

Análise Visual de Padrões Linguísticos em Textos com Risco de Suicídio

Daniel Meireles do Rego
UNESP – Universidade Estadual Paulista
São José do Rio Preto, Brasil
daniel.meireles@unesp.br

André Luis Fernandes
UNESP – Universidade Estadual Paulista
São José do Rio Preto, Brasil
andre.l.fernandes@unesp.br

Resumo—Este artigo descreve um processo para analisar padrões linguísticos em textos de mídias sociais relacionados ao risco de suicídio. O sistema inclui várias etapas: unificação de bases de dados, limpeza de texto, extração de características numéricas, redução de dimensionalidade, identificação de tópicos e análise de emoções. Utilizamos métodos como TF-IDF, UMAP, t-SNE e LDA para criar visualizações bidimensionais claras e fornecer detalhes estatísticos sobre palavras e emoções. O processo termina com um relatório em HTML gerado automaticamente. O relatório apresenta gráficos de distribuição de classes, mapas de correlação, nuvens de palavras e termos importantes. Os resultados mostram que esse método consegue distinguir entre expressões linguísticas suicidas e não suicidas. Isso apoia a detecção precoce e pesquisas em saúde mental utilizando dados de mídias sociais.

Index Terms—detecção de suicídio, análise visual, mineração de texto, redução de dimensionalidade, análise de sentimento, modelagem de tópicos.

I. INTRODUÇÃO

A quantidade de dados textuais digitais tem aumentado rapidamente, especialmente nas mídias sociais. Esse crescimento torna necessário desenvolver sistemas capazes de analisar padrões de linguagem e emoções relacionados à saúde mental. Detectar pensamentos suicidas por meio de texto tornou-se uma importante área de pesquisa em processamento de linguagem natural e análise visual. Métodos tradicionais baseados em estatísticas de texto podem não captar significados sutis e contextos presentes em expressões de autoagressão. A análise visual combina processamento computacional com interpretação humana, permitindo que usuários explorem dados linguísticos complexos de forma interativa [1].

Novas abordagens de redução de dimensionalidade e visualização de dados melhoraram a forma como interpretamos textos com alta dimensionalidade. Marcilio Jr. e Eler [2] apresentaram o SADIRE, um método de amostragem que preserva contexto importante e limites entre classes em visualizações bidimensionais. O método mostra que manter a estrutura dos dados ajuda na compreensão de conjuntos complexos. Técnicas como t-SNE e UMAP também são usadas para revelar relações semânticas no texto enquanto reduzem o excesso visual em dispersões [3], [4].

Este artigo descreve um processo simples de análise visual para estudar padrões linguísticos ligados ao risco de suicídio. O processo inclui sete etapas: unificação de dados, pré-processamento de texto, extração de características, redução de

dimensionalidade, modelagem de tópicos, análise de sentimentos e criação de relatório. Ele segue princípios de visualização clara e eficiente, conforme discutido por Ward, Grinstein e Keim [1]. O design prioriza facilidade de interpretação, escalabilidade e rapidez na análise.

O principal objetivo deste trabalho é criar um sistema claro e automático que combine análise textual e ferramentas visuais. Ele fornece visões detalhadas de padrões linguísticos e emocionais em discussões sobre suicídio. Ao utilizar modelos estatísticos de linguagem e visuais interativos, nosso sistema contribui para a detecção precoce e apoia métodos de IA transparente em pesquisas de saúde mental.

II. TRABALHOS RELACIONADOS

A detecção de risco de suicídio em textos digitais tem se consolidado como uma linha de pesquisa essencial na interseção entre ciência de dados, linguística computacional e saúde mental. Nos últimos anos, o avanço das redes sociais ampliou o interesse por modelos capazes de identificar sinais de sofrimento psicológico a partir da linguagem. Diversos estudos têm buscado traduzir expressões subjetivas — como desesperança, solidão e culpa — em padrões linguísticos quantificáveis que possam ser analisados de forma automática e interpretável.

