

MMLU 基準測試

本文評估了一個語言模型在 MMLU（大規模多任務語言理解）基準測試上的表現。

MMLU 基準測試全面檢驗了模型在多種學科領域內執行各項任務的能力。它包含多項選擇題，涵蓋數學、歷史、法律和醫學等多樣化的主題。

數據集鏈接：

- Papers with Code
- Hugging Face Datasets

```
import torch
from datasets import load_dataset
import requests
import json

# 加載 MMLU 數據集
subject = "abstract_algebra" # 選擇你的科目
dataset = load_dataset("cais/mmlu", subject, split="test")

# 格式化提示，包含少量示例
def format_mmlu_prompt(example, few_shot_examples=5):
    prompt = " 以下是關於{}的多項選擇題（含答案）。\n\n".format(subject.replace("_", " "))
    prompt += f" 問題: {example['question']}\n"
    prompt += " 選項:\nA. {} \nB. {} \nC. {} \nD. {} \n".format(*example['choices'])
    prompt += f" 答案: {example['answer']}\n\n"

    # 添加當前問題
    prompt += f" 問題: {example['question']}\n"
    prompt += " 選項:\nA. {} \nB. {} \nC. {} \nD. {} \n".format(*example['choices'])
    prompt += " 答案:"

    return prompt

# 評估循環
correct = 0
total = 0
```

```

for example in dataset:
    prompt = format_mmlu_prompt(example)

# 向 llama-server 發送請求

url = "http://localhost:8080/v1/chat/completions"
headers = {"Content-Type": "application/json"}
data = {
    "messages": [{"role": "user", "content": prompt}],
    "max_tokens": 5,
    "temperature": 0,
}

response = requests.post(url, headers=headers, data=json.dumps(data))

if response.status_code == 200:
    output_text = response.json()["choices"][0]["message"]["content"]
    predicted_answer = output_text.strip()[0] if len(output_text.strip()) > 0 else ""
else:
    predicted_answer = ""
    print(f" 錯誤: {response.status_code} - {response.text}")

# 與真實答案比較

if predicted_answer.upper() == example["answer"]:
    correct += 1
total += 1

# 計算準確率

accuracy = correct / total
print(f" 科目: {subject}")
print(f" 準確率: {accuracy:.2%} ({correct}/{total})")

```