

GENAI LAB 1

SRN- pes2ug23cs129

Section – B

Name – bhuvan hebbar

Inference: Learning Generative AI & NLP Using Google Colab

Working through this Google Colab notebook was a helpful introduction to the fundamentals of Generative AI and Natural Language Processing. Instead of focusing only on theory, the notebook allowed hands-on interaction with real models, which made the concepts much easier to understand and relate to real-world applications.

The initial setup introduced the Hugging Face ecosystem and the transformers library, which clearly showed how modern NLP development is done today. Hugging Face felt like a central hub for AI models, where powerful pre-trained models can be used directly without needing to train them from scratch. The use of the pipeline() function was especially useful, as it simplified complex tasks like text generation, summarization, and question answering into a few lines of code.

One of the most interesting parts of the notebook was comparing different text generation models. Using the same prompt with a smaller model (distilgpt2) and a slightly larger one (gpt2) clearly showed how model size affects output quality. The smaller model was fast but produced shallow and repetitive text, while the larger model generated more coherent and meaningful responses. This comparison helped in understanding the trade-off between efficiency and quality when selecting models.

The notebook also explained what happens behind the scenes in NLP. Concepts like tokenization, part-of-speech tagging, and named entity recognition were demonstrated using simple examples. Seeing how sentences are broken into subwords, how grammatical roles are assigned, and how important entities are extracted from text gave a better understanding of how language models process human language.

More advanced tasks such as summarization, question answering, and masked language modeling further strengthened this understanding. Different summarization models showed noticeable differences in output quality, again reinforcing the idea that better results often require more computational resources. The question-answering task highlighted how models extract answers strictly from the given context, rather than reasoning beyond it.

Overall, this notebook provided a clear and practical learning experience. It helped connect theoretical NLP concepts with real implementations and emphasized the importance of choosing the right model based on the task. The hands-on approach made Generative AI feel less abstract and more applicable, making it a valuable starting point for anyone learning NLP.

Task	Model	Classification	Observation (What happened)	Why this happened
Text Generation	BERT	Failure	The model couldn't really generate text properly and either failed or gave meaningless output.	BERT is an encoder-only model and isn't meant to generate the next word in a sentence.
	RoBERTa	Failure	Similar to BERT, text generation didn't work as expected.	RoBERTa is also encoder-only and is designed for understanding text, not generating it.
	BART	Success	BART was able to continue the sentence and generate readable text.	BART has a decoder and is trained for sequence-to-sequence generation.
Fill-Mask	BERT	Success	It predicted sensible words like "create" or "generate" for the masked token.	BERT is trained using masked language modeling, so this task fits it well.
	RoBERTa	Success	RoBERTa gave accurate and fluent predictions for the missing word.	RoBERTa is an improved version of BERT and performs very well on MLM tasks.
	BART	Partial Success	It gave predictions, but they were not as clean or consistent as BERT or RoBERTa.	BART supports masking but is mainly optimized for generation tasks.
Question Answering	BERT	Poor / Partial	The answers were incomplete or sometimes incorrect.	The base BERT model is not fine-tuned for question answering.
	RoBERTa	Poor / Partial	Similar behavior to BERT, with weak or unreliable answers.	RoBERTa-base is also not trained specifically for QA tasks.
	BART	Poor / Failure	The output was mostly irrelevant or incorrect.	BART-base is not designed for extractive QA without fine-tuning.