

A Web of Hate: Tackling Hateful Speech in Online Social Spaces

Processamento de Linguagem Natural





Integrantes

Caique de Camargo

11091312

caique.camargo@aluno.ufabc.edu.br

Jean Augusto

21025614

jean.a@aluno.ufabc.edu.br

July Anne Pinheiro

11094013

july.pinheiro@aluno.ufabc.edu.br

Marcela Yamashita

21083913

marcela.a@aluno.ufabc.edu.br



Introdução

Artigo: A Web of Hate: Tackling Hateful Speech in Online Social Spaces

Autores: H. M. Saleem, K. P. Dillon, S. Benesch, D. Ruths

46 citações. Apresentado no First Workshop on Text Analytics for Cybersecurity and Online Safety at LREC (2016)

Proposta do artigo:

- Comparar abordagens diferentes para definição de discurso de ódio
- Propor uma nova abordagem para detecção de discurso de ódio

Community based x Keyword based



Community based x Keyword based

Community based: utilização da linguagem que emerge de grupos que se auto-organizam (como Reddit). Grupos se formam através de práticas linguísticas comuns. **“the group is defined by speech and the speech comes to define the group”**

Keyword based: classificação baseada em palavras-chave



Proposta do grupo: keyword-based

Inabilidade de encontrar datasets de grupos fechados

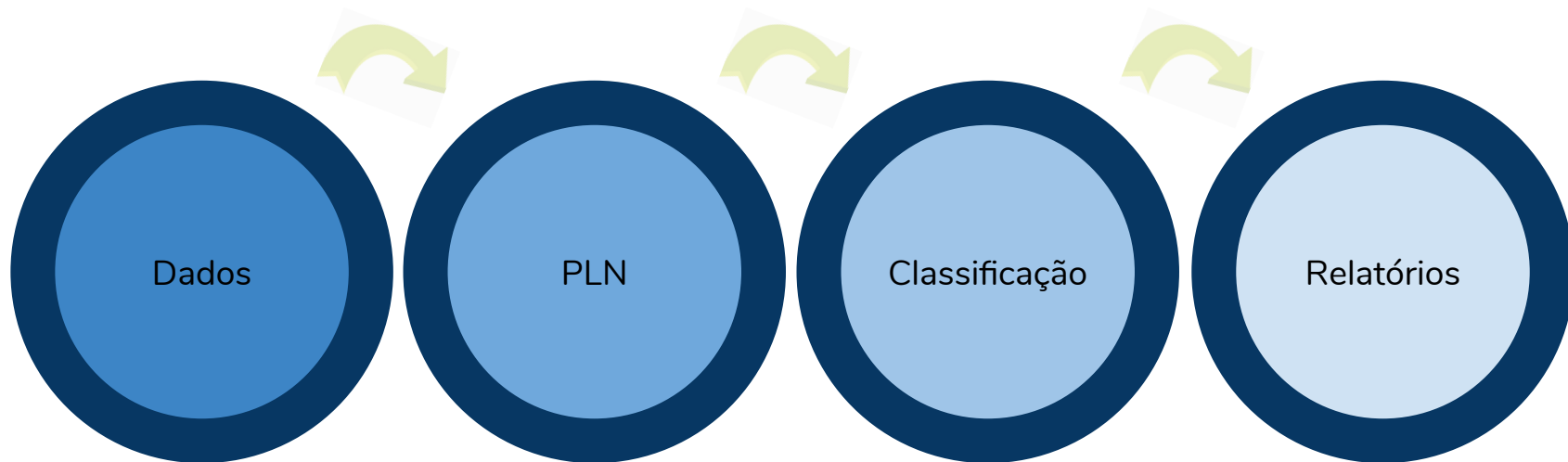
Testar a performance de diferentes classificadores (Naive Bayes, SVM e LR) no dataset do Twitter

Twitter: ~8 milhões de usuários no Brasil (julho/2019)

Grande influência no cenário social e político.