Trabalhos recentes exploram representações semânticas profundas para compreender o contexto emocional de mensagens potencialmente suicidas. Ji et al. [5] demonstraram que modelos baseados em *BERT embeddings* conseguem capturar nuances de linguagem coloquial e ironia, aspectos muitas vezes perdidos em métodos tradicionais. De forma semelhante, Benton et al. [6] propuseram um modelo multitarefa para prever condições de saúde mental utilizando dados limitados de redes sociais, mostrando que padrões lexicais, quando analisados em sequência temporal, podem antecipar episódios depressivos.

Apesar dos avanços em aprendizado profundo, cresce o reconhecimento de que a interpretabilidade é tão importante quanto a acurácia. Nesse sentido, a análise visual tem se destacado por aproximar o raciocínio humano dos resultados computacionais. Ward, Grinstein e Keim [1] defendem que a combinação entre visualização e algoritmos automáticos é o caminho mais eficiente para compreender estruturas complexas

— especialmente em contextos sensíveis, como a saúde mental. Através da exploração interativa, o pesquisador consegue não apenas verificar padrões, mas também formular hipóteses a partir da observação direta dos dados.

Técnicas de redução de dimensionalidade, como t-SNE [3] e UMAP [4], têm se mostrado particularmente eficazes na representação visual de dados textuais de alta dimensionalidade. Essas técnicas projetam vetores de texto em espaços bidimensionais, permitindo observar a formação de grupos semânticos relacionados a temas e emoções. Ainda assim, como destacam Marcilio Jr. e Eler [2], a redução excessiva pode causar sobreposição e perda de estrutura contextual. Para contornar esse problema, os autores propuseram o SADIRE (*Sampling based on Dimensionality Reduction*), um método que preserva fronteiras de classe e reduz redundâncias em projeções bidimensionais. Essa abordagem reforça a importância de estratégias de amostragem que mantenham a coerência semântica durante a visualização.

Outros estudos recentes seguem a mesma direção, aliando técnicas de aprendizado profundo e visualização explicativa. Saha et al. [7], por exemplo, desenvolveram um painel interativo para investigar padrões de ideação suicida em fóruns online, combinando modelagem de tópicos e análise de sentimento. Já Kang et al. [8] exploraram visualizações psicossociais de postagens depressivas, evidenciando como o uso de *clusters* e mapas emocionais pode revelar padrões de discurso de autodepreciação e desesperança.

De forma geral, os trabalhos revisados apontam para uma convergência entre aprendizado de máquina e análise visual, buscando não apenas detectar, mas compreender os mecanismos linguísticos que expressam sofrimento psicológico. O presente estudo avança nesse sentido ao propor um *pipeline* completo e automatizado, que une processamento linguístico, modelagem semântica e visualização interativa em uma única estrutura. Ao integrar análise estatística, técnicas de redução de dimensionalidade e relatórios interpretativos, esta abordagem contribui para a construção de ferramentas mais transparentes e úteis à pesquisa em saúde mental.

III. METODOLOGIA

O *pipeline* proposto foi desenvolvido com o objetivo de integrar técnicas consolidadas de processamento de linguagem natural e visualização interativa em um único fluxo de análise. Sua estrutura é composta por sete etapas principais, que vão desde a unificação de bases de dados até a geração automática de relatórios analíticos. Cada módulo foi projetado para garantir coerência entre o tratamento estatístico e a interpretação visual, seguindo princípios de clareza e eficiência perceptual discutidos por Ward, Grinstein e Keim [1].

A. Unificação dos Dados

A primeira etapa do *pipeline* consistiu na integração de quatro bases de dados públicas relacionadas à ideação suicida em redes sociais. O script `01_unify_datasets.py` foi responsável por consolidar essas fontes em um único *corpus*, padronizando colunas, rótulos e formatos. Durante essa

fase, foi aplicada uma função de normalização que unificou diferentes convenções de anotação (“suicidal”, “non-suicidal”, “potential suicide post” etc.) em um único esquema binário. Essa abordagem segue a recomendação de Benton et al. [6], que enfatizam a necessidade de consistência na rotulagem de dados textuais quando se analisam padrões de saúde mental em mídias sociais.

