# Cheat Sheet: Generative AI Overview and Data Preparation

| Package/Method | Description | Code example |
|---|---|---|
| NLTK | NLTK is a Python library used in natural language processing (NLP) for tasks such as tokenization and text processing.<br><br>The code example shows how you can tokenize text using the NLTK word-based tokenizer. | ```python<br>import nltk<br>nltk.download("punkt")<br>from nltk.tokenize import word_tokenize<br><br>text = "Unicorns are real. I saw a unicorn yesterday. I couldn't see it today."<br>token = word_tokenize(text)<br>print(token)<br>``` |
| spaCy | spaCy is an open-source library used in NLP. It provides tools for tasks such as tokenization and word embeddings.<br><br>The code example shows how you can tokenize text using spaCy word-based tokenizer. | ```python<br>import spacy<br>text = "Unicorns are real. I saw a unicorn yesterday. I couldn't see it today."<br>nlp = spacy.load("en_core_web_sm")<br>doc = nlp(text)<br>token_list = [token.text for token in doc]<br>print("Tokens:", token_list)<br>``` |
| BertTokenizer | BertTokenizer is a subword-based tokenizer that uses the WordPiece algorithm.<br><br>The code example shows how you can tokenize text using BertTokenizer. | ```python<br>from transformers import BertTokenizer<br><br>tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")<br>tokenizer.tokenize("IBM taught me tokenization.")<br>``` |
| XLNetTokenizer | XLNetTokenizer tokenizes text using Unigram and SentencePiece algorithms.<br><br>The code example shows how you can tokenize text using XLNetTokenizer. | ```python<br>from transformers import XLNetTokenizer<br><br>tokenizer = XLNetTokenizer.from_pretrained("xlnet-base-cased")<br>tokenizer.tokenize("IBM taught me tokenization.")<br>``` |
| torchtext | The torchtext library is part of the PyTorch ecosystem and provides the tools and functionalities required for NLP. | ```python<br>from torchtext.vocab import build_vocab_from_iterator<br><br># Defines a dataset<br>dataset = [<br>``` |

The code example shows how you can use torchtext to generate tokens and convert them to indices.

```
    (1,"Introduction to NLP"),
    (2,"Basics of PyTorch"),
    (1,"NLP Techniques for Text Classification"),
    (3,"Named Entity Recognition with PyTorch"),
    (3,"Sentiment Analysis using PyTorch"),
    (3,"Machine Translation with PyTorch"),
    (1,"NLP Named Entity,Sentiment Analysis, Machine
Translation"),
    (1,"Machine Translation with NLP"),
    (1,"Named Entity vs Sentiment Analysis NLP")]

# Applies the tokenizer to the text to get the tokens as a
list
from torchtext.data.utils import get_tokenizer

tokenizer = get_tokenizer("basic_english")
tokenizer(dataset[0][1])

# Takes a data iterator as input, processes text from the
iterator, and yields the tokenized output individually
def yield_tokens(data_iter):
    for _,text in data_iter:
        yield tokenizer(text)

# Creates an iterator
my_iterator = yield_tokens(dataset)

# Fetches the next set of tokens from the data set
next(my_iterator)

# Converts tokens to indices and sets <unk> as the default
word if a word is not found in the vocabulary
vocab = build_vocab_from_iterator(yield_tokens(dataset),
specials=["<unk>"])
vocab.set_default_index(vocab["<unk>"])
```

| | | ```
# Gives a dictionary that maps words to their corresponding
numerical indices
vocab.get_stoi()
``` |
|---|---|---|
| vocab | The vocab object is part of the PyTorch torchtext library. It maps tokens to indices.<br><br>The code example shows how you can apply the vocab object to tokens directly. | ```python
# Takes an iterator as input and extracts the next tokenized
# sentence. Creates a list of token indices using the vocab
# dictionary for each token.

def get_tokenized_sentence_and_indices(iterator):
    tokenized_sentence = next(iterator)
    token_indices = [vocab[token] for token in
tokenized_sentence]
    return tokenized_sentence, token_indices

# Returns the tokenized sentences and the corresponding
# token indices. Repeats the process.

tokenized_sentence, token_indices = \
get_tokenized_sentence_and_indices(my_iterator)
next(my_iterator)

# Prints the tokenized sentence and its corresponding token
# indices.

print("Tokenized Sentence:", tokenized_sentence)
print("Token Indices:", token_indices)
``` |
| Special tokens in PyTorch:<br><eos> and <bos> | Special tokens are tokens introduced to input sequences to convey specific information or serve a particular purpose during training.<br><br>The code example shows the use of <bos> and <eos> during tokenization. The <bos> token denotes the beginning of the input sequence, and the <eos> token denotes the end. | ```python
# Appends <bos> at the beginning and <eos> at the end of the
# tokenized sentences using a loop that iterates over the
# sentences in the input data

