

Protótipo de sistema de processamento de linguagem natural para inteligência de mercado baseada em notícias no porto de Santos

Prototype of a system for news-based market intelligence at the port of Santos

Julianna Lerner 

Fatec Santos
julerner12@gmail.com

Natália Freitas 

Fatec Santos
natalia1208cfreitas@gmail.com

Guilherme Amorim 

Fatec Santos
tqamorim@gmail.com

RESUMO

Este artigo propõe um método inovador para a sumarização de textos relacionados a notícias portuárias, indicadores de desempenho e produtos. O método se baseia em técnicas avançadas de processamento de linguagem natural, visando identificar as informações mais relevantes e cruciais nos textos analisados. A avaliação do método foi conduzida utilizando um conjunto de dados de textos com notícias portuárias. Os resultados obtidos revelaram a capacidade do método em gerar sumários precisos e concisos, destacando sua eficácia na extração de informações cruciais. Além disso, o modelo demonstrou habilidade em identificar discrepâncias e imprecisões nos textos, sugerindo uma utilidade potencial em sistemas de verificação de notícias. A aplicação prática deste protótipo promete aprimorar significativamente a eficiência na análise e compreensão de textos relacionados ao setor portuário, proporcionando informações mais condensadas e relevantes. Sua capacidade de detecção de imprecisões também destaca sua utilidade em promover a precisão e confiabilidade das informações veiculadas, contribuindo para a integridade de sistemas de verificação de notícias.

PALAVRAS-CHAVE: Sumarização de texto; Processamento de linguagem natural; Inteligência de mercado; Notícias portuárias.

ABSTRACT

This article proposes an innovative method for summarizing texts related to port news, performance indicators, and products. The method is based on advanced natural language processing techniques, aiming to identify the most relevant and crucial information in the analyzed texts. The method was evaluated using a dataset of texts containing port news. The results revealed the method's ability to generate accurate and concise summaries, highlighting its effectiveness in extracting crucial information. Additionally, the model demonstrated the ability to identify discrepancies and inaccuracies in the texts, suggesting potential utility in news verification systems. The practical application of this prototype promises to significantly enhance efficiency in the analysis and understanding of texts related to the port sector, providing more condensed and relevant information. Its ability to detect inaccuracies also underscores its usefulness in promoting the accuracy and reliability of conveyed information, contributing to the integrity of news verification systems.

KEY-WORDS: *Text summarization; Natural language processing; Market intelligence; Port news.*

INTRODUÇÃO

Os portos são importantes instalações para a logística nacional e internacional de um país. De acordo com dados da ANTAQ (2023), há hoje, no Brasil, 175 instalações portuárias de cargas, que incluem portos, terminais marítimos e instalações aquaviárias. O Porto de Santos, na Baixada Santista, em São Paulo, é o mais importante para a economia brasileira.

A globalização dos mercados e a intensificação das transações comerciais e financeiras entre as diversas economias exigem que as empresas estabeleçam mudanças significativas na forma como lidamos com os dados. Dessa maneira, a coleta de grandes volumes de dados necessita de novas soluções, pois, se por um lado temos um volume cada vez maior de dados aumentando vertiginosamente, do outro temos cada vez menos tempo para transformá-los em informações úteis.

De acordo com Probst (2019), apesar das mudanças tecnológicas desenvolvidas nas últimas décadas para lidar com dados e informações, os usuários gastam mais tempo para recuperar informações existentes do que analisando e gerando novos conhecimentos.

A sumarização de textos refere-se ao processo de resumir um texto longo ou um conjunto de textos em um resumo conciso e coerente, preservando as informações mais importantes e relevantes (HUTCHINS, 1987).

Ao relatar um evento a alguém, é comum oferecer um resumo do ocorrido em vez de narrar todos os detalhes minuciosos. De forma inconsciente, todos nós praticamos a arte da sumarização regularmente.