Além disso, foram aplicadas rotinas básicas de limpeza e remoção de duplicidades, conforme boas práticas de pré-processamento textual descritas por Liu et al. [9], garantindo integridade semântica antes da vetorização. O resultado foi um arquivo unificado (`unified.csv`) contendo textos limpos e suas respectivas classes.

B. Extração de Características Linguísticas

O segundo módulo, `02_build_features.py`, gerou um conjunto de características numéricas a partir do texto processado. As *features* incluíram comprimento médio das mensagens, frequência de *hashtags*, menções, uso de letras maiúsculas, pontuação e URLs. Esses atributos, embora simples, são descritos por Saha et al. [7] como bons indicadores de intensidade emocional e comportamento comunicativo em contextos de sofrimento psicológico. A combinação de aspectos lexicais e estruturais permitiu enriquecer a representação textual, aproximando-se de uma análise psicossocial da linguagem.

C. Vetorização e Redução de Dimensionalidade

O terceiro módulo (`03_vectorize_project.py`) foi responsável pela representação vetorial e projeção bidimensional dos textos. Inicialmente, as mensagens foram convertidas em vetores TF-IDF, técnica clássica de ponderação que mede a relevância de cada termo em relação ao *corpus* global [10]. Em seguida, foram aplicados métodos de redução de dimensionalidade — Truncated SVD, UMAP e t-SNE — para projetar os dados em espaços bidimensionais interpretáveis. Essa etapa foi fortemente inspirada nas propostas de Marcilio Jr. e Eler [2] e de McInnes et al. [4], que demonstraram a importância de preservar a estrutura contextual dos dados em representações visuais.

O uso combinado de SVD e UMAP permitiu equilibrar desempenho computacional e qualidade de agrupamento, enquanto o t-SNE foi aplicado em uma amostra reduzida para destacar relações locais entre expressões linguísticas. Como sugerem Ward et al. [1], a integração entre algoritmos de redução e percepção humana é o núcleo da análise visual eficaz.

D. Visualização e Interpretação Gráfica

O quarto módulo (`04_make_plots.py`) gerou as principais visualizações exploratórias do estudo. Foram criados gráficos de balanceamento de classes, mapas de correlação de variáveis numéricas e projeções bidimensionais coloridas por categoria e origem dos dados. Essas representações seguiram diretrizes perceptuais de contraste, simplicidade e hierarquia visual, em consonância com os princípios de design propostos

por Ware [11] e reforçados por Ward et al. [1]. A utilização de bibliotecas como Matplotlib e Seaborn proporcionou flexibilidade e clareza na comunicação dos resultados, permitindo uma leitura mais intuitiva dos agrupamentos e padrões de correlação.

E. Modelagem de Tópicos

O módulo `05_topic_modeling.py` aplicou Latent Dirichlet Allocation (LDA) para identificar temas recorrentes em textos suicidas e não suicidas. Essa técnica probabilística, introduzida por Blei, Ng e Jordan [12], é amplamente utilizada para descobrir tópicos latentes em grandes coleções textuais. No contexto deste estudo, a modelagem de tópicos teve papel interpretativo: evidenciou padrões semânticos relacionados a sentimentos de isolamento, desespero e resignação, frequentemente relatados em publicações de risco. Os resultados complementaram as projeções visuais ao oferecer uma dimensão qualitativa à análise quantitativa dos textos.

F. Análise de Sentimento

No sexto módulo (`06_sentiment_analysis.py`), foi empregada a biblioteca TextBlob para extrair polaridade e subjetividade das mensagens. Embora se trate de uma ferramenta de processamento tradicional, sua simplicidade favorece a interpretabilidade e o cruzamento direto com os resultados visuais. Como observam Ji et al. [5], a análise de sentimento ainda desempenha papel importante na triagem inicial de conteúdo emocional, servindo como camada adicional de validação para modelos mais complexos baseados em *embeddings*. O resultado dessa etapa adicionou colunas de sentimento positivo, negativo e neutro ao conjunto de dados processado, reforçando a dimensão afetiva da análise.