tokenizer_en = get_tokenizer('spacy',
language='en_core_web_sm')
tokens = []
max_length = 0

for line in lines:
``` |

| | | |
|---|---|---|
| | | ```
tokenized_line = tokenizer_en(line)
tokenized_line = ['<bos>'] + tokenized_line + ['<eos>']
tokens.append(tokenized_line)
max_length = max(max_length, len(tokenized_line))
``` |
| Special tokens in PyTorch: <pad> | The code example shows the use of <pad> token to ensure all sentences have the same length. | ```
# Pads the tokenized lines
for i in range(len(tokens)):
    tokens[i] = tokens[i] + ['<pad>'] * (max_length - len(tokens[i]))
``` |
| Dataset class in PyTorch | The Dataset class enables accessing and retrieving individual samples from a data set.<br><br>The code example shows how you can create a custom data set and access samples. | ```
# Imports the Dataset class and defines a list of sentences
from torch.utils.data import Dataset

sentences = ["If you want to know what a man's like, take a good look at how he treats his inferiors, not his equals.", "Fae's a fickle friend, Harry."]

# Downloads and reads data
class CustomDataset(Dataset):
    def __init__(self, sentences):
        self.sentences = sentences

    # Returns the data length
    def __len__(self):
        return len(self.sentences)

    # Returns one item on the index
    def __getitem__(self, idx):
        return self.sentences[idx]

# Creates a dataset object
dataset=CustomDataset(sentences)

# Accesses samples like in a list
E.g., dataset[0]
``` |

| | | |
|---|---|---|
| DataLoader class in PyTorch | A DataLoader class enables efficient loading and iteration over data sets for training deep learning models.<br><br>The code example shows how you can use the DataLoader class to generate batches of sentences for further processing, such as training a neural network model | ```python<br># Creates an iterator object<br>data_iter = iter(dataloader)<br><br># Calls the next function to return new batches of samples<br>next(data_iter)<br># Creates an instance of the custom data set<br>from torch.utils.data import DataLoader<br>custom_dataset = CustomDataset(sentences)<br><br># Specifies a batch size<br>batch_size = 2<br><br># Creates a data loader<br>dataloader = DataLoader(custom_dataset,<br>batch_size=batch_size, shuffle=True)<br><br># Prints the sentences in each batch<br>for batch in dataloader:<br>    print(batch)<br>``` |
| Custom collate function in PyTorch | The custom collate function is a user-defined function that defines how individual samples are collated or batched together. You can utilize the collate function for tasks such as tokenization, converting tokenized indices, and transforming the result into a tensor.<br><br>The code example shows how you can use a custom collate function in a data loader. | ```python<br># Defines a custom collate function<br>def collate_fn(batch):<br>    tensor_batch = []<br><br># Tokenizes each sample in the batch<br>    for sample in batch:<br>        tokens = tokenizer(sample)<br><br># Maps tokens to numbers using the vocab<br>        tensor_batch.append(torch.tensor([vocab[token] for<br>token in tokens]))<br><br># Pads the sequences within the batch to have equal lengths<br>    padded_batch =<br>pad_sequence(tensor_batch,batch_first=True)<br>    return padded_batch<br>``` |

| | | ```python
# Creates a data loader using the collate function and the custom dataset
dataloader = DataLoader(custom_dataset, batch_size=batch_size, shuffle=True, collate_fn=collate_fn)
``` |
| --- | --- | --- |