Essa técnica é amplamente utilizada em diversas formas escritas, como em notícias veiculadas em jornais, artigos de revistas e resumos de textos científicos, entre outros contextos. Desse modo, esses textos são posteriormente analisados, proporcionando insights valiosos. Essa abordagem desempenha um papel fundamental ao conseguir identificar padrões e embasar decisões com maior precisão.

Ao incorporar esse avanço tecnológico nos sistemas de verificação de notícias, é possível realizar uma rápida análise do conteúdo informativo, checar a credibilidade das fontes, examinar a estrutura e contexto do texto, além de identificar discrepâncias ou imprecisões. Isso agilizaria a detecção de notícias, possibilitando sua propagação com a certeza de que as informações são confiáveis e precisas.

Para este trabalho, considera-se como objetivo geral apresentar o processo de coleta de dados em sites da web relacionados a notícias portuárias, indicadores de desempenho e produtos e garantir que as informações compartilhadas e consumidas sejam confiáveis e precisas, permitindo que as pessoas tomem decisões informadas e baseadas em fatos com a sumarização dos textos.

2 FUNDAMENTAÇÃO TEÓRICA

Este segmento trata da pesquisa bibliográfica realizada para estabelecer a base teórica que sustenta cada decisão tomada durante a elaboração deste trabalho.

2.1 SUMÁRIOS, ÍNDICES E EXTRATOS

Hutchins (1987) categoriza os sumários científicos em três tipos distintos: indicativos, informativos e críticos.

Os sumários indicativos destacam os pontos essenciais de um texto, sem entrar em detalhes sobre resultados, argumentos ou conclusões. Enquanto isso, os sumários informativos são considerados substitutos do texto original, abarcando todos os aspectos principais. Se o texto original está organizado em seções como dados, métodos, hipóteses e conclusões, um sumário informativo deve conter as informações principais de cada uma dessas seções. Por fim, os sumários críticos atuam como avaliadores, fornecendo, por exemplo, uma análise comparativa entre o conteúdo do texto original e o contexto de outros trabalhos relacionados na mesma área específica.

Hutchins argumenta que a produção automática de sumários indicativos é mais simples do que a modelagem adequada da sumarização humana para os outros tipos de sumários, devido à complexidade envolvida.

Os índices, por sua vez, têm aplicação na classificação de documentos bibliográficos de forma geral, oferecendo uma indicação do seu conteúdo e agilizando o acesso às informações relevantes. Eles permitem que um leitor de uma enciclopédia, por exemplo, vá diretamente ao volume ou página que aborda o tema de seu interesse. Assim, a utilidade dos índices é mais direta e sua função mais limitada do que a dos sumários, o que facilita uma avaliação mais sólida de sua eficiência e qualidade.

Os extratos são outra forma de sumários, sendo uma composição de segmentos relevantes de um texto. Na sumarização automática, sentenças inteiras podem ser extraídas de um texto e agrupadas sem alterações para formar um sumário (PAICE, 1981).

2.2 A ESTRUTURA E CONTEÚDO TEXTUAL

É fundamental que o sumarizador identifique e reproduza as ideias centrais do texto para se criar um resumo eficaz. Da mesma forma, um sistema automático de sumarização precisa compreender o conteúdo apresentado. A compreensão do discurso envolve a análise das estruturas textuais, o que torna essencial investigar essa estrutura antes da sumarização propriamente dita.

Além dos sinais utilizados pelo autor - sejam eles estruturais ou linguísticos - para determinar o conteúdo relevante de um texto-fonte, outros fatores devem ser considerados pelo sumarizador humano. Isso inclui o domínio do assunto específico e o conhecimento prévio como alguém familiarizado com o campo em questão (MUSHAKOJI, 1993).

Embora essas considerações sejam cruciais para a compreensão do processo de sumarização, são desafiadoras de modelar computacionalmente devido à sua subjetividade. Elas demandam uma representação complexa do conhecimento de mundo e do domínio em questão, assim como um modelo sofisticado do escritor eleitor para lidar com as decisões variáveis relacionadas ao conteúdo textual.

