

✓ Predict Product Prices

And now, to evaluate our fine-tuned open source model

```
# pip installs
```

```
!pip install -q --upgrade torch==2.5.1+cu124 torchvision==0.20.1+cu124 torchaudio==2.5.1+cu124 --index-url https://download
!pip install -q --upgrade requests==2.32.3 bitsandbytes==0.46.0 transformers==4.48.3 accelerate==1.3.0 datasets==3.2.0 peft
```

```

_____ 908.2/908.2 MB 1.5 MB/s eta 0:00:00
_____ 7.3/7.3 MB 36.8 MB/s eta 0:00:00
_____ 3.4/3.4 MB 19.9 MB/s eta 0:00:00
_____ 24.6/24.6 MB 40.7 MB/s eta 0:00:00
_____ 883.7/883.7 kB 53.1 MB/s eta 0:00:00
_____ 13.8/13.8 MB 91.9 MB/s eta 0:00:00
_____ 664.8/664.8 MB 1.2 MB/s eta 0:00:00
_____ 363.4/363.4 MB 1.5 MB/s eta 0:00:00
_____ 211.5/211.5 MB 2.3 MB/s eta 0:00:00
_____ 56.3/56.3 MB 15.3 MB/s eta 0:00:00
_____ 127.9/127.9 MB 7.8 MB/s eta 0:00:00
_____ 207.5/207.5 MB 1.5 MB/s eta 0:00:00
_____ 188.7/188.7 MB 7.1 MB/s eta 0:00:00
_____ 99.1/99.1 kB 9.5 MB/s eta 0:00:00
_____ 21.1/21.1 MB 107.8 MB/s eta 0:00:00
_____ 209.6/209.6 MB 6.4 MB/s eta 0:00:00
_____ 6.2/6.2 MB 103.7 MB/s eta 0:00:00
_____ 44.4/44.4 kB 4.0 MB/s eta 0:00:00
_____ 64.9/64.9 kB 6.4 MB/s eta 0:00:00
_____ 67.0/67.0 MB 11.5 MB/s eta 0:00:00
_____ 9.7/9.7 MB 111.4 MB/s eta 0:00:00
_____ 336.6/336.6 kB 31.7 MB/s eta 0:00:00
_____ 480.6/480.6 kB 30.2 MB/s eta 0:00:00
_____ 374.8/374.8 kB 19.0 MB/s eta 0:00:00
_____ 313.9/313.9 kB 24.9 MB/s eta 0:00:00
_____ 8.7/8.7 MB 94.2 MB/s eta 0:00:00
_____ 20.0/20.0 MB 83.2 MB/s eta 0:00:00
_____ 179.3/179.3 kB 14.2 MB/s eta 0:00:00
_____ 3.1/3.1 MB 62.6 MB/s eta 0:00:00
```

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is google-colab 1.0.0 requires requests==2.32.4, but you have requests 2.32.3 which is incompatible.
 gcfs 2025.3.0 requires fsspec==2025.3.0, but you have fsspec 2024.9.0 which is incompatible.
 google-adk 1.17.0 requires requests<3.0.0,>=2.32.4, but you have requests 2.32.3 which is incompatible.

```
# imports
```

```

import os
import re
import math
from tqdm import tqdm
from google.colab import userdata
from huggingface_hub import login
import torch
import torch.nn.functional as F
import transformers
from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig, set_seed
from datasets import load_dataset, Dataset, DatasetDict
from datetime import datetime
from peft import PeftModel
import matplotlib.pyplot as plt
```

```
# Constants
```

```

BASE_MODEL = "Qwen/Qwen2.5-7B"
PROJECT_NAME = "pricer"
HF_USER = "harish-anandaramanujam"
```

```
# The run itself
```

```

RUN_NAME = "2025-10-29_05.43.51"
PROJECT_RUN_NAME = f"{PROJECT_NAME}-{RUN_NAME}"
REVISION = None
FINETUNED_MODEL = f"{HF_USER}/{PROJECT_RUN_NAME}"
```

```
# Data
```

```
DATASET_NAME = f"{HF_USER}/amzn-appliances-price-lite-data"
```

```
# Hyperparameters for QLoRA
```

```

QUANT_4_BIT = True

%matplotlib inline

# Used for writing to output in color

GREEN = "\033[92m"
YELLOW = "\033[93m"
RED = "\033[91m"
RESET = "\033[0m"
COLOR_MAP = {"red":RED, "orange": YELLOW, "green": GREEN}

```

Log in to HuggingFace

```

# Log in to HuggingFace

hf_token = userdata.get('HF_TOKEN')
login(hf_token, add_to_git_credential=True)

```

```

dataset = load_dataset(DATASET_NAME)
train = dataset['train']
test = dataset['test']