Organização





Base de dados

Original

- Stuck_In_The_Matrix Reddit Dataset (250GB compress.)
- Voat
- Web forums

Adaptação

- “Hate Speech Identification” (Twitter Dataset) - Data World (cerca de 15mil tweets)



Algoritmos

Implementação e comparação de resultados de três algoritmos de classificação:

- Naive Bayes
- Support Vector Machine (SVM)
- Logistic Regression (LR)

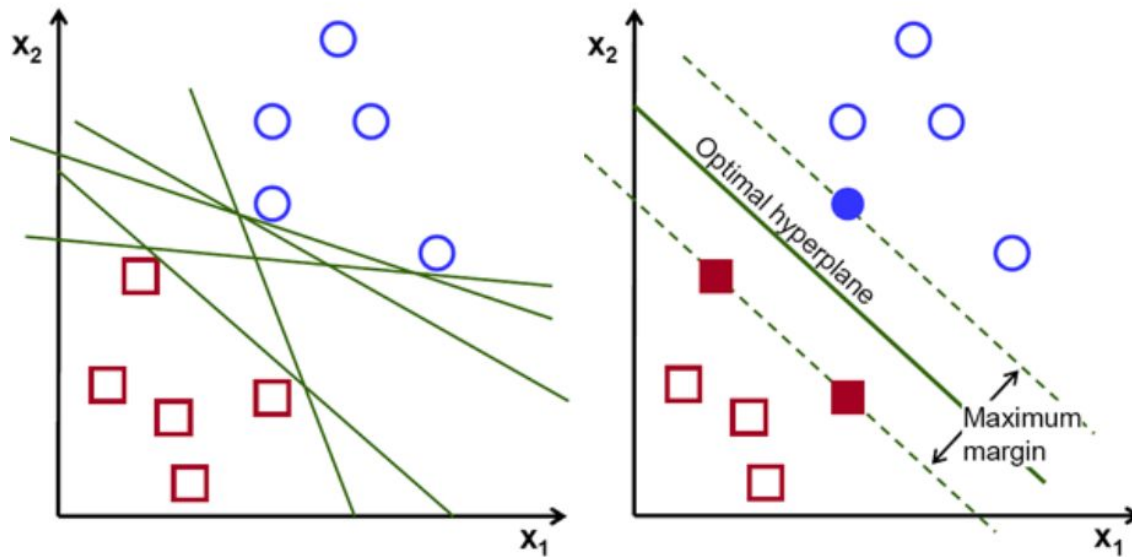


Naive Bayes

$$p(C_k \mid \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} \mid C_k)}{p(\mathbf{x})}$$



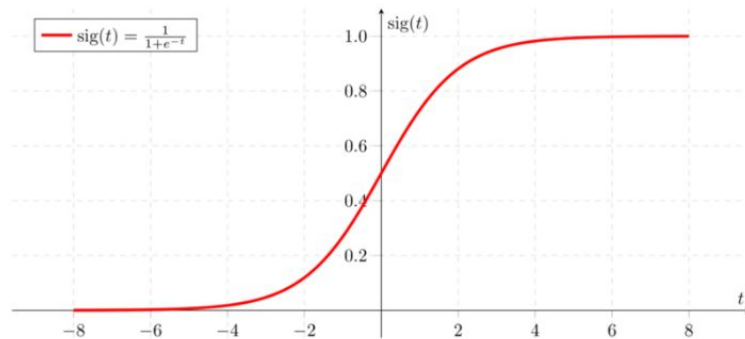
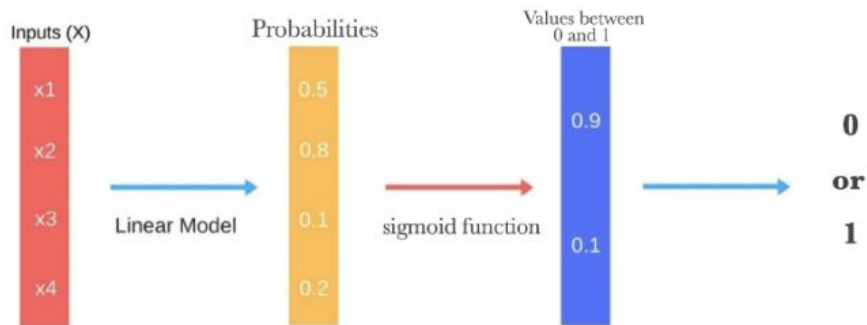
SVM

