

REINFORCEMENT LEARNING

Reinforcement Learning (RL) is a type of machine learning where an agent learns to make decisions by interacting with an environment. The agent takes actions based on the current state of the environment and receives feedback in the form of rewards or penalties. The goal of reinforcement learning is for the agent to learn a strategy, or policy, that maximizes the cumulative reward over time by exploring and exploiting its environment.

COMPONENTS OF REINFORCEMENT LEARNING

- **Agent:** The learner or decision-maker.
- **Environment:** Everything the agent interacts with.
- **State:** A specific situation in which the agent finds itself.
- **Action:** All possible moves the agent can make.
- **Reward:** Feedback from the environment based on the action taken.

WORKING OF REINFORCEMENT LEARNING

RL operates on the principle of learning optimal behavior through trial and error. The agent takes actions within the environment, receives rewards or penalties, and adjusts its behavior to maximize the cumulative reward. This learning process is characterized by the following elements:

Policy: A strategy used by the agent to determine the next action based on the current state.

Reward Function: A function that provides a scalar feedback signal based on the state and action.

Value Function: A function that estimates the expected cumulative reward from a given state.

Model of the Environment: A representation of the environment that helps in planning by predicting future states and rewards.

EXAMPLE OF REINFORCEMENT LEARNING

Teaching a Dog to Sit

Imagine you're training a dog to sit on command. This is an example of **Reinforcement Learning** because the dog (the **agent**) learns through interaction with the environment (the process of being trained).

1. **State:** The dog is either standing or sitting.
2. **Action:** The dog can either **sit** or **not sit** when you give the command.

3. **Reward:** When the dog sits on command, you give it a treat (positive reward). If the dog doesn't sit, there is no reward (or perhaps a small negative consequence, like ignoring the dog).
4. **Goal:** The dog's goal is to learn to sit as quickly as possible to get treats. Over time, the dog learns that sitting leads to rewards, and not sitting leads to no rewards.

Learning Process:

- At the beginning, the dog doesn't know what "sit" means. It might try random actions like standing, walking, or even jumping.
- When it accidentally sits, you reward it with a treat, telling it that "sitting" is a good action.
- After several repetitions, the dog learns that sitting is associated with a reward and starts sitting whenever you give the command.

In RL, this is a process of trial and error: - The **dog (agent)** takes actions (sitting or not sitting). - The **environment** provides feedback (the reward of a treat or no treat). - The dog **learns** to associate sitting with rewards and figures out the best action to maximize its rewards.

This process mirrors the way RL agents learn in more complex environments like video games, robotics, or business strategies. The agent explores its options, tries different actions, and learns the optimal behavior over time to maximize rewards.

[]: *CartPole Environment in OpenAI Gym*

The CartPole environment is a classic reinforcement learning problem where the goal is to balance a pole on a cart by applying forces to the left or right.

```
import gym
import numpy as np
import warnings
```

Suppress specific deprecation warnings

```
warnings.filterwarnings("ignore", category=DeprecationWarning)
```

Load the environment with render mode specified

```
env = gym.make('CartPole-v1', render_mode="human")
```

Initialize the environment to get the initial state

```
state = env.reset()
```

Print the state space and action space

```
print("State space:", env.observation_space)
print("Action space:", env.action_space)
```

Run a few steps in the environment with random actions

```
for _ in range(10):
    env.render() Render the environment for visualization
    action = env.action_space.sample() Take a random action
```

Take a step in the environment

```
step_result = env.step(action)
```

Check the number of values returned and unpack accordingly

```
if len(step_result) == 4:
```

```
    next_state, reward, done, info = step_result
```

```
    terminated = False
```

```
else:
```

```
    next_state, reward, done, truncated, info = step_result
```

```
    terminated = done or truncated
```

```
print(f"Action: {action}, Reward: {reward}, Next State: {next_state}, Done: {done}, Info: {info}")
```

```
if terminated:
```

```
    state = env.reset() Reset the environment if the episode is finished
```

```
env.close() Close the environment when done
```

State space: Box([-4.8000002e+00 -3.4028235e+38 -4.1887903e-01 -3.4028235e+38], [4.8000002e+00 3.4028235e+38 4.1887903e-01 3.4028235e+38], (4,), float32)

Action space: Discrete(2)

Action: 1, Reward: 1.0, Next State: [0.01660546 0.18100183 -0.04156917 -0.34862286], Done: False, Info: {}

Action: 0, Reward: 1.0, Next State: [0.02022549 -0.013505 -0.04854162 -0.06933248], Done: False, Info: {}

Action: 1, Reward: 1.0, Next State: [0.01995539 0.18227807 -0.04992827 -0.37692678], Done: False, Info: {}

Action: 0, Reward: 1.0, Next State: [0.02360095 -0.01210052 -0.0574668 -0.10039505], Done: False, Info: {}

Action: 0, Reward: 1.0, Next State: [0.02335894 -0.20635381 -0.05947471 0.17361793], Done: False, Info: {}

Action: 1, Reward: 1.0, Next State: [0.01923187 -0.01043331 -0.05600235 -0.13721846], Done: False, Info: {}

Action: 0, Reward: 1.0, Next State: [0.0190232 -0.20471027 -0.05874672 0.13728404], Done: False, Info: {}

Action: 0, Reward: 1.0, Next State: [0.01492899 -0.39894372 -0.05600104 0.41087002], Done: False, Info: {}

Action: 1, Reward: 1.0, Next State: [0.00695012 -0.20307444 -0.04778364 0.10107096], Done: False, Info: {}

Action: 0, Reward: 1.0, Next State: [0.00288863 -0.39748016 -0.04576222 0.37830368], Done: False, Info: {}

NATURAL LANGUAGE PROCESSING

INTRODUCTION TO NLP

NLP stands for Natural Language Processing. It is the branch of Artificial Intelligence that gives

the ability to machine understand and process human languages. Human languages can be in the form of text or audio format.

- **Core Areas:** NLP involves various fields such as linguistics, computer science, and cognitive psychology to process language data.
- **Key Components:** Common aspects of NLP include syntax (sentence structure), semantics (meaning), and pragmatics (contextual meaning).
- **Applications:** NLP powers applications like machine translation, sentiment analysis, chatbots, speech recognition, and language generation.
- **Techniques:** Techniques used in NLP range from rule-based models to statistical models, and more recently, deep learning models like Transformers.
- **Machine Learning Models:** NLP utilizes machine learning models such as Bag of Words, Word2Vec, and advanced deep learning models like BERT and GPT.
- **Challenges:** Key challenges include understanding context, handling ambiguity, processing idiomatic expressions, and managing large datasets.
- **Tools and Libraries:** Popular NLP libraries include NLTK, SpaCy, Hugging Face Transformers, and OpenNLP.
- **Future of NLP:** With advancements in AI, NLP is expected to become more accurate and capable of understanding complex human interactions.

NLP LIBRARIES

NLTK

The Natural Language Toolkit (NLTK) is a Python library designed for natural language processing, offering tools for text analysis, preprocessing, and linguistic research. It includes functions for tokenization, stemming, lemmatization, and syntactic parsing, along with access to extensive corpora and lexical resources. NLTK is widely used for educational and research purposes, especially suitable for small to medium-scale NLP projects.

- **Comprehensive NLP Toolkit:** NLTK is a widely used Python library for natural language processing (NLP) tasks, offering tools for text manipulation and analysis.
- **Text Processing Functions:** Supports essential NLP functions such as:
 - **Tokenization** (splitting text into words or sentences)
 - **Stemming and Lemmatization** (reducing words to their base forms)
 - **Part-of-Speech Tagging** (identifying nouns, verbs, etc.)
 - **Named Entity Recognition** (identifying people, places, etc.)
 - **Syntactic Parsing** (analyzing sentence structure)
- **Extensive Corpus and Lexicon Resources:** Includes access to well-known text corpora (e.g., Brown, Reuters, WordNet) for training and evaluation.
- **Statistical and Machine Learning Tools:** Provides tools for text classification, sentiment analysis, and other statistical NLP tasks.

- **Educational Focus:** NLTK is ideal for beginners and academic purposes due to its intuitive design, modular structure, and extensive documentation.
- **Prototyping and Research:** Suitable for prototyping NLP models and linguistic research; allows quick experimentation with various NLP techniques.
- **Scalability Consideration:** More appropriate for small to medium-sized projects; libraries like spaCy or Hugging Face Transformers may be better for high-speed or large-scale NLP applications.
- **Community and Resources:** Supported by a robust community, with a wealth of tutorials and guides for learning NLP fundamentals.

[]: `pip install nltk` *nltk installation*

```
import nltk to download some NLTK resources (corpora and models)
nltk.download('punkt') For tokenization
nltk.download('wordnet') For lemmatization
nltk.download('averaged_perceptron_tagger') For part-of-speech tagging
nltk.download('punkt') For tokenization
nltk.download('wordnet') For lemmatization
nltk.download('averaged_perceptron_tagger') For part-of-speech tagging
```

Requirement already satisfied: nltk in /usr/local/lib/python3.10/dist-packages (3.9.1)

Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from nltk) (8.1.7)

Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from nltk) (1.4.2)

Requirement already satisfied: regex<=2021.8.3 in /usr/local/lib/python3.10/dist-packages (from nltk) (2024.9.11)

Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.6)

[nltk_data] Downloading package punkt to /root/nltk_data...

