

Biscoitos Tia Gertrudes

github.com/ImLusca/Biscoitos-Tia-Gertrudes

Tia Gertrudes é uma senhora muito doce que possui o dom da culinária. Porém, nos tempos modernos, seus quitutes vendidos em sua bela lojinha vêm perdendo espaço para as vendas do uaiComida, aplicativo de celular que permite que as pessoas comprem açúcares e gorduras saturadas sem sair de casa. E como essa nova geração é preguiçosa e acredita que um bom doce feito à mão com muito amor e carinho não vale uma caminhada até sua lojinha, tia Gertrudes vem vendendo cada vez menos.

Então, para superar esse problema, tia Gertrudes pede para que seu muito prestativo sobrinho Lucas crie um desses sites online para que sua lojinha também possa fazer parte dessa nova era. Mas tia Gertrudes não conseguia navegar no próprio site. E para consertar isso, pediu para que Lucas fizesse com que ela pudesse conversar com o computador para ele levá-la para as páginas corretas. E assim começamos essa epopeia:

Do que se trata?

Este projeto tem como permitir busca semântica para encontrar recursos dentro de um site.

Como é feito?

Para isso, geramos um Json com todos os recursos disponíveis no site. Geramos palavras-chave que podem ser relacionadas a esses recursos e criamos embeddings delas. Depois disso, quando o usuário faz uma solicitação, pegamos o seu request, e comparamos o embedding da sua solicitação com os embeddings disponíveis, e assim, geramos o contexto que alimentará um prompt para um modelo de geração de texto que irá construir a response correta baseada no contexto que lhe foi fornecido.

Mas oque são Embeddings?

Embeddings são uma técnica fundamental em processamento de linguagem natural (PLN) e aprendizado de máquina que transforma dados categóricos, como palavras, frases ou até mesmo imagens, em vetores numéricos de dimensão fixa. Esses vetores capturam informações semânticas e relacionamentos entre os dados, permitindo que algoritmos de machine learning trabalhem com eles de forma eficiente.

Exemplo:

The screenshot shows a REST client interface. On the left, a POST request is configured to `127.0.0.1:3000/embed`. The body is a JSON object: `{ "text": "me dá a lista dos cookies mais baratos" }`. On the right, the response is shown with a status of 200 OK, size of 18.77 KB, and time of 1.91 s. The response body is a JSON object with an `embedding` array containing 28 numerical values.

```
{
  "embedding": [
    -0.03652685,
    -0.038854048,
    -0.009643949,
    -0.03146773,
    0.002622830,
    -0.021716274,
    -0.04267368,
    0.042876046,
    -0.015392374,
    -0.047479846,
    0.04290134,
    -0.04034649,
    0.0030623488,
    0.0037058059,
    0.03915759,
    -0.005631434,
    -0.050995935,
    -0.01771957,
    -0.00039583666,
    0.015139419,
    0.01259721,
    0.0008011592,
    -0.02658568,
    0.009232895,
    0.028301666,
    0.0006521523,
    0.0009106417,
    -0.0053721536,
    0.012135565,
    -0.019249953,
    -0.020059412,
    -0.037412196,
    -0.016075356,
    0.016745688,
    -0.0075254417
  ]
}
```

