

file title: **Solution for HW 1**  
student: **FINICHIU Eduard - Adelin**  
group: **1241 EB (FILS)**

!OBS: THE REPORT (EX3) IS WRITTEN AFTER THE SOLUTIONS FOR EX 1 AND EX 2

## EXERCISE 1

1. Try dif. temp. values and see which one gives best results!

TEMP = 0.7

```
init.py > ...
1 from transformers import AutoTokenizer, AutoModelForCausalLM
2 model_name = "gpt2"
3 tokenizer = AutoTokenizer.from_pretrained(model_name)
4 model = AutoModelForCausalLM.from_pretrained(model_name)
5
6 #prepare input
7 inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
8 #generate output
9 outputs = model.generate(**inputs, max_new_tokens=50)
10
11 #add some extra arguments temp, top k, top_p, do_sample, repetition_penalty etc.
12 outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.7, top_k=50, top_p=0.95, do_sample=True, repetition_penalty=1.2)
13
14
15 #decode and print the output
16 print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

PS C:\Users\PC\Music\an4-sem1\softEng\lab1> & C:/Users/PC/Music/an4-sem1/softEng/lab1/llm\_lab/Scripts/Activate.ps1  
(llm\_lab) PS C:\Users\PC\Music\an4-sem1\softEng\lab1> & C:/Users/PC/Music/an4-sem1/softEng/lab1/llm\_lab/Scripts/python.exe c:/Users/PC/Music/an4-sem1/softEng/lab1/init.py  
Setting 'pad\_token\_id' to 'eos\_token\_id':50256 for open-end generation.  
Hello, my dog is cute and he's been doing really well. It was great to see him get some exercise in the morning."  
2) A friend of his told me that she recently had a heart attack during training at her gym because there were "a lot more people

TEMP = 0.3

```
init.py > ...
1 from transformers import AutoTokenizer, AutoModelForCausalLM
2 model_name = "gpt2"
3 tokenizer = AutoTokenizer.from_pretrained(model_name)
4 model = AutoModelForCausalLM.from_pretrained(model_name)
5
6 #prepare input
7 inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
8 #generate output
9 outputs = model.generate(**inputs, max_new_tokens=50)
10
11 #add some extra arguments temp, top k, top_p, do_sample, repetition_penalty etc.
12 outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.3, top_k=50, top_p=0.95, do_sample=True, repetition_penalty=1.2)
13
14
15 #decode and print the output
16 print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

(llm\_lab) PS C:\Users\PC\Music\an4-sem1\softEng\lab1> & C:/Users/PC/Music/an4-sem1/softEng/lab1/llm\_lab/Scripts/python.exe c:/Users/PC/Music/an4-sem1/softEng/lab1/init.py  
Setting 'pad\_token\_id' to 'eos\_token\_id':50256 for open-end generation.  
Hello, my dog is cute. I love him so much."  
The next day he was given a bottle of water and told to go home for the night before going back in with his family at 7am on Friday morning after  
having been sleeping all afternoon when they heard noises coming

TEMP = 0.95

```
init.py> ...
1 from transformers import AutoTokenizer, AutoModelForCausalLM
2 model_name = "gpt2"
3 tokenizer = AutoTokenizer.from_pretrained(model_name)
4 model = AutoModelForCausalLM.from_pretrained(model_name)
5
6 #prepare input
7 inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
8 #generate output
9 #outputs = model.generate(**inputs, max_new_tokens=50)
10
11 #add some extra arguments temp, top k, top p, do_sample, repetition_penalty etc.
12 outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.95, top_k=50, top_p=0.95, do_sample=True, repetition_penalty=1.2)
13
14
15 #decode and print the output
16 print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

Conclusion: In my case, higher temperature gave more “normal” results, even though higher temperature usually means higher “creativity.”

## EXERCISE 2

**2. Run the code with 3 diff. models of your choice (keep them small)!**

I will run 3 versions:

## # 1. "gpt2" - done in class

```

1 from transformers import AutoTokenizer, AutoModelForCausalLM
2 model_name = "gpt2"
3 tokenizer = AutoTokenizer.from_pretrained(model_name)
4 model = AutoModelForCausalLM.from_pretrained(model_name)
5
6 #prepare input
7 inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
8 #generate output
9 #outputs = model.generate(**inputs, max_new_tokens=50)
10
11 #add some extra arguments temp, top k, top_p, do_sample, repetition_penalty etc.
12 outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.95, top_k=50, top_p=0.95, do_sample=True, repetition_penalty=1.2)
13
14
15 #decode and print the output
16 print(tokenizer.decode(outputs[0], skip_special_tokens=True))

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

(llm\_lab) PS C:\Users\PC\Music\an4-sem1\softEng\lab1> & C:\Users\PC\Music\an4-sem1\softEng\lab1\llm\_lab\Scripts\python.exe c:\Users\PC\Music\an4-sem1\softEng\lab1\init.py

Setting 'pad\_token\_id' to 'eos\_token\_id':50256 for open-end generation.