2.3 PROCESSAMENTO DE LINGUAGEM NATURAL E SUMARIZAÇÃO DE TEXTOS

Quando se fala em sumarização, é importante lembrar que cada sumário envolve pressuposições e características diversas, assim como conteúdos e correspondência com suas fontes de teores variados.

Os sumários derivados de textos desempenham papéis específicos, podendo atuar como guias ou oferecer informações completas por si só. No primeiro cenário, os sumários são consultados para identificar o tema do texto original; caso interesse, o leitor acessa o texto completo para se aprofundar. No segundo caso, os sumários são suficientemente informativos, dispensando a leitura do texto original, mas ainda oferecendo seus pontos principais.

Dada sua utilidade, frequência e os avanços na área de Processamento de Linguagem Natural (PLN), a automação da sumarização tem despertado grande interesse. A partir do final dos anos 50, métodos estatísticos começaram a surgir para extrair as principais sentenças de um texto. As pesquisas prosseguiram nas décadas seguintes, trazendo avanços significativos à área, como discutido neste relatório, abordando duas perspectivas principais do PLN: a superficial e a profunda, representando métodos distintos de sumarização automática. A abordagem superficial se baseia, em sua maioria, em métodos experimentais e estatísticos, enquanto a perspectiva profunda está associada a teorias formais e linguísticas. (MARTINS, 2001).

Para a automação da sumarização, essas perspectivas apresentam um desafio significativo: modelá-las de maneira apropriada para garantir que os resultados automáticos capturem a variedade de sumários sem perder sua ligação essencial com os textos originais.

É importante destacar que os sumários estão diretamente ligados aos eventos ou textos de origem, e sua elaboração deve garantir a preservação do significado original, embora apresentem menos detalhes e possam adotar estruturas distintas em comparação com as fontes originais. Por isso, é fundamental estabelecer claramente os conceitos fundamentais relacionados aos sumários antes de explorar os princípios que possibilitam a modelagem automatizada.

Sabendo da necessidade da sumarização, o objetivo deste estudo é de exibir o processo de sumarização de textos relacionados a notícias portuárias, indicadores de desempenho e produtos. O método se baseia em técnicas avançadas de processamento de linguagem natural, visando a organização da informação e identificar as informações mais relevantes e cruciais nos textos analisados, permitindo que as pessoas e empresas tomem decisões informadas e baseadas em fato.

3. PROCEDIMENTOS METODOLÓGICOS

A metodologia de pesquisa visa classificar e especificar os métodos empregados na condução do estudo. Seu propósito é explicar as decisões tomadas e descrever os procedimentos realizados, com o intuito de fundamentar os resultados obtidos e possibilitar a replicação do experimento.

3.1 PROCESSAMENTO DE LINGUAGEM NATURAL

Para Alcarde (2023), a Linguagem de Processamento Natural (PLN) é uma área multidisciplinar que engloba conhecimentos de ciência da computação, linguística e estatística, visando desenvolver algoritmos e modelos capazes de entender, interpretar e gerar texto ou fala de maneira similar à forma como os humanos se comunicam. O objetivo do PLN é permitir que os computadores compreendam, interpretem e processem a linguagem humana em seus diversos aspectos, como texto escrito, fala e diálogo.

Conforme Alcarde (2023), alguns exemplos de aplicações de PLN conhecidos são: Assistentes virtuais, como Siri da Apple e a Alexa da Amazon, tradutores automáticos, como Google Tradutor e os chatbots, utilizados para atendimento de clientes.

3.2 PYTHON

A principal ferramenta usada neste estudo é o Python, que é uma linguagem de programação de alto nível. Isso significa que ela se aproxima mais da linguagem humana, sendo composta por tipagem dinâmica, orientada a objetos e multiparadigma, além dos diversos recursos que podem ser encontrados nas bibliotecas e frameworks frequentemente desenvolvidos por toda a comunidade. O código Python é aberto e sua utilização é gratuita, podendo ser instalado em praticamente todos os sistemas operacionais (PYTHON, 2020).

