

Figura 2: Dez palavras com maiores e menores frequência.

Fonte: autor

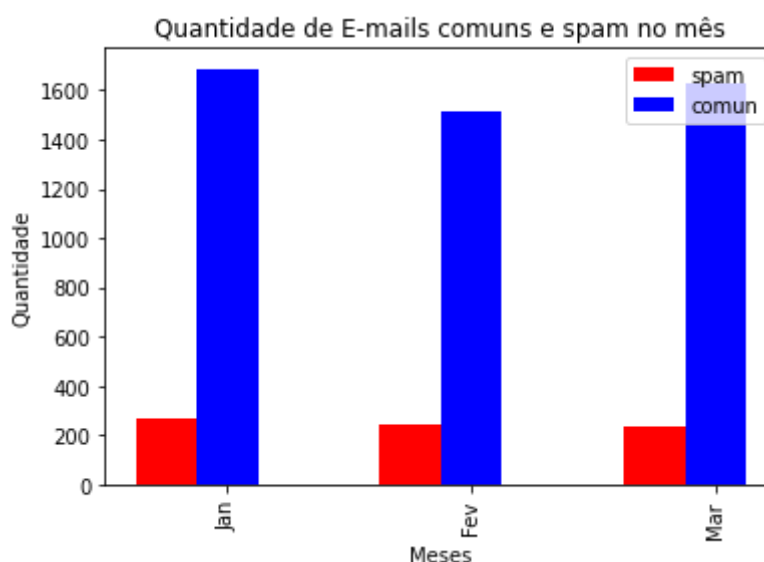


Figura 3: Quantidade de *sms* comuns e spam mensais.

Fonte: autor

Com esses gráficos é possível observar as palavras com maior repetição, e através da Figura 3 é possível observar que a quantidade de mensagens spam é menor do que as comuns.

Também foram extraídas as estatísticas de cada mês da contagem de palavras da coluna *Word_Count* apresentada na Tabela 1.

Tabela 1: Estatísticas

Mês	Média	Mediana	Desvio	Variância
Janeiro	16	13	12,52	157
Fevereiro	16	13	11,00	121
Março	16	12	11,53	133

Fonte: autor

Após isso se deu início a implementação dos modelos de classificação e foram escolhidos três modelos que são: Regressão Logística(LR), Floresta Aleatória(FA) e uma Rede Neural(RN).

A Regressão Logística foi escolhida por ser boa em classificação binária sendo um algoritmo baseado no conceito de probabilidade ela utiliza a função sigmóide para classificação. A Floresta Aleatória foi utilizada por ser um conjunto de árvores de decisão que tendem a ter respostas melhores que apenas uma árvore de decisão. Já a Rede Neural foi escolhida por ter um *dataset* ligeiramente desbalanceado com apenas 13% de sms de spam.

Antes da implementação foram realizadas algumas operações no dataset, foi realizada a *tokenização* da coluna que continha a mensagem inteira para que os modelos pudessem interpretar. Também foram realizadas normalizações de colunas para que não houvesse valores discrepantes na matriz ou que o modelo não conseguisse interpretar.

Os modelos implementados foram avaliados através das métricas observadas na Tabela 2.

Tabela 2: Métricas utilizadas nos modelos.

Modelo	Acurácia (%)	Precisão (%)	Recall (%)	F1 (%)	ROC (%)
RL	98.38	98.42	94.37	96.27	98.74
FA	93.72	96.66	75.86	82.34	97.43
RN	98.83	97.14	93.79	95.43	96.69

Fonte: autor

3. CONCLUSÃO

Todos os modelos implementados tiveram respostas satisfatórias, o modelo de Floresta Aleatória apresentou o pior desempenho de acordo com as métricas utilizadas. Já os modelos de regressão logística e rede neural apresentaram resultados bons, como essa aplicação soft que caso não classifique corretamente não causará danos físicos aos usuários. O modelo de Regressão Logística tem sua implementação mais simples oque facilita validação do projeto.

Para sugestão de melhoria do projeto poderia ser explorado melhor as colunas e diminuí-las quando possível. Também explorar as configurações nos

modelos mudando ativadores, quantidades e também a implementação de um comitê formado por alguns modelos. Por último examinar mais adequadamente as métricas que se encaixam de forma mais correta e quais seriam as que mediriam qual seria o melhor produto final.

4. REFERÊNCIAS

JAIN, Tarun. SMS Spam Classification Using Machine Learning Techniques. **IEEE Xplore**, Noida, v. 213, n. 4, p. 213-214, dez. 2022.

KUNUMI, Blog. **Métricas de Avaliação em Machine Learning: Classificação**. Disponível em: <https://medium.com/kunumi/m%C3%A9tricas-de-avalia%C3%A7%C3%A3o-em-machine-learning-classifica%C3%A7%C3%A3o-49340dcdb198>. Acesso em: 27 dez. 2023.

POVO, Correio do. **Mais de 155 milhões de brasileiros possuem celular para uso pessoal, aponta IBGE**. 2022. Disponível em: <https://www.correiodopovo.com.br/jornalcomtecnologia/mais-de-155-milh%C3%B5es-de-brasileiros-possuem-celular-para-uso-pessoal-aponta-ibge-1.891007>. Acesso em: 27 jan. 2023.