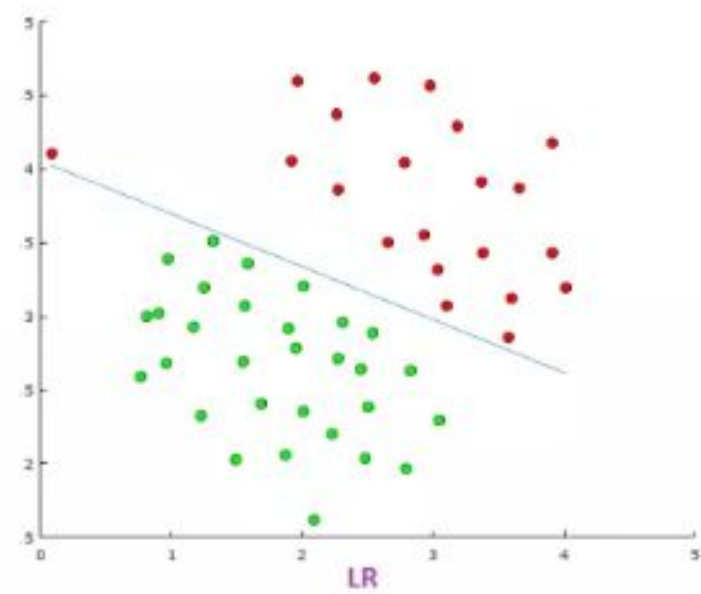


Obs: Utilização do kernel linear do SVM.



