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

| | | |
|---|---|---|
| Special tokens in PyTorch: <eos> and <bos> | input sequences to convey specific information or serve a particular purpose during training. 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. | ```\n4    tokens = []\n5    max_length = 0\n6    for line in lines:\n7        tokenized_line = tokenizer_en(line)\n8        tokenized_line = ['<bos>'] + tokenized_line + ['<eos>']\n9        tokens.append(tokenized_line)\n10       max_length = max(max_length, len(tokenized_line))\n``` |
| Special tokens in PyTorch: <pad> | The code example shows the use of <pad> token to ensure all sentences have the same length. | ```\n1    # Pads the tokenized lines\n2    for i in range(len(tokens)):\n3        tokens[i] = tokens[i] + ['<pad>'] * (max_length - len(tokens[i]))\n``` |
| Dataset class in PyTorch | The Dataset class enables accessing and retrieving individual samples from a data set. The code example shows how you can create a custom data set and access samples. | ```\n1    # Imports the Dataset class and defines a list of sentences\n2    from torch.utils.data import Dataset\n3    sentences = ["If you want to know what a man's like, take a\n4    good look at how he treats his inferiors, not his equals.",\n5    "Fae's a fickle friend, Harry."]\n6    # Downloads and reads data\n7    class CustomDataset(Dataset):\n8        def __init__(self, sentences):\n9            self.sentences = sentences\n10       # Returns the data length\n11       def __len__(self):\n12           return len(self.sentences)\n13       # Returns one item on the index\n14       def __getitem__(self, idx):\n15           return self.sentences[idx]\n16   # Creates a dataset object\n17   dataset=CustomDataset(sentences)\n18   # Accesses samples like in a list\n19   E.g., dataset[0]\n``` |
| DataLoader class in PyTorch | A DataLoader class enables efficient loading and iteration over data sets for training deep learning models. 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 | ```\n1    # Creates an iterator object\n2    data_iter = iter(dataloader)\n3    # Calls the next function to return new batches of samples\n4    next(data_iter)\n5    # Creates an instance of the custom data set\n6    from torch.utils.data import DataLoader\n7    custom_dataset = CustomDataset(sentences)\n8    # Specifies a batch size\n9    batch_size = 2\n10   # Creates a data loader\n11   dataloader = DataLoader(custom_dataset, batch_size=batch_size, shuffle=True)\n12   # Prints the sentences in each batch\n13   for batch in dataloader:\n14       print(batch)\n``` |
| 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. The code example shows how you can use a custom collate function in a data loader. | ```\n1    # Defines a custom collate function\n2    def collate_fn(batch):\n3        tensor_batch = []\n4    # Tokenizes each sample in the batch\n5        for sample in batch:\n6            tokens = tokenizer(sample)\n7    # Maps tokens to numbers using the vocab\n8            tensor_batch.append(torch.tensor([vocab[token] for token in tokens]))\n9    # Pads the sequences within the batch to have equal lengths\n10       padded_batch = pad_sequence(tensor_batch,batch_first=True)\n11       return padded_batch\n12   # Creates a data loader using the collate function and the custom dataset\n13   dataloader = DataLoader(custom_dataset, batch_size=batch_size, shuffle=True, collate_fn=collate_`,\n``` |