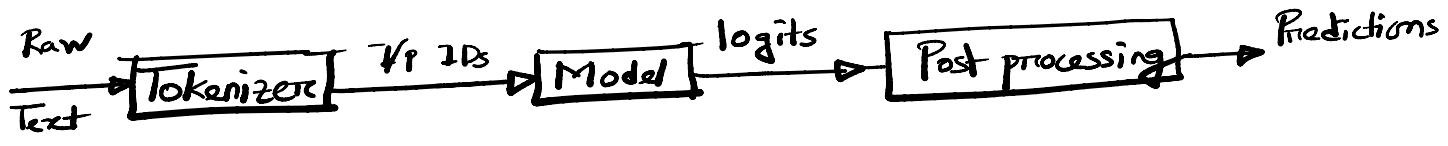


Main features of huggingface's Xmerc library:-

Ease of use:-

Flexibility

Simplicity



Raw text

This course is amazing!

↓  
Tokens

[this, course, is, amazing, !]

↓  
Special tokens

[CLS], this, course, is, amazing, !, SEP

↓  
IP IDs

[ 101, 2023, - - - , - - - , 102 ]

Tokenizer

Preprocessing with a tokenizer:-

Splitting IP into (sub) words, symbols → tokens

Mapping each token into integers

+ additional tokens like [CLS] [SEP]

Should be done exactly same way as the model PT.

Xmerc models only use tensors as IP.

For specifying the type of tensors to get back, we use

return\_tensors = "pt"

An argument for all from list of lists return -2(0)

Argument is ~~a file~~ an input list of lists return -~~(o)~~ 20)

Going through the model:-

downloaded & cached.

O/P:- hidden states/ features.

hidden states or head 20 i/p from the list - 21)

diff. tasks 21 uses same archi. but diff. head.

there 20 O/P large 21 3D:-

Batch size: No. of seq. processed at a time.

Seq. length: length of the numerical rep. of seq.

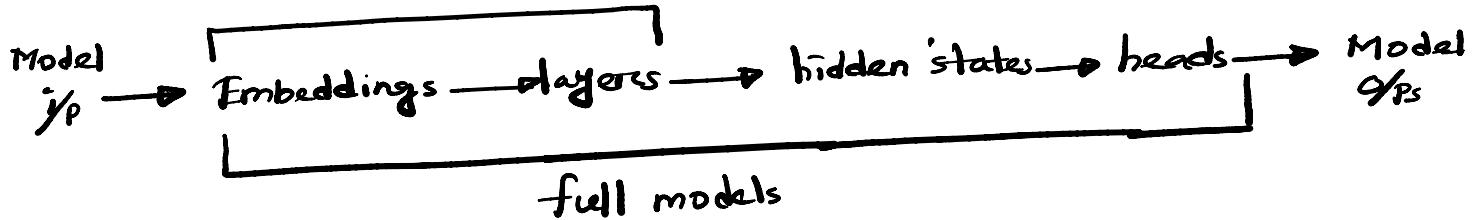
hidden size:- Vector rep. of each model i/p.

Model heads:- Making sense out of numbers.



composed of one/few linear layers.

Transformer Net.



Postprocessing the O/P:-

O/P of the model x make sense

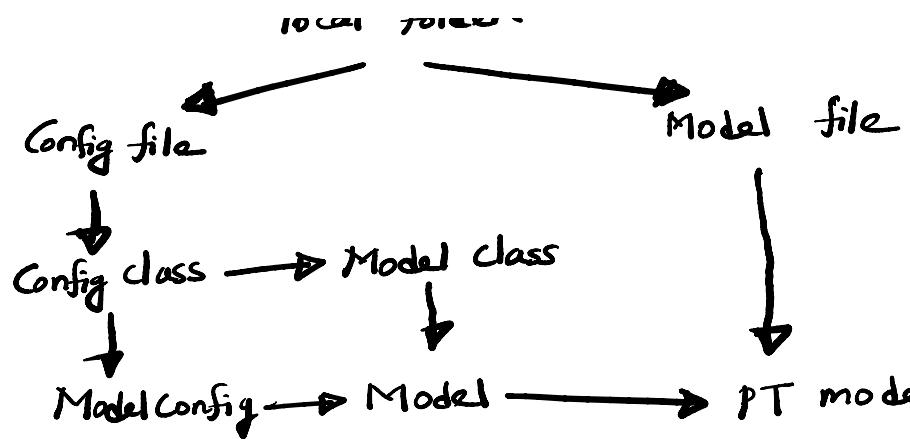
Outputs logits

model.config.id2label

Instantiate a Transformers model:-

checkpoint/  
local folder

← → ... ↘ ↗



**Wrappers:-** prog. components simplifies usage/interaction with other S/W lib./API/System.

### Section 2.3 Models:-

AutoModel & Sim. wrappers automatically guess appropriate model archi.