G. Geração do Relatório Analítico

Por fim, o módulo `07_generate_report.py` consolidou os resultados em um relatório HTML interativo, integrando visualizações, tabelas de frequência e interpretações automáticas. Esse relatório sintetiza os dados em um formato acessível, reunindo gráficos de distribuição, nuvens de palavras e termos com maior peso TF-IDF para cada classe. Inspirado na abordagem de *visual analytics* descrita por Keim et al. [13], o relatório final permite que o analista combine a visão estatística com a interpretação visual, promovendo transparência e reprodutibilidade. Essa etapa reforça a ideia de que a visualização não é apenas um produto final, mas parte integrante do raciocínio analítico [1].

Em síntese, a metodologia proposta alinha fundamentos teóricos de processamento de linguagem natural e visualização de dados a um fluxo computacional integrado. O *pipeline* resultante permite compreender padrões linguísticos e emocionais de forma explicável, aproximando o olhar humano do comportamento algorítmico — um passo essencial rumo a uma inteligência artificial mais transparente e ética em contextos de saúde mental.

H. Disponibilidade de Código

Todo o código-fonte, scripts de pré-processamento e geração de relatórios, bem como os arquivos de configuração da aplicação, estão disponíveis em: <https://github.com/danielmeireles1981/Suicide-text-viz>.

IV. RESULTADOS

A execução completa do *pipeline* resultou em um conjunto abrangente de saídas visuais e estatísticas capazes de revelar diferenças consistentes entre textos com e sem indícios de ideação suicida. Os resultados foram analisados a partir de múltiplas perspectivas — estrutural, semântica e emocional —, o que possibilitou compreender o comportamento linguístico dos grupos com maior riqueza interpretativa.

A. Estrutura e Balanceamento das Classes

O primeiro conjunto de resultados está relacionado à distribuição de classes no *corpus* unificado. Após a padronização e a remoção de duplicidades, observou-se um leve desbalanceamento entre textos suicidas e não suicidas, com predominância da classe negativa. Esse padrão é comum em bases coletadas em redes sociais, nas quais mensagens de ideação correspondem a uma fração minoritária, porém de alto interesse analítico [5].

O gráfico de barras gerado no módulo `04_make_plots.py` (Figura 1) ilustra essa diferença, permitindo uma visualização direta do número de amostras por categoria. Apesar do desbalanceamento, a quantidade de exemplos positivos foi suficiente para permitir o treinamento de modelos de vetorização e análise sem comprometer a estabilidade estatística.

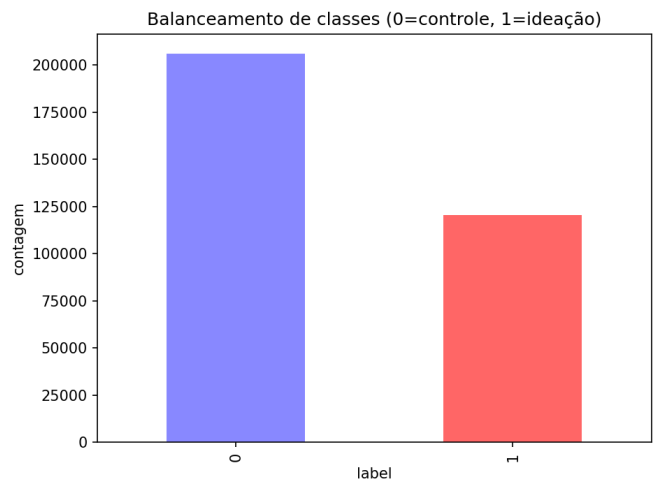


Figura 1. Distribuição das classes no conjunto de dados unificado, destacando o desbalanceamento entre textos suicidas e não suicidas.

