

2.1 Model scratch ~~corpus~~ train major brand new tokenizer train ~~corpus~~
एवं ।

HF tokenizer library
fast tokenizer

2.2 Training a new tokenizer from old one:-

model x available

corpus is different from PT corpus

retrain model from scratch. need a new tokenizer

Most Xmerc model use subword tokenization algo.

Training a tokenizer:-

x same as training a model

randomized by nature

set some seeds to get same results when doing same training

Statistical process

which subwords to pick depend on tokenization algo.

Deterministic

Assembling a Corpus:-

Autotokenizer.train_new_from_iterators()

GPT-2 from scratch on Python code.

Code-search-net

Dissimilarities arise from:-

New lang.

" characters

" domain

" style

6.3 Fast Tokenizer's Special powers:-

6.3 Fast Tokenizer's Special powers:-

Slow tokenizers written in Python
fast " " " Rust.

Special tokens



Model



Pre-tokenization

Splits



Normalization

lower case, removing accents

Word IDs application

whole word masking
token classification

Offset mapping application

token classification

QA

O/P of tokenizer

x simple python dictionary

Special batch encoding object.

Subclass of dict.

Fast tokenizers

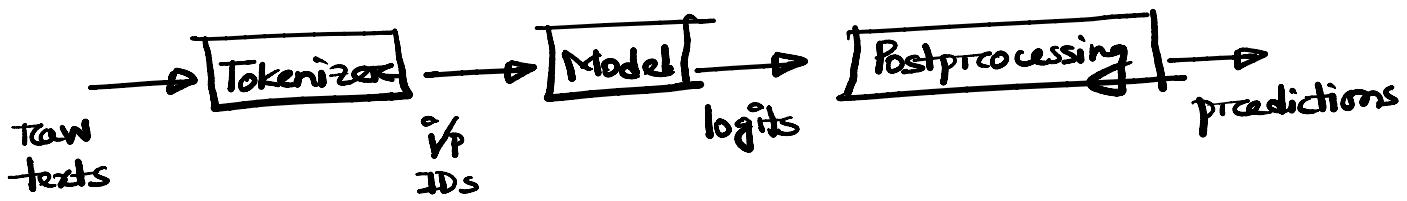
parallelization capabilities

keep track of original span of texts → offset mapping

Token classification pipeline steps:-

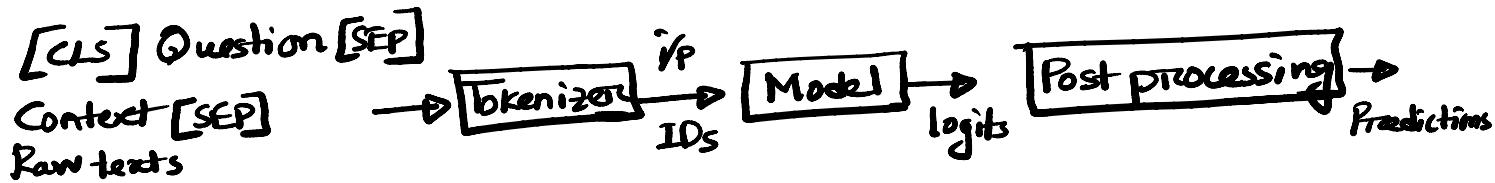


Token classification pipeline steps:-



→ 1 2 3 4
→ 5 6 7 8

6.4 fast tokenizers in QA-pipeline:-



6.5 Normalization & pretokenization:-

Hello, how are U + day?

↓
hello how are u today?

↓
[hello, how, are, u, today, ?]

↓
[hello, how, are, u, td, #day, ?]