Loading a PT Xmerc:-

~~BertModel.from\_pretrained("bert-base-uncased")~~

Weights downloaded & cached into  
~/.cache/huggingface/transformers.

Saving Methods:-

model.save\_pretrained("dir...")

dir... → ~~py~~ file save into

Config.json

Attributes necessary to build model archi.

Contains metadata

where checkpoints originated, Xmerc version

when last saved the checkpoint

pytorch\_model.bin

State dir.

Contains all weights

} model archi.

Using a Xmerc for inference:-

Xmerc model only process numbers that tokenizer generates.  
... commands increase.

- Xmer model only processes numbers that tokenizer generates.
- Tokenizer casts i/p to appropriate framework tensors.
- i/p Segs.  $\rightarrow$  tokenizer  $\rightarrow$  vocab. indices (input IDs), list of lists
- Tensors accept rectangular shapes (Matrices)
- input IDs  $\xrightarrow{\text{torch.tensor}(\quad)}$  tensors, i/p to Xmer models.

## Section 2.1 Tokenizers:-

Core components of NLP  
translate text into numbers (data that Xmer model can understand).  
goal:- find most meaningful & smallest rep.

### Word-based tokenizers:-

Easy to set up

use only with few rules

decent performance often

different ways to split texts.

by whitespace using Python's `split()` fn.

Disadv:- may result very large vocab.

inflections are identified as diff. words:-

dog, dogs

run, running

unknown / [UNK]

$\uparrow$  [UNK],  $\uparrow$  losses of info.

### Character based tokenizers:-

Vocab. smaller (256)

fewer OOV tokens

charac. x convey meaning

large no. tokens (1 word  $\rightarrow$  10+ tokens)

whole english words (1,70,000) with just 256 char.

Whole English words (1,70,000) with just 256 char.

Subword tokenization:-

frequently used words X Split  
rare " ✓ "

good coverage + small vocab.

~ no unknown tokens

helpful in Turkish lang. where long complex words can be formed

loading & saving of tokenizer:-

from\_pretrained()

Save\_pretrained()

Encoding:- 2 steps

tokenization

Conversion to i/p IDs

Instantiate some tokenizer as the model → same rules used to PT.  
" Vocab.

Tokenization:- Auto/BertTokenizer. from\_pretrained()

Tokens to i/p IDs:- tokenizer.convert\_tokens\_to\_ids()

Decoding:- tokenizer.decode()

## Section 2.5 Handling multiple seq.

How to handle multiple seq. of variable length?

Vocab. only i/p to allow model to work well?

Too long seq.?

Model expects a batch of i/p's:-

Seq. → list of no. → tensor → i/p to the model.

Xmore expects multiple sentences by default

Batching is the act of sending multiple sentences through the model.

1 sentence input, you can batch with a single sentence.

= [ids, ids]

~

multiple sentences connect with

= [ids, ids]

Issue 1:- If you have multiple sentence expect राहे

if 2 sentences having diff. length.

Tensors need to rectangular shape.  
pad the inputs.

Issue 3:- padding कराये तो 2 दूरे logit same ताकि ना model each token को contextualize कराये जाए।

Soln:- Attention mask use करें padding को ignore कराएं।

Attention Mask:- Is indicating → attended into

0s " → ignore.

Longer seq:-

limit to the lengths of seq. → 512/1024.

seq. > 512/1024

Crash

Soln:- use a model with a longer supported seq. length, Longformer LED  
truncate your seq.

## Section 2.6:- Putting it all together:-

tokenizer( ) fn directly use करा।

can handle everything related to tokenization.

multiple seq.

padding : longest , max\_length

Conversion to specific framework tensors

pt = Pytorch tensors

tf = Tensorflow "

np = Numpy arrays.

2.7

Basic usage completed:-

2.8 Quiz:-