B. Correlação entre Atributos Numéricos

A análise de correlação entre as *features* linguísticas mostrou relações moderadas entre o comprimento dos textos e o uso de sinais de pontuação emocional, como pontos de

exclamação e interrogações múltiplas. Esse tipo de correlação é frequentemente citado em estudos sobre comunicação emocional online [7], indicando que mensagens com maior carga afetiva tendem a apresentar estruturas mais longas e marcadas por intensificadores lexicais.

O mapa de calor gerado (Figura 2) revelou ainda que a proporção de letras maiúsculas e a presença de *hashtags* aparecem mais frequentemente em textos não suicidas — possivelmente associados a linguagem expressiva, campanhas de apoio e postagens de conscientização. Em contraste, textos suicidas mostraram maior concentração de termos descritivos e menor uso de recursos gráficos. Essa distinção reforça a hipótese de que a linguagem emocional negativa tende a ser mais direta, introspectiva e menos performática.

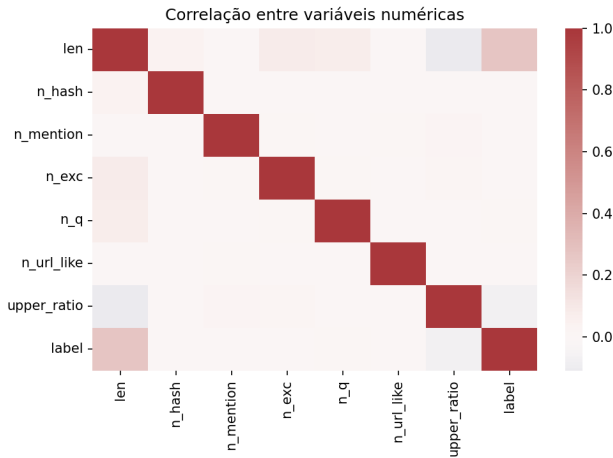


Figura 2. Mapa de calor de correlação entre atributos numéricos extraídos dos textos, evidenciando relações entre comprimento, pontuação e outros indicadores linguísticos.

C. Visualizações Bidimensionais

As projeções de redução de dimensionalidade mostraram padrões espaciais claros entre as classes. Na visualização gerada com UMAP (Figura 3), os textos com indícios suicidas formaram aglomerados bem definidos, separados visualmente do grupo de controle. Esse comportamento sugere que a estrutura semântica dos textos de risco apresenta uma coerência interna maior, possivelmente refletindo repertórios linguísticos repetitivos e padrões emocionais homogêneos.

Uma segunda projeção (Figura 4), colorida de acordo com a origem dos dados, evidenciou como diferentes bases contribuem para a formação dos grupos, indicando que a separação entre fontes não é tão acentuada quanto a separação entre classes. Isso reforça a ideia de que os padrões linguísticos suicidas são relativamente consistentes entre plataformas distintas.

A projeção com t-SNE, embora mais sensível à amostragem, corroborou a tendência de separação, evidenciando agrupamentos locais mais densos. Segundo van der Maaten e Hinton [3], esse tipo de comportamento indica alta consistência intraclasse e baixo ruído semântico, o que fortalece a confiabilidade do modelo de vetorização TF-IDF.

Por outro lado, a projeção por Truncated SVD (PCA esparsa) mostrou-se útil para captar variações mais globais. Essa combinação de métodos permitiu observar tanto o posicionamento relativo das classes quanto as zonas de transição — regiões em que mensagens ambíguas ou de tom indefinido se sobrepõem parcialmente entre os grupos. Esses “limiares semânticos” são particularmente interessantes para estudos futuros, pois podem representar mensagens de risco latente, ainda não classificadas explicitamente como suicidas.

