

EXTRAINDO TEMA DE DOCUMENTOS DIPLOMÁTICOS

Alice Duarte Scarpa

29 de Outubro de 2015

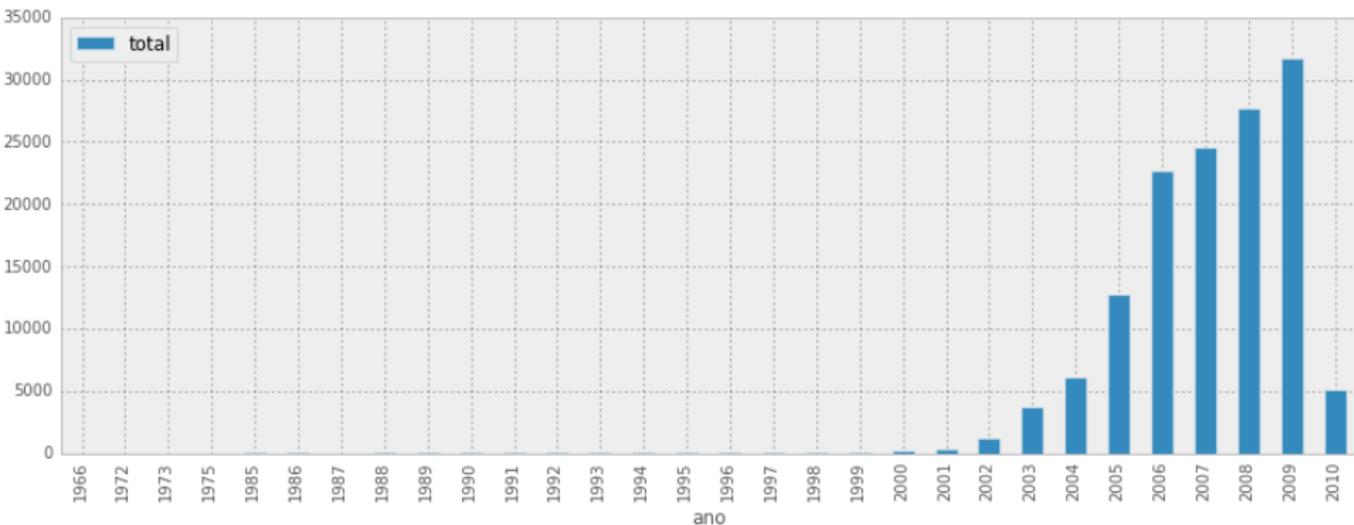
INTRODUÇÃO

Neste trabalho vamos explorar dados do [cablegate](#), que são cabos diplomáticos entre embaixadas dos Estados Unidos que foram divulgados como parte do [WikiLeaks](#).

