Recuperação de Informação 2022.1





Jefferson Rodrigues

Ciência da computação jrm@cin.ufpe.br

Igor Eduardo

Ciência da computação iem@cin.ufpe.br

Victor Vasconcelos

Ciência da computação vvb2@cin.ufpe.br







Domínio: Receitas de Bolo

Sites Escolhidos:

- https://www.tudogostoso.com.br/
- https://www.comidaereceitas.com.br/
- https://receitaculinaria.com.br/
- https://receitas.globo.com/
- https://vovopalmirinha.com.br/
- https://www.receiteria.com.br/
- https://www.receitasnestle.com.br/
- https://cybercook.com.br/
- https://www.panelinha.com.br/
- https://br.recepedia.com/

cin.ufpe.br



Crawler

Victor Borges

Ciência da computação ivvb2@cin.ufpe.br



Crawler:

```
def decodeSite(url):
    if url == "#primary":
        ua = UserAgent()
        header = {'User-Agent': str(ua.chrome)}
        retorno =""
        for t in range(tentativasSite):
        result = req.get(url, headers=header)
        if result.status_code == 200:
            retorno = result.content.decode()
            break
        return retorno
def getSoup(url):
    return BeautifulSoup(decodeSite(url), "html.parser")
def salvarArquivo(url,dominio,texto):
    guid = str(uuid.uuid4())
    path = os.path.join(config['PastaRaiz'],"Crawler",dominio)
    bancoSite.append((guid,url))
    f = open(os.path.join(path, guid + ".txt"), "w",encoding='utf8')
    f.write(str(texto.encode('utf-8')))
    f.close()
```

O Crawler decodifica a página até encontrar o status 200.

Depois do decoder, ele vai para o salvar arquivo onde gera o Guid e Path para salvar o arquivo na pasta do seu domínio.



Crawling:

```
def navegacao2(paginapai,pagina):
    if paginapai == pagina or pagina == "#primary" or len(bancoSite) > 100 or existeUrl(pagina):
        return

paginaText = decodeSite(pagina)
    salvarArquivo(pagina,dominio,paginaText)

print(str(len(bancoSite)) + " - " + paginapai + " - " + pagina)

soup1 = BeautifulSoup(paginaText, "html.parser")
    listAncora1 = soup1.find_all('a', href=True)

for site in listAncora1:
    url = site["href"]
    if pegarDominio(url) == pegarDominio(paginapai):
        navegacao2(pagina,url)
```

Aqui o Crawler pega as âncoras que foram salvas e faz a recursão baixando os sites Positivos.





Crawler:

Dificuldades:

- Links que levam para fora do site (aplicativo, redes sociais, loja, etc) Sites com alta "indisposição" (verificação de robôs) a serem minerados
- Links com prefixo
- Links com redirecionadores
- Especificidades de links de determinadas páginas
- Parâmetros causando comportamento obsessivo na busca heurística



Igor Eduardo

Ciência da computação iem@cin.ufpe.br



Treino e teste de rótulos

- → Usamos 10 links positivos e 10 negativos rotulados previamente de cada site para avaliar os classificadores;
- → Um total de 200 HTMLs;
- → As páginas selecionadas como exemplos positivos são paginas que são claramente
 uma receita de bolo.
- As páginas selecionadas como **exemplos negativos são paginas que não são uma**receita de bolo , porém ainda podem ser outros tipos de receitas ou paginas dos sites escolhidos que não são receitas como dicas de cuilinaria, blog, loja, etc.



TF-IDF:

O TF-IDF significa frequência de termos - documentos inversos de frequência.

TF-IDF é na verdade a combinação de duas métricas a tf e a idf com o intuito de se obter um peso frequentemente usado nas áreas de recuperação de informações, mineração de texto, processamento de linguagem natural, sistemas de recomendação, etc. E esse peso é uma medida estatística normalmente usada para avaliar a importância de uma palavra para um documento a partir de uma coleção de documentos, chamada de corpus.

cin.ufpe.br



TF:

O TF:

Frequência de termo, que mede com que frequência um termo ocorre em um documento.