[nltk_data] Unzipping tokenizers/punkt.zip.

[nltk_data] Downloading package wordnet to /root/nltk_data...

[nltk_data] Package wordnet is already up-to-date!

[nltk_data] Downloading package averaged_perceptron_tagger to

[nltk_data] /root/nltk_data...

[nltk_data] Unzipping taggers/averaged_perceptron_tagger.zip.

[nltk_data] Downloading package punkt to /root/nltk_data...

[nltk_data] Package punkt is already up-to-date!

[nltk_data] Downloading package wordnet to /root/nltk_data...

[nltk_data] Package wordnet is already up-to-date!

[nltk_data] Downloading package averaged_perceptron_tagger to

[nltk_data] /root/nltk_data...

[nltk_data] Package averaged_perceptron_tagger is already up-to-

[nltk_data] date!

[]: True

SPACY

spaCy is a powerful, modern Python library specifically designed for high-performance natural language processing (NLP). Unlike NLTK, which is more academically focused, spaCy is optimized for real-world applications and is built with efficiency and scalability in mind. It is particularly popular in industry settings, as it offers a fast, production-ready pipeline that handles various NLP tasks with state-of-the-art accuracy.

- **Efficient NLP Pipeline:** Provides an optimized pipeline for text processing, with components for **tokenization**, **lemmatization**, **part-of-speech tagging**, **named entity recognition (NER)**, and **dependency parsing**.
- **Pre-trained Models:** Comes with pre-trained models for several languages, supporting multi-language NLP applications and offering fast deployment without needing extensive training data.
- **Word Vectors and Similarity:** Uses word vectors for semantic similarity, allowing it to capture the meaning of words and their relationships, which is crucial for tasks like similarity matching and entity linking.
- **Rule-Based and Statistical Matching:** Includes flexible rule-based matching and dependency parsing, which are useful for identifying specific patterns and structures in text.
- **Integration and Customization:** Easily integrates with other ML libraries and offers customization options, including training custom components and adding new languages or entity types.
- **Performance and Scalability:** Built on Cython for speed and efficiency, making it ideal for large datasets and real-time applications.
- **Advanced NLP Applications:** Supports deep learning integrations with libraries like PyTorch and TensorFlow, making it well-suited for advanced NLP tasks such as text classification and sequence labeling.

[]: `pip install spacy` *spacy installation*

`python -m spacy download en_core_web_sm` *need to download a language model (e.g., English)*

Requirement already satisfied: spacy in /usr/local/lib/python3.10/dist-packages (3.7.5)

Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /usr/local/lib/python3.10/dist-packages (from spacy) (3.0.12)

Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.10/dist-packages (from spacy) (1.0.5)

Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.10/dist-packages (from spacy) (1.0.10)

Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.10/dist-packages (from spacy) (2.0.8)

Requirement already satisfied: preshed<3.1.0,>=3.0.2 in

/usr/local/lib/python3.10/dist-packages (from spacy) (3.0.9)
 Requirement already satisfied: thinc<8.3.0,>=8.2.2 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (8.2.5)
 Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (1.1.3)
 Requirement already satisfied: srsly<3.0.0,>=2.4.3 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (2.4.8)
 Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (2.0.10)
 Requirement already satisfied: weasel<0.5.0,>=0.1.0 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (0.4.1)
 Requirement already satisfied: typer<1.0.0,>=0.3.0 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (0.13.0)
 Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (4.66.6)
 Requirement already satisfied: requests<3.0.0,>=2.13.0 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (2.32.3)
 Requirement already satisfied: pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (2.9.2)
 Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages
 (from spacy) (3.1.4)
 Requirement already satisfied: setuptools in /usr/local/lib/python3.10/dist-
 packages (from spacy) (75.1.0)
 Requirement already satisfied: packaging>=20.0 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (24.2)
 Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in
 /usr/local/lib/python3.10/dist-packages (from spacy) (3.4.1)
 Requirement already satisfied: numpy>=1.19.0 in /usr/local/lib/python3.10/dist-
 packages (from spacy) (1.26.4)
 Requirement already satisfied: language-data>=1.2 in
 /usr/local/lib/python3.10/dist-packages (from langcodes<4.0.0,>=3.2.0->spacy)
 (1.2.0)
 Requirement already satisfied: annotated-types>=0.6.0 in
 /usr/local/lib/python3.10/dist-packages (from
 pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4->spacy) (0.7.0)
 Requirement already satisfied: pydantic-core==2.23.4 in
 /usr/local/lib/python3.10/dist-packages (from
 pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4->spacy) (2.23.4)
 Requirement already satisfied: typing-extensions>=4.6.1 in
 /usr/local/lib/python3.10/dist-packages (from
 pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4->spacy) (4.12.2)
 Requirement already satisfied: charset-normalizer<4,>=2 in
 /usr/local/lib/python3.10/dist-packages (from requests<3.0.0,>=2.13.0->spacy)
 (3.4.0)
 Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-
 packages (from requests<3.0.0,>=2.13.0->spacy) (3.10)
 Requirement already satisfied: urllib3<3,>=1.21.1 in
 /usr/local/lib/python3.10/dist-packages (from requests<3.0.0,>=2.13.0->spacy)

(2.2.3)

Requirement already satisfied: certifi>=2017.4.17 in
/usr/local/lib/python3.10/dist-packages (from requests<3.0.0,>=2.13.0->spacy)
(2024.8.30)

Requirement already satisfied: blis<0.8.0,>=0.7.8 in
/usr/local/lib/python3.10/dist-packages (from thinc<8.3.0,>=8.2.2->spacy)
(0.7.11)

Requirement already satisfied: confection<1.0.0,>=0.0.1 in
/usr/local/lib/python3.10/dist-packages (from thinc<8.3.0,>=8.2.2->spacy)
(0.1.5)

Requirement already satisfied: click>=8.0.0 in /usr/local/lib/python3.10/dist-
packages (from typer<1.0.0,>=0.3.0->spacy) (8.1.7)

Requirement already satisfied: shellingham>=1.3.0 in
/usr/local/lib/python3.10/dist-packages (from typer<1.0.0,>=0.3.0->spacy)
(1.5.4)

Requirement already satisfied: rich>=10.11.0 in /usr/local/lib/python3.10/dist-
packages (from typer<1.0.0,>=0.3.0->spacy) (13.9.4)

Requirement already satisfied: cloudpathlib<1.0.0,>=0.7.0 in
/usr/local/lib/python3.10/dist-packages (from weasel<0.5.0,>=0.1.0->spacy)
(0.20.0)

Requirement already satisfied: smart-open<8.0.0,>=5.2.1 in
/usr/local/lib/python3.10/dist-packages (from weasel<0.5.0,>=0.1.0->spacy)
(7.0.5)

Requirement already satisfied: MarkupSafe>=2.0 in
/usr/local/lib/python3.10/dist-packages (from jinja2->spacy) (3.0.2)

Requirement already satisfied: marisa-trie>=0.7.7 in
/usr/local/lib/python3.10/dist-packages (from language-
data>=1.2->langcodes<4.0.0,>=3.2.0->spacy) (1.2.1)

Requirement already satisfied: markdown-it-py>=2.2.0 in
/usr/local/lib/python3.10/dist-packages (from
rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy) (3.0.0)

Requirement already satisfied: pygments<3.0.0,>=2.13.0 in
/usr/local/lib/python3.10/dist-packages (from
rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy) (2.18.0)

Requirement already satisfied: wrapt in /usr/local/lib/python3.10/dist-packages
(from smart-open<8.0.0,>=5.2.1->weasel<0.5.0,>=0.1.0->spacy) (1.16.0)

Requirement already satisfied: mdurl~=0.1 in /usr/local/lib/python3.10/dist-
packages (from markdown-it-py>=2.2.0->rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy)
(0.1.2)

Collecting en-core-web-sm==3.7.1

Downloading [https://github.com/explosion/spacy-](https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.7.1/en_core_web_sm-3.7.1-py3-none-any.whl)
models/releases/download/en_core_web_sm-3.7.1/en_core_web_sm-3.7.1-py3-none-
any.whl (12.8 MB)

12.8/12.8 MB

49.5 MB/s eta 0:00:00

Requirement already satisfied: spacy<3.8.0,>=3.7.2 in
/usr/local/lib/python3.10/dist-packages (from en-core-web-sm==3.7.1) (3.7.5)

Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in

/usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.0.12)
 Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.0.5)
 Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.0.10)
 Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.0.8)
 Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.0.9)
 Requirement already satisfied: thinc<8.3.0,>=8.2.2 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (8.2.5)
 Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.1.3)
 Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.4.8)
 Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.0.10)
 Requirement already satisfied: weasel<0.5.0,>=0.1.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.4.1)
 Requirement already satisfied: typer<1.0.0,>=0.3.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.13.0)
 Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (4.66.6)
 Requirement already satisfied: requests<3.0.0,>=2.13.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.32.3)
 Requirement already satisfied: pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.9.2)
 Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.1.4)
 Requirement already satisfied: setuptools in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (75.1.0)
 Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (24.2)

Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.4.1)

Requirement already satisfied: numpy>=1.19.0 in /usr/local/lib/python3.10/dist-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.26.4)

Requirement already satisfied: language-data>=1.2 in /usr/local/lib/python3.10/dist-packages (from langcodes<4.0.0,>=3.2.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.2.0)

Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.10/dist-packages (from pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.7.0)

Requirement already satisfied: pydantic-core==2.23.4 in /usr/local/lib/python3.10/dist-packages (from pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.23.4)

Requirement already satisfied: typing-extensions>=4.6.1 in /usr/local/lib/python3.10/dist-packages (from pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (4.12.2)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.4.0)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.10)

Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.2.3)

Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2024.8.30)

Requirement already satisfied: blis<0.8.0,>=0.7.8 in /usr/local/lib/python3.10/dist-packages (from thinc<8.3.0,>=8.2.2->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.7.11)

Requirement already satisfied: confection<1.0.0,>=0.0.1 in /usr/local/lib/python3.10/dist-packages (from thinc<8.3.0,>=8.2.2->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.1.5)

Requirement already satisfied: click>=8.0.0 in /usr/local/lib/python3.10/dist-packages (from typer<1.0.0,>=0.3.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (8.1.7)

Requirement already satisfied: shellingham>=1.3.0 in /usr/local/lib/python3.10/dist-packages (from typer<1.0.0,>=0.3.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.5.4)

Requirement already satisfied: rich>=10.11.0 in /usr/local/lib/python3.10/dist-packages (from typer<1.0.0,>=0.3.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (13.9.4)

Requirement already satisfied: cloudpathlib<1.0.0,>=0.7.0 in

/usr/local/lib/python3.10/dist-packages (from
 weasel<0.5.0,>=0.1.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.20.0)
 Requirement already satisfied: smart-open<8.0.0,>=5.2.1 in
 /usr/local/lib/python3.10/dist-packages (from
 weasel<0.5.0,>=0.1.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (7.0.5)
 Requirement already satisfied: MarkupSafe>=2.0 in
 /usr/local/lib/python3.10/dist-packages (from jinja2->spacy<3.8.0,>=3.7.2->en-
 core-web-sm==3.7.1) (3.0.2)
 Requirement already satisfied: marisa-trie>=0.7.7 in
 /usr/local/lib/python3.10/dist-packages (from language-
 data>=1.2->langcodes<4.0.0,>=3.2.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1)
 (1.2.1)
 Requirement already satisfied: markdown-it-py>=2.2.0 in
 /usr/local/lib/python3.10/dist-packages (from
 rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1)
 (3.0.0)
 Requirement already satisfied: pygments<3.0.0,>=2.13.0 in
 /usr/local/lib/python3.10/dist-packages (from
 rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1)
 (2.18.0)
 Requirement already satisfied: wrapt in /usr/local/lib/python3.10/dist-packages
 (from smart-open<8.0.0,>=5.2.1->weasel<0.5.0,>=0.1.0->spacy<3.8.0,>=3.7.2->en-
 core-web-sm==3.7.1) (1.16.0)
 Requirement already satisfied: mdurl~=0.1 in /usr/local/lib/python3.10/dist-
 packages (from markdown-it-
 py>=2.2.0->rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy<3.8.0,>=3.7.2->en-core-web-
 sm==3.7.1) (0.1.2)

Download and installation successful

You can now load the package via `spacy.load('en_core_web_sm')`

Restart to reload dependencies

If you are in a Jupyter or Colab notebook, you may need to restart Python in order to load all the package's dependencies. You can do this by selecting the 'Restart kernel' or 'Restart runtime' option.

TEXT PREPROCESSING

LOWERCASING

Lowercasing converts all characters in the text to lowercase. This step is crucial because in most NLP tasks, “Apple” and “apple” are considered the same word. By making everything lowercase, you ensure that the analysis is case-insensitive and doesn’t count uppercase and lowercase forms as different tokens, reducing data sparsity and improving model performance.

```
[ ]: text = "This is a Sample TEXT."
text_lower = text.lower()
print(text_lower) Output: "this is a sample text."
```

this is a sample text.

REMOVING PUNCTUATIONS

Punctuation, like commas, periods, and exclamation points, often doesn't carry meaningful information in many NLP tasks, so it is usually removed. This step helps clean up text for tasks like sentiment analysis and text classification where punctuation may create noise rather than add value. However, for some tasks, punctuation may be retained if it carries specific meaning (e.g., in sentiment analysis where "!" might denote stronger emotions).

```
[ ]: import string

text = "Hello, world! Welcome to NLP."
text_no_punctuation = text.translate(str.maketrans("", "", string.punctuation))
print(text_no_punctuation) Output: "Hello world Welcome to NLP"
```

Hello world Welcome to NLP

TOKENIZATION

Tokenization splits the text into smaller units, typically words or sentences, called tokens. These tokens are the basis for analysis in NLP. Tokenization allows the model to treat individual words as separate entities, making it easier to perform word-level processing, like frequency counting, stemming, or vectorization. In languages with complex grammar, specialized tokenizers may be required to handle the nuances correctly.

```
[ ]: import nltk
nltk.download('punkt_tab')
from nltk.tokenize import word_tokenize

text = "Tokenization is crucial for text processing."
tokens = word_tokenize(text)
print(tokens) Output: ['Tokenization', 'is', 'crucial', 'for', 'text', 'processing', '.']
```

[nltk_data] Downloading package punkt_tab to /root/nltk_data...

[nltk_data] Unzipping tokenizers/punkt_tab.zip.

['Tokenization', 'is', 'crucial', 'for', 'text', 'processing', '.']

```
[ ]: import spacy

Load spaCy's English model
nlp = spacy.load("en_core_web_sm")

Sample text
text = "Tokenization is crucial for text processing."

Process the text using spaCy
doc = nlp(text)

Tokenize the text
tokens = [token.text for token in doc]
```

```
print(tokens) Output: ['Tokenization', 'is', 'crucial', 'for', 'text', 'processing', '.']
```

```
['Tokenization', 'is', 'crucial', 'for', 'text', 'processing', '.']
```

STOP WORDS REMOVAL

Stop words are common words like “the,” “is,” and “in” that generally don’t add significant meaning in tasks like topic modeling or text classification. Removing stop words can reduce the size of the dataset, improve processing speed, and reduce noise, though in some cases, keeping stop words can help capture nuances in meaning, especially in sentiment analysis.

```
[ ]: import nltk

Download 'punkt_tab' before calling word_tokenize
nltk.download('punkt_tab')

from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

text = "This is an example sentence demonstrating stop words removal."
stop_words = set(stopwords.words('english'))
tokens = word_tokenize(text)
filtered_text = [word for word in tokens if word.lower() not in stop_words]
print(filtered_text) Output: ['example', 'sentence', 'demonstrating', 'stop', 'words', 'removal', '.']
```

```
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
```

```
['example', 'sentence', 'demonstrating', 'stop', 'words', 'removal', '.']
```

```
[nltk_data] Unzipping tokenizers/punkt_tab.zip.
```

```
[ ]: import spacy

Load spaCy's English model
nlp = spacy.load("en_core_web_sm")

Sample text
text = "This is an example sentence demonstrating stop words removal."

Process text using spaCy
doc = nlp(text)

Remove stop words using spaCy's stop word list
filtered_text = [token.text for token in doc if not token.is_stop]
print(filtered_text) Output: ['example', 'sentence', 'demonstrating', 'stop', 'words', 'removal', '.']
```

```
['example', 'sentence', 'demonstrating', 'stop', 'words', 'removal', '.']
```

POS TAGGING

POS tagging involves identifying the grammatical group (noun, verb, adjective, etc.) of each word in a sentence. This is an essential task in NLP that helps in understanding the structure and meaning of text. POS tagging can be done using predefined models or using libraries like spaCy or nltk.

```
[22]: import nltk
      from nltk.tokenize import word_tokenize
      from nltk import pos_tag

      Download necessary NLTK resources
      nltk.download('punkt')
      Download the 'averaged_perceptron_tagger_eng' data package.
      nltk.download('averaged_perceptron_tagger_eng')
      Download the 'punkt_tab' data package.
      nltk.download('punkt_tab')

      Example text
      text = "The quick brown fox jumps over the lazy dog."

      Tokenize the text
      tokens = word_tokenize(text)

      Perform POS tagging
      tagged_tokens = pos_tag(tokens)

      Print POS tags for each word
      for token, tag in tagged_tokens:
          print(f'{token}: {tag}')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger_eng to
[nltk_data]   /root/nltk_data...
[nltk_data]   Unzipping taggers/averaged_perceptron_tagger_eng.zip.
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]   Package punkt_tab is already up-to-date!
```

```
The: DT
quick: JJ
brown: NN
fox: NN
jumps: VBZ
over: IN
the: DT
lazy: JJ
dog: NN
```

∴

```
[19]: import spacy
```

Load the English language model

```
nlp = spacy.load("en_core_web_sm")
```

Example text

```
text = "The quick brown fox jumps over the lazy dog."
```

Process the text using spaCy

```
doc = nlp(text)
```

Print POS tags for each word

```
for token in doc:
```

```
    print(f'{token.text}: {token.pos_}')
```