```

```

README.md: 100%                                412/412 [00:00<00:00, 14.1kB/s]

data/train-00000-of-00001.parquet: 100%          9.50M/9.50M [00:01<00:00, 6.60MB/s]

data/test-00000-of-00001.parquet: 100%           759k/759k [00:00<00:00, 1.43MB/s]

Generating train split: 100%                     25000/25000 [00:00<00:00, 95758.03 examples/s]

Generating test split: 100%                      2000/2000 [00:00<00:00, 32776.19 examples/s]

```

```
test[0]
```

```

{'text': 'How much does this cost to the nearest dollar?\n\nCONTINENTAL REFRIGERATOR 2-705 Gasket, SNAP in\nSPECIFICATIONS LENGTH 22 3/8, 568 mm WIDTH 25 3/8, 644 mm WEIGHT 1.225 lb GASKET TYPE 4 SIDED, MAGNETIC, DART PART REFERENCE INFO CONTINENTAL REFRIGERATOR 2-705 MODEL REFERENCE INFO CONTINENTAL REFRIGERATOR WORK TOP REFRIGERATOR MODELS SW48 CONTINENTAL REFRIGERATOR REFRIGERATOR DOOR GASKET 2-705 High Quality Gasket Brand Name Continental Refrigerator, Model Info Weight 1.2\n\nPrice is $',
 'price': 36.22}

```

Now load the Tokenizer and Model

```
# pick the right quantization (thank you Robert M. for spotting the bug with the 8 bit version!)
```

```

if QUANT_4_BIT:
    quant_config = BitsAndBytesConfig(
        load_in_4bit=True,
        bnb_4bit_use_double_quant=True,
        bnb_4bit_compute_dtype=torch.bfloat16,
        bnb_4bit_quant_type="nf4"
    )
else:
    quant_config = BitsAndBytesConfig(
        load_in_8bit=True,
        bnb_8bit_compute_dtype=torch.bfloat16
    )

```

```
# Load the Tokenizer and the Model
```

```

tokenizer = AutoTokenizer.from_pretrained(BASE_MODEL, trust_remote_code=True)
tokenizer.pad_token = tokenizer.eos_token
tokenizer.padding_side = "right"

```

```

base_model = AutoModelForCausalLM.from_pretrained(
    BASE_MODEL,
    quantization_config=quant_config,
    device_map="auto",
)
base_model.generation_config.pad_token_id = tokenizer.pad_token_id

```

```
# Load the fine-tuned model with PEFT
```

```

if REVISION:
    fine_tuned_model = PeftModel.from_pretrained(base_model, FINETUNED_MODEL, revision=REVISION)
else:
    fine_tuned_model = PeftModel.from_pretrained(base_model, FINETUNED_MODEL)

```

```
print(f"Memory footprint: {fine_tuned_model.get_memory_footprint() / 1e6:.1f} MB")
```

```
tokenizer_config.json:      7.23k/? [00:00<00:00, 272kB/s]
vocab.json:                2.78M/? [00:00<00:00, 45.8MB/s]
merges.txt:                1.67M/? [00:00<00:00, 22.4MB/s]
tokenizer.json:            7.03M/? [00:00<00:00, 69.7MB/s]
config.json: 100%          686/686 [00:00<00:00, 40.4kB/s]
model.safetensors.index.json: 27.8k/? [00:00<00:00, 3.08MB/s]
Downloading shards: 100%   4/4 [02:33<00:00, 40.64s/it]
model-00001-of-00004.safetensors: 100%      3.95G/3.95G [00:31<00:00, 255MB/s]
model-00002-of-00004.safetensors: 100%      3.86G/3.86G [00:33<00:00, 167MB/s]
model-00003-of-00004.safetensors: 100%      3.86G/3.86G [00:38<00:00, 77.4MB/s]
model-00004-of-00004.safetensors: 100%      3.56G/3.56G [00:47<00:00, 110MB/s]
Loading checkpoint shards: 100%              4/4 [01:05<00:00, 16.08s/it]
generation_config.json: 100%              138/138 [00:00<00:00, 13.6kB/s]
adapter_config.json: 100%                 739/739 [00:00<00:00, 83.4kB/s]
adapter_model.safetensors: 100%          80.8M/80.8M [00:06<00:00, 14.9MB/s]
Memory footprint: 5524.0 MB
```

```
fine_tuned_model
```

```
(lora_A): ModuleDict(
  (default): Linear(in_features=3584, out_features=32, bias=False)
)
(lora_B): ModuleDict(
  (default): Linear(in_features=32, out_features=512, bias=False)
)
(lora_embedding_A): ParameterDict()
(lora_embedding_B): ParameterDict()
(lora_magnitude_vector): ModuleDict()
)
(v_proj): lora.Linear4bit(
  (base_layer): Linear4bit(in_features=3584, out_features=512, bias=True)
  (lora_dropout): ModuleDict(
    (default): Dropout(p=0.1, inplace=False)
  )
  (lora_A): ModuleDict(
    (default): Linear(in_features=3584, out_features=32, bias=False)
  )
  (lora_B): ModuleDict(
    (default): Linear(in_features=32, out_features=512, bias=False)
  )
  (lora_embedding_A): ParameterDict()
  (lora_embedding_B): ParameterDict()
  (lora_magnitude_vector): ModuleDict()
)
)
```

```

    )
    (lm_head): Linear(in_features=3584, out_features=152064, bias=False)
  )
)
)

