

Guia pratico: uso da OpenAI em um projeto Python

1) Chaves e segurança

- Guarde a OPENAI_API_KEY em variável de ambiente (.env) e nunca comite no git.
- Em runtime, valide se a chave existe e, se não existir, solicite o input no terminal.

2) Modelos e custo

- Para embeddings, use um modelo de embeddings e mantenha o mesmo na ingestão e na busca.
- Para respostas (LLM), escolha um modelo adequado ao tamanho do contexto e latência desejada.
- Controle custos limitando o tamanho do contexto e o número de documentos recuperados.

3) Embeddings e RAG

- Extraia texto do PDF, divida em chunks (por exemplo 1000 caracteres com overlap 150).
- Gere embeddings de cada chunk e armazene em um banco vetorial (por exemplo Postgres com pgvector).
- Na pergunta, gere embedding da query e recupere os top-k mais similares (k=10).
- Monte o prompt com CONTEXTO e REGRAS para responder somente com base no contexto recuperado.

4) Qualidade e confiabilidade

- Use temperatura baixa para reduzir alucinação.
- Logue k, scores e páginas para depurar recuperação.
- Se o contexto não contiver a informação, responda com uma mensagem padrão de falta de dados.

5) Operação

- Defina timeouts de rede e trate erros de API e de conexão com o banco.
- Rode a ingestão antes do chat e mantenha uma rotina para resetar a coleção em desenvolvimento.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.

Boas práticas adicionais: mantenha prompts versionados, use validação de entrada, e monitore latência e taxa de erros para evoluir o sistema com segurança.