

Treinamento de modelo LLM fundacional para a descrição de produtos da Amazon

1.Quanto ao objetivo

O objetivo do finetuning é treinar um modelo LLM fundacional com um dataset de produtos da Amazon, com a finalidade de que o modelo possa descrever os produtos a partir do título do mesmo.

2.Quanto à estrutura do projeto

Além deste relatório, o projeto contém os seguintes arquivos:

- dataset_preparation.ipynb = Notebook responsável pela preparação do dataset para posterior uso no fine tuning e no RAG.
- model_import.ipynb = Notebook responsável por realizar o download do modelo base e armazená-lo em disco.
- fine_tuning.ipynb = Notebook responsável pelo processo de fine tuning do modelo base.
- model_testing.ipynb = Notebook utilizado para comparar o modelo base com o modelo treinado, realizando um conjunto de perguntas para cada um deles.
- RAG.ipynb = Notebook responsável pelo processo de RAG, realizado utilizando o modelo treinado.
- evaluation_history.pdf = Arquivo contendo a comparação das respostas em cada etapa do processo de fine tuning e, depois, com o modelo final utilizando RAG.

3.Quanto ao modelo base escolhido

O foundational model escolhido para treinamento foi o GOOGLE/FLAN-T5, uma variante mais recente e aprimorada do T5, disponível inclusive para uso comercial. O modelo foi pré treinado com o objetivo de mapear sequências de texto, com bases supervisionadas e não supervisionadas.

O modelo está disponível em 5 variações de tamanho, sendo elas em ordem crescente: small, base, large, xl e xxl. Os modelos variam de 80 milhões de parâmetros à 11 bilhões. Por limitações de hardware para treinamento, o modelo escolhido para uso foi o

base, pois a partir do large o hardware disponível se mostrou inviável. A versão base do modelo possui 248 milhões de parâmetros.

4.Quanto ao dataset de treinamento

Os dados foram tratados de forma a buscar a otimização do treinamento, a limpeza foi realizada removendo as linhas que possuíssem alguma das colunas (title ou content) sem conteúdo. As demais colunas foram também removidas. O texto foi tratado removendo entidades html, que alguns registros possuíam, e removendo caracteres especiais. Para preservar o contexto do texto, letras maiúsculas e pontuação convencional (. , ! ?) foram mantidas.

O dataset tratado, com aproximadamente 1.3 milhões de registros, foi dividido em datasets menores com 20 mil registros cada, com a finalidade de facilitar o treinamento, visto que o tempo e consumo de hardware para treinar com o dataset inteiro de uma única vez seria muito grande e inviável para o hardware disponível.

5.Quanto ao processo de fine tuning (treinamento)

O treinamento foi realizado iterativamente sobre os subdatasets de 20 mil registros, cada um com 3 épocas de treinamento. Cada subdataset treinado gerou um novo modelo, que foi usado para treinar com o próximo dataset.

O input foi estruturado da seguinte maneira: Please describe this product: {título do produto}, com a descrição do produto fornecida como resposta.

O Google Colab foi utilizado para treinamento, com GPU Nvidia L4, disponível de forma paga na plataforma.

5.1.Quanto aos detalhes de implementação

O fine tuning foi realizado utilizando a biblioteca transformers (huggingface), utilizando T5Tokenizer para tokenização, e T5ForConditionalGeneration para manusear o modelo. Para o treinamento foi utilizado a classe Seq2SeqTrainer, a qual é focada em atividades sequence-to-sequence, como sumarização e tradução.

Os hiperparâmetros mais relevantes optados durante o processo de treinamento são descritos abaixo:

- **LEARNING_RATE**: Aqui buscou-se no valor de $1e-4$ (0.0001) um equilíbrio. Como o modelo é relativamente pequeno e não possui uma quantidade grande de parâmetros, ele deve aprender bem com os novos dados, mas a taxa não deve ser grande ao ponto de bagunçar o que o modelo já sabe.

- PER_DEVICE_TRAIN_BATCH_SIZE: Optou-se pelo valor de 8, para que o treinamento não ficasse lento demais, e pudesse ser realizado com a memória disponível.
- GRADIENT_ACCUMULATION_STEPS: Utilizado para simular um batch maior, interessante quando se possui limitação de memória.
- WEIGHT_DECAY: Utilizado o valor de 0.01 como medida para prevenção de overfitting do modelo, aplicando penalidade aos pesos do modelo.
- NUM_TRAIN_EPOCHS: Para o caso foi utilizado 3 épocas de treinamento, visto que a progressão é lenta, e um maior número de épocas demandaria maior capacidade computacional.

5.2.Quanto à execução

Cada subdataset de 20.000 registros foi treinado em 3 épocas, sendo 5% (1.000) para teste e o restante para treinamento(19.000).

Para cada treinamento foi exibido o training loss, validation loss e o rouge score (Recall-Oriented Understudy for Gisting Evaluation), que calcula a eficiência do modelo para tarefas de question/answer, quanto maior o score, melhor a performance do modelo.

Por motivos de limitação de hardware/custo, não foi possível realizar o treinamento do dataset de forma total, foram treinados 10 subdatasets, correspondendo à um total de 200.000 registros.