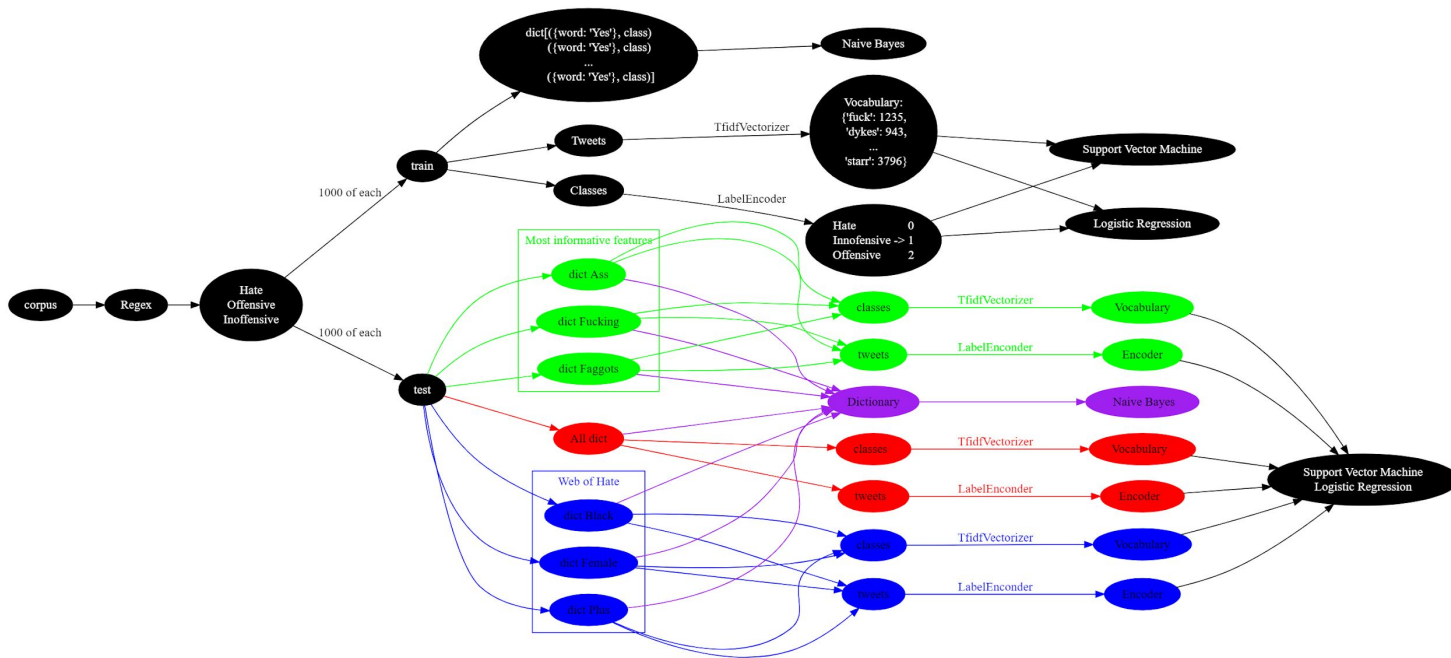
Logistic Regression





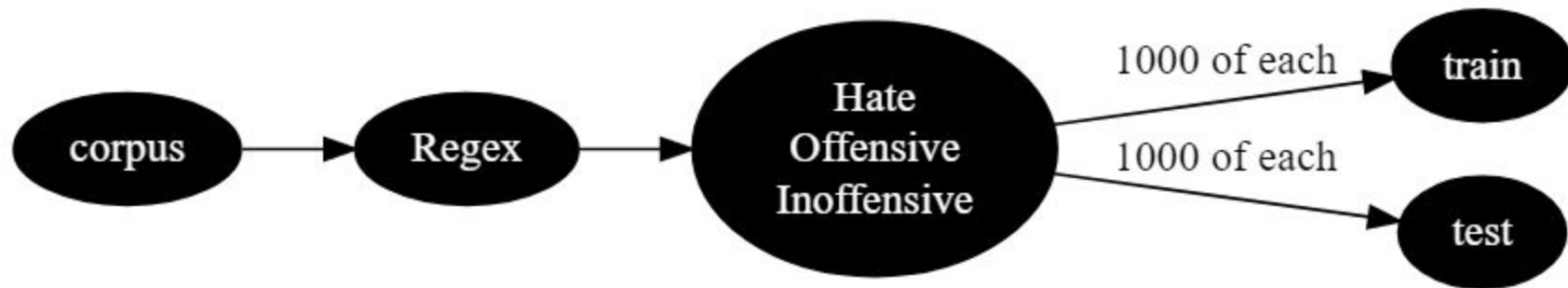


Algoritmo



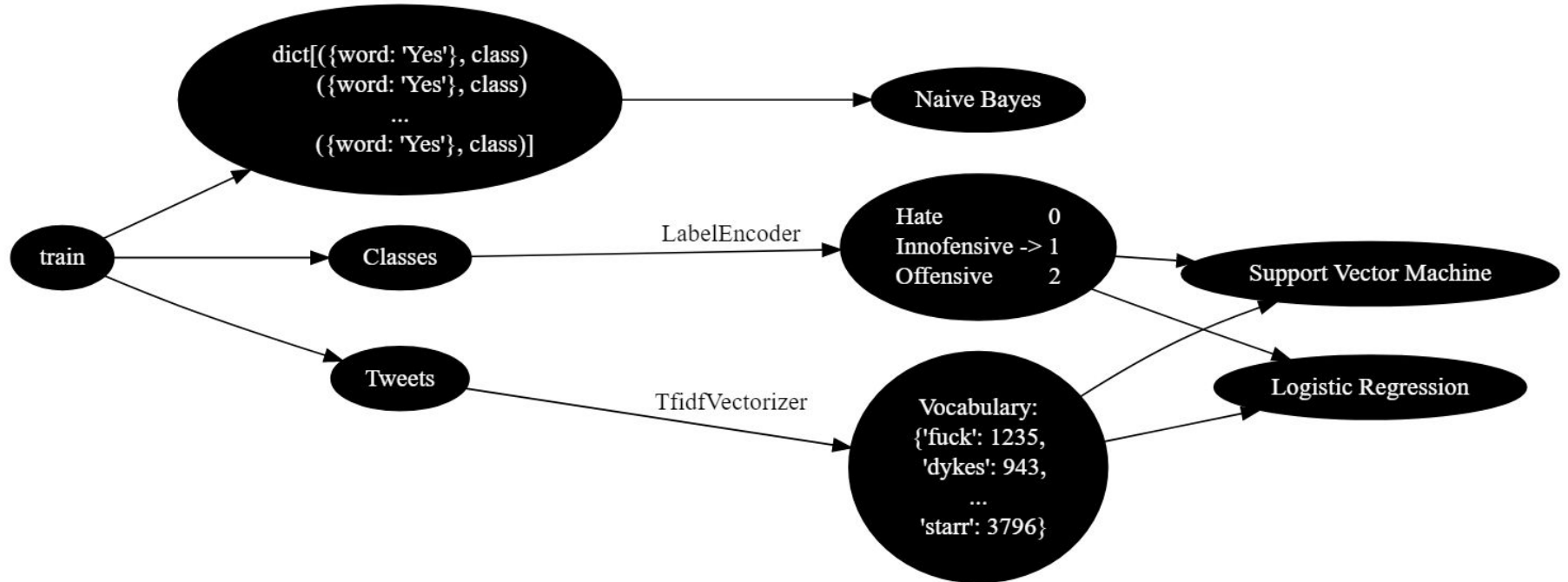