A biblioteca Pandas funciona como um bloco de construção de alto nível para elaborar análises de dados em Python, sendo uma das ferramentas mais utilizadas no mundo por ser altamente adaptável (Python, 2022). "Pandas é uma biblioteca de código aberto, licenciada pelo BSD, que fornece estruturas de dados de alto desempenho e fáceis de usar, além de ferramentas de análise de dados para a linguagem de programação Python" (PYTHON, 2022).

O *Natural Language Toolkit*, ou NLTK, trata-se de uma biblioteca desenvolvida em Python que funciona como uma caixa de ferramentas computacionais voltadas tanto para o Processamento de Língua Natural quanto para a análise computacional da linguagem, oferecendo uma ampla gama de recursos para processamento de linguagem natural (PLN), incluindo *tokenização*, *stemming*, *lemmatização*, *POS tagging*, *parsing* e análise de sentimento. Esses recursos podem ser utilizados para diversas aplicações de PLN, incluindo sumarização de texto.

A sumarização de texto é uma tarefa de PLN que consiste em gerar um resumo conciso de um texto original. Existem dois principais métodos de sumarização automática de texto, a sumarização extrativa: Este método consiste em selecionar as sentenças mais relevantes do texto original para formar o resumo e a sumarização abstractiva, que consiste em gerar um novo texto que resume o conteúdo do texto original. O método utilizado nesse trabalho foi o de sumarização.

3.3 COLETA DE DADOS

A coleta de dados para a realização deste trabalho foi feita por meio da extração manual de notícias portuárias durante o período de janeiro de 2023 até outubro de 2023, através dos sites: NovaCana, Conab-Companhia Nacional de Abastecimento, Forbes, Broadcast, Canal Rural, Globo Rural, Portal Celulose, Investing, Logweb, Datamar News a fim de se criar uma base de dados.

As notícias foram coletadas de acordo com os seguintes critérios: a notícia deve ser sobre o Porto de Santos ou sobre o setor portuário, estar ligada às cargas movimentadas no porto referido ou foi publicada em um dos sites referenciais.

Os requisitos desta base foram definidos como a data da notícia, o título, a notícia e um link para acesso, totalizando cerca de 165 registros armazenados em um arquivo de extensão CSV.

Após a coleta de dados, a base foi transferida para o Python, utilizando as bibliotecas Pandas, para a manipulação dos dados, como a leitura e escrita dos dados, a criação de tabelas e a análise dos dados e NLTK foi utilizada para o processamento de linguagem natural, como a tokenização, a remoção de stop words e a criação de um dicionário para sentenças mais importantes. A tokenização divide o texto em partes menores, como palavras ou frases. A remoção de stopwords remove palavras comuns e sem significado. A criação de um dicionário

para sentenças mais importantes identifica as sentenças mais importantes e as armazena em um dicionário.

A base de dados em Python permitiu realizar as seguintes operações: armazenamento das notícias para facilitar o acesso e a consulta e sumarizar das notícias para obter um resumo do conteúdo.

3.4 TOKENIZAÇÃO

A tokenização, no âmbito do pré-processamento em Processamento de Linguagem Natural (PLN), consiste em fragmentar o texto em unidades menores, conhecidas como tokens. Esses tokens podem abranger palavras, números, símbolos, letras ou até mesmo frases (ALCARDE, 2023). Uma abordagem comum e simples de tokenização é dividir o texto em palavras baseando-se nos espaços existentes ou fragmentar as palavras em partes menores. O propósito da tokenização é viabilizar uma compreensão e processamento mais eficazes do texto por parte dos modelos de PLN.

A sumarização de notícias segue uma metodologia rigorosa, dividida em várias etapas para extrair os dados mais cruciais de um texto informativo. Primeiramente, a notícia é tokenizada usando o método `tokenize`, que segmenta cada palavra na sentença e as dispõe em uma lista indexada para organização, conforme a figura 1.