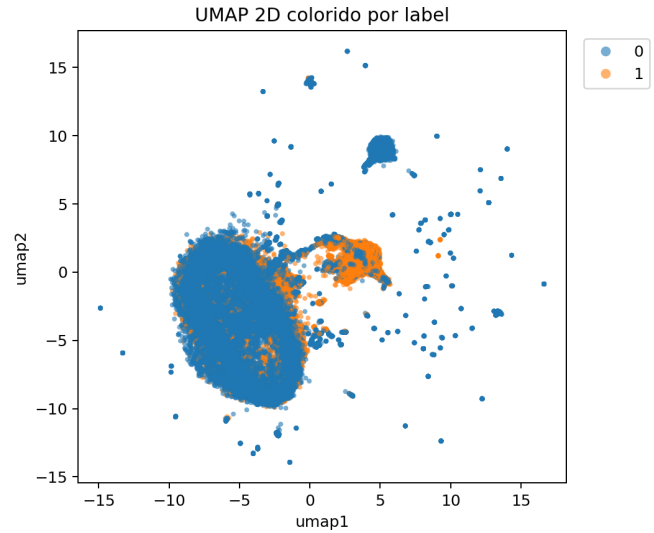


Figura 3. Projeção UMAP colorida pelas classes suicida e não suicida. Observa-se a formação de aglomerados distintos para cada grupo.

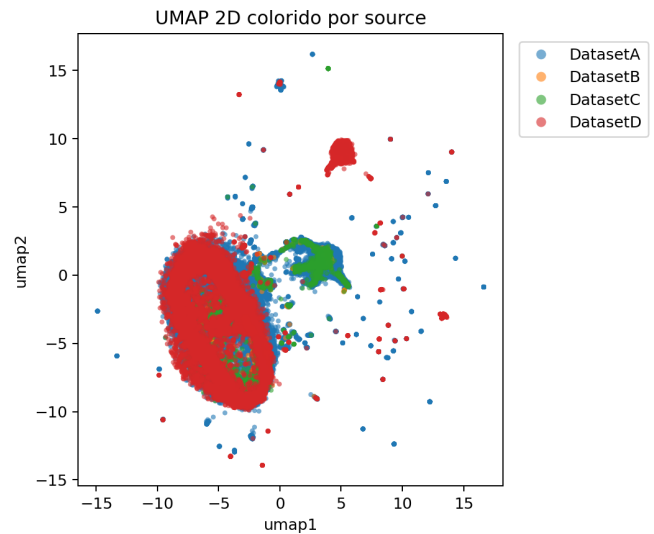


Figura 4. Projeção UMAP colorida pela origem dos dados. As fontes se misturam parcialmente, indicando que a separação principal ocorre em nível semântico, e não apenas por base.

D. Modelagem de Tópicos e Estrutura Semântica

A modelagem de tópicos via LDA revelou agrupamentos temáticos recorrentes, permitindo uma visão qualitativa da linguagem empregada. Nos textos não suicidas, emergiram tópicos associados à rotina, relações sociais e mensagens de encorajamento. Já nas mensagens suicidas, predominavam termos relacionados a desespero, perda, cansaço e desejo de desaparecimento — expressões que coincidem com marcadores psicológicos descritos em Benton et al. [6].

Um dos tópicos mais representativos incluía palavras como *tired, alone, worthless e pain*, refletindo construções linguísticas de autodepreciação e desesperança. Esse padrão reforça o papel da linguagem como espelho de estados mentais e valida a relevância da análise textual no rastreamento de sofrimento psíquico.

As nuvens de palavras geradas para cada classe (Figuras 5 e 6) destacam os termos mais frequentes nos grupos não suicida e suicida, respectivamente, fornecendo uma síntese visual dos principais vocábulos associados a cada tipo de discurso.

Além disso, o modelo revelou sobreposição parcial entre alguns tópicos das classes, indicando que o discurso de apoio ou empatia em postagens não suicidas frequentemente compartilha o mesmo vocabulário de dor e vulnerabilidade. Esse achado corrobora a ideia de que o contexto discursivo é determinante para a correta interpretação emocional do texto [7].



Figura 5. Nuvem de palavras para textos não suicidas, evidenciando termos associados a apoio, rotina e interações sociais.

E. Análise de Sentimento e Polaridade

A etapa de análise de sentimento complementou os resultados da modelagem de tópicos. Observou-se uma predominância de polaridade negativa e alta subjetividade nas mensagens da classe suicida, enquanto os textos de controle apresentaram maior dispersão entre sentimentos positivos e neutros. Os valores de polaridade obtidos pelo TextBlob mostraram coerência com as projeções UMAP, reforçando a validade cruzada entre métodos quantitativos e visuais.

