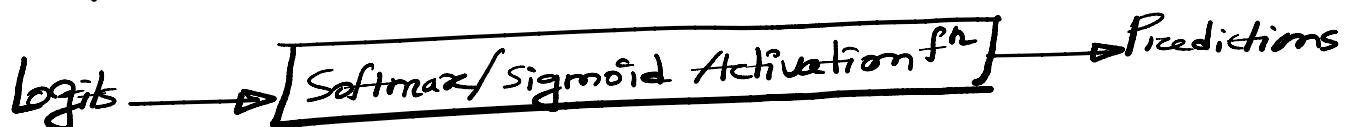
3 stages



Attn. mask = 0 याने padding, दूसरे Attn. fn के लिए
= 1 " Text. Attn. fn के लिए 1

Automodel loads a model w/o pretraining head.

Logits:- Few, unnormalized predictions of model



Preprocessing with a Tokenizer:-

Can't process Text directly } problem

Convert text i/p into no. }

Soln → use Tokenizer

Tokenizer responsible for:-

Splitting i/p into tokens

Tokens Mapping → Integers

+ extra i/p

Model pretraining के लिए यह preprocessing करता है जैसे exactly
यहाँके preprocessing करता है।

1. Download info. from Hub.

2. ... from - pretrained(model_name)

1. Download info. from Hub.

AutoTokenizer.from_pretrained(model_name)

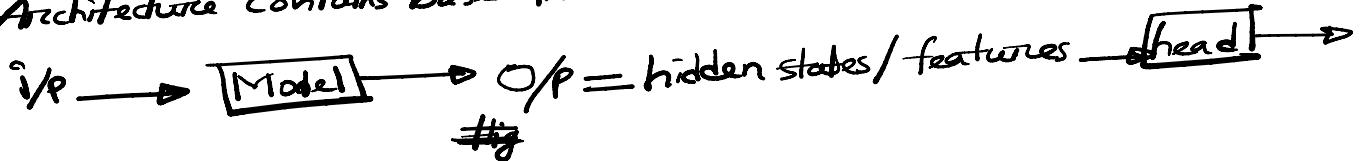
Tokenizer returns dict. (I/P ID)

need to convert into Tensor

Model Download:-

AutoModel.from_pretrained(checkpoint)

Architecture contains base Transformer Module



Features = high dim. vector rep. contextual understanding of I/P.

Same architecture use करते Diff. Task करते लिए head diff. होते हैं।

Vectors O/p of Transformer Model:-

Large dim. vectors as O/p. Three dim.

Batch size:- At a time करने से sentence process करते

seq. len.:- Padding करते हैं वाले तक Sentence की len. का

hidden size:- Vec. dim. of each model I/P 768 होती है।

O/p behave like named tuples / dictionaries

Access the elements by

by attributes op.last_hidden_state.shape

or key op["last"]

or index op[0]