Figura 1 - Processo de criação

```
# Criar loop para iterar sobre a coluna "Noticia"
corpus = df_base_dados['Noticia'][0]
sentencas = nltk.tokenize.sent_tokenize(corpus)
palavras = nltk.tokenize.word_tokenize(corpus.lower())
```

Fonte: Autores (2023)

A tokenização é executada utilizando a classe *RegexTokenizer* do *Natural Language Toolkit* (NLTK). Essa classe viabiliza a segmentação baseada em expressões regulares, adotando `r'[A-z]\w*'` como o padrão. Esta expressão define o padrão para reconhecer palavras, abrangendo letras maiúsculas e minúsculas, bem como caracteres alfanuméricos e sublinhados, conforme exemplificado na figura 2.

Figura 2- Tokenizer

```
# Regex para tirar acentuacao
tokenizer = RegexpTokenizer(r'[A-z]\w*')
palavras = tokenizer.tokenize(corpus)
```

Fonte: Autores (2023)

Após a tokenização, o próximo passo é lidar com as stop words, palavras comuns como "o", "a", "e", "de", que geralmente não carregam significado essencial na frase. O módulo stopwords do NLTK é usado para obter uma lista dessas palavras, permitindo sua remoção do texto e concentrando-se nas palavras mais relevantes.

Segue-se a criação de uma distribuição de frequência das palavras. Nesse estágio, cada palavra na lista tokenizada é contada para determinar sua ocorrência. Um dicionário é então construído, associando as palavras às suas frequências. Esse dicionário desempenha um papel crucial na formulação do resumo da notícia.

Figura 3 - Stopwords

```
# Remocao das stopwords
stopwords = nltk.corpus.stopwords.words('portuguese')
palavras_sem_stopwords = [w.lower() for w in palavras if w not in stopwords]
frequencia = nltk.probability.FreqDist(palavras_sem_stopwords)

# Criando dicionario para frequencia
sentencas_importantes = defaultdict(int)
```

Fonte: Autores (2023)

O resumo é construído iterando sobre o dicionário de frequências por meio de um loop for. Durante esse processo, as sentenças são avaliadas conforme suas frequências no texto original, e o resumo é composto, dando prioridade às sentenças mais frequentes. As sentenças escolhidas são então agregadas a uma lista, agindo como um repositório temporário para armazenar o resumo da notícia, conforme mostra as figuras 4 e 5.

Figura 4 – Criação de repositório temporário

```
for i, sentenca in enumerate(sentencas):
    for palavra in word_tokenize(sentenca.lower()):
        if palavra in frequencia:
            sentencas_importantes[i] += frequencia[palavra]

idx_setencas_importantes = nlargest(4, sentencas_importantes, sentencas_importantes.get)

# Criar um vetor para armazenar os resumos da noticias
resumo_noticia = []

for i in sorted(idx_setencas_importantes):
    resumo_noticia.append(sentencas[i])
```

Fonte: Autores (2023)

Figura 5 – Exibição do repositório

```
print(df_base_dados['Titulo'][0])
print(resumo_noticia[0])
print(df_base_dados['Link'][0])
```

