

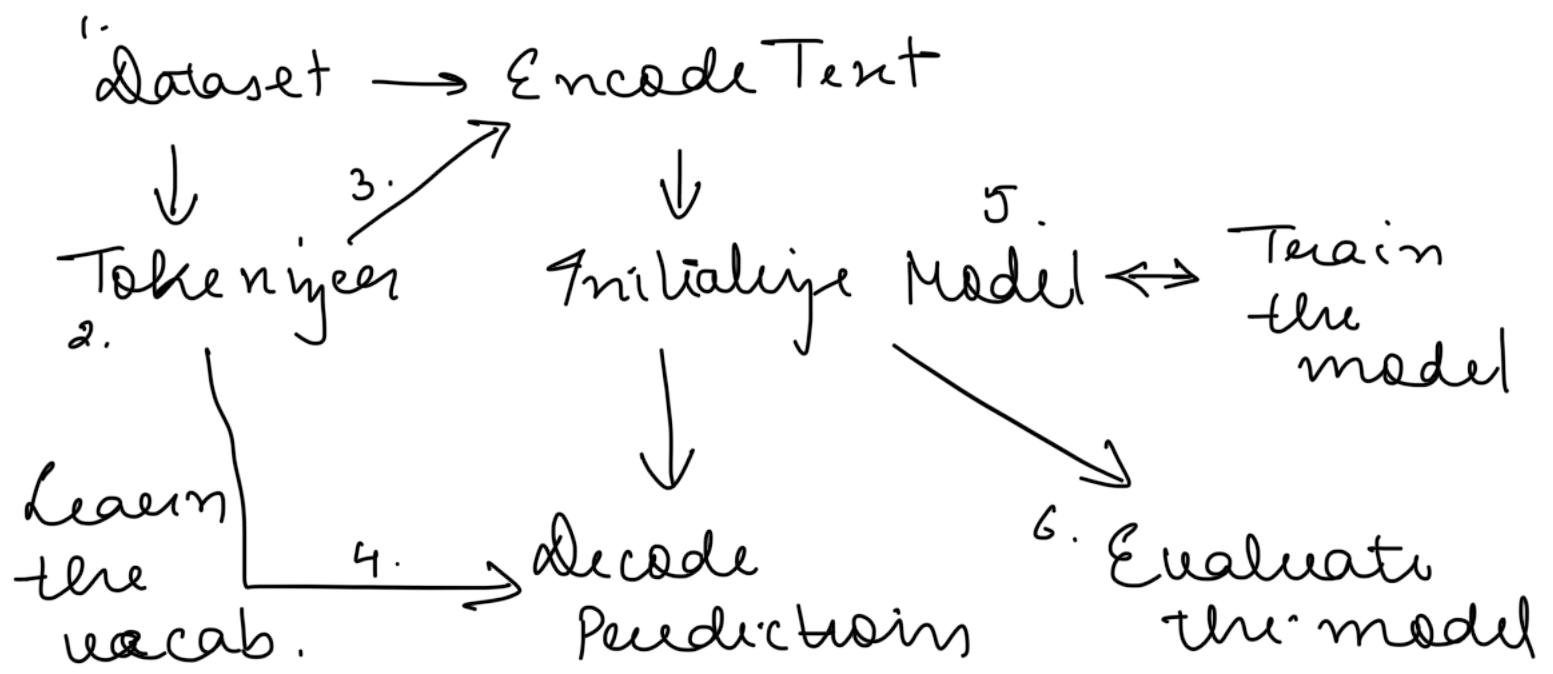
## Week - 1

### ↳ Intro to Course & Recap of NN Models

- NLP enables computers to understand, interpret & gen. human lang.
- Input Text  
↓  
Translate      Sentiment      NER  
(Named entity  
recognition)
- n gram models → GPT 4, Flama  
(with emerging abilities)
- 4 types of architectures - MLP (backprop. alg. for learning the params. feed-forward fully connected neural network), CNN (CV, wt. sharing concept where the wts. in the kernels are shared across the input), RNN (recurrent LSTMs, transformers, etc. use the same wt. for all elements in the input seq.).
- Now, it is easy to build these architectures using PyTorch & Tensorflow.
- Each layer can be derived from

- 'torch.nn.module' is PyTorch.
- MLP → The model is composed of linear layers & non-linear activation (assume ReLU). ∴ nn.linear and torch.relu.
- CNN → The model is composed of nn.linear, nn.Conv2d, MaxPool 2D & torch.relu.
- RNN → The model is composed of nn.linear, nn.Embedding, nn.RNN. m.RNN has nn.linear, torch.relu and other required modules.
- Transformer → Model is composed of nn.embedding, nn.Transformer, nn.Linear. m.Transformer has nn.MultiHeadAttention and other required modules.