Abaixo segue a execução de cada um dos subdatasets treinados:

Subdataset 1:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeLsum
0	3.578700	3.193809	0.10753 1	0.02698 8	0.08866 3	0.096138
1	3.432500	3.143379	0.10825 7	0.03147 7	0.09043 5	0.097658
2	3.377900	3.127922	0.10988 3	0.03181 0	0.09168 6	0.099655

Subdataset 2:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeLsum
0	3.478000	3.166302	0.10207 1	0.01846 0	0.07964 0	0.089572

1	3.409200	3.139682	0.10592 7	0.02322 2	0.08416 4	0.093236
2	3.370900	3.130713	0.10581 8	0.02315 8	0.08352 3	0.093628

Subdataset 3:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.473200	3.161193	0.09445 4	0.01812 0	0.07500 1	0.083629
1	3.399500	3.127282	0.09535 2	0.01861 0	0.07533 7	0.084209
2	3.354300	3.118323	0.09541 2	0.01874 8	0.07574 5	0.084342

Subdataset 4:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.487300	3.223335	0.08346 2	0.01341 6	0.06793 2	0.075505
1	3.429400	3.199837	0.08138 1	0.01302 6	0.06629 8	0.073516
2	3.401500	3.192307	0.08065 8	0.01290 0	0.06547 3	0.072030

Subdataset 5:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.404000	3.131171	0.09871 2	0.01780 9	0.07866 3	0.087920
1	3.338900	3.106844	0.10043 7	0.01913 2	0.07997 0	0.089308
2	3.309100	3.099408	0.09914 8	0.01876 8	0.07957 5	0.088745

Subdataset 6:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.349000	3.096937	0.12074 1	0.02684 4	0.09586 2	0.106956

1	3.299900	3.083470	0.11214 2	0.02590 9	0.08988 8	0.099633
2	3.260500	3.077748	0.11519 7	0.02748 1	0.09189 3	0.102308

Subdataset 7:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.366300	3.109919	0.11067 4	0.02194 6	0.08717 2	0.097535
1	3.309600	3.089842	0.10861 1	0.02408 5	0.08739 4	0.097108
2	3.268400	3.083262	0.10963 4	0.02421 1	0.08831 2	0.097732

Subdataset 8:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.349800	3.077527	0.11007 5	0.02085 4	0.08652 5	0.096420
1	3.264400	3.048553	0.10425 6	0.01821 1	0.08243 8	0.091196
2	3.228000	3.040243	0.10395 5	0.01810 1	0.08190 2	0.091138

Subdataset 9:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.344600	3.076648	0.10012 9	0.02250 1	0.08041 4	0.089185
1	3.288300	3.057238	0.10047 3	0.02236 7	0.08170 8	0.090020
2	3.261500	3.050043	0.10304 8	0.02370 4	0.08329 9	0.092229

Subdataset 10:

Epoch	Training Loss	Validation Loss	Rouge1	Rouge2	RougeL	RougeSum
0	3.313600	3.065063	0.10651 3	0.02031 8	0.08361 3	0.093346

1	3.257800	3.047984	0.105555	0.020653	0.082840	0.092829
2	3.231800	3.043123	0.105747	0.019800	0.082561	0.092483

6.Quanto à conclusão e aos resultados

A cada etapa de treinamento(20.000 registros) foi realizado a comparação entre o modelo treinado e o modelo base, realizando uma série de perguntas solicitando a descrição dos produtos. Estes resultados estão no arquivo evaluation_history.pdf.

É notável que antes do treinamento o modelo não tinha capacidade nenhuma para descrever os produtos, na maioria das vezes retornando o nome do próprio produto ou alucinando. Após o treinamento com 200.000 registros, o modelo consegue descrever de forma simples e curta, as vezes incompleta, alguns dos produtos solicitados, para outros produtos a resposta ficou confusa e também houve alucinação. Com isso, para melhorar a performance, foi realizado o RAG sobre o modelo treinado, para buscar descrições mais precisas dos produtos solicitados. Vale ressaltar ainda que os produtos solicitados nas perguntas testes foram escolhidos aleatoriamente entre os 200.000 que foram utilizados para o treinamento e para o RAG.

Após a execução do RAG, foram realizadas novamente as mesmas perguntas e, claramente, houve um grande ganho de performance nas respostas. Apesar de fornecer respostas legais, o maior problema do modelo foi a repetição indiscriminada de texto para algumas das respostas, repetindo palavras e trechos de forma confusa.

Apesar de algumas respostas confusas, o modelo final com a utilização do RAG forneceu descrições legais para alguns produtos, não de forma perfeita, mas para um modelo base relativamente pequeno, que não conseguia descrever nenhum dos produtos solicitados, houve considerável progresso na tarefa.

Fontes utilizadas:

<https://www.datacamp.com/tutorial/flan-t5-tutorial>

<https://huggingface.co/google/flan-t5-base>

<https://medium.com/georgian-impact-blog/the-practical-guide-to-lms-flan-t5-6d26cc5f14c0>