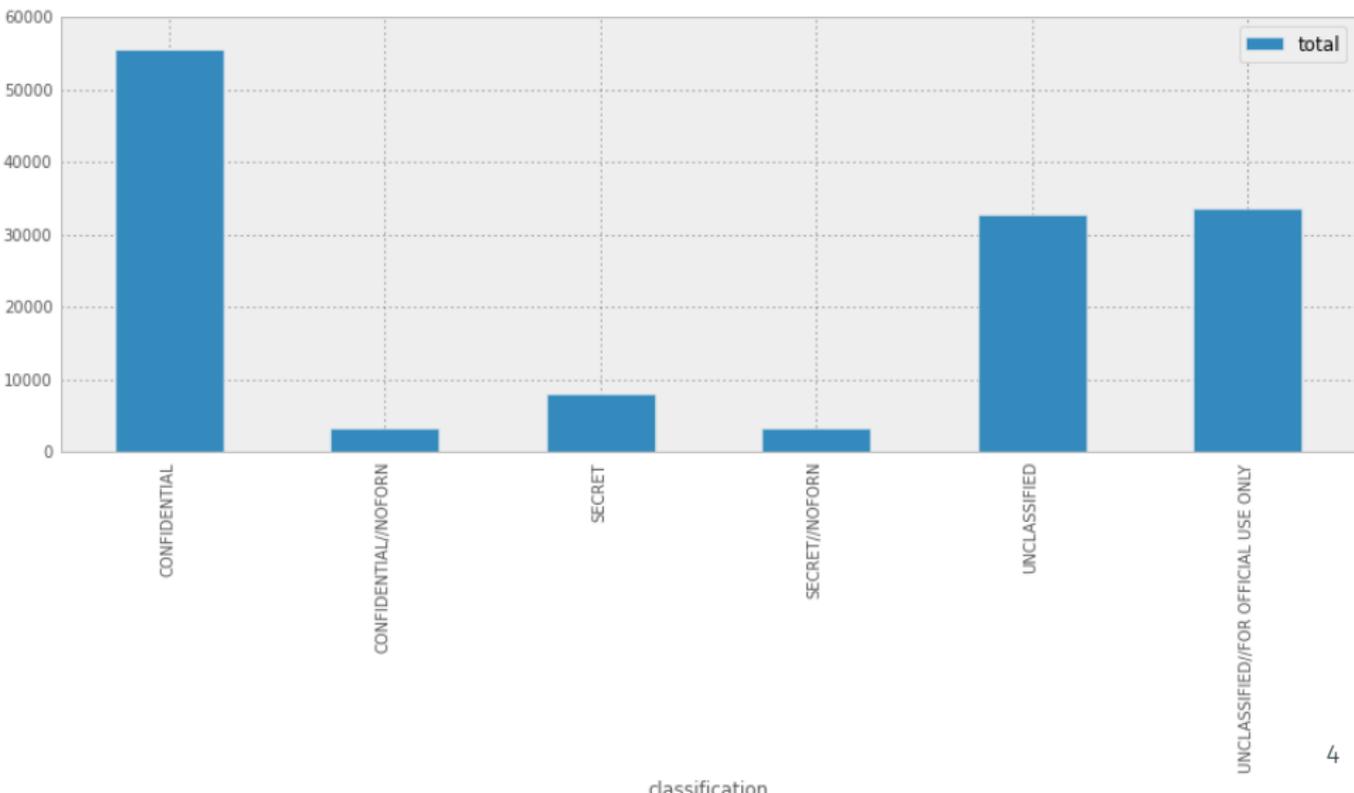
Neste trabalho vamos explorar dados do [cablegate](#), que são cabos diplomáticos entre embaixadas dos Estados Unidos que foram divulgados como parte do [WikiLeaks](#).

São 250 mil documentos, de 1966 a 2010 e diferentes níveis de confidencialidade. A grande maioria dos documentos é de 2003 a 2010.