/usr/local/lib/python3.10/dist-packages/spacy/util.py:1740: UserWarning: [W111] Jupyter notebook detected: if using `prefer_gpu()` or `require_gpu()`, include it in the same cell right before `spacy.load()` to ensure that the model is loaded on the correct device. More information:

<http://spacy.io/usage/v3jupyter-notebook-gpu>

```
warnings.warn(Warnings.W111)
```

The: DET

quick: ADJ

brown: ADJ

fox: NOUN

jumps: VERB

over: ADP

the: DET

lazy: ADJ

dog: NOUN

∴ PUNCT

STEMMING

Reduces words to their root form by removing suffixes, which helps to group related words together. However, it may produce non-standard words (e.g., “running” becomes “run”).

```
[ ]: from nltk.stem import PorterStemmer
```

```
text = "running runner runs"
```

```
stemmer = PorterStemmer()
```

```
stems = [stemmer.stem(word) for word in text.split()]
```

```
print(stems)    Output: ['run', 'runner', 'run']
```

```
['run', 'runner', 'run']
```

```
[ ]: import spacy

    Load spaCy's English model
nlp = spacy.load("en_core_web_sm")

    Sample text
text = "running runner runs"

    Process text and apply custom stemming (simple suffix removal)
stems = [word.text[:-3] if word.text.endswith('ing') else word.text[:-2] if_
sword.text.endswith(('ed', 'es')) else word.text[:-1] if word.text.
sendswith('s') else word.text for word in nlp(text)]

print(stems)  Output: ['run', 'runner', 'run']
```

['runn', 'runner', 'run']

LEMMATIZATION

Converts words to their base form (lemma), typically producing a more linguistically accurate root word based on context (e.g., “better” to “good”). Lemmatization tends to be more accurate than stemming and is useful in applications where grammatical correctness is important.

```
[ ]: import nltk
nltk.download('wordnet')  Download the wordnet dataset

from nltk.stem import WordNetLemmatizer

text = "running better ran"
lemmatizer = WordNetLemmatizer()
lemmas = [lemmatizer.lemmatize(word, pos="v") for word in text.split()]
print(lemmas)  Output: ['run', 'good', 'run']
```

[nltk_data] Downloading package wordnet to /root/nltk_data...

['run', 'better', 'run']

```
[ ]: import spacy

    Load spaCy's English model
nlp = spacy.load("en_core_web_sm")

    Sample text
text = "running better ran"

    Process the text using spaCy
doc = nlp(text)

    Lemmatization (using token.lemma_ instead of stemming)
```



```
lemmas = [token.lemma_ for token in doc]
print(lemmas)    Output: ['run', 'well', 'run']
```

['run', 'well', 'run']

N-GRAMS

N-grams are contiguous sequences of ‘n’ items from a given text. They are used to analyze text patterns and help in tasks like text classification, language modeling, and sentiment analysis. Common types are unigrams (1-gram), bigrams (2-grams), and trigrams (3-grams).

```
[23]: from sklearn.feature_extraction.text import CountVectorizer
```

Example text

```
text = ["I love programming", "Programming is fun"]
```

Initialize the CountVectorizer with n-gram range

```
vectorizer = CountVectorizer(ngram_range=(2, 2))    Bigrams
```

Fit and transform the text

```
X = vectorizer.fit_transform(text)
```

Print the bigrams

```
print(vectorizer.get_feature_names_out())
```

['is fun' 'love programming' 'programming is']

SMOOTHING

Smoothing techniques are applied in N-gram models to handle the issue of zero probabilities for unseen n-grams (sequences that are not present in the training data). Laplace smoothing is a common technique that adds a small value (usually 1) to each probability to ensure no probability is zero.

```
[24]: import random
import nltk
from nltk import bigrams
from nltk.probability import FreqDist, MLEProbDist
```

Example text

```
text = "I love programming and programming is fun".split()
```

Generate bigrams

```
bi_grams = list(bigrams(text))
```

Create frequency distribution for bigrams

```
fdist = FreqDist(bi_grams)
```

Apply MLE ProbDist for bigrams with smoothing

```
prob_dist = MLEProbDist(fdist)
```

```
Generate probability for a bigram (with smoothing)
prob = prob_dist.prob(('programming', 'is'))
print(f'Probability of "programming is": {prob}')
```

Probability of "programming is": 0.16666666666666666

```
[25]: import nltk
      from nltk import bigrams
      from nltk.probability import FreqDist, MLEProbDist

      Example text
      text = "I love programming and programming is fun".split()

      Generate bigrams
      bi_grams = list(bigrams(text))

      Create frequency distribution for bigrams
      fdist = FreqDist(bi_grams)

      Apply MLE ProbDist for bigrams with smoothing
      prob_dist = MLEProbDist(fdist)

      Generate probability for a bigram (with smoothing)
      prob = prob_dist.prob(('programming', 'is'))
      print(f'Probability of "programming is": {prob}')
```

Probability of "programming is": 0.16666666666666666

REMOVING NUMBERS

In many NLP applications, numbers don't carry meaningful information, especially in text analysis tasks where the primary focus is on words. For example, in sentiment analysis or text classification, numbers may not add much value. However, if numbers are important for a specific task, such as in financial or scientific documents, they might be retained.

```
[ ]: import re

      text = "There are 3 cats and 7 dogs."
      text_no_numbers = re.sub(r'\d+', "", text)
      print(text_no_numbers) Output: "There are cats and dogs."
```

There are cats and dogs.

HANDLING CONTRACTIONS

Expanding contractions, like changing "can't" to "cannot," helps in text normalization, ensuring that different forms of a word or phrase are treated the same. Contractions can lead to inconsistencies in analysis if not expanded, especially in sentiment analysis or keyword extraction tasks.

```
[ ]: !pip install contractions
```

Requirement already satisfied: contractions in /usr/local/lib/python3.10/dist-packages (0.1.73)
Requirement already satisfied: textsearch>=0.0.21 in /usr/local/lib/python3.10/dist-packages (from contractions) (0.0.24)
Requirement already satisfied: anyascii in /usr/local/lib/python3.10/dist-packages (from textsearch>=0.0.21->contractions) (0.3.2)
Requirement already satisfied: pyahocorasick in /usr/local/lib/python3.10/dist-packages (from textsearch>=0.0.21->contractions) (2.1.0)

```
[ ]: from contractions import fix
```

```
text = "I can't do this because it isn't working."  
expanded_text = fix(text)  
print(expanded_text) Output: "I cannot do this because it is not working."
```

I cannot do this because it is not working.

REMOVING SPECIAL CHARACTERS

Special characters (e.g., @, , \$, ^) are usually removed because they often don't carry semantic meaning in natural language. However, in social media analysis, special characters like “” for hashtags may be retained as they provide contextual information. Removing these characters is beneficial for tasks where only plain text is analyzed.

```
[ ]: import re
```

```
text = "Hello @world! NLP is ^awesome." text_no_special_chars  
= re.sub(r'^A-Za-z0-9\s', "", text)  
print(text_no_special_chars) Output: "Hello world NLP is awesome"
```

Hello world NLP is awesome

HANDLING MISSPELLINGS

Correcting misspelled words improves data quality, making the text more standardized and reducing noise. Misspellings can cause similar words to be treated as different tokens, which reduces the effectiveness of text-based machine learning models, so correcting them helps models understand the text better.

```
[ ]: !pip install pyspellchecker
```

Collecting pyspellchecker

Downloading pyspellchecker-0.8.1-py3-none-any.whl.metadata (9.4 kB)

Downloading pyspellchecker-0.8.1-py3-none-any.whl (6.8 MB)

6.8/6.8 MB

45.5 MB/s eta 0:00:00

Installing collected packages: pyspellchecker

Successfully installed pyspellchecker-0.8.1

```
[ ]: from spellchecker import SpellChecker

spell = SpellChecker()
text = "I am lerning NLP and data scince."
corrected_text = " ".join([spell.correction(word) for word in text.split()])
print(corrected_text)  Output: "I am learning nap and data since."
```

I am learning nap and data since

NORMALIZATION

Normalization standardizes various elements, such as abbreviations, slang, or alternative spellings, to a common format. For example, converting “gonna” to “going to” or “u” to “you” makes the text more consistent, enhancing model interpretability and improving the effectiveness of NLP tasks.

```
[ ]: text = "I'm gonna use NLP in my proj."
normalized_text = text.replace("gonna", "going to").replace("proj", "project")
print(normalized_text)  Output: "I'm going to use NLP in my project."
```

I'm going to use NLP in my project.

