

Classificação de texto

com Hadoop + Spark Mlib

Erik Perillo, RA135582

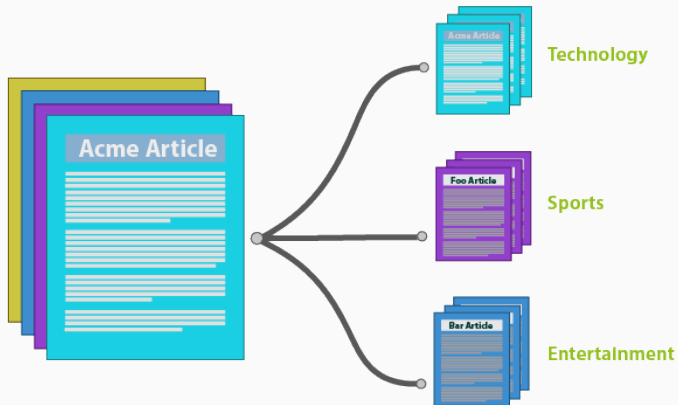
Universidade Estadual de Campinas

Objetivos

Um classificador de textos binário:

- Com uma das melhores técnicas atuais
- Com grande volume de dados

Classificação de Texto



$$p(C_k|x_1,\dots,x_n) = \frac{1}{Z}p(C_k)\prod_{i=1}^np(x_i|C_k)$$

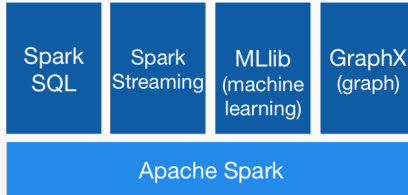
O projeto



- Computação em *batches* de alto desempenho
- Análise de dados avançada
- Processamento de *stream* em tempo real
- Interfaces para Java, Scala, Python



Spark: componentes



- Diversas técnicas de *Machine Learning*
- Desenvolvimento rápido e escalável



Juntando tudo

Juntando tudo

- Tudo foi feito em Python :)
- Foi implementada visualização pela linha da comando

- Arquivos de uma database entram (20_newsgroup)
- Pré-processamento do texto é feito (*stemming, stopwords...*)
- Conversão de texto para vetor de *features* (if-idf)
- Textos positivos/negativos são convertidos para .csv
- Arquivo .csv é convertido em formato para *Spark*
- Treinamento é feito com *Naive Bayes*

Resultados