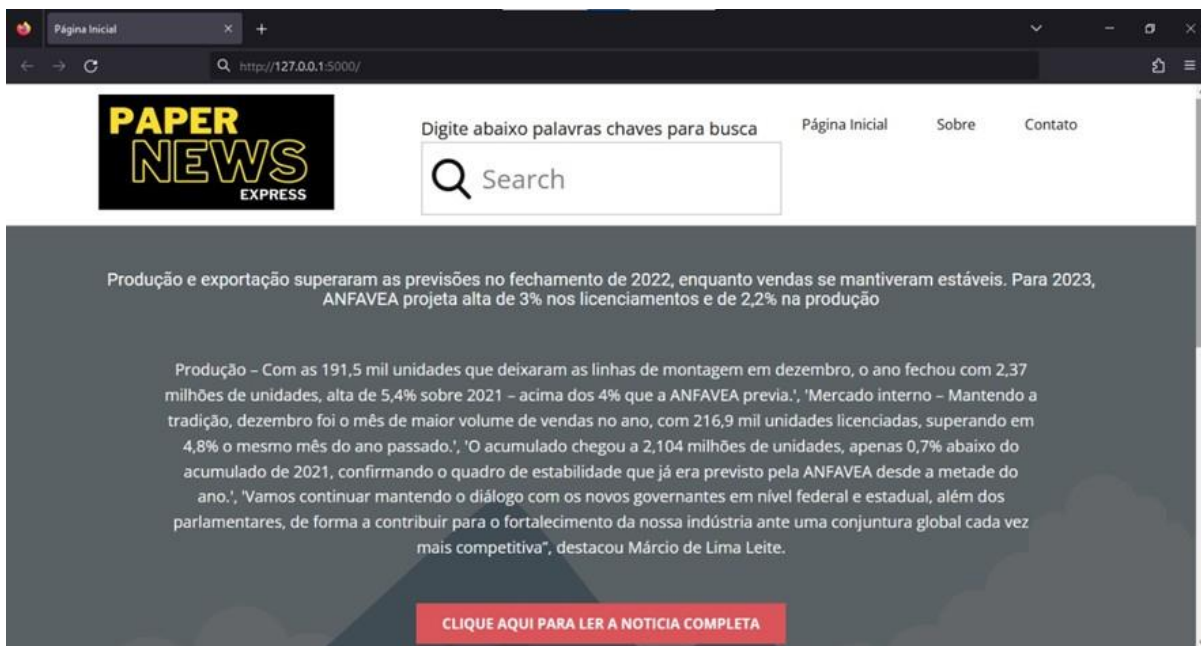
Produção e exportação superaram as previsões no fechamento de 2022, enquanto vendas se mantiveram estáveis. Para 2023, ANFAVEA projeta alta de 3% nos licenciamentos e de 2,2% na produção
Produção - Com as 191,5 mil unidades que deixaram as linhas de montagem em dezembro, o ano fechou com 2,37 milhões de unidades, alta de 5,4% sobre 2021 - acima dos 4% que a ANFAVEA previa.
<https://anfavea.com.br/site/wp-content/uploads/2023/01/RELEASE-JANEIRO-2.pdf>

Fonte: Autores (2023)

Adicionalmente, para tornar a implementação mais acessível ao usuário final, foi criado um protótipo de um aplicativo web. Esse aplicativo possibilita aos usuários a busca livre por notícias e a criação de seus próprios resumos de maneira automatizada.

O resultado final do código gera um formato que inclui o título da notícia, acompanhado pelo resumo gerado automaticamente e um link para acessar a notícia completa, conforme a figura 6. Essa abordagem tem como objetivo oferecer aos usuários uma visão rápida e eficaz do conteúdo das notícias, facilitando a compreensão e o acesso às informações cruciais. Esses métodos metodológicos estabelecem a base essencial para a bem-sucedida implementação do processo de sumarização de notícias.

Figura 6- Página Web



Fonte: Autores (2023)

4. RESULTADOS E DISCUSSÕES

O método de sumarização descrito ainda está em fase de implementação e teste, apesar dos avanços já alcançados, ainda existem algumas dificuldades a serem superadas.

No entanto, foi possível atingir os resultados da criação de um banco de dados real, onde pode ser feito o teste do algoritmo. O resultado é um relatório com uma interação pré-estabelecida para que o usuário possa fazer suas próprias consultas.

Uma das principais dificuldades encontradas foi a automatização da coleta de dados. As APIs disponíveis apresentam problemas de confiabilidade, e o método de web scraping acabou não sendo viável devido à estrutura de cada site, pois os sites utilizados para o teste possuem estruturas diferentes, o que dificulta a criação de um *script* genérico para a coleta dos dados.

Ainda é necessário realizar mais testes para avaliar a qualidade dos resumos gerados. Além disso, é preciso melhorar a automatização da coleta de dados para que o método possa ser utilizado em escala.

Com o desenvolvimento de novas técnicas e a melhoria da infraestrutura tecnológica, é possível que o método de sumarização descrito se torne uma ferramenta valiosa para a pesquisa e auxílio para tomada de decisões.

Desse modo, com o desenvolvimento dessas melhorias, o método de sumarização descrito poderá se tornar uma ferramenta ainda mais útil e eficiente, visando a organização da

informação e permitindo que as pessoas e empresas tomem decisões informadas e baseadas em fatos.