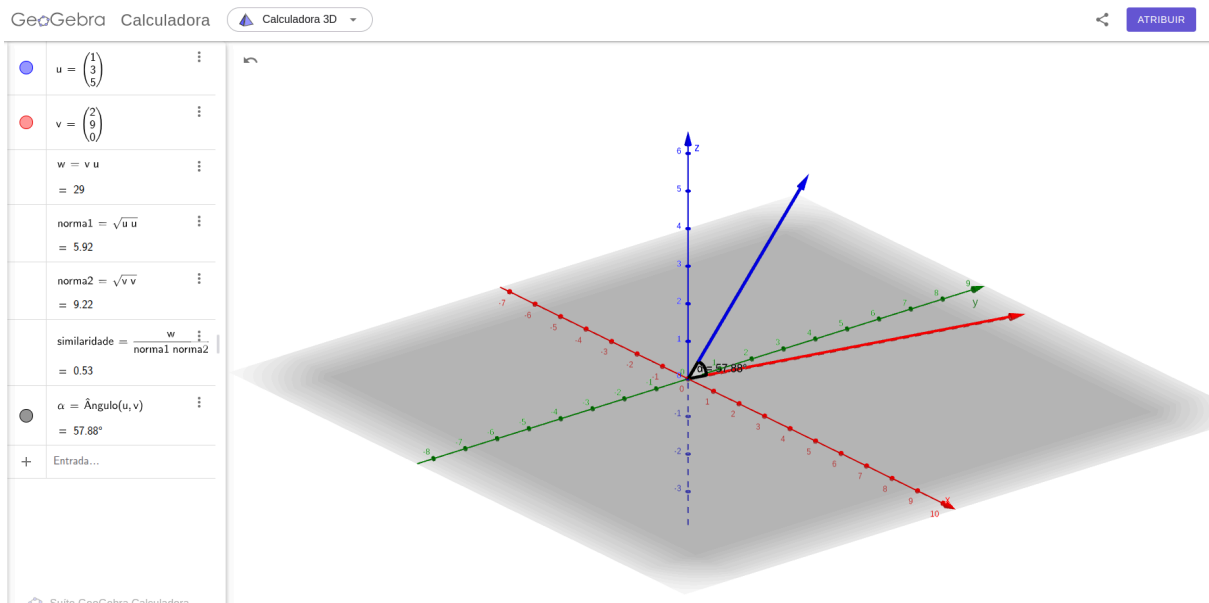
A frase "me dá a lista dos cookies mais baratos" após ser convertida, se torna uma lista de números entre -1 e 1 que podem ser representados como um vetor em um espaço multidimensional.

Se também gerarmos o embedding para a frase "me fala quais são os doces menos caros", receberemos outra lista de números, com o mesmo número de dimensões, que no espaço multidimensional, é muito próximo do primeiro vetor gerado.

Para calcular o quão similar um vetor em N dimensões é de outro vetor, podemos usar similaridade de cossenos, que é uma técnica que divide o produto vetorial dos dois vetores pela magnitude desses dois vetores, resultando em um número entre 0 e 1 que indica o quão próximos esses vetores são!!

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

Exemplo em 3 dimensões:



Com tudo isso, conseguimos filtrar quais são os dados mais relevantes para entregarmos ao nosso modelo de geração de texto como contexto, reduzindo o tempo de processamento, custo de tokens e as chances de alucinações.

```
const generateResponse = async (
  userInput: string,
  references: Referencia[]
): Promise<Resposta> => {
  const prompt = `
    - Dadas as instruções abaixo, retorne exatamente e apenas um objeto JSON seguindo a interface:
    Resposta {
      page: string;
      sort: boolean;
      filter: boolean;
      sort_value?: string;
      filter_value?: string;
    }
    - Você receberá um input vindo do usuário e uma lista de referências extraídas de uma base maior.
    Caso não consiga encontrar uma resposta, retorne um JSON vazio.
    - O input do usuário é: ${userInput}
    - As referências são: ${JSON.stringify(references)}
  `;

  try {
    const completion = await openai.chat.completions.create({
      model: "gpt-3.5-turbo",
      messages: [
        {
          role: "user",
          content: prompt,
        },
      ],
    });
  }
};
```

Depois disso, por segurança, validamos se a resposta do gerador de texto realmente é um json no formato que solicitamos, e se for, Retornamos para o usuário!

Pontos de melhoria:

Usar banco de vetores em vez de jsons,
Migrar para uma linguagem que permita paralelização real mais facilmente,
Aumentar a base de dados e palavras chaves,
monitorar logs de users para alimentar a base de dados.