## L2 Why HF? Read Map



- HF does the heavy lifting
  1. import datasets
  2. import tokenizers
  3. tokenzier.encode()
  4. tokenzier.decode()
  5. import transformers
  6. import evaluate
- HF has 150k datasets, 350k models, & 170k spaces.
- Modern NLP
  - Prompt → Text (sentiment, summarizing, fill in the blank, generate story)
  - Lang. Modelling ↓

Output response conditioned on  
Prompt.

- wide range of NLP tasks - sentiment, machine transl.", NER, question-answers, textual entailment, summary and generation.
- For each of these tasks we have 100s of datasets with 1000s of samples.
- Dataset can be varied sizes & formats understand the format & write a script to parse & load the text data.
- If memory is limited, we can stream the samples from the dataset.

### L3 Datasets Module HF, Loading

- access to datasets → diff. format, diff. size for various tasks.  
↳ datasets module
- inspect datasets (movie reviews in SL)  
with `huggingface.co/datasets/`
- store it in cache memory. The cached data is stored in a memory-mapped columnar format, e.g. less memory

- now load the dataset  $\rightarrow$  data dict.  
train, test, unsupervised.
- each split has features & no. of rows  
`imbd_data['train']`  
↳ now this is the train split  
↳ key
- remove unsupervised (using pop on dict)  
`dataset.pop('unsupervised')`
- load ('.csv', split = "train")  
↳ to download only train split.
- we can train on certain dataset,  
while show the evalua" on some other  
dataset
- train  $\rightarrow$  can also be split  
`train.train-test_split(test_size = 0.2)`
- data
  - |- train.csv
  - |- test.csvsame no. of columns  
and same col. name
- `data_files = [-, -]`  
`load_dataset("csc", data_files = [ ])`
- convert dataset from any format to  
pyarrow format by saving it to the  
disk.
- load it back from `load_from_disk`  
func".

## L4 Accessing Samples

- `index = 1000  
ex = dataset['train'][1000]` } or send a list  
`print(ex)`  
`label — }` } dictionary obj.  
`text — }`  
[ ] . select ([idx])
- Now explore "translate" dataset  
wmt/wmt14  
`get_dataset_config_name("—")`  
→ cs-en, — —  
`get_dataset_split_names("—", "hi-en")`  
→ train, validation, test
- if you want to get the entire dataset, we can combine the splits
- features → define the int. structure.  
It specifies the underlying serialized task.  
`print(dataset['train'].features)`
- Glue: data → sub-config: mapc

## L5 Common Methods

- `dataset dict` → filter, map, concatenate  
↳ operates on all samples across the splits

- filtering  $\rightarrow$  no. of words in a sentence should be  $\geq 100$ .

$$m = 1000$$

data.filter(lambda eg: len(eg['text']) - split['']) > m)

- map  $\rightarrow$  add prefix IMDB: to all the samples.

(eg['text'] = "IMDB: " + eg['text'])  
dataset.map(addprefix)

- concatenate  $\rightarrow$  same features & same no. of splits

$\hookrightarrow$  eg. IMDB + rotten tomatoes

datasets.concatenate\_datasets

([—, —], axis=0)

- interleaving datasets  $\rightarrow$  n-skewed datasets, that is no. of samples in the dataset might differ drastically

interleave\_datasets([—, —],  
probabs. = [0.6, 0.4])

stopping-strategy = first\_exhausted  
 $\hookrightarrow$  default

- iterable dataset  $\rightarrow$  use on the fly dataset.

load(—, streaming=True)

the other methods like map, filter  
also does its work on the fly