```

✓ THE MOMENT OF TRUTH!

```

def extract_price(s):
    if "Price is $" in s:
        contents = s.split("Price is $")[1]
        contents = contents.replace(',', '')
        match = re.search(r"[-+]?[d*]\d+|\d+", contents)
        return float(match.group()) if match else 0
    return 0

```

```
extract_price("Price is $a fabulous 899.99 or so")
```

```
899.99
```

Original prediction function takes the most likely next token

```

def model_predict(prompt):
    set_seed(42)
    inputs = tokenizer.encode(prompt, return_tensors="pt").to("cuda")
    attention_mask = torch.ones(inputs.shape, device="cuda")
    outputs = fine_tuned_model.generate(inputs, attention_mask=attention_mask, max_new_tokens=3, num_return_sequences=1)
    response = tokenizer.decode(outputs[0])
    return extract_price(response)

```

An improved prediction function takes a weighted average of the top 3 choices
 # This code would be more complex if we couldn't take advantage of the fact
 # That Llama generates 1 token for any 3 digit number

```
top_K = 3
```

```

def improved_model_predict(prompt, device="cuda"):
    set_seed(42)
    inputs = tokenizer.encode(prompt, return_tensors="pt").to(device)
    attention_mask = torch.ones(inputs.shape, device=device)

    with torch.no_grad():
        outputs = fine_tuned_model(inputs, attention_mask=attention_mask)
        next_token_logits = outputs.logits[:, -1, :].to('cpu')

    next_token_probs = F.softmax(next_token_logits, dim=-1)
    top_prob, top_token_id = next_token_probs.topk(top_K)
    prices, weights = [], []
    for i in range(top_K):
        predicted_token = tokenizer.decode(top_token_id[0][i])
        probability = top_prob[0][i]
        try:
            result = float(predicted_token)
        except ValueError as e:
            result = 0.0
        if result > 0:
            prices.append(result)
            weights.append(probability)
    if not prices:
        return 0.0, 0.0
    total = sum(weights)
    weighted_prices = [price * weight / total for price, weight in zip(prices, weights)]
    return sum(weighted_prices).item()

```

Start coding or [generate](#) with AI.

```

class Tester:

    def __init__(self, predictor, data, title=None, size=250):
        self.predictor = predictor
        self.data = data
        self.title = title or predictor.__name__.replace("_", " ").title()
        self.size = size
        self.guesses = []
        self.truths = []
        self.errors = []
        self.sles = []

```

```

        self.colors = []

    def color_for(self, error, truth):
        if error<40 or error/truth < 0.2:
            return "green"
        elif error<80 or error/truth < 0.4:
            return "orange"
        else:
            return "red"

    def run_datapoint(self, i):
        datapoint = self.data[i]
        guess = self.predictor(datapoint["text"])
        truth = datapoint["price"]
        error = abs(guess - truth)
        log_error = math.log(truth+1) - math.log(guess+1)
        sle = log_error ** 2
        color = self.color_for(error, truth)
        title = datapoint["text"].split("\n\n")[1][:20] + "..."
        self.guesses.append(guess)
        self.truths.append(truth)
        self.errors.append(error)
        self.sles.append(sle)
        self.colors.append(color)
        print(f"{COLOR_MAP[color]}{i+1}: Guess: ${guess:,.2f} Truth: ${truth:,.2f} Error: ${error:,.2f} SLE: {sle:,.2f} It")

    def chart(self, title):
        max_error = max(self.errors)
        plt.figure(figsize=(12, 8))
        max_val = max(max(self.truths), max(self.guesses))
        plt.plot([0, max_val], [0, max_val], color='deepskyblue', lw=2, alpha=0.6)
        plt.scatter(self.truths, self.guesses, s=3, c=self.colors)
        plt.xlabel('Ground Truth')
        plt.ylabel('Model Estimate')
        plt.xlim(0, max_val)
        plt.ylim(0, max_val)
        plt.title(title)
        plt.show()

    def report(self):
        average_error = sum(self.errors) / self.size
        rmsle = math.sqrt(sum(self.sles) / self.size)
        hits = sum(1 for color in self.colors if color=="green")
        title = f"{self.title} Error=${average_error:,.2f} RMSLE={rmsle:,.2f} Hits={hits/self.size*100:.1f}%"
        self.chart(title)

    def run(self):
        self.error = 0
        for i in range(self.size):
            self.run_datapoint(i)
        self.report()

    @classmethod
    def test(cls, function, data):
        cls(function, data).run()

```

```
Tester.test(improved_model_predict, test)
```

