

GENAI UNIT 1 DOCUMENTATION

NAME:ADITYA CS

SRN:PESU2UG23CS030

1) What I learnt from the HANDS_ON file

- ➔ what is transformer and how it can be used by importing it from hugging face and also what is hugging face which is a GitHub for all pre-trained models which you can import and use in your code
- ➔ Then generating text using different models like distilgpt and gpt and also the difference between how one of them is slow and dumb while the other is fast and smart
- ➔ Also learnt about seeding that is set_seed, which is used so that we don't get random output every time we run the code again and again but I also noticed that seed works only if you run the cell which has the code again if you restart the session and run again with the same seed the output is different
- ➔ Then also what is NLP and how sentences are broken down into tokens and we can use different tokenizers according to your model and then each token is converted to IDs and stored
- ➔ Also worked on different experiments like summarization in which encoder-decoder models like bart are very good and then filling in missing words in which only encoder models like bert and good and then at last question answers in this only encoder-decoder perform well since we have the context and also the output

→ FIXING THE ERROR IN THE CODE

There was a bug/error in the POS , I discovered that I nlk packages that were installed were outdated and we had to import new packages for it to work so I installed

```
nltk.download('punkt_tab')
nltk.download('averaged_perceptron_tagger_eng')
```

these 2 new packages and the code worked

3.2 POS Tagging (Part-of-Speech)

Why? To understand grammar. Is 'book' a noun (the object) or a verb (to book a flight)? **What?** We label each word as Noun (NN), Verb (VB), Adj

```
[15] # Download necessary NLTK data
      nltk.download('punkt', quiet=True)
      nltk.download('punkt')
      |nltk.download('punkt_tab')
      |nltk.download('averaged_perceptron_tagger_eng')
      |nltk.download('averaged_perceptron_tagger')

... [nltk_data] Downloading package punkt to C:\Users\Aditya
[nltk_data]   CS\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package punkt_tab to C:\Users\Aditya
[nltk_data]   CS\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt_tab is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger_eng to
[nltk_data]   C:\Users\Aditya CS\AppData\Roaming\nltk_data...
[nltk_data]   Package averaged_perceptron_tagger_eng is already up-to-
[nltk_data]   date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   C:\Users\Aditya CS\AppData\Roaming\nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]   date!

... True
```

Let's tag our sentence.

```
[16] tokens = nltk.word_tokenize(sample_sentence)
      pos_tags = nltk.pos_tag(tokens)
      print(f"POS Tags: {pos_tags}")
```

2) What I learnt from the UNIT1 QUESTIONS IPYNB

→ SO in this we had to perform the 3 different experiments by using 3 different models

BERT – encoder only

ROBERTA – encoder only

BART- encoder and decoder

1) TEXT GENERATION

```
EXPERIMENT 1 - TEXT GENERATION

[44] set_seed(69)
prompt="The future of Artificial Intelligence is"
Python

1. BERT MODEL

[45] bert_model = pipeline("text-generation",model="bert-base-uncased")
bert_output = bert_model(prompt,max_length=50,num_return_sequences=1)
print(bert_output[0]['generated_text'])
Python

... If you want to use 'BertLMHeadModel' as a standalone, add 'is_decoder=True'.
Device set to use 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 Both 'max_new_tokens' (=256) and 'max_length' (=50) seem to have been set. 'max_new_tokens' will take precedence. Please refer to the documentation for more information. (https://huggingface.co/docs/transformers/main/en/main\_classes/text\_generation.html#transformers.BertLMHeadModel)
The future of Artificial Intelligence is....
```



```
2. ROBERTA MODEL

[46] roberta_model = pipeline("text-generation",model="roberta-base")
roberta_output = roberta_model(prompt,max_length=50,num_return_sequences=1)
print(roberta_output[0]['generated_text'])
Python

... Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: 'pip install huggingface_hub[hf_xet]' or 'pip install hf_xet'.
If you want to use 'RobertaLMHeadModel' as a standalone, add 'is_decoder=True'.
Device set to use 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 Both 'max_new_tokens' (=256) and 'max_length' (=50) seem to have been set. 'max_new_tokens' will take precedence. Please refer to the documentation for more information. (https://huggingface.co/docs/transformers/main/en/main\_classes/text\_generation.html#transformers.RobertaLMHeadModel)
The future of Artificial Intelligence is....
```



```
3. BART MODEL

[47] bart_model = pipeline("text-generation",model="facebook/bart-base")
bart_output = bart_model(prompt,max_length=50,num_return_sequences=1)
print(bart_output[0]['generated_text'])
Python

... Some weights of BartForCausalLM were not initialized from the model checkpoint at facebook/bart-base and are newly initialized: ['lm_head.weight', 'model.decoder.embed_tokens.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
Device set to use 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 Both 'max_new_tokens' (=256) and 'max_length' (=50) seem to have been set. 'max_new_tokens' will take precedence. Please refer to the documentation for more information. (https://huggingface.co/docs/transformers/main/en/main\_classes/text\_generation.html#transformers.BartForCausalLM)
The future of Artificial Intelligence is govern 215 215 Ripple51 globally Bitcoinasurredudud Tankkudud Severalud aeros Tank Tank dest summons Tank Tank dest Disintensiveintensiverunner Stupid 215 Stupid 215 Collect Tank aeros res....
```

