

5.2 Dataset with HF Hub & APIload_dataset

flexible

Single path

list of paths

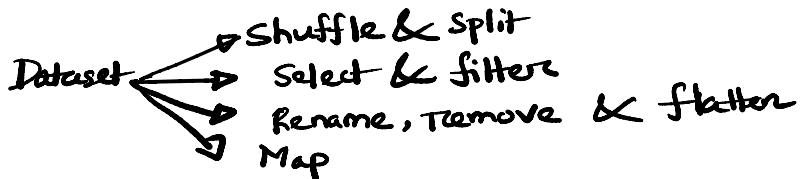
dict. maps split names to file path.

Automatic decompression of zip files.

GZIP, ZIP, TAR

5.3 Slice & dice a dataset

Dataset won't perfectly prepare for training

5.4 Big data:-

Millions of doc

GB of size

HF Datasets library

frees memory management problems

memory-mapped files.

GPT-J corpus → The pile 1.2 TB

T-5 u → CC9 ≈ 3 TB

Datasets library uses

Apache arrow Traditional Row/Column format

Streaming API. Iterable dataset.

Arrow's Memory mapped format

enables access to bigger than RAM datasets

multiple processes to work with same dataset w/o moving/copying

Pile

English text corpus.

File

English text corpus.

Electra-PTI

Diverse range of topics

14.3 GB

2 columns :- meta, text

num_rows: 15518009

RAM usage:-

psutil.Process().memory_info().rss

RSS = resident set size

Dataset size

name_of_dataset["train"].dataset_size

size \approx 20 GB

→ RAM used only 5600 MB

Wes McKinney's famous rule of thumb for pandas library

5x/10x RAM \leq dataset size असे देवता !

Datasets library treat each datasets as a memory-mapped file
mapping betw RAM & file storage w/o fully load into memory

done by apache arrow memory format & pyarrow lib.

Parquet:- compress करें file & store करें रास्ते data

80% Compression Ratio.

→ if dataset असे size आणि येतो ट्रॅक, disk तरीके रुपात

वॅल Streamming = True argument नंतर dataset एक element

तरीके access करा शक्य on-the-fly, w/o need to download whole dataset.

access first element

next(iter(name_of_dataset))

elements can be processed by

IterableDataset.map()

shuffle करा शक्य

Shuffle করা কৈ

```
sh_ds = name_of_dataset.shuffle(buffer_size=10,000, seed=42)
next(iterator(sh_ds))
```

Select করা কৈ

```
d = dataset.take(5)
list(d)
```

Skip করা কৈ . skip()

অনেকগুলো dataset রে combine করা কৈ : -

```
interleave_datasets()
```

```
combined_dataset = interleave_datasets([dataset_1, dataset_2])
list(islice(combined_dataset, 2))
```

5.5 Creating your own dataset :-

Creating a corpus of GitHub issues :-

to track bugs/features in GitHub repo.

use cases :-

time to take to close open issues/pull requests.

training a multilabel classifier to tag metadata based on issue's description

Creating semantic search engines to find issues match user query.

Getting data:-

all issues of the datasets in repo's issues tab

to download all repo's issues, use GitHub REST API :-

5.6 Semantic search with FAISS :-

Sentence \rightarrow Transformer Encoders \rightarrow embed. vectors.

Cosine sim. for measuring how close two embeddings are : -

$$\cos\theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|}$$

FAISS (Facebook AI similarity Search)

FAISS (Facebook AI Similarity Search)

Sim. search based on embed.

convert long passages into single embed.

" Question " "

Compare 2 embed. using FAISS

5.6 Use embedding for semantic search:-

rep. each token as embed. vector

pool individual embed. for vector rep. of whole sentences.

Loading & preparing