Como cada documento tem um comprimento diferente, é possível que um termo apareça muito mais vezes em documentos longos do que em documentos mais curtos.

Por isso a frequência de um termo é frequentemente dividida pelo comprimento do documento, como forma de normalização.

supondo que um documento tenha 20 palavras e 5 delas é a palavra "açucar". o TF será calculado como: TF: 5/20 = 0.25

cin.ufpe.br



TF:

```
def pegaTFPalavrasDoc(path):
   dicTF = dict()
   f = open(path, mode="r", encoding="utf-8")
   soup = BeautifulSoup(f.read())
   for script in soup(["script", "style"]):
       script.extract()
   text = soup.get text()
  lines = (line.strip() for line in text.splitlines())
   chunks = (phrase.strip() for line in lines for phrase in line.split(" "))
   text = '\\n'.join(chunk for chunk in chunks if chunk)
  listText = text.split("\\n")
   for linha in listText:
       for palavra in linha.split(" "):
           palavra = palavra.lower()
          if dicTF.get(palavra) is None :
               dicTF[palavra] = 1
              dicTF[palavra] = dicTF[palavra] + 1
   return dicTF
```





IDF:

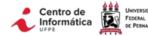
Precisamos ponderar os termos frequentes enquanto aumentamos os termos raros.

Devido a isso o IDF é calculado como logaritmo do numero dos documentos no corpus dividido pelo numero de documentos onde o termo especifico aparece.

suponha que temos **5** documentos no total e a palavra "açucar" aparece em **2** deles o IDF será calculado como: IDF: **log(5/2)=0.398**

Portanto, o tf-idf da palavra açucar será o produto dessas duas metricas.

0,25(tf)x(idf)0,398 = 0.0995



IDF:

```
def computaIDF(palavra,tipo):
    if(tipo == "positivo"):
        qtdOcorrenciaPositiva = qtdOcorrencia(listaTFDicpositivo,palavra)
        if qtdOcorrenciaPositiva is None or qtdOcorrenciaPositiva == 0:
            return 0
        return math.log10(200 / qtdOcorrenciaPositiva)
    else:
        qtdOcorrenciaNegativa = qtdOcorrencia(listaTFDicNegativo,palavra)
        if qtdOcorrenciaNegativa is None or qtdOcorrenciaNegativa == 0:
            return 0
        return math.log10(200 / qtdOcorrenciaNegativa)
```





Dificuldades:

- Muitos dos htmls baixados estavam muitos "sujos", o tratamento dos dados foram

uma das dificuldades.

```
for key in dictPAth:
   for file in os.listdir(dictPAth[key]):
       textoDict = dict()
       fullpath = os.path.join(dictPAth[key],file)
       f = open(fullpath, mode="r", encoding="utf-8")
       soup = BeautifulSoup(f.read())
       for script in soup(["script", "style"]):
           script.extract()
       text = soup.get_text()
       lines = (line.strip() for line in text.splitlines())
       chunks = (phrase.strip() for line in lines for phrase in line.split(" "))
       text = '\\n'.join(chunk for chunk in chunks if chunk)
       listText = text.split("\\n")
       for linha in listText:
           for palavra in linha.split(" "):
               palavra = palavra.lower()
               if textoDict.get(palavra) is None :
                   textoDict[palavra] = 1
                   textoDict[palavra] = textoDict[palavra] + 1
       if key == "positivo":
           listaTFDicpositivo.append(textoDict)
           listaTFDicNegativo.append(textoDict)
```