CONCLUSION FROM TEXT GENERATION

Bert and Roberta are only encoder models so they can only read text and cannot generate text as we want , the the outputs for both them are none(no words generated)

But Bart even though its a encoder and decoder model which is capable of generating text , the output it gave was gibberish and not proper because it is not a decoder only model like gpt , so bart can be used only for sequence to sequence sentences like summarisation or translation and not suitable for text generation

2) MISSING WORDS

```
EXPERIMENT 2 - MASKED LANGUAGE MODELLING (MISSING WORDS)

1.BERT MODEL

mask_filler= pipeline("fill-mask",model="bert-base-uncased")
masked_sentence = "The goal of Generative AI is to [MASK] new content."
preds = mask_filler(masked_sentence)

for p in preds:
    print(f'{p["token_str"]}:{p["score"]:.2f}')


Some weights of the model checkpoint at bert-base-uncased were not used when initializing BertForMaskedLM: ['bert.pooler.dense.bias', 'bert.pooler.dense.weight', 'cls.seq_relationship.bias', 'cls.seq_relationship.weight']
- This IS expected if you are initializing BertForMaskedLM 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 BertForMaskedLM from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

Device set to use cpu
create: 0.54
generate: 0.16
produce: 0.05
develop: 0.04
add: 0.02

2.ROBERTA MODEL

mask_filler= pipeline("fill-mask",model="roberta-base")
masked_sentence = "The goal of Generative AI is to <mask> new content."
preds = mask_filler(masked_sentence)

for p in preds:
    print(f'{p["token_str"]}:{p["score"]:.2f}')


Device set to use cpu
generate: 0.37
create: 0.37
discover: 0.08
find: 0.02
provide: 0.02

3.BART MODEL

mask_filler=pipeline("fill-mask",model="facebook/bart-base")
masked_sentence = f"The goal of Generative AI is to {mask_filler.tokenizer.mask_token} new content."
preds = mask_filler(masked_sentence)

for p in preds:
    print(f'{p["token_str"]}:{p["score"]:.2f}')


Device set to use cpu
create:0.07
help:0.07
provide:0.06
enable:0.04
improve:0.03
```

CONCLUSION OF MASKED LANGUAGE MODELLING

Bert and Roberta are extremely good at predicting masked words with high score/probability but in the case of Bart it predicted the masked word but with less confidence and less score , at this point i dont know in what Bart is good , didnt perform well in text generation nor predicting missed words ,lets see what it does in the next one

Also i was copy pasting the code of bert for roberta and it was giving me error everytime , then i got to know its not [MASK] in roberta its <mask> so had to make that change and also its diffrent for each model so you can also use a general method which i used in Bart that is mask.filler.tokeniser.mask_token

