

Profa Daniela Souza Moreira da Silva



Mineração de Textos

- É baseada na busca por padrões em textos digitais em linguagem natural
- Texto digital visto como uma das principais fontes de informação.
- Dados textuais estão, de modo geral, não estruturados.
- Mineração de Dados x Mineração de Textos
- Etapas da Mineração de Textos



- Também chamada de PLN
- Coleta de dados textuais
- Normalização textual

Mineração de Textos

O que linguagem natural?

É aquela que foi desenvolvida e evoluída por seres humanos a partir do <u>uso natural do dom de se comunicar</u>, isto é, não foi criada artificialmente como uma linguagem de programação (SARKAR, 2016).

- Etapas mais comuns da Normalização Textual:
 - Limpeza textual (Retirada de marcadores textuais, tags, comentários ..)
 - Uniformização maiúscula/minúsculas
 - Remoção de símbolos e pontuação
 - Tokenização do texto (unidade mínima, chamada de token, normalmente são palavras)
 - Expansão de contrações (tromba d'agua -> tromba de agua)
 - Remoção de stopwords (artigos, pronomes, advérbios, conjunções)
 - Correção da grafia
 - Lematização/ Radicalização (stemming)
 - Lematização: busca a forma canônica de um conjunto de palavras (Será classificada como verbo ou substantivo.
 - Radicalização: busca o radical em comum de um conjunto de palavras.

Lematização

- Estudo
- Estudas
- Estudamos

Estudo (N)

Ou Estudar (V)

Radicalização

- Estudo
- Estudas
- Estudamos
 Estud (radical)

Sobre WEKA



WEKA (Waikato Environment for Knowledge Analysis) é um software para mineração de dados, desenvolvido pela Universidade de Waikato, Nova Zelândia.

Weka é uma coleção de algoritmos de aprendizado de máquina para tarefas de mineração de dados. Ele contém ferramentas para preparação de dados, classificação, regressão, agrupamento, mineração de regras de associação e visualização

Link para download: https://waikato.github.io/weka-wiki/downloading_weka/

Sobre WEKA



- WEKA implementado na linguagem de programação JAVA, com a característica principal de sua portabilidade funcionando em diferentes sistemas operacionais.
- É um software de código aberto sob a licença GPL (General Public License).
- Algoritmos (Métodos) implementados no WEKA:
 - Métodos de classificação (SVM, Árvore de Decisão, Naive Bayes, Regressão Lógica...)
 - Métodos de predição numérica (Regressão Linear, Percepton multicamadas...)
 - Métodos de agrupamento (SimpleKMeans, Clope, Cobweb...)
 - Métodos de associação (Apriori, FPGrowth...)

Sobre WEKA

Possui vários

 algoritmos com a
 finalidade de realizar o
 pré-processamento
 dos dados.

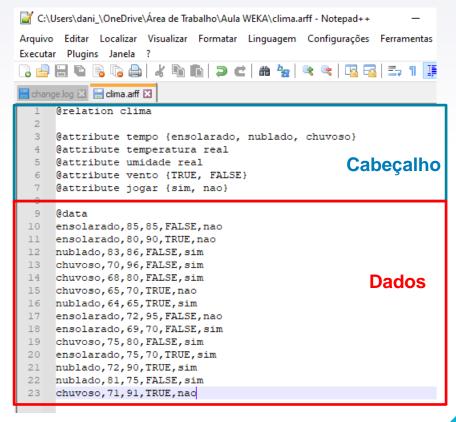
 Os dados precisam estar organizados para realizar a mineração. Os arquivos podem estar em um formato específico, planilha ou banco de dados.

O WEKA possui um formato para a organização dos dados ARFF(Attribute-Relation File Format).