VECTORIZATION

Vectorization transforms text into numerical form for machine learning models. Common methods include: - **Bag of Words**: Represents text as word frequencies. - **TF-IDF (Term Frequency-Inverse Document Frequency)**: Adjusts word frequency based on its commonness across all documents, giving less weight to common words. - **Word Embeddings**: Word2Vec and GloVe capture contextual meaning, improving the representation of words in relation to others.

```
[ ]: (TF-IDF)

from sklearn.feature_extraction.text import TfidfVectorizer

texts = ["I love NLP", "NLP is amazing"]
vectorizer = TfidfVectorizer()
vectorized_texts = vectorizer.fit_transform(texts)
print(vectorized_texts.toarray())  Output: [[0.707 0.707 0.          ], [0.          0.
0.707 0.707]]
```

```
[[0.          0.          0.81480247 0.57973867]
 [0.6316672 0.6316672 0.          0.44943642]]
```

COSINE SIMILARITY

Cosine Similarity is a metric used to measure how similar two vectors are in terms of their angle (cosine of the angle between them). It's commonly used in text mining and natural language processing (NLP) to measure the similarity between two documents or texts by comparing their vector representations.

The cosine similarity score ranges from -1 to 1:

- 1 indicates that the vectors are identical.

- 0 indicates that the vectors are orthogonal (no similarity).
- -1 indicates that the vectors are completely opposite.

$$\text{Cosine Similarity} = (\mathbf{A} \cdot \mathbf{B}) / (||\mathbf{A}|| \ ||\mathbf{B}||)$$

Where: - $\mathbf{A} \cdot \mathbf{B}$ is the dot product of vectors A and B. - $||\mathbf{A}||$ and $||\mathbf{B}||$ are the magnitudes (norms) of vectors A and B.

```
[27]: from sklearn.feature_extraction.text import TfidfVectorizer
      from sklearn.metrics.pairwise import cosine_similarity

      Example text documents
      text1 = "I love programming in Python"
      text2 = "Python programming is amazing"

      Create the TF-IDF Vectorizer
      vectorizer = TfidfVectorizer()

      Convert the text documents to a matrix of TF-IDF features
      tfidf_matrix = vectorizer.fit_transform([text1, text2])

      Compute cosine similarity
      cosine_sim = cosine_similarity(tfidf_matrix[0:1], tfidf_matrix[1:2])

      Print the cosine similarity
      print(f"Cosine Similarity between the two documents: {cosine_sim[0][0]}")
```

Cosine Similarity between the two documents: 0.3360969272762575

WORD CLOUD

A word cloud is a visual representation of the most frequent words in a text, with word size proportional to frequency. It is often used to visualize text data and show key themes. Popular libraries like WordCloud in Python can help generate such visualizations.

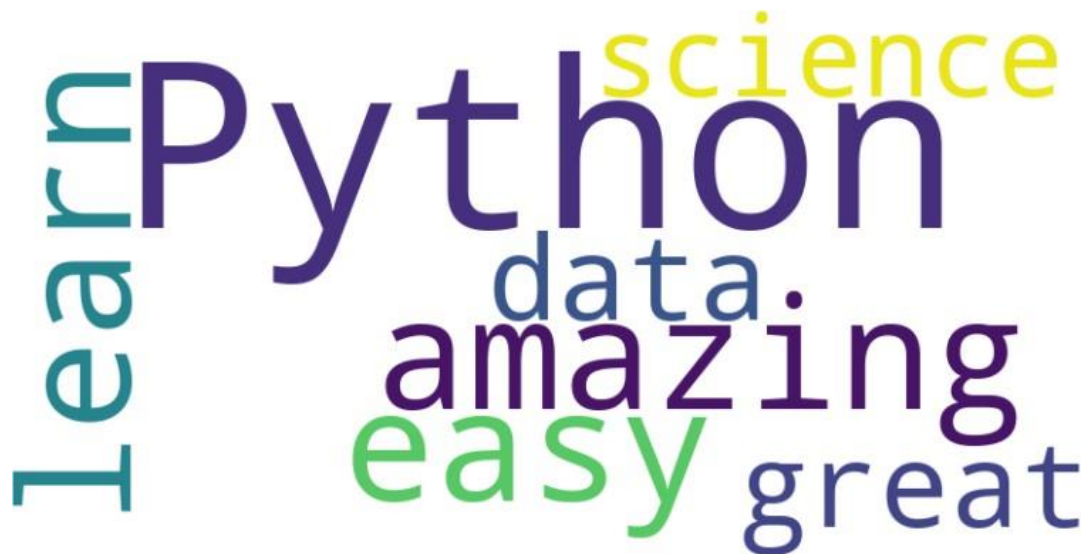
```
[26]: from wordcloud import WordCloud
      import matplotlib.pyplot as plt

      Example text
      text = "Python is amazing, Python is easy to learn, and Python is great for_
             ,data science."

      Generate word cloud
      wordcloud = WordCloud(width=800, height=400, background_color="white").
                      .generate(text)

      Display the word cloud plt.figure(figsize=(10,
      5)) plt.imshow(wordcloud,
      interpolation="bilinear")
```

```
plt.axis("off")  
plt.show()
```



TRANSFORMERS

Transformer is a deep learning model architecture used primarily in Natural Language Processing (NLP). Introduced in 2017, it relies on self-attention mechanisms to process entire sequences of data in parallel, instead of sequentially like RNNs or LSTMs. This allows Transformers to efficiently capture long-range dependencies and relationships between words in a sentence. The model consists of two main parts: the encoder, which processes input data, and the decoder, which generates output. Transformers have become the foundation of powerful NLP models like BERT and GPT, enabling significant advancements in tasks such as language translation, text generation, and classification.

NEED OF TRANSFORMER

1. **Handling Long-Range Dependencies:** Traditional models like RNNs and LSTMs struggled with long-range dependencies (i.e., understanding relationships between distant words in a sentence). Transformers, through their self-attention mechanism, can capture these long-range dependencies more effectively, making them ideal for tasks such as language translation and text generation.
2. **Parallelization:** Unlike RNNs, which process sequences step-by-step, Transformers process the entire sequence simultaneously. This parallelization allows for much faster training and better scalability, especially on large datasets.
3. **Scalability:** Transformers are highly scalable, allowing them to be trained on large datasets with billions of parameters, leading to more accurate and powerful models. This scalability has enabled the development of large pre-trained models like BERT, GPT, and T5.
4. **Flexibility:** The Transformer architecture can be adapted for a wide range of tasks, from

machine translation to text classification, question answering, and even image processing. Its modular nature allows for easy customization and fine-tuning for specific tasks.

5. **Improved Performance:** Transformers consistently outperform traditional models in many NLP tasks, setting new benchmarks for performance in areas like sentiment analysis, summarization, and machine translation.
6. **Self-Attention:** The self-attention mechanism in Transformers allows the model to weigh the importance of different parts of the input, which makes it particularly useful for understanding context and semantics in natural language, leading to better decision-making and prediction accuracy.

EXAMPLE

There is a word - 'Point', and we use it in two different contexts given below

- The needle has a sharp point.
- It is not polite to point at people.

Here, the word 'Point' has two different contexts in both of the sentences, but when embedding is done the context is not taken into consideration. Therefore, there was a need for a different architecture — Transformer.

COMPONENTS OF TRANSFORMER

1. **Self-Attention Mechanism:**
 - Helps the model understand relationships between all words in a sentence, regardless of their position.
 - It allows the model to focus on the most important words when processing each word.
2. **Multi-Head Attention:**
 - Multiple self-attention operations are done at once, allowing the model to focus on different aspects of the input simultaneously.
3. **Positional Encoding:**
 - Adds information about the order of words in the sequence since Transformers don't process words one by one like RNNs. It helps the model understand word positions.
4. **Feedforward Neural Networks:**
 - After the attention layers, each word's representation is passed through a simple neural network to process it further and refine its understanding.
5. **Layer Normalization:**
 - Normalizes the output of each layer to improve training speed and stability.
6. **Residual Connections:**
 - Shortcuts between layers that help gradients flow easily during training, preventing problems like vanishing gradients.
7. **Encoder:**
 - The part of the Transformer that processes the input data, understanding its meaning.
8. **Decoder:**
 - The part of the Transformer that generates the output (like translating a sentence or generating text). It uses information from the encoder to do this.

WORKING OF TRANSFORMERS

1. Input Embedding:

- The first step involves converting the input words or tokens (e.g., from a sentence) into numerical vectors using embeddings. These embeddings represent each word as a dense vector of numbers.

2. Positional Encoding:

- Since Transformers process all tokens in parallel (instead of sequentially like RNNs), they need information about the order of words in the sequence. **Positional encoding** adds this order information to the word embeddings, ensuring the model knows the position of each word in the sequence.

3. Self-Attention Mechanism:

- This is the core of how Transformers work. For each word in the input, the self-attention mechanism calculates how much focus or attention it should give to all other words in the sequence.
- The model computes three vectors for each word: **Query (Q)**, **Key (K)**, and **Value (V)**. These vectors help determine the relationship between words:
 - **Query**: Represents the current word.
 - **Key**: Represents all the other words in the sequence.
 - **Value**: Contains the actual information from the other words.
- The attention score between words is computed by taking the dot product of the Query and Key vectors. This score tells the model how much attention to give to other words when processing the current word.
- The result is a weighted sum of all words, where more relevant words have higher attention weights.

4. Multi-Head Attention:

- Instead of having a single attention mechanism, **multi-head attention** performs multiple self-attention operations in parallel, each with different parameters (or “heads”). This allows the model to focus on different parts of the sentence at the same time, capturing various relationships between words.