3) QUESTION AND ANSWERS

```
EXPERIMENT 3 QUESTION ANSWERING

1.BERT MODEL

[38]
qa_pipeline = pipeline("question-answering", model="bert-base-uncased")
question = "What are the risks?"
context="Generative AI poses significant risks such as hallucinations, bias, and deepfakes."
res = qa_pipeline(question=context)
print(f"\nQ: {question}")
print(f"A: {res['answer']}")

...
Some weights of BertForQuestionAnswering were not initialized from the model checkpoint at bert-base-uncased and are newly initialized: ['qa_outputs.bias', 'qa_outputs.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
Device set to use cpu

Q: What are the risks?
A: risks such as hallucinations

2.ROBERTA MODEL

[40]
qa_pipeline = pipeline("question-answering",model="roberta-base")
question = "What are the risks?"
context="Generative AI poses significant risks such as hallucinations, bias, and deepfakes."
res = qa_pipeline(question=context)
print(f"\nQ: {question}")
print(f"A: {res['answer']}")

...
Some weights of RobertaForQuestionAnswering were not initialized from the model checkpoint at roberta-base and are newly initialized: ['qa_outputs.bias', 'qa_outputs.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
Device set to use cpu

Q: What are the risks?
A: deepfakes.

3.BART MODEL

[41]
qa_pipeline = pipeline("question-answering",model="facebook/bart-base")
question = "What are the risks?"
context="Generative AI poses significant risks such as hallucinations, bias, and deepfakes."
res = qa_pipeline(question=context)
print(f"\nQ: {question}")
print(f"A: {res['answer']}")

...
Some weights of BartForQuestionAnswering were not initialized from the model checkpoint at facebook/Bart-base and are newly initialized: ['qa_outputs.bias', 'qa_outputs.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
Device set to use cpu

Q: What are the risks?
A: risks such as hallucinations, bias, and deepfakes
```

CONCLUSION FROM QUESTION ANSWERING

as we can see bert and roberta both gave answers to the question unlike experiment 1 text generation where they didnt give any answer but here also the answer the gave was not full , it was just a partial form only reason is its only encoder model

if we talk about Bart this where it shines , because since its a decoder and encoder model and also we give it a proper context and input the answer it gave was complete and correct , finally i am happy for Bart it this something successfully in life

→ FINAL OUTPUT

Deliverable: Observation Table				
Deliverable: Observation Table				
Task	Model	Classification (Success/Failure)	Observation (What actually happened?)	Why did this happen? (Architectural Reason)
Generation	BERT	Failure	Failed to generate coherent text and returned minimal or noisy output .	BERT is an encoder-only model and is not trained for autoregressive text generation.
	RoBERTa	Failure	Failed to generate coherent text and returned minimal or noisy output similar to BERT.	RoBERTa is also an encoder-only model optimized for understanding rather than generation.
	BART	Somewhat Success(generated words but wrong)	Generated incoherent or unstable text when used with the text-generation pipeline.	BART requires task-specific, conditional generation and is not trained for free-form continuation.
Fill-Mask	BERT	Success	Correctly predicted suitable words such as "generate" and "create" for the masked token.	BERT is explicitly trained using Masked Language Modeling (MLM).
	RoBERTa	Success	Produced highly relevant and context-aware predictions with strong confidence scores.	RoBERTa improves MLM performance through larger training data and dynamic masking.
	BART	Partial Success	Predicted reasonable masked tokens but with lower confidence than BERT and RoBERTa.	BART is trained as a denoising autoencoder rather than a token-level MLM model.
QA	BERT	Success(not full answer)	Returned short and precise answers extracted directly from the context.	BERT is trained for extractive question answering by predicting answer spans.
	RoBERTa	Success(not full answer)	Provided concise and accurate span-based answers similar to BERT.	RoBERTa's optimized encoder architecture supports effective extractive QA.
	BART	Success(full answer)	Generated complete, fluent, and human-readable answers.	BART's encoder-decoder architecture enables conditional generative question answering.

&&&&&&&&&&&&&THANK YOU&&&&&&&&&&&