Hello, my dog is cute. And I'm sure there's something to this story we're going through tonight," he said during the phone call from his local veterinarian in North Carolina (I've only seen one person so far).

"It could be that she got bitten by

## # 2. "distilgpt2" - found it online lighter than gpt 2

```
 initpy > ...  
1 from transformers import AutoTokenizer, AutoModelForCausalLM  
2 model_name = "distilgpt2"  
3 tokenizer = AutoTokenizer.from_pretrained(model_name)  
4 model = AutoModelForCausalLM.from_pretrained(model_name)  
5  
6 #prepare input  
7 inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")  
8 #generate output  
9 outputs = model.generate(**inputs, max_new_tokens=50)  
10  
11 #add some extra arguments temp, top k, top_p, do_sample, repetition_penalty etc.  
12 outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.95, top_k=50, top_p=0.95, do_sample=True, repetition_penalty=1.2)  
13  
14  
15 #decode and print the output  
16 print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

### # 3. "facebook/opt-125m" - for sure the best results yet

```
init.py > ...  
1 from transformers import AutoTokenizer, AutoModelForCausalLM  
2 model_name = "facebook/opt-125m"  
3 tokenizer = AutoTokenizer.from_pretrained(model_name)  
4 model = AutoModelForCausalLM.from_pretrained(model_name)  
5  
6 #prepare input  
7 inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")  
8 #generate output  
9 outputs = model.generate(**inputs, max_new_tokens=50)  
10  
11 #add some extra arguments temp, top_k, top_p, do_sample, repetition_penalty etc.  
12 outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.95, top_k=50, top_p=0.95, do_sample=True, repetition_penalty=1.2)  
13  
14  
15 #decode and print the output  
16 print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

#### # 4. "EleutherAI/gpt-neo-125M" - results seem to be ok

PHOTO OF RESULTS ON NEXT PAGE

[illegible]

more serious and simple. I think OPT-125M is my favourite. GPT-Neo-125M is creative but not always coherent.

Changing the temperature really affects how the models behave: low values make them safe but boring, and high values make them imaginative but less logical. As I have previously said, I like when models behave in a more creative manner.

TEMP = 0.7

```
from transformers import AutoTokenizer, AutoModelForCausalLM
model_name = "gpt2"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(model_name)

# Generate text
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50)

# Decode and print the output
print(tokenizer.decode(outputs[0], skip_special_tokens=True))

# Run the model on a custom prompt
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.7, top_k=50, top_p=0.9, do_sample=True, repetition_penalty=1.2)
print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

TEMP = 0.3

```
from transformers import AutoTokenizer, AutoModelForCausalLM
model_name = "gpt2"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(model_name)

# Generate text
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50)

# Decode and print the output
print(tokenizer.decode(outputs[0], skip_special_tokens=True))

# Run the model on a custom prompt
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.3, top_k=50, top_p=0.9, do_sample=True, repetition_penalty=1.2)
print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

TEMP = 0.95

```
from transformers import AutoTokenizer, AutoModelForCausalLM
model_name = "gpt2"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(model_name)

# Generate text
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50)

# Decode and print the output
print(tokenizer.decode(outputs[0], skip_special_tokens=True))

# Run the model on a custom prompt
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.95, top_k=50, top_p=0.9, do_sample=True, repetition_penalty=1.2)
print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

distilGPT vs Meta

```
from transformers import AutoTokenizer, AutoModelForCausalLM
model_name = "gpt2"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(model_name)

# Generate text
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50)

# Decode and print the output
print(tokenizer.decode(outputs[0], skip_special_tokens=True))

# Run the model on a custom prompt
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.95, top_k=50, top_p=0.9, do_sample=True, repetition_penalty=1.2)
print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

```
from transformers import AutoTokenizer, AutoModelForCausalLM
model_name = "gpt2"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(model_name)

# Generate text
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50)

# Decode and print the output
print(tokenizer.decode(outputs[0], skip_special_tokens=True))

# Run the model on a custom prompt
inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50, temperature=0.95, top_k=50, top_p=0.9, do_sample=True, repetition_penalty=1.2)
print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```