



Instituto Infnet

GRADUAÇÃO EM BANCO DE DADOS

Projeto de Bloco: Ciência de Dados Aplicada [24E3_5]

WANDERSON RAFAEL MENDONÇA BATISTA

TESTE DE PERFORMANCE – TP4

PROF. DIEGO DA SILVA RODRIGUES

RIO DE JANEIRO, 2024

Configuração do Ambiente de Desenvolvimento:

Configure seu ambiente de desenvolvimento, incluindo Git para controle de versão e preparação para deploy. Lembre-se de seguir a estrutura do CRISP-DM para organizar seu projeto de forma eficiente e escalável.

<https://github.com/wanderaf/ProjetoDeBloco>

1. Identificação e Escolha do Modelo LLM (Local):

- **Critérios de Seleção:** Pesquise e selecione o modelo de linguagem natural mais adequado para a sua aplicação, considerando os seguintes critérios:
- **Desempenho:** Avalie a precisão, a capacidade de resposta e o tipo de tarefa para a qual o modelo foi treinado (ex: GPT, BERT, T5).
- **Custo Computacional Local:** Verifique os recursos de hardware necessários para rodar o modelo no ambiente local (como GPU, memória RAM, etc.).
- **Acessibilidade Local:** Garanta que o modelo escolhido pode ser carregado localmente, utilizando bibliotecas como Transformers da HuggingFace para download e execução do modelo em sua própria máquina.

O modelo escolhido para a aplicação foi o t5-small, da arquitetura T5, devido à sua capacidade de realizar tarefas de transformação de texto com precisão e eficiência, como previsão baseada em padrões históricos. Esse modelo apresenta requisitos computacionais moderados, sendo adequado para execução em CPUs com pelo menos 8 GB de RAM ou em GPUs para maior desempenho, garantindo viabilidade em ambientes locais com recursos limitados. Além disso, sua acessibilidade é facilitada pela integração com a biblioteca Hugging Face Transformers, permitindo carregamento e execução local de forma independente, com suporte a otimizações como cache de recursos para maior eficiência operacional.

2. Integração de LLM com FastAPI no Ambiente Local:

Execução Local: Utilize um modelo da HuggingFace ou um modelo treinado localmente para realizar tarefas de processamento de linguagem natural (como geração de texto, resumo automático, classificação de sentimento, etc.) e integre-o ao seu backend FastAPI, sem necessidade de conexão com a nuvem.

Rota FastAPI: Implemente uma rota em FastAPI que se conecte ao modelo LLM rodando localmente para processar dados textuais fornecidos pela aplicação.

Exemplo de rota:

POST /processar_texto: Rota que recebe um texto enviado pelo usuário e retorna a análise ou processamento realizado pelo modelo de linguagem, como um resumo ou análise de sentimento.

Execução Local do HuggingFace: Baixe o modelo diretamente em seu ambiente de desenvolvimento local e utilize as funções da biblioteca Transformers para carregar e executar o modelo sem depender de serviços externos.