5. Feedforward Neural Networks:

- After the attention mechanism, the output is passed through a simple **feedforward neural network** (a series of linear transformations and activation functions). This step helps refine the representation of each word by applying non-linear transformations.

6. Residual Connections and Layer Normalization:

- To stabilize training, **residual connections** (shortcuts) are used around both the attention and feedforward layers. These help the model avoid problems like vanishing gradients during training.
- **Layer normalization** is applied after each step to standardize the output, ensuring that the model converges efficiently during training.

7. Stacking Layers:

- The Transformer model consists of multiple layers of attention and feedforward neural networks, stacked on top of each other. Each layer refines the understanding of the input, helping the model build more complex representations as the data passes through.

8. Encoder and Decoder:

- The **encoder** processes the input sequence and generates a set of contextualized word representations.
- The **decoder** generates the output sequence, based on the encoder’s output and previous

words generated (in tasks like translation or text generation). The decoder has an additional attention layer that attends to the encoder's output.

9. Output:

- Finally, after passing through the encoder and decoder (in the case of models like machine translation), the output is a sequence of predictions (like translated words or generated text). The output sequence is produced using the softmax function, which assigns probabilities to different possible words or tokens.

HUGGING FACE

Hugging Face is a company and an open-source platform that is widely recognized for its contributions to Natural Language Processing (NLP) and machine learning. Hugging Face is particularly known for its Transformers library, which provides tools and pre-trained models that make it easier to work with advanced machine learning models, especially those based on the transformer architecture (such as BERT, GPT, T5, and RoBERTa).

Key Features of Hugging Face:

1. Transformers Library:

- **Hugging Face's Transformers library** is one of the most popular tools in the NLP community. It offers a large collection of pre-trained models that can be fine-tuned for a wide variety of NLP tasks, such as:
 - Text classification
 - Named entity recognition (NER)
 - Question answering
 - Text generation
 - Machine translation
- It supports models like BERT, GPT, T5, and more, which can be used with both **PyTorch** and **TensorFlow**.

2. Pre-trained Models:

- Hugging Face hosts thousands of pre-trained models that can be used out-of-the-box for a variety of NLP tasks. These models are trained on massive datasets and have achieved state-of-the-art performance on tasks like sentiment analysis, language understanding, and summarization.
- These pre-trained models are easily accessible via the Hugging Face Model Hub, where developers and researchers can find, upload, and share models.

3. Ease of Use:

- The **Transformers library** abstracts much of the complexity of working with transformers. It provides simple APIs to load pre-trained models, tokenize text, and perform inference or fine-tune models on custom datasets.

4. Integration with Other Tools:

- Hugging Face integrates well with other tools, such as **Datasets**, **Accelerate**, and **Optimum**, allowing users to work seamlessly across different stages of model training, optimization, and deployment.

5. Community and Open Source:

- Hugging Face fosters a large and active **open-source community**. Their models and tools are available to everyone, and many researchers and companies contribute to the development and improvement of models.
- They also provide various tutorials, documentation, and resources to help developers get started with NLP tasks.

6. Hub for Models and Datasets:

- Hugging Face provides a centralized **Model Hub**, which allows users to explore, upload, and share pre-trained models.
- The **Datasets Hub** offers access to a wide range of datasets for training and evaluation.

[8]: *Install Hugging Face's transformers library (if not already installed)*

```
!pip install transformers
```

Requirement already satisfied: transformers in /usr/local/lib/python3.10/dist-packages (4.46.2)

Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from transformers) (3.16.1)

Requirement already satisfied: huggingface-hub<1.0,>=0.23.2 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.26.2)

Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (1.26.4)

Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from transformers) (24.2)

Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (6.0.2)

Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (2024.9.11)

Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from transformers) (2.32.3)

Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.4.5)

Requirement already satisfied: tokenizers<0.21,>=0.20 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.20.3)

Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.10/dist-packages (from transformers) (4.66.6)

Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub<1.0,>=0.23.2->transformers) (2024.10.0)

Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub<1.0,>=0.23.2->transformers) (4.12.2)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.4.0)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.10)

Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (2.2.3)

Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (2024.8.30)

```
[9]: from transformers import BertTokenizer, BertForSequenceClassification, pipeline
```

Step 1: Load the tokenizer and model from Hugging Face's pre-trained models

```
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
```

```
model = BertForSequenceClassification.from_pretrained('bert-base-uncased')
```

Step 2: Create a sentiment analysis pipeline

```
classifier = pipeline('sentiment-analysis', model=model, tokenizer=tokenizer)
```

Step 3: Input text for sentiment analysis

```
texts = [  
    "I love using Hugging Face models for natural language processing tasks!",  
    "I am so frustrated with this error, it just won't work!",  
    "This is an amazing experience, I am really enjoying it."  
]
```

Step 4: Perform sentiment analysis on each text

```
results = classifier(texts)
```

Step 5: Display the results

```
for i, result in enumerate(results):  
    print(f"Text: {texts[i]}")  
    print(f"Sentiment: {result['label']}, Confidence: {result['score']:.4f}\n")
```

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initialized:

['classifier.bias', 'classifier.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Hardware accelerator e.g. GPU is available in the environment, but no `device` argument is passed to the `Pipeline` object. Model will be on CPU.

Text: I love using Hugging Face models for natural language processing tasks!

Sentiment: LABEL_1, Confidence: 0.5491

Text: I am so frustrated with this error, it just won't work!

Sentiment: LABEL_1, Confidence: 0.5806

Text: This is an amazing experience, I am really enjoying it.

Sentiment: LABEL_1, Confidence: 0.5502

For further information refer, <https://huggingface.co/docs>

VARIANTS OF TRANSFORMERS

GPT

GPT is a deep learning model designed for natural language understanding and generation. It uses the Transformer architecture, specifically the **decoder** part, and is pretrained on a large corpus

of text data. GPT generates coherent and contextually relevant text based on a given prompt by predicting the next word in a sequence, which makes it powerful for tasks like text generation, conversational AI, and more.

Key Points about GPT:

1. Architecture:

- **Decoder-only Model:** GPT uses only the decoder part of the Transformer architecture, processing text in a **left-to-right** manner (sequentially).
- It generates text by predicting the next word based on the previous context, making it particularly good for generating coherent, flowing text.

2. Pretraining:

- **Causal Language Modeling:** GPT is pretrained using a causal language model, where it learns to predict the next word in a sentence from the previous words.
- The pretraining is done on massive text datasets from various sources, helping the model to learn patterns, syntax, and grammar across a wide range of topics.

3. Fine-Tuning:

- After pretraining, GPT can be **fine-tuned** for specific tasks by adjusting the model on smaller, task-specific datasets (e.g., summarization, translation, question answering).
- Fine-tuning is usually done with supervised learning, where the model adapts to a particular type of text generation or NLP task.

4. Text Generation:

- GPT excels in **generating human-like text**. Given a prompt or a starting sentence, it can generate coherent, contextually appropriate text that follows the style and content of the input.
- It is used in applications like automated storytelling, chatbots, code generation, and content creation.

5. Autoregressive Nature:

- The model generates text **one token at a time** by considering the previously generated tokens and predicting the next one. This makes it highly effective for tasks like autocomplete and dialogue generation.

6. Scalability:

- GPT models (especially GPT-2 and GPT-3) are highly **scalable** with millions or even billions of parameters, allowing them to learn rich language representations and handle complex tasks.

7. Applications:

- **Chatbots:** GPT is widely used in building conversational agents like chatbots.
- **Creative Writing:** It is used for generating creative content, including stories, poems, and articles.
- **Code Generation:** GPT models like Codex are capable of generating programming code based on natural language descriptions.
- **Text Completion:** GPT is used for applications where text needs to be completed or predicted, such as email composition, content generation, and more.

8. Versions:

- **GPT-1:** The original version, introduced in 2018, demonstrated the potential of generative pretraining but had limited capabilities.
- **GPT-2:** Released in 2019, it was much larger, generating more coherent text but still prone to mistakes.
- **GPT-3:** Released in 2020, it is one of the largest language models ever created with 175

billion parameters, capable of producing highly realistic text and performing a variety of tasks without specific fine-tuning.