Extrator

Jefferson Rodrigues

Ciência da computação jrm@cin.ufpe.br



Alguns Exemplos de Extratores Abaixo



Extrator:(comida e receitas)

```
path = "C:\\Users\\jefferson-pc\\Documents\\ProjetoReceitas\\Classificador\\comidasereceitas"

def criaCsvInformacoes(informacoes,path):
    with open(path, mode="w", encoding="utf-8") as csvfile:
    fieldnames = ["NomeReceita", "Ingredientes", "Passos"]
    writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
    writer.writeheader()
    for informacao in informacoes:
        writer.writerow({"NomeReceita": informacao[0], "Ingredientes": informacao[1], "Passos": informacao[2]})

25
```

```
for file in os.listdir(path):
    fullPaht = os.path.join(path,file)
    f = open(fullPaht, oode="n", encoding="utf-8")
    soup = BeautifulSoup(f.read())

ingredientes = soup.find("p").find_next_siblings("ul")

textoIngredientes = " ".join(ingrediente.get_text() for ingrediente in ingredientes)

modoPreparo = soup.find("p",("class":"no-print")).find_next_siblings("ul")[0]

textoModoPreparo = " ".join(modoPreparo.get_text() for modoPreparo in modoPreparo)

nomeReceita = soup.find("h1").get_text()

informacoes = [(nomeReceita, textoIngredientes, textoModoPreparo)]

print(informacoes)
file= "informacoes"
    criaCsvInformacoes(informacoes, "C:\\Users\\jefferson-pc\\Documents\\ProjetoReceitas\\Classificador\\comidasereceitas\\" + file + ".csv")

f.close()
```





Extrator:(cybercook)

```
path = "C:\\Users\\jefferson-pc\\Documents\\ProjetoReceitas\\Classificador\\cybercook"

def criaCsvInformacoes(informacoes,path):
    with open(path, mode="w", encoding="utf-8") as csvfile:
    fieldnames = ["NomeReceita", "Ingredientes", "Passos"]
    writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
    writer.writeheader()
    for informacao in informacoes:
        writer.writerow({"NomeReceita": informacao[0], "Ingredientes": informacao[1], "Passos": informacao[2]})

25
```

```
for file in os.listdir(path):
    fullPaht = os.path.join(path,file)
    f = open(fullPaht, mode="r", encoding="utf-8")
soup = BeautifulSoup(f.read())

ingredientes = soup.find("div",{"class":"mt-4"}).find_all("p")

textoIngredientes = " ".join(ingrediente.get_text() for ingrediente in ingredientes)

modoPreparo = soup.find("div",{"class":"mt-0"}).find_all("p")

textoModoPreparo = " ".join(modoPreparo.get_text() for modoPreparo in modoPreparo)

nomeReceita = soup.find("h1").get_text()

informacoes = [(nomeReceita, textoIngredientes, textoModoPreparo)]
```





Extrator:(receiteria)

```
path = "C:\\Users\\jefferson-pc\\Documents\\ProjetoReceitas\\Classificador\\Receiteria"

def criaCsvInformacoes(informacoes,path):
    with open(path, mode="w", encoding="utf-8") as csvfile:
        fieldnames = ["NomeReceita", "Ingredientes", "Passos"]
        writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
        writer.writeheader()
    for informacoa in informacoes:
        writer.writerow({"NomeReceita": informacao[0], "Ingredientes": informacao[1], "Passos": informacao[2]})

26
```

```
for file in os.listdir(path):
    fullPaht = os.path.join(path,file)
    f = open(fullPaht, mode="r", encoding="utf-8")
    soup = BeautifulSoup(f.read())

ingredientes = soup.find("h2").find_next_siblings("ul")[0].find_all("li")

textoIngredientes = " ".join(ingrediente.get_text() for ingrediente in ingredientes)

print(textoIngredientes)

passos = soup.find("h2").find_next_siblings("ol")[0].find_all("li")

textoPassos = " ".join(passo.get_text() for passo in passos)

print(textoPassos)

nomeReceita = soup.find("h1").get_text()

informacoes = [(nomeReceita, textoIngredientes, textoPassos)]

file= "informacoes"

criaCsvInformacoes(informacoes, "C:\\Users\\jefferson-pc\\Documents\\ProjetoReceitas\\Classificador\\Receiteria\\" + file + ".csv")
```





Extrator:

Dificuldades:

- Fazer um extrator pra cada site

- ✓ Extrator
 - comidasEreceitas.py
- cybercook.py
- extratorCompleto.py
- ExtratorReceiteria.py
- globo.py
- panelinha.py
- receitaculinaria.py

Fazer um extrator que funcione para todos os sites é dificil, pois cada um tem sua particularidade