Modelo HuggingFace

```
7 # Carregar o modelo HuggingFace local
8 tokenizer = AutoTokenizer.from_pretrained("t5-small")
9 model = AutoModelForSeq2SeqLM.from_pretrained("t5-small")
10
11 # Função para gerar texto com o modelo
12 def generate_response(input_text: str) -> str:
13     inputs = tokenizer(input_text, return_tensors="pt", max_length=512, truncation=True)
14     outputs = model.generate(inputs["input_ids"], max_length=512, num_beams=5, early_stopping=True)
15     return tokenizer.decode(outputs[0], skip_special_tokens=True)
```

```
50 # Configurar o pipeline do HuggingFace para geração de texto
51 try:
52     tokenizer = AutoTokenizer.from_pretrained("distilgpt2")
53     model = AutoModelForCausalLM.from_pretrained("distilgpt2")
54     text_generator = pipeline("text-generation", model=model, tokenizer=tokenizer)
55     print("Modelo de geração de texto carregado com sucesso.")
56 except Exception as e:
57     print(f"Erro ao carregar o modelo: {e}")
58     text_generator = None
```

```

152 @app.post("/chat/converse", response_model=ChatResponse)
153 async def converse(chat_request: ChatRequest):
154     """
155     Endpoint para responder perguntas específicas com base nos dados do MongoDB.
156     """
157     try:
158         pergunta = chat_request.question.lower().strip()
159
160         # Processar a pergunta com o modelo HuggingFace
161         processed_question = generate_response(f"Interprete a pergunta: {pergunta}")
162
163         resposta = ""
164
165         # Normalização para nomes de unidades e procedimentos
166         def normalize(text):
167             return text.lower().strip()
168
169         # Verificar qual pergunta está sendo feita e buscar os dados correspondentes
170         if "procedimentos apresentados" in processed_question:
171             unidade_nome = normalize(processed_question.split("na unidade")[-1])
172             query = {"FANTASIA": {"$regex": f".*{unidade_nome}.*", "$options": "i"}}
173             print(f"Consulta MongoDB: {query}") # Log da consulta
174             dados = list(collection_sia_analise.find(query, {"_id": 0}))
175             print(f"Resultados encontrados: {dados}") # Log dos resultados
176
177             if not dados:
178                 return ChatResponse(
179                     question=chat_request.question,
180                     response=f"Nenhuma unidade encontrada com o nome '{unidade_nome}'."
181                 )
182
183             total_procedimentos = sum(item.get("TOTAL_PA_QTDPRO", 0) for item in dados)
184             resposta = f"A unidade '{unidade_nome}' apresentou um total de {total_procedimentos} procedimentos."
185
186         elif "procedimentos aprovados" in processed_question:
187             unidade_nome = normalize(processed_question.split("na unidade")[-1])
188             query = {"FANTASIA": {"$regex": f".*{unidade_nome}.*", "$options": "i"}}
189             print(f"Consulta MongoDB: {query}") # Log da consulta
190             dados = list(collection_sia_analise.find(query, {"_id": 0}))
191             print(f"Resultados encontrados: {dados}") # Log dos resultados

```

```

192
193         if not dados:
194             return ChatResponse(
195                 question=chat_request.question,
196                 response=f"Nenhuma unidade encontrada com o nome '{unidade_nome}'."
197             )
198
199         total_aprovados = sum(item.get("TOTAL_PA_QTDAPR", 0) for item in dados)
200         resposta = f"A unidade '{unidade_nome}' aprovou um total de {total_aprovados} procedimentos."
201
202         elif "procedimentos do tipo" in processed_question and "foram realizados na unidade" in processed_question:
203             partes = processed_question.split("foram realizados na unidade")
204             procedimento_nome = normalize(partes[0].split("do tipo")[-1])
205             unidade_nome = normalize(partes[-1])
206             query = {
207                 "FANTASIA": {"$regex": f".*{unidade_nome}.*", "$options": "i"},
208                 "IP_DSCR": {"$regex": f".*{procedimento_nome}.*", "$options": "i"}
209             }
210             print(f"Consulta MongoDB: {query}") # Log da consulta
211             dados = list(collection_sia_analise.find(query, {"_id": 0}))
212             print(f"Resultados encontrados: {dados}") # Log dos resultados
213
214             if not dados:
215                 return ChatResponse(
216                     question=chat_request.question,
217                     response=f"Nenhuma unidade ou procedimento encontrado com os critérios fornecidos."
218                 )
219
220             total_procedimentos = sum(item.get("TOTAL_PA_QTDPRO", 0) for item in dados)
221             resposta = (
222                 f"Foram realizados {total_procedimentos} procedimentos do tipo '{procedimento_nome}' na unidade '{unidade_nome}'."
223             )
224
225         elif "historico de procedimentos realizados na unidade" in processed_question:
226             unidade_nome = normalize(processed_question.split("na unidade")[-1])
227             query = {"FANTASIA": {"$regex": f".*{unidade_nome}.*", "$options": "i"}}
228             print(f"Consulta MongoDB: {query}") # Log da consulta
229             dados = list(collection_sia_analise.find(query, {"_id": 0}))
230             print(f"Resultados encontrados: {dados}") # Log dos resultados

```

```

232 ~         if not dados:
233 ~             return ChatResponse(
234 ~                 question=chat_request.question,
235 ~                 response=f"Nenhuma unidade encontrada com o nome '{unidade_nome}'."
236 ~             )
237 ~
238 ~         historico = {}
239 ~         for item in dados:
240 ~             ano_mes = item.get("PA_MVMR", "Desconhecido")
241 ~             historico[ano_mes] = historico.get(ano_mes, 0) + item.get("TOTAL_PA_QTDPRO", 0)
242 ~
243 ~         resposta = "Histórico de procedimentos realizados:\n"
244 ~         for periodo, total in sorted(historico.items()):
245 ~             resposta += f"- {periodo}: {total} procedimentos\n"
246 ~
247 ~         else:
248 ~             resposta = "Desculpe, não entendi sua pergunta. Por favor, reformule."
249 ~
250 ~         return ChatResponse(
251 ~             question=chat_request.question,
252 ~             response=resposta
253 ~         )
254 ~
255 ~     except Exception as e:
256 ~         print(f"Erro no chat: {e}")
257 ~         raise HTTPException(status_code=500, detail="Erro interno no servidor.")