Algoritmo

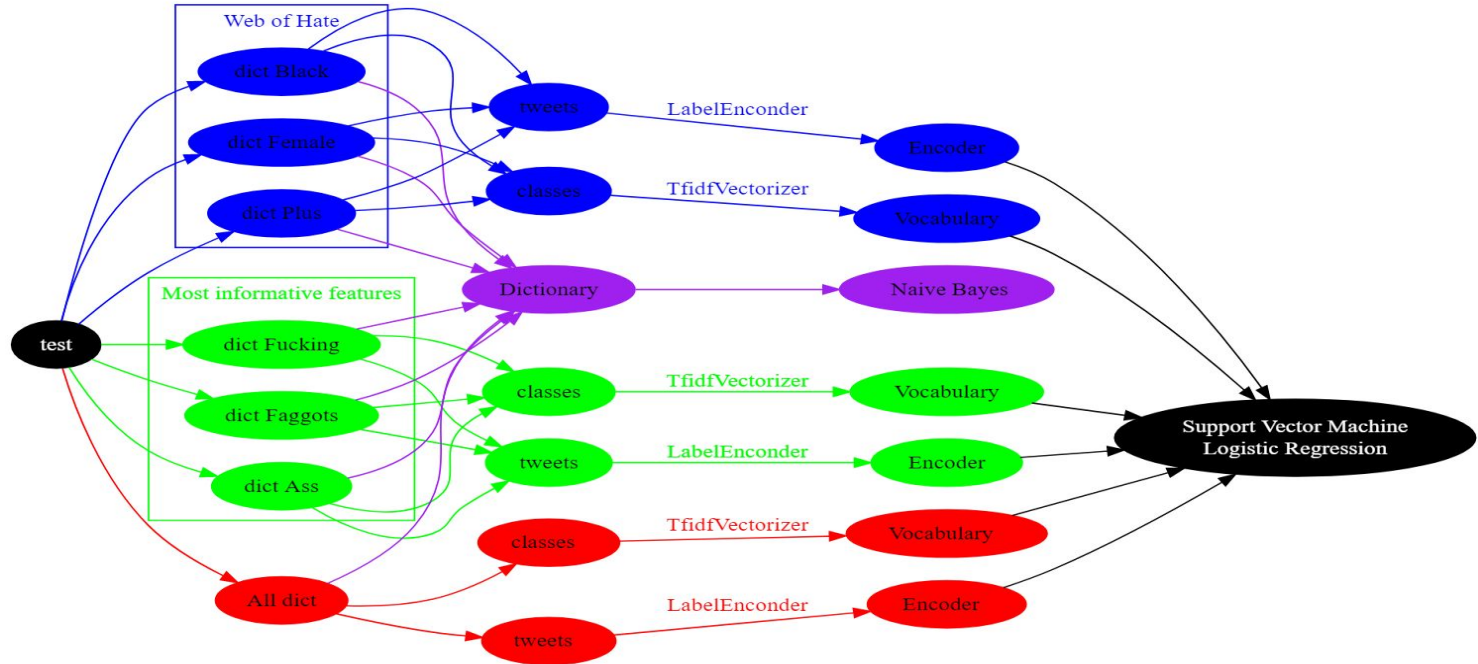




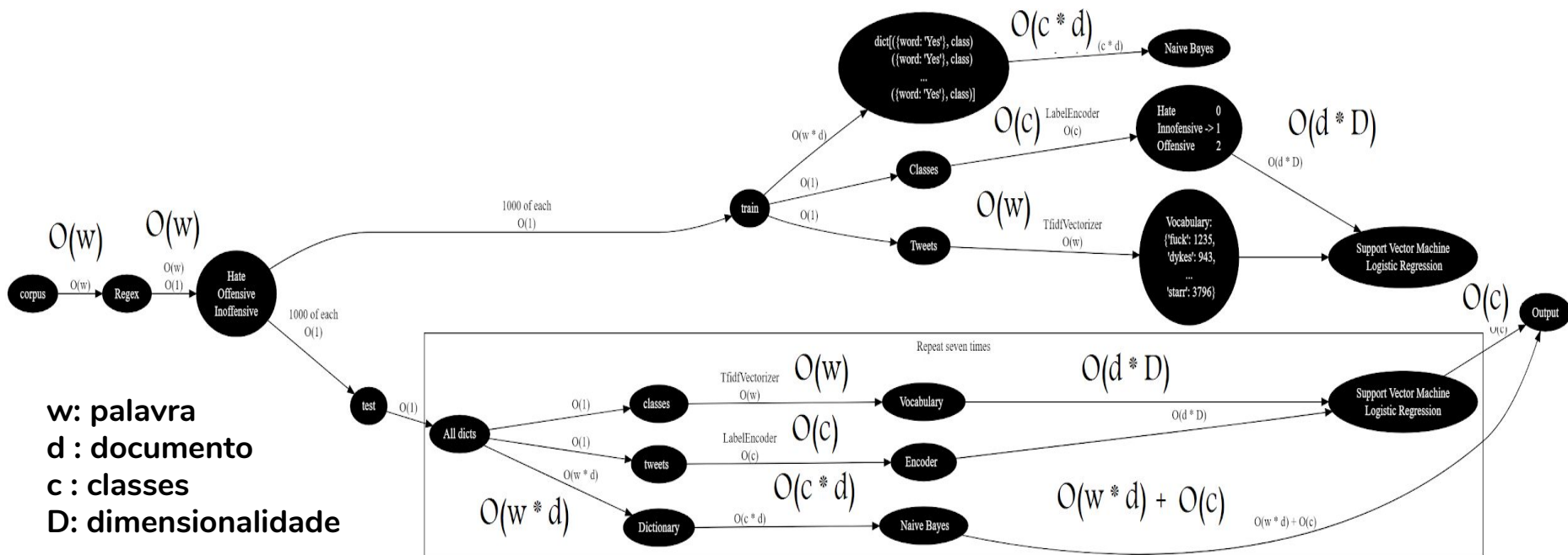
Algoritmo



Algoritmo



Complexidade