5. CONSIDERAÇÕES FINAIS

O PLN é um campo importante para viabilizar a extração de informação automática de notícias, o que se inicia como pré-processamento dos códigos. A partir desse contexto, esta pesquisa investigou um conjunto de algoritmos de PLN para o pré-processamento de cláusulas de notícias portuárias brasileiras, visando a criação de um produto para automatizar a transformação e a sumarização das informações.

A principal contribuição deste estudo é permitir uma maior compreensão no processo de sumarização e contribui também para uma maior representatividade de trabalhos que demonstram o processo de sumarização em língua portuguesa.

No futuro, os autores têm a intenção de desenvolver uma base de dados com extração automática, substituindo o processo manual. Isso busca automatizar a coleta de dados para garantir a confiabilidade e eficiência do processo, elevando a qualidade dos resumos gerados para assegurar precisão, concisão e clareza. Com isso, almejam criar uma interface mais amigável para o usuário, simplificando a interação com o sistema.

REFERÊNCIAS

ANTAQ. Disponível em: <https://www.gov.br/antag/pt-br>. Acesso em: 28 set. 2023.

BARBOSA, A.; CAVALCANTI, A. **Web Scraping e Análise de dados**. Disponível em: https://www.editorarealize.com.br/editora/anais/conapesc/2020/TRABALHO_EV138_MD4_SA_24_ID1284_24112020001516.pdf. Acesso em: 2 out. 2023.

Broadcast. Disponível em: <http://broadcast.com.br/cadernos/agro/>. Acesso em: 20 set. 2023.

Canal Rural. Disponível em: <https://www.canalrural.com.br/>. Acesso em: 20 set. 2023.

Conab. Disponível em: <https://www.conab.gov.br/>. Acesso em: 20 out. 2023.

Datamar News. Disponível em: <https://datamarnews.com/pt/pt/home-pt/>. Acesso em: 20 set. 2023.

Forbes. Disponível em: <https://forbes.com.br/forbesagro/>. Acesso em: 21 out. 2023.

Globo Rural. Disponível em: <https://globorural.globo.com/>. Acesso em: 20 set. 2023.

GRACIANO, H.; RAMALHO, R. SCRAPERCI. **Um *web scraper* para coleta de dados científicos**. 2023. Disponível em: <https://www.scielo.br/j/eb/a/9QYwtw5kgByRpDFFQB778Tj/?format=pdf&lang=pt>. Acesso em: 1 out. 2023.

HUTCHINS, J. **Summarization: Some problems and Methods**. In: JONES, P. (Org.). **Meaning: The frontier of informatics**. Cambridge: London, 1987. p. 151-173.

Investing. Disponível em: <https://br.investing.com/news/commodities-news>. Acesso em: 20 set. 2023.

MARTINS, C.B.; Pardo, T.A.S.; ESPINA, A.P.; RINO, L.H.M. (2001). “**Introdução à Sumarização Automática**”. Relatório Técnico RT-DC 002/2001, Departamento de Computação, Universidade Federal de São Carlos.

MUSHAKOJI, S. **Constructing ‘Identity’ and ‘Differences’ in Original Scientific Texts and Their Summaries: Its Problems and Solutions**. Seminar Report of Summarizing Text for Intelligent Communication Seminar. Dagstuhl, Germany, 1993.

PAICE, C. D. **The automatic generation of literature abstracts: an approach based on the identification of self-indicating phrases**. In: Information Retrieval Research. Butterworth & Co. (Publishers), 1981.

PROBSTEIN, S. **Reality check: still spending more time gathering instead of analyzing**. Forbes Technology Council, 2019. Disponível em: <https://www.forbes.com/sites/forbestechcouncil/2019/12/17/reality-check-still-spending-more-time-gathering-instead-of-analyzing>. Acesso em: 25 set. 2023.

PYTHON. **Documentação Python**. 2022. Disponível em: <https://www.python.org/about>. Acesso em: 31 set. 2023.