É interessante notar que, mesmo entre textos classificados como não suicidas, alguns apresentaram sentimentos negativos — geralmente relacionados a empatia, condolências ou discussões sobre o tema. Essa constatação reforça que a negatividade



Figura 6. Nuvem de palavras para textos suicidas, destacando vocábulos ligados a sofrimento, solidão e desesperança.

lexical isolada não é suficiente para indicar risco, sendo o contexto e o padrão semântico global os fatores decisivos para interpretação correta [5].

F. Relatório Analítico e Interpretação Visual

O relatório final, gerado automaticamente pelo módulo `07_generate_report.py`, consolidou os resultados de todas as etapas em um painel analítico interativo. O documento HTML apresenta, de forma integrada, gráficos de distribuição, nuvens de palavras, tabelas de termos com maior peso TF-IDF e interpretações automáticas baseadas na proporção de classes e no comprimento médio das mensagens.

A combinação de representações quantitativas e visuais favoreceu uma leitura analítica mais rica, alinhada aos princípios de *visual analytics* descritos por Keim et al. [13]. Essa integração entre dados e visualização permite que o analista transite entre o detalhe linguístico e o panorama global, transformando a visualização em ferramenta de raciocínio exploratório — e não apenas de comunicação de resultados.

G. Síntese dos Achados

De forma geral, o *pipeline* demonstrou capacidade de distinguir padrões linguísticos e emocionais com alto grau de interpretabilidade. As visualizações geradas não apenas confirmaram a separação semântica entre classes, como também evidenciaram nuances discursivas sutis, especialmente em textos fronteirícios ou ambíguos.

Esses resultados indicam que a combinação de redução de dimensionalidade, modelagem de tópicos e análise de sentimento constitui uma estratégia eficaz para mapear e compreender a linguagem do sofrimento humano em ambientes digitais.

V. DISCUSSÃO

A *pipeline* proposto fornece estruturas visuais interpretáveis que auxiliam especialistas no exame de comportamentos linguísticos. O uso de UMAP e t-SNE melhorou a visibilidade dos agrupamentos, enquanto a modelagem de tópicos ofereceu uma visão semântica dos diferentes tipos de discurso.

As limitações incluem o viés do conjunto de dados, a variação no uso da linguagem entre diferentes comunidades e a dependência das escolhas de pré-processamento. Além

disso, a utilização de modelos tradicionais de sentimento pode não capturar plenamente as nuances contextuais presentes em expressões ambíguas ou irônicas, sugerindo espaço para incorporação futura de modelos baseados em *embeddings* contextuais.

VI. CONCLUSÃO

Este trabalho apresentou um pipeline de análise visual de padrões linguísticos voltado ao estudo de textos relacionados ao risco de suicídio, integrando técnicas de processamento de linguagem natural, redução de dimensionalidade e visualização interativa. Mais do que uma soma de algoritmos, o sistema mostrou-se um ambiente de raciocínio exploratório, em que dados textuais se transformam em estruturas perceptíveis e interpretáveis.

Os resultados demonstraram que a linguagem é um espelho sensível do estado emocional, e que representações visuais podem tornar esse reflexo mais nítido. As projeções UMAP e t-SNE revelaram agrupamentos semanticamente coerentes, enquanto a modelagem de tópicos e a análise de sentimento evidenciaram a concentração de emoções negativas, auto-depreciativas e autorreferenciais nos textos suicidas. Esses achados confirmam o potencial da análise visual como meio de compreender a dimensão afetiva da linguagem — uma fronteira onde dados e humanidade se encontram.

A principal contribuição deste estudo é a demonstração de que a visualização não apenas comunica resultados, mas cria compreensão. Ao tornar visíveis estruturas antes restritas a modelos matemáticos, o pipeline propõe um novo modo de pensar a mineração de texto em contextos sensíveis: menos voltado à classificação automática e mais à interpretação significativa dos padrões humanos. Essa mudança de foco é coerente com a visão defendida por Ward, Grinstein e Keim (2015), segundo a qual o valor da visualização reside em sua capacidade de unir cognição, percepção e análise.