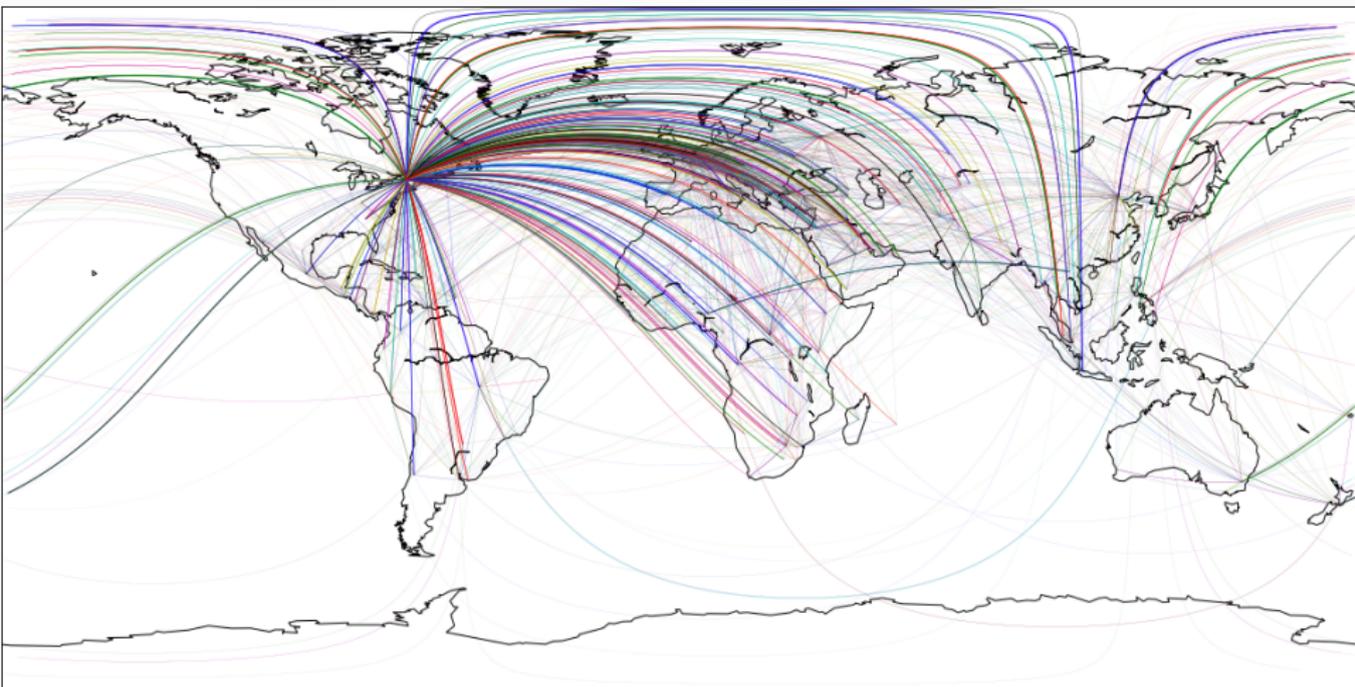
Os DADOS



Os DADOS



Os DADOS



OBJETIVOS

- Criar uma base de dados explorável

OBJETIVOS

- Criar uma base de dados explorável
- Comparar diferentes técnicas existentes de extração de tema

OBJETIVOS

- Criar uma base de dados explorável
- Comparar diferentes técnicas existentes de extração de tema
- Explorar características dos documentos para adaptar alguma técnica existente

TÉCNICAS

Bag of Words Considera o texto apenas como o conjunto de palavras que nele aparece e a quantidade de vezes que cada palavra aparece.

- Bag of Words** Considera o texto apenas como o conjunto de palavras que nele aparece e a quantidade de vezes que cada palavra aparece.
- N-gramas** Como o Bag of Words, mas olha para expressões de n palavras.

- Bag of Words** Considera o texto apenas como o conjunto de palavras que nele aparece e a quantidade de vezes que cada palavra aparece.
- N-gramas** Como o Bag of Words, mas olha para expressões de n palavras.
- Noun Phrases** Procura o sujeito e o objeto de cada oração.

N.E.E. (Named Entities Extraction) Tenta responder: Quem?
Como? Onde?

N.E.E. (Named Entities Extraction) Tenta responder: Quem?
Como? Onde?

Lexical Chaining Procura cadeias semânticas no texto usando
sinônimos, antônimos e repetições

N.E.E. (Named Entities Extraction) Tenta responder: Quem?
Como? Onde?

Lexical Chaining Procura cadeias semânticas no texto usando
sinônimos, antônimos e repetições

Candidatos a tema A partir das Part-Of-Speech Tags, determina
alguns candidatos a tema e dá uma nota para cada
um deles

APLICAÇÕES

- Extrair conhecimento a partir de tendências nos temas dos textos

- Extrair conhecimento a partir de tendências nos temas dos textos
- Facilitar o estudo dessa base de dados por pesquisadores de Ciências Sociais

Experiments with Theme Extraction in Explanatory Texts (1996),
Francois Paradis , Catherine Berrut

Extracting Key Terms From Noisy and Multi-theme Documents (2009),
Maria Grineva, Maxim Grinev, Dmitry Lizorkin

PERGUNTAS?