

به نام خدا

امیرمحمد کمیجانی ۹۹۵۲۲۰۳۲

توضیحات تکلیف سری ۷

ابتدا دیتاست را با استفاده از لینک گیت‌هاب داده شده لود میکنیم و دیتای تست و آموزش را به دیتافریم های کتابخانه pandas تبدیل کردیم.

سپس مدل داده شده را طبق لینک هاگینگ فیس تکمیل میکنیم.

```
model_id = "meta-llama/Meta-Llama-3-8B-Instruct"

pipeline = transformers.pipeline(
    "text-generation",
    model=model_id,
    model_kwargs={"torch_dtype": torch.bfloat16},
    device_map="auto",
)
```

در مرحله بعد توابع zero shot, one shot, five shot را تعریف نموده ایم و از تعریفی که در داکيومنت سوال آمده است استفاده مینماییم.

و در نهایت تابعی برای بهبود کد برای اینکه هر یک از n_shot ها را کنترل کنیم نوشتیم.

برای نمایش خروجی از تابعی استفاده کردیم که کد آن در هاگینگ فیس نیز موجود بود.

```

messages = [
    {"role": "system", "content": "You are a pirate chatbot who always responds in pirate speak!"},
    {"role": "user", "content": "Who are you?"},
]

prompt = pipeline.tokenizer.apply_chat_template(
    messages,
    tokenize=False,
    add_generation_prompt=True
)

terminators = [
    pipeline.tokenizer.eos_token_id,
    pipeline.tokenizer.convert_tokens_to_ids("<|eot_id|>")
]

outputs = pipeline(
    prompt,
    max_new_tokens=256,
    eos_token_id=terminators,
    do_sample=True,
    temperature=0.6,
    top_p=0.9,
)

print(outputs[0]['generated_text'][len(prompt):])

```

سپس در قسمت برای تست کردن مدل تابعی را نوشتیم که ابتدا تابع prompt را کال می‌کنیم تا هر یک از zero,one,five shot انجام شوند و هر یک از سناریو ها انجام میشوند و برای مقدار مشخصی از دیتای تست که آنرا مشخص می‌کنیم مقدار خروجی پرامپت و خروجی واقعی را می‌سنجیم در انتها برای درصد موفقیت هر یک از سناریو های نامبرده شده تعداد پیش بینی های درست را بر تعداد پیش بینی ها تقسیم می‌کنیم.

در انتها برای اینکه تعداد دیتای تست زیاد بود تعداد ۵۰ تا از دیتا ها را به صورت رندوم انتخاب کردیم.(پارامتر دوم در تابع sample همان seed میباشد)

سپس برای سناریو های مختلف تابع تست را فراخوانی میکنیم و در انتها نتایج را نمایش میدهیم.

نکته ای که وجود دارد این است که طبق توضیحات سوال یک فرض برای این بود که مسئله را **text classification** در نظر بگیریم چون به نوعی میتوانیم در نظر بگیریم که دو جمله از نظر کلاس اگر یکسان باشند ممکن است شباهت داشته باشند که خوب مشخصا برای این تسک مناسب نخواهد بود و دقت بسیار پایینی خواهد داشت چون خیلی از جملاتی که در یک کلاس یا دسته هستند نمیتوانند مترادف باشند و ما به دقت بیشتری برای این کار نیاز داریم.

همچنین با توجه به لینکهای کمکی ارسال شده تسک مدل را **Text generation** انتخاب کردیم.

- به دلیل محدودیت دیسک و **gpu** نتوانستم تمام مدل ها را ران بگیرم ولی کد برای تمام مدل ها یکسان است و کدی که زدم را برای همه ران میگیرم.
- برای مدل دوم میتوانیم از **load_in_8bit** استفاده کنیم که باز به دلیل محدودیت نمیتوانستم کار بیشتری انجام دهم.