[37]: `pip install transformers torch`

Requirement already satisfied: transformers in /usr/local/lib/python3.10/dist-packages (4.46.2)

Requirement already satisfied: torch in /usr/local/lib/python3.10/dist-packages (2.5.0+cu121)

Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from transformers) (3.16.1)

Requirement already satisfied: huggingface-hub<1.0,>=0.23.2 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.26.2)

Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (1.26.4)

Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from transformers) (24.2)

Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (6.0.2)

Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transformers) (2024.9.11)

Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from transformers) (2.32.3)

Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.4.5)

Requirement already satisfied: tokenizers<0.21,>=0.20 in /usr/local/lib/python3.10/dist-packages (from transformers) (0.20.3)

Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.10/dist-packages (from transformers) (4.66.6)

Requirement already satisfied: typing-extensions>=4.8.0 in /usr/local/lib/python3.10/dist-packages (from torch) (4.12.2)

Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch) (3.4.2)

Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch) (3.1.4)

Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packages (from torch) (2024.10.0)

Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.10/dist-packages (from torch) (1.13.1)

Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from sympy==1.13.1->torch) (1.3.0)

Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch) (3.0.2)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.4.0)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->transformers) (3.10)

Requirement already satisfied: urllib3<3,>=1.21.1 in

/usr/local/lib/python3.10/dist-packages (from requests->transformers) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in
/usr/local/lib/python3.10/dist-packages (from requests->transformers)
(2024.8.30)

```
[5]: from transformers import GPT2LMHeadModel, GPT2Tokenizer

    Load pre-trained model and tokenizer
model_name = "gpt2"    You can use 'gpt2-medium', 'gpt2-large', or 'gpt2-xl'
                        for larger models
model = GPT2LMHeadModel.from_pretrained(model_name)
tokenizer = GPT2Tokenizer.from_pretrained(model_name)

    Encode input text (prompt) to generate from
input_text = "Once upon a time"
input_ids = tokenizer.encode(input_text, return_tensors='pt')

    Generate text from the model
output = model.generate(input_ids, max_length=100, num_return_sequences=1,
                        no_repeat_ngram_size=2, temperature=0.7)

    Decode the generated output to readable text
generated_text = tokenizer.decode(output[0], skip_special_tokens=True)

    Print the generated text
print(generated_text)
```

The attention mask and the pad token id were not set. As a consequence, you may observe unexpected behavior. Please pass your input's `attention_mask` to obtain reliable results.

Setting `pad_token_id` to `eos_token_id`:None for open-end generation.

Once upon a time, the world was a place of great beauty and great danger. The world of the gods was the place where the great gods were born, and where they were to live.

The world that was created was not the same as the one that is now. It was an endless, endless world. And the Gods were not born of nothing. They were created of a single, single thing. That was why the universe was so beautiful. Because the cosmos was made of two

For further information refer, <https://www.manning.com/books/how-gpt-works>

BERT

BERT is a deep learning model developed for natural language understanding tasks. Unlike traditional models, BERT processes text bidirectionally, meaning it considers the context from both the left and the right of each word in a sentence. This bidirectional approach allows BERT to generate better word representations, making it highly effective for tasks like text classification, question answering, and named entity recognition.

Key Points about BERT:

1. Architecture:

- **Encoder-only Model:** BERT uses only the **encoder** part of the Transformer architecture. It processes the entire input sequence at once and generates contextualized embeddings for each word by looking at both its left and right context.
- This bidirectional approach is a key feature, allowing BERT to understand a word in relation to all other words in the sentence.

2. Pretraining:

- **Masked Language Model (MLM):** During pretraining, BERT uses the MLM task, where it randomly masks some of the words in a sentence and trains the model to predict them. This enables BERT to learn deeper, more context-sensitive word representations.
- **Next Sentence Prediction (NSP):** BERT is also trained with NSP, where it learns to predict if one sentence logically follows another. This helps the model understand sentence relationships and context beyond individual words.

3. Fine-Tuning:

- After pretraining, BERT can be **fine-tuned** for specific downstream tasks, such as text classification, sentiment analysis, question answering, and named entity recognition. Fine-tuning involves training the model with labeled data for the specific task, with minimal changes to the architecture.

4. Bidirectional Context Understanding:

- BERT processes input bidirectionally, considering both the left and right context of each word. This helps it understand the full meaning of a word in context, unlike earlier models like GPT, which process text unidirectionally.

5. Transfer Learning:

- One of BERT's strengths is **transfer learning**. Once pretrained, BERT can be fine-tuned on a relatively small dataset for a variety of NLP tasks, making it highly efficient and adaptable to different applications.

6. Applications:

- **Question Answering:** BERT has been highly effective in tasks like question answering, where it can understand the context of both the question and the passage of text to provide accurate answers.
- **Sentiment Analysis:** BERT is commonly used for analyzing the sentiment of a piece of text, helping in applications like customer feedback analysis and social media monitoring.
- **Text Classification:** BERT can be used for classifying text into categories, such as spam detection or topic categorization.
- **Named Entity Recognition (NER):** BERT can identify and classify named entities (such as people, organizations, or locations) in text.

7. Version Variants:

- **BERT-Base:** The original BERT model with 12 layers and 110 million parameters.
- **BERT-Large:** A larger version with 24 layers and 340 million parameters, providing better performance on more complex tasks.

8. Performance:

- BERT set new benchmarks in multiple NLP tasks when it was released, outperforming traditional models in tasks like the **SQuAD** (Stanford Question Answering Dataset) and **GLUE** (General Language Understanding Evaluation) benchmarks.

9. Multilingual BERT:

- A multilingual version of BERT (mBERT) was also introduced, enabling the model to

work with multiple languages. It was trained on text from 104 languages, making it useful for cross-lingual tasks like translation and multilingual text classification.

```
[6]: from transformers import BertTokenizer, BertForSequenceClassification
import torch

Load pre-trained model and tokenizer
model_name = "bert-base-uncased" You can use other variants like_
s 'bert-large-uncased'
model = BertForSequenceClassification.from_pretrained(model_name, num_labels=2)
s num_labels=2 for binary classification
tokenizer = BertTokenizer.from_pretrained(model_name)

Example input text
input_text = "I love this movie!"

Tokenize input text
inputs = tokenizer(input_text, return_tensors="pt", truncation=True,
s padding=True, max_length=512)

Forward pass through the model
with torch.no_grad():
    outputs = model(**inputs)

Get the prediction (logits)
logits = outputs.logits

Apply softmax to get probabilities
probabilities = torch.nn.functional.softmax(logits, dim=-1)

Print the probabilities for each class
print("Class Probabilities:", probabilities)

Get the predicted class (0 or 1 in this case)
predicted_class = torch.argmax(probabilities, dim=-1).item()
print("Predicted Class:", predicted_class)
```

config.json: 0%| | 0.00/570 [00:00<?, ?B/s]

model.safetensors: 0%| | 0.00/440M [00:00<?, ?B/s]

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initialized:

['classifier.bias', 'classifier.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

tokenizer_config.json: 0%| | 0.00/48.0 [00:00<?, ?B/s]

vocab.txt: 0%| | 0.00/232k [00:00<?, ?B/s]


```
tokenizer.json: 0%|          | 0.00/466k [00:00<?, ?B/s]
```

```
Class Probabilities: tensor([[0.6700, 0.3300]])
```

```
Predicted Class: 0
```

For further information refer, https://huggingface.co/transformers/v3.0.2/model_doc/bert.html

USE CASE OF PRE-TRAINED MODELS

TEXT CLASSIFICATION

Text classification involves categorizing text into predefined labels, such as spam detection or sentiment analysis. Pre-trained models like BERT or RoBERTa can be fine-tuned to classify text into categories.

```
[10]: from transformers import pipeline

      Load pre-trained BERT model for text classification
classifier = pipeline("text-classification", model="bert-base-uncased")

      Sample text for sentiment analysis
text = "I love using Hugging Face transformers!"

      Get sentiment result
result = classifier(text)

      print(result)
```

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initialized:

```
['classifier.bias', 'classifier.weight']
```

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Hardware accelerator e.g. GPU is available in the environment, but no `device` argument is passed to the `Pipeline` object. Model will be on CPU.

```
[{'label': 'LABEL_0', 'score': 0.5556365847587585}]
```

NAMED ENTITY RECOGNITION

NER is used to identify and classify entities like names, dates, and locations within a text. Pre-trained models like BERT and spaCy's NER model can extract these entities.

```
[11]: import spacy

      Load spaCy's pre-trained NER model
nlp = spacy.load("en_core_web_sm")

      Sample text for NER
text = "Apple was founded by Steve Jobs in Cupertino, California in 1976."

      Process the text with spaCy
```

```
doc = nlp(text)

Print named entities
for ent in doc.ents:
    print(ent.text, ent.label_)
```

```
Apple ORG
Steve Jobs PERSON
Cupertino GPE
California GPE
1976 DATE
```

MACHINE TRANSLATION

Machine translation converts text from one language to another. Models like mBERT, XLM-R, and T5 can handle translation tasks between languages by leveraging their multilingual training.

```
[12]: from transformers import pipeline

Load pre-trained model for translation (English to French)
translator = pipeline("translation_en_to_fr", model="t5-base")

Sample text for translation
text = "Hello, how are you?"

Get translation result
result = translator(text)

print(result)
```

```
config.json: 0%|          | 0.00/1.21k [00:00<?, ?B/s]
model.safetensors: 0%|          | 0.00/892M [00:00<?, ?B/s]
generation_config.json: 0%|          | 0.00/147 [00:00<?, ?B/s]
spiece.model: 0%|          | 0.00/792k [00:00<?, ?B/s]
tokenizer.json: 0%|          | 0.00/1.39M [00:00<?, ?B/s]
```

Hardware accelerator e.g. GPU is available in the environment, but no `device` argument is passed to the `Pipeline` object. Model will be on CPU.

```
[{'translation_text': 'Bonjour, comment êtes-vous?'}]
```

SUMMARIZATION

Summarization involves shortening a long document while retaining key information. Pre-trained models like BART and T5 can generate concise summaries of lengthy texts.

```
[13]: from transformers import pipeline

Load pre-trained BART model for summarization
```

```
summarizer = pipeline("summarization", model="facebook/bart-large-cnn")
```

Sample long text for summarization

```
text = """
```

Hugging Face provides an open-source library for Natural Language Processing (NLP).