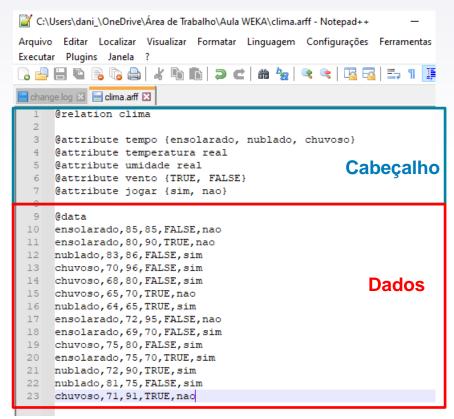
Arquivo ARFF (Attribute-Relation File Format)

- Arquivo dividido entre 2 seções: Cabeçalho e Dados.
- Cabeçalho:
 - Primeira linha deve apresentar a relação do arquivo "@relation"
 - Se o nome dessa relação contiver espaços deve estar entre aspas
 - Os atributos devem ser únicos, terão um nome e um tipo precedidos da palavra reservada @attribute.
 - O nome do atributo deve começar com letra e, se contiver espaços, deve estar entre aspas
 - Tipos de Dados:
 - Números (reais ou inteiros): Numeric ou Real
 - ► Texto "livre": String
 - Atributos categóricos (lista de valores)
 - Data: Date [< date-format >]



Arquivo ARFF (Attribute-Relation File Format)

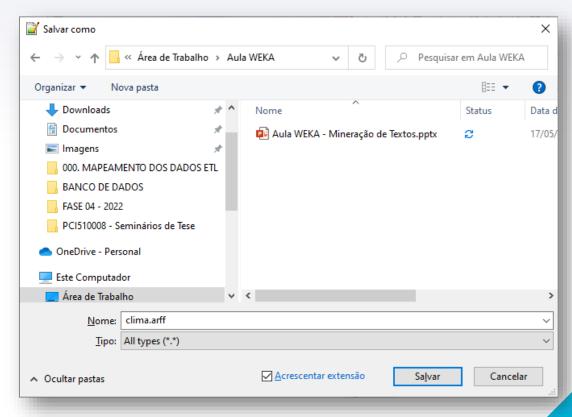
- Arquivo dividido entre 2 seções: Cabeçalho e Dados.
- Dados:
 - A base de dados vem logo após o @data
 - É a lista de todas as instancias com os valores dos atributos separados por vírgulas. Cada instância (registro) é representada em uma única linha.
 - Os atributos devem aparecer na ordem em que são declarados no cabeçalho.
 - Por default, o WEKA trata o último atributo especificado no cabeçalho como o atributo classe e os demais como atributos preditivos.
 - No caso de dados textuais, na área de dados, os valores devem estar entre aspas.



Arquivo ARFF (Attribute-Relation File Format)

Para gerar um arquivo arff utilizar o notepad++, e salvar o arquivo com o nome.arff, e o tipo (all types (*.*)

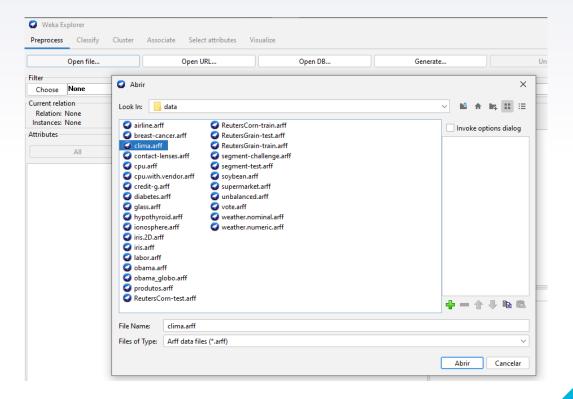
Outra opção é o notepad, salvando o arquivo com o nome.arff e o tipo "Todos os arquivos (*.*)



A guia Explorer será utilizada para as atividades de mineração (dados ou textos).

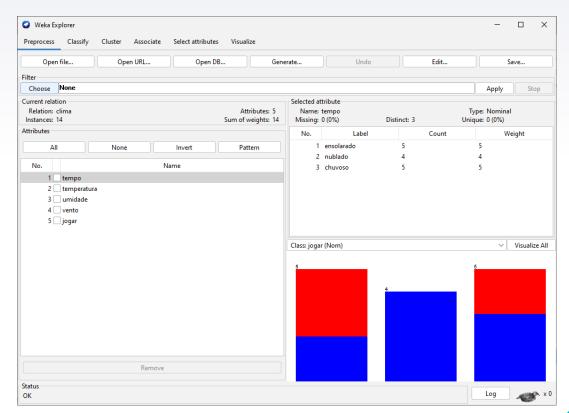


 Após clicar no Explorer, clicar na opção "Open File" e selecione o arquivo desejado (na extensão arff) e o programa carregará os dados na interface.