↓
Els, hello, how, are , u, td, #day, ?, SEP

Normalization

Pretokenization

Model

Postprocessor

Normalization operation automatically included when tokenizing
ন্যো এটা tokenize করার পথে অনেক <unk>/<oov>

| bytes | 01100011 1100100 10100111 11000011 10100111 | Decoding (UTF-8) |
|---------------|---|------------------|
| code point(s) | U+0063 U+0327 | unicode mapping |
| characters | S | |

characters

S

G & "mapping"

unicode normalization standards:-

NFC, NFD, NFKC, NFKD

Sentencepiece algorithm:-

Consider text as a seq. of unicode char.

replaces spaces with special char. —

with unigram algo.

x req. pre-tokenization

easily reversible and invertible

replace — with space

| Model | Byte pair Encoding | Wordpiece | Unigram |
|---------------|---|--|---|
| Training | Starts with small vocab. & learn rules to merge tokens. | → Same | with large vocab. learn rules to remove tokens. |
| Training Step | Merge tokens correspondingly to most common pair. | on best score | remove tokens that min. loss |
| Learn | Merge rules & vocab. | Just a vocab. | a vocab. with a slot for each token. |
| Encoding | Splits a char. & apply merges during training | finds longest subword then same for rest | |

Byte-Pair Encoding (BPE) Algorithm:-

for compression of texts

OpenAI for tokenization during training GPT

Open AI for tokenization during training (or)

(Ro) DeBERTa.

Normalize → Pretokenize → ~~Assign word as unk~~
word counter add गाड़ा रेत । → Split into characters. →
तुम्हीं तुम्हीं वह चर्के फ्रेम पैटर्न लगूनो frequency नज़र ।
Highest freq. pair in vocab. & splits ज़ add गाड़ा ।

BPE training algo. details

unique sets of words in the corpus (after Normalization & pretokenization)

Building the vocab. by taking all symbols to write all words

Contain all ASCII char. & some unicode char.

If char. x in corpus, during tokenizing texts unknown token
very bad with emojis.

GPT-2, RoBERTa

x · unicode char

✓ u Bytes

256 size

every char. included

Byte-level BPE

base vocab. अलगी नियम

add new tokens

until desired vocab. size reached.

Start learning merges

merges 2 char. into 1

Search most freq. pair of existing tokens

-Tokenization Algorithm:-

search most freq. pair of words

Tokenization Algorithm:-

Normalization

Pre-tokenization

Splitting words into char.

6.6 Word Piece Tokenization:-

Google

PT BERT

Sim. to BPE but tokenization is diff.

X open sourced.

Training Algorithm:-

Small vocab.

special tokens

initial alphabets

word

w # # o # # r # # d

learns merge rules

Score = $\frac{\text{freq of pair}}{\text{freq. of first element} \times \text{freq. of second element}}$

Tokenization Algo:-

WP only saves final vocab.

X merge rules learned.

finds longest subword

Split it hugs \rightarrow hug, # # s

bugs \rightarrow b, # # u, # # g s

bun \rightarrow [UNK] BPE

\rightarrow b, # # u, [UNK] WP

→ b, #HU, [UNK] WP

6.8 Unigram Tokenization:-

Starts with big vocab.

removes tokens until reaches desired vocab. size.

big vocab build करने वाले way

most common substrings in Protokenized

apply BPE on corpus with large vocab.

Training के लिए unigram algo. loss का कैसे

overall loss ↑ हो जाता है if one symbol removed.

find out the symbol will have lower effect on overall loss.

one symbol remove करना अपेक्षा costly

per symbol → काफ़ी 2% symbols associated with lowest loss increase.

process is repeated until desired vocab. size is reached.

X remove base char.

Tokenization algo.

independent of context

predict most common token.

$$\text{probability} = \frac{\text{freq. in original corpus}}{\text{sum of all tokens freq.}} = \frac{20}{210}$$

tokenization with least tokens will have highest probability

Split a token into least no. of tokens possible.

pug → p u, g / P, ug

depends on which segmentation encountered first

Viterbi algo.

• 1 In 1-best possible commentaries

Viterbi algo.

find build a graph to detect possible segmentations

Training:-

loss computed by tokenizing every word "in corpus"

$$\text{loss} = -\log(P(\text{word}))$$

hug: "hug" score: 0.071428

pug: "pu", "g" u: 0.007710

:

$$10 \times (-\log(0.071428)) + 5 \times (-\log(0.007710)) + 12 \dots$$

remove वाक्य loss वा \uparrow होती है remove वाक्य!