```

Get -> Busca nome das unidades através de palavras chaves

```

62 ~ @app.get("/unidades/", response_model=List[dict])
63 ~ async def consultar_unidades(pa_coduni: str = Query(..., description="Texto para busca (PA_CODUNI ou parte de FANTASIA)")):
64 ~     """
65 ~     Consulta unidades de saúde na coleção SIA_ANALISE pelo PA_CODUNI ou parte do nome FANTASIA.
66 ~     - Se 'pa_coduni' corresponder a parte do PA_CODUNI, retorna todos os registros correspondentes.
67 ~     - Se 'pa_coduni' corresponder a parte do FANTASIA, retorna todos os registros correspondentes.
68 ~     """
69 ~     try:
70 ~         # Construir o filtro de busca dinâmico
71 ~         query_filter = {
72 ~             "$or": [
73 ~                 {"PA_CODUNI": {"$regex": pa_coduni, "$options": "i"}}, # Busca parcial no PA_CODUNI
74 ~                 {"FANTASIA": {"$regex": pa_coduni, "$options": "i"}} # Busca parcial no FANTASIA
75 ~             ]
76 ~         }
77 ~
78 ~         # Buscar dados no MongoDB
79 ~         resultados = collection_sia_analise.find(query_filter, {"_id": 0, "PA_CODUNI": 1, "FANTASIA": 1})
80 ~
81 ~         # Garantir unicidade dos resultados
82 ~         dados_unicos = list({res["PA_CODUNI"], res["FANTASIA"]} for res in resultados if "PA_CODUNI" in res and "FANTASIA" in res)
83 ~
84 ~         # Verificar se há registros encontrados
85 ~         if not dados_unicos:
86 ~             return {"message": "Nenhum registro encontrado para a consulta."}
87 ~
88 ~         # Formatar como lista de dicionários
89 ~         resposta_formatada = [{"PA_CODUNI": item[0], "FANTASIA": item[1]} for item in dados_unicos]
90 ~
91 ~         return resposta_formatada
92 ~
93 ~     except Exception as e:
94 ~         # Log detalhado do erro para depuração
95 ~         print(f"Erro ao consultar unidades: {e}")
96 ~         raise HTTPException(status_code=500, detail="Erro interno no servidor.")

```

Post -> Acrescenta novo registro no banco

```
98 @app.post("/unidades/")
99 v async def adicionar_unidade(unidade: UnidadeData):
100     """
101     Endpoint para adicionar manualmente um registro na coleção SIA_ANALISE.
102     """
103     try:
104         # Converter o objeto Pydantic para um dicionário
105         unidade_dict = unidade.dict()
106
107         # Inserir o registro no MongoDB
108         resultado = collection_sia_analise.insert_one(unidade_dict)
109
110         # Verificar se a inserção foi bem-sucedida
111         if not resultado.inserted_id:
112             raise HTTPException(status_code=500, detail="Erro ao inserir o registro no banco de dados.")
113
114         return {"message": "Registro inserido com sucesso!", "id": str(resultado.inserted_id)}
115
```

3. Manipulação das Respostas da API:

Response Models: Crie Response Models em FastAPI para garantir que as respostas da API sejam consistentes e estruturadas, utilizando modelos de resposta claros e bem definidos.

Exemplo: A rota `/processar_texto` deve retornar um JSON estruturado com informações detalhadas, como o texto processado, o tipo de análise realizada (resumo, classificação, etc.) e os resultados gerados.

Validação das Respostas: Assegure-se de que todos os retornos da API estejam validados, garantindo que os dados enviados e recebidos estejam no formato correto e sejam compreensíveis para o cliente.