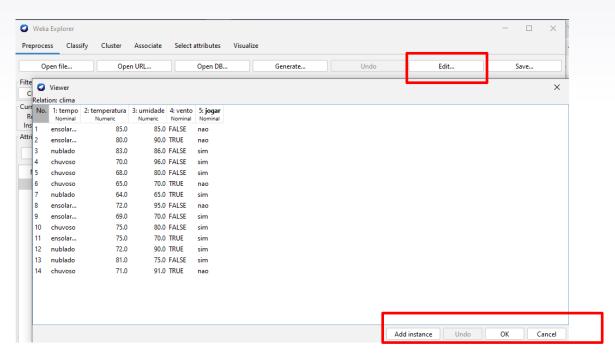


Na aba "Preprocess" é onde será realizado o préprocessamento dos dados por meio dos filtros(Filter->Choose).

Em "Current relation" podese ver o nome do arquivo(Relation), a quantidade de atributos(Attributes) e de instâncias(Instances)..

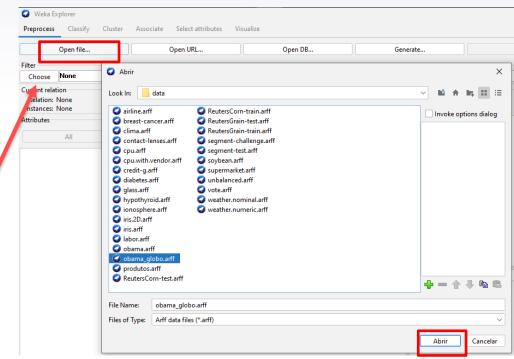


Na opção "Edit" pode ser editado o arquivo de dados.

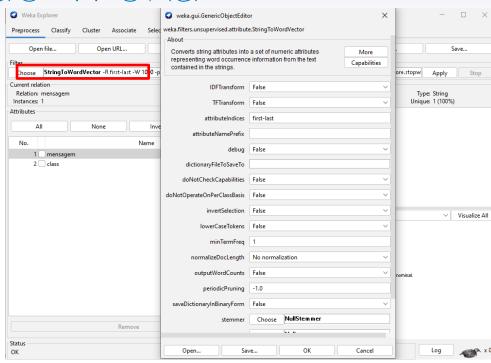


- Análise de textos utilizando WEKA: discurso do ex-presidente dos USA Barack Obama, que ocorreu no Theatro Municipal da cidade do Rio de Janeiro no dia 20/03/2011 (discurso disponível no G1).
- Objetivo: Minerar o texto para verificar o número de ocorrências das palavras para identificar qual foi o enfoque do discurso.
- Metodologia: Normalizar o texto (remoção de stopwords e uniformização maiúscula/minúscula) e utilizar um algoritmo que fracione o texto em palavras.
 - Algoritmo: StringToWordVector. Cria um vetor de palavras separando o texto.
 - Ajustar os parâmetros.

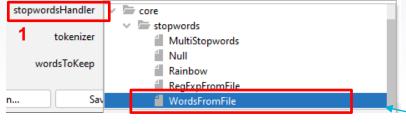
- № 1º passo) Criar o arquivo arff com o discurso;
- 2º passo) Carregar o arquivo no WEKA, clicando em Open File, selecionando o diretório onde o arquivo foi salvo.
- ▶ 3º passo) Após abrir o arquivo, acessar a opção "Filter/Choose", e escolher a opção "filters->unsupervised->attribute->StringToWordVector".



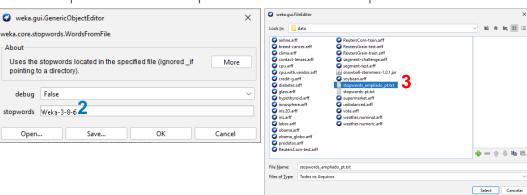
- 4º passo) Para ajustar os parâmetros do algoritmo, clicar na guia choose do filter e abrirá uma janela para configuração:
 - LowerCaseTokens: Transforma todas as palavras para minúscula
 - MinTermFreq: Foi estabelecido o valor mínimo de 5 palavras por ocorrência, a depender do texto a média de palavras por ocorrência pode ser aumentada.
 - OutputWordCounts: Conta a quantidade de ocorrências de cada palavra.
 - WordsToKeep: Estabelece o número máximo de palavras a serem analisadas.