It has a large number of pre-trained models for various NLP tasks, including sentiment analysis, text classification, and more.

These models help developers to quickly integrate state-of-the-art NLP functionality into their applications.

```
"""
```

Get summary result

```
result = summarizer(text, max_length=50, min_length=25, do_sample=False)
```

```
print(result)
```

```
config.json: 0%|          | 0.00/1.58k [00:00<?, ?B/s]
```

```
model.safetensors: 0%|          | 0.00/1.63G [00:00<?, ?B/s]
```

```
generation_config.json: 0%|          | 0.00/363 [00:00<?, ?B/s]
```

```
vocab.json: 0%|          | 0.00/899k [00:00<?, ?B/s]
```

```
merges.txt: 0%|          | 0.00/456k [00:00<?, ?B/s]
```

```
tokenizer.json: 0%|          | 0.00/1.36M [00:00<?, ?B/s]
```

Hardware accelerator e.g. GPU is available in the environment, but no `device` argument is passed to the `Pipeline` object. Model will be on CPU.

```
[{'summary_text': 'Hugging Face provides an open-source library for Natural Language Processing (NLP) It has a large number of pre-trained models for various NLP tasks.'}]
```

TEXT GENERATION

Text generation creates coherent and contextually relevant text based on an input prompt. Models like GPT-2, GPT-3, and DialoGPT are trained for such tasks and can generate creative text for applications like chatbots or content creation.

```
[14]: from transformers import pipeline
```

Load pre-trained GPT-2 model for text generation

```
generator = pipeline("text-generation", model="gpt2")
```

Sample text prompt

```
text = "Once upon a time, in a land far, far away,"
```

Generate continuation of the text

```
result = generator(text, max_length=50, num_return_sequences=1)
```

```
print(result)
```

Hardware accelerator e.g. GPU is available in the environment, but no ``device`` argument is passed to the ``Pipeline`` object. Model will be on CPU.
Truncation was not explicitly activated but ``max_length`` is provided a specific value, please use ``truncation=True`` to explicitly truncate examples to max length. Defaulting to 'longest_first' truncation strategy. If you encode pairs of sequences (GLUE-style) with the tokenizer you can select this strategy more precisely by providing a specific strategy to ``truncation``.
Setting ``pad_token_id`` to ``eos_token_id``:None for open-end generation.

```
{'generated_text': 'Once upon a time, in a land far, far away, and in the middle of the ocean so distant that by its very entrance into the world a man who was in a state so far away could see only the smallest part of the world,'}]
```

QUESTION ANSWERING

Question answering models are designed to extract answers from a given context based on the posed question. Models like BERT and ALBERT are fine-tuned to answer questions effectively.

```
[15]: from transformers import pipeline
```

```
Load pre-trained BERT model for question answering
```

```
qa_pipeline = pipeline("question-answering",  
                        model="bert-large-uncased-whole-word-masking-finetuned-squad")
```

```
Sample context and question
```

```
context = "Hugging Face is a company that provides state-of-the-art NLP models."  
question = "What does Hugging Face provide?"
```

```
Get answer from the context
```

```
result = qa_pipeline(question=question, context=context)
```

```
print(result)
```

```
config.json: 0%|          | 0.00/443 [00:00<?, ?B/s]
```

```
model.safetensors: 0%|          | 0.00/1.34G [00:00<?, ?B/s]
```

Some weights of the model checkpoint at bert-large-uncased-whole-word-masking-finetuned-squad were not used when initializing BertForQuestionAnswering:

```
['bert.pooler.dense.bias', 'bert.pooler.dense.weight']
```

- This IS expected if you are initializing BertForQuestionAnswering from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPreTraining model).

- This IS NOT expected if you are initializing BertForQuestionAnswering from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

tokenizer_config.json: 0%| | 0.00/48.0 [00:00<?, ?B/s]

vocab.txt: 0%| | 0.00/232k [00:00<?, ?B/s]

tokenizer.json: 0%| | 0.00/466k [00:00<?, ?B/s]

Hardware accelerator e.g. GPU is available in the environment, but no `device` argument is passed to the `Pipeline` object. Model will be on CPU.

```
{'score': 0.65113765001297, 'start': 40, 'end': 67, 'answer': 'state-of-the-art NLP models'}
```

SENTIMENT ANALYSIS

Sentiment Analysis is a key task in Natural Language Processing (NLP) that involves determining the emotional tone or sentiment expressed in a piece of text. This can include classifying text into categories like positive, negative, or neutral sentiments or even detecting emotions like anger, joy, or sadness.

Purpose:

To understand the feelings or emotions behind a text, commonly applied in social media monitoring, customer feedback, reviews, etc.

Types:

- **Binary Sentiment Analysis:** Classifying text as either positive or negative.
- **Multi-class Sentiment Analysis:** Classifying into multiple categories, such as positive, negative, and neutral.
- **Emotion Detection:** Identifying emotions like happiness, anger, sadness, etc.

WORKING OF SENTIMENT ANALYSIS

Tokenization: The input text is broken down into smaller units called tokens.

Model Inference: The tokens are passed through a pre-trained sentiment analysis model (like BERT or RoBERTa).

Prediction: The model classifies the text into one or more sentiment categories (e.g., positive, negative, neutral).

Output: The model provides the sentiment label along with the associated confidence score.

```
[18]: from transformers import pipeline, BertTokenizer, BertForSequenceClassification
import torch
```

*Step 1: Load a pre-trained model and tokenizer for sentiment analysis
Using a multilingual sentiment analysis model for broader language support*

```
model_name = "nlp-town/bert-base-multilingual-uncased-sentiment"
model = BertForSequenceClassification.from_pretrained(model_name)
tokenizer = BertTokenizer.from_pretrained(model_name)
```

Step 2: Initialize a sentiment analysis pipeline

```
sentiment_analyzer = pipeline("sentiment-analysis", model=model,
                               tokenizer=tokenizer)
```

Step 3: Define a function to analyze multiple texts for sentiment

```
def analyze_sentiment(texts):
    results = []
    for text in texts:
        Tokenize the text and get model output
        result = sentiment_analyzer(text)
        Collect sentiment results
        results.append({
            'text': text,
            'sentiment': result[0]['label'],
            'confidence': result[0]['score']
        })
    return results
```

Step 4: Example usage with multiple sentences

```
texts = [
    "I absolutely love this product, it works amazing!",
    "This is the worst experience I have ever had.",
    "The movie was okay, but not as good as I expected.",
    "Hugging Face makes it so easy to use state-of-the-art models!"
]
```

Step 5: Analyze sentiment for each text

```
results = analyze_sentiment(texts)
```

Step 6: Print out the results

```
for result in results:
    print(f"Text: {result['text']}")
    print(f"Sentiment: {result['sentiment']}, Confidence: {result['confidence']:.4f}")
    print("-" * 50)
```

Hardware accelerator e.g. GPU is available in the environment, but no `device` argument is passed to the `Pipeline` object. Model will be on CPU.

Text: I absolutely love this product, it works amazing!

Sentiment: 5 stars, Confidence: 0.9541

Text: This is the worst experience I have ever had.

Sentiment: 1 star, Confidence: 0.9376

Text: The movie was okay, but not as good as I expected.

Sentiment: 3 stars, Confidence: 0.8514

Text: Hugging Face makes it so easy to use state-of-the-art models!

Sentiment: 5 stars, Confidence: 0.7232

HANDS-ON

1. How would you build a sentiment analysis model for classifying text into positive or negative categories using traditional machine learning algorithms?

Hint: Preprocess the text data with tokenization and vectorization (TF-IDF), then train a classifier like SVM or Logistic Regression.

2. Explain how you would implement Named Entity Recognition (NER) using Conditional Random Fields (CRF) and extract entities such as names, dates, and locations from a given text.

Hint: Extract text features like word, POS tags, and character-level information, and train a CRF model on a labeled NER dataset.

3. Describe the process of using pre-trained word embeddings like Word2Vec or GloVe for text classification tasks. How would you apply these embeddings to a text classification model?

Hint: Load pre-trained embeddings, average them for each document, and use a classifier like Logistic Regression or SVM for text classification.

4. How would you implement an extractive text summarization system to select key sentences from a document for generating a summary?

Hint: Use sentence importance based on term frequency (TF-IDF) or TextRank to select the most relevant sentences for summarization.

5. How would you design and implement a simple rule-based chatbot that responds to user queries based on predefined patterns?

Hint: Use regular expressions to detect user input patterns and map them to predefined responses.

6. Explain how you would build an n-gram language model to predict the next word in a sequence, and how would you handle unseen n-grams?

Hint: Generate n-grams from text, calculate probabilities, and use smoothing techniques like Laplace to handle unseen sequences.

7. How would you perform topic modeling on a corpus using Latent Dirichlet Allocation (LDA) and evaluate the coherence of the extracted topics?

Hint: Preprocess the text, apply LDA to extract topics, and evaluate coherence using metrics like coherence score or visualization tools.