Foram realizadas modificações em relação ao código apresentado na questão anterior afim de aplicar o Response Models:

```

135 @app.post("/chat/converse", response_model=ChatResponse)
136 async def converse(chat_request: ChatRequest):
137     """
138     Endpoint para responder perguntas específicas com base nos dados do MongoDB.
139     """
140     try:
141         pergunta = chat_request.question.lower()
142
143         # Processar a pergunta com o modelo HuggingFace
144         processed_question = generate_response(f"Interprete a pergunta: {pergunta}")
145
146         print(f"Processed Question: {processed_question}") # Log para depuração
147
148         # Mapear tipo de análise baseado na pergunta
149         analysis_type = ""
150         results = []
151
152         if "procedimentos apresentados" in pergunta or "total apresentados" in pergunta:
153             analysis_type = "procedimentos_apresentados"
154         elif "procedimentos aprovados" in pergunta or "total aprovados" in pergunta:
155             analysis_type = "procedimentos_aprovados"
156         elif "procedimentos do tipo" in pergunta:
157             analysis_type = "procedimentos_tipo"
158         elif "histórico de procedimentos" in pergunta:
159             analysis_type = "historico_procedimentos"
160         else:
161             return ChatResponse(
162                 question=chat_request.question,
163                 processed_question=ProcessedQuestionResponse(
164                     processed_question=processed_question,
165                     analysis_type="unknown"
166                 ),
167                 results=[]
168             )
169
170         # Extrair nome da unidade de saúde
171         unidade_match = re.search(r"unidade\s(.+)", pergunta)
172         unidade_nome = unidade_match.group(1).strip() if unidade_match else None
173
174         if not unidade_nome:
175             return ChatResponse(
176                 question=chat_request.question,
177                 processed_question=ProcessedQuestionResponse(
178                     processed_question=processed_question,
179                     analysis_type=analysis_type
180                 ),
181                 results=[]
182             )

```

```

183
184     # Buscar no MongoDB com base no tipo de análise
185     query = {"FANTASIA": {"$regex": unidade_nome, "$options": "i"}}
186     dados = list(collection_sia_analise.find(query, {"_id": 0}))
187
188     if not dados:
189         return ChatResponse(
190             question=chat_request.question,
191             processed_question=ProcessedQuestionResponse(
192                 processed_question=processed_question,
193                 analysis_type=analysis_type
194             ),
195             results=[]
196         )
197
198     if analysis_type == "procedimentos_apresentados":
199         total_procedimentos = sum(item.get("TOTAL_PA_QTDPRO", 0) for item in dados)
200         results.append(QueryResult(unidade=unidade_nome, total=total_procedimentos))
201     elif analysis_type == "procedimentos_aprovados":
202         total_aprovados = sum(item.get("TOTAL_PA_QTDAPR", 0) for item in dados)
203         results.append(QueryResult(unidade=unidade_nome, total=total_aprovados))
204
205     return ChatResponse(
206         question=chat_request.question,
207         processed_question=ProcessedQuestionResponse(
208             processed_question=processed_question,
209             analysis_type=analysis_type
210         ),
211         results=results
212     )
213
214     except Exception as e:
215         print(f"Erro no chat: {e}")
216         raise HTTPException(status_code=500, detail="Erro interno no servidor.")

```

4. Tratamento Robusto de Erros na API:

Exceções HTTP: Implemente um tratamento robusto de erros na API, utilizando exceções HTTP específicas em FastAPI para lidar com problemas que possam ocorrer durante as requisições (como falha ao carregar o modelo, problemas com os dados de entrada, ou timeouts).

Exemplo: Se o modelo LLM não estiver carregado corretamente ou ocorrer um erro na manipulação dos dados, a aplicação deve retornar um código de erro adequado (ex: 503 Serviço Indisponível) e uma mensagem explicativa.

Melhoria da Experiência do Usuário: Garanta que os erros sejam tratados de forma clara e informativa, proporcionando feedback útil para o usuário final e facilitando o diagnóstico de possíveis problemas.

Exemplos de tratamento

```
166         except Exception as e:
167             raise HTTPException(
168                 status_code=503,
169                 detail=f"Erro ao processar a pergunta no modelo LLM: {e}"
170             )
```

```
189         try:
190             unidade_match = re.search(r"unidade\s(.+)", pergunta)
191             unidade_nome = unidade_match.group(1).strip() if unidade_match else None
192             if not unidade_nome:
193                 raise ValueError("Nome da unidade não encontrado na pergunta.")
194         except Exception as e:
195             raise HTTPException(
196                 status_code=400,
197                 detail=f"Erro ao identificar a unidade de saúde: {e}"
198             )
```

```
213         except Exception as e:
214             raise HTTPException(
215                 status_code=500,
216                 detail=f"Erro ao consultar o banco de dados: {e}"
217             )
```