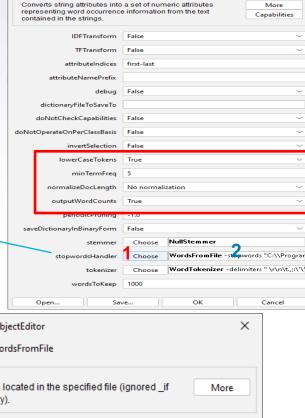


5º passo) Para utilizar o arquivo das stopwords, clicar na opção stopwordsHandler/Choose/WordsFromFile (ler palavras do arquivo)



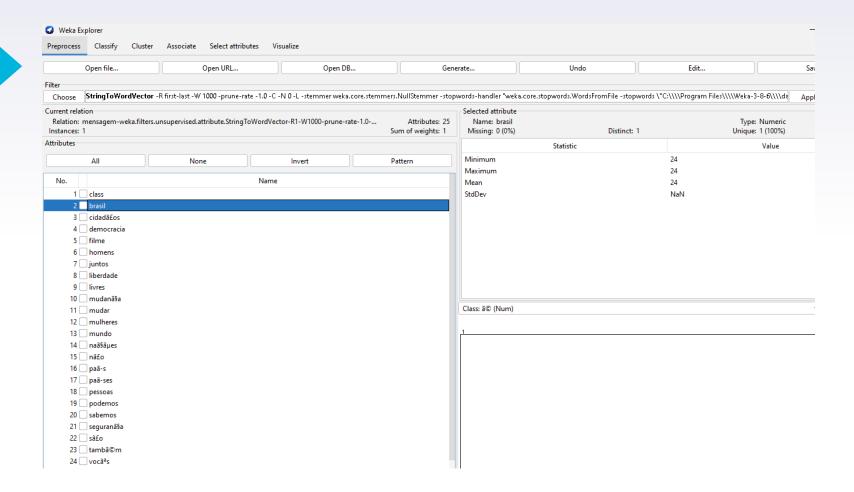
6º passo) Clicar na área wordFromFile e indicar o local onde está o arquivo com as stopwords. Coloca-lo na pasta data do WEKA.

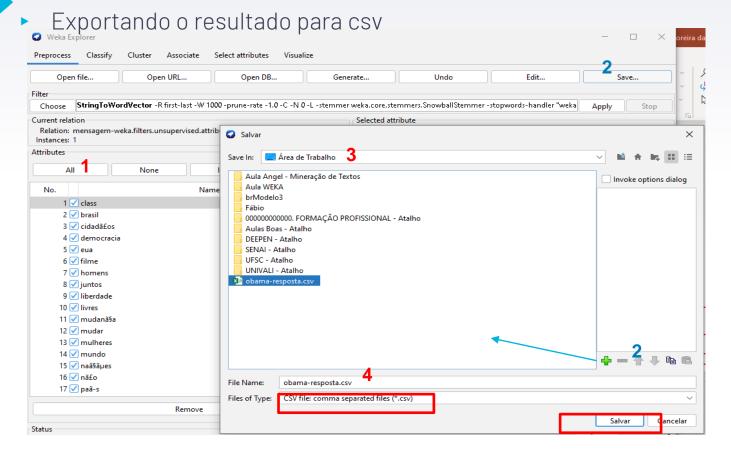




weka.gui.GenericObjectEditor

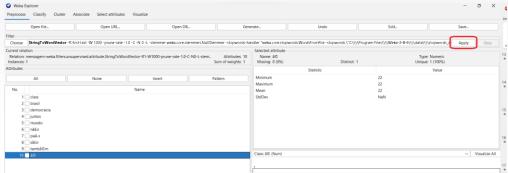
weka.filters.unsupervised.attribute.StringToWordVector