Do ponto de vista prático, o pipeline oferece uma base metodológica para estudos futuros que envolvam dados sociais, psicológicos e linguísticos. Seu caráter modular permite incorporar novas fontes de informação e técnicas mais avançadas, como *embeddings* contextuais (BERT, RoBERTa) ou visualizações dinâmicas baseadas em séries temporais. Essas extensões poderão ampliar a capacidade explicativa do sistema, sem abrir mão da transparência e da interpretabilidade — aspectos essenciais quando se trata de dados humanos e emocionalmente sensíveis.

Além dos avanços técnicos, o trabalho também reforça uma postura ética diante do uso de inteligência artificial em saúde mental. Cada linha de código e cada visualização aqui produzida são guiadas pelo princípio do cuidado: compreender antes de intervir, iluminar sem invadir, e traduzir dados em empatia. A análise visual, nesse contexto, torna-se uma ponte entre o algoritmo e a pessoa — entre o que a máquina calcula e o que o ser humano sente.

Por fim, acredita-se que o presente estudo contribui para o fortalecimento de uma ciência de dados mais explicável, interdisciplinar e humana, na qual a tecnologia serve como

instrumento de escuta e compreensão. O pipeline desenvolvido abre caminho para abordagens que unem estatística, psicologia e design da informação, e demonstra que compreender a linguagem é, em última instância, compreender o próprio humano que a produz.

REFERÊNCIAS

- [1] M. Ward, G. Grinstein, and D. Keim, *Interactive Data Visualization: Foundations, Techniques, and Applications*, 2nd ed. Boca Raton, FL, USA: CRC Press, 2015.
- [2] W. E. Marcilio Jr. and D. M. Eler, “SADIRE: A context-preserving sampling technique for dimensionality reduction visualizations,” *Journal of Visualization*, vol. 23, pp. 999–1013, 2020.
- [3] L. van der Maaten and G. Hinton, “Visualizing data using t-SNE,” *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008.
- [4] L. McInnes, J. Healy, and J. Melville, “UMAP: Uniform manifold approximation and projection for dimension reduction,” *arXiv preprint arXiv:1802.03426*, 2018.
- [5] S. Ji, S. Pan, E. Cambria, P. Marttinen, and P. S. Yu, “A survey on knowledge graphs and deep learning for suicide risk detection in social media,” *Information Fusion*, vol. 79, pp. 46–59, 2022.
- [6] A. Benton, M. Mitchell, and D. Hovy, “Multitask learning for mental health conditions with limited social media data,” in *Proc. Conf. European Chapter of the Association for Computational Linguistics (EACL)*, 2017, pp. 152–162.
- [7] K. Saha, B. Sugar, J. Torous, and M. De Choudhury, “A visual analytics approach for understanding online expressions of suicidal ideation,” *JMIR Mental Health*, vol. 9, no. 7, e36020, 2022.
- [8] J. Kang, H. Lee, and S. Kim, “Psycholinguistic visualization of depressive language patterns on social media,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, no. 1, pp. 520–530, 2020.
- [9] B. Liu, *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*, 2nd ed. Cambridge, U.K.: Cambridge Univ. Press, 2019.
- [10] G. Salton and C. Buckley, “Term-weighting approaches in automatic text retrieval,” *Information Processing & Management*, vol. 24, no. 5, pp. 513–523, 1988.
- [11] C. Ware, *Information Visualization: Perception for Design*, 3rd ed. San Francisco, CA, USA: Morgan Kaufmann, 2013.
- [12] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent Dirichlet allocation,” *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.
- [13] D. A. Keim, J. Kohlhammer, G. Ellis, and F. Mansmann, Eds., *Mastering the Information Age: Solving Problems with Visual Analytics*. Goslar, Germany: Eurographics Association, 2010.