$$O(11w + 14c + 7d + 14(w * d) + 8(c * d) + 9(d * D))$$



Resultados

Medidas:

- Accuracy
- Precision
- Recall
- F_measure

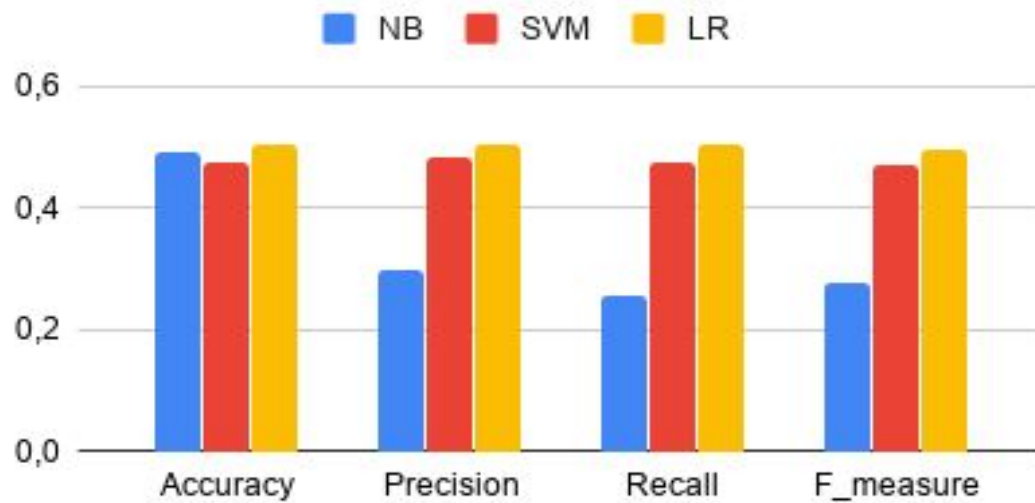
Termos utilizados:

- Black
- Plus
- Female



Resultados

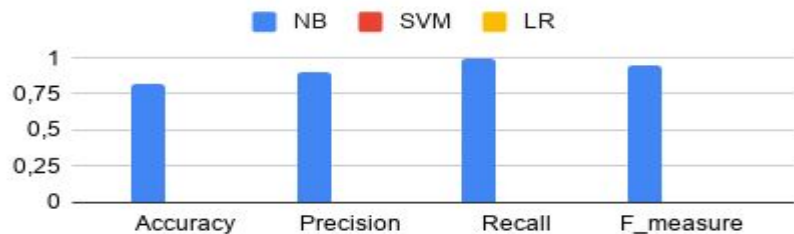
TOTAL



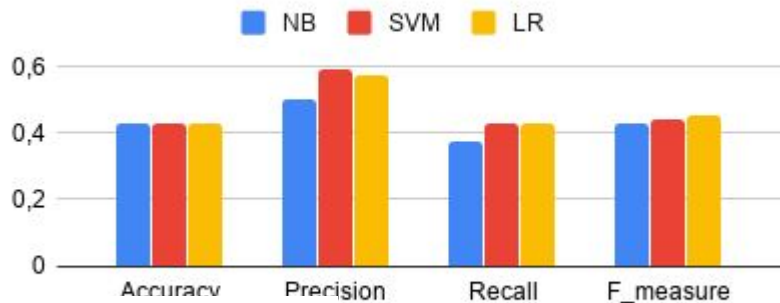


Resultados

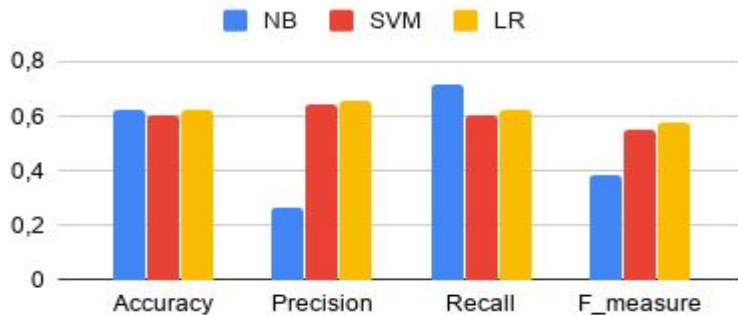
PLUS



FEMALE



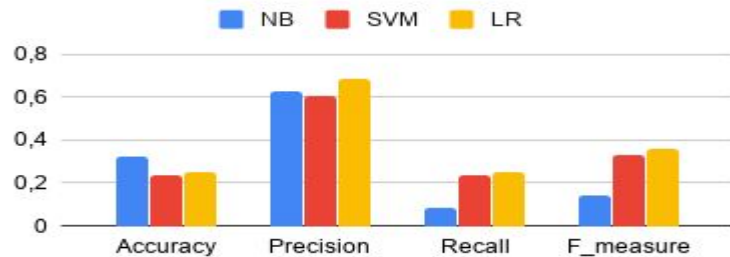
BLACK



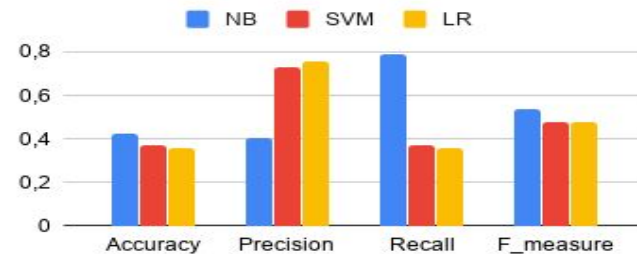


Resultados

FUCKING



FAGGOTS



ASS





Obrigado!