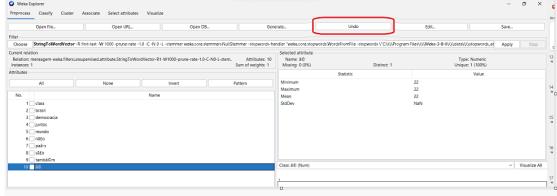


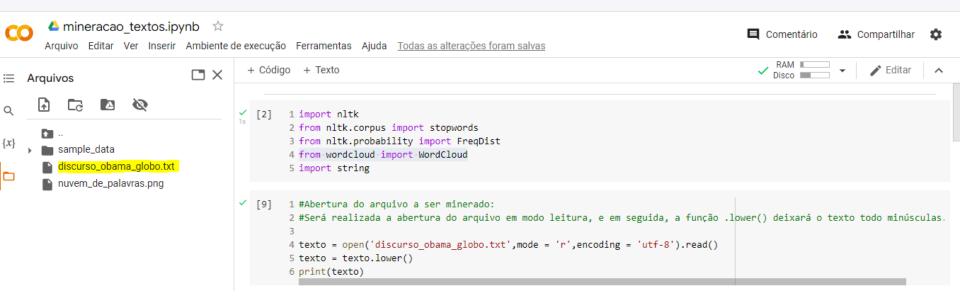
Executando os testes

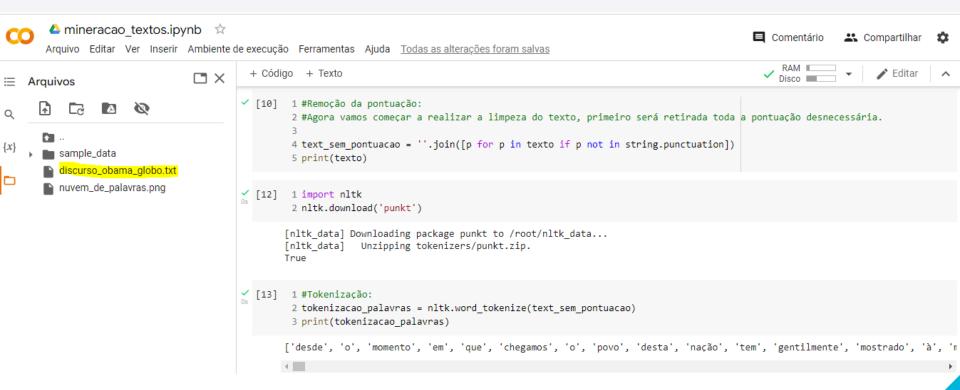
Após realizar a parametrização do filtro, clicar em "Apply" para rodar o algoritmo.



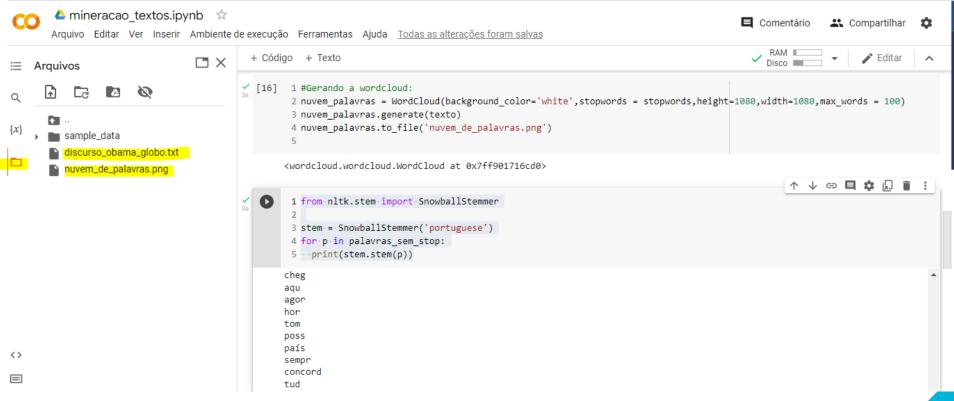
Para realizar novos testes utilizar a opção "undo". Ela remove as configurações













Referências

- https://www.cs.waikato.ac.nz/ml/weka/index.html
- https://www.ranks.nl/stopwords/portuguese
- https://g1.globo.com/obama-no-brasil/noticia/2011/03/leiaintegra-do-discurso-de-barack-obama-no-theatromunicipal.html
- https://www.ufsm.br/pet/sistemas-deinformacao/2021/07/12/introducao-a-mineracao-de-textoscom-python/
- https://github.com/Prof-Rodrigo-Silva/Text-Mining
- https://www.linkedin.com/pulse/minera%C3%A7%C3%A3ode-texto-usando-stringtowordvector-do-weka-rosa-da-silva/



Referências

https://code.google.com/archive/p/ptstemmer/



Obrigada! Perguntas?

Dani.smoreira@gmail.com







Free templates for all your presentation needs



For PowerPoint and Google Slides



100% free for personal or commercial use

Ready to use, professional and customizable Blow your audience away with attractive visuals