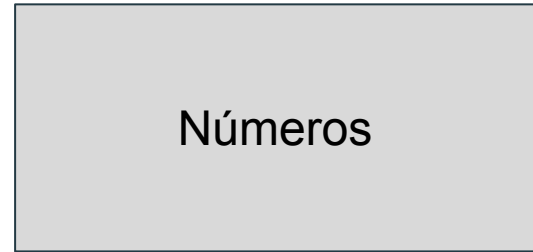




Como visualizar Word Embeddings

Feature extraction

Problema



Soluções

- Bag of Words;
- Bag of n-gram;
- Word2vec; e
- Paragraph2vec

Bag of Words

Frase	Brasil	acima	de	todos	Deus	tudo
Brasil acima de todos	1	1	1	1	0	0
Deus acima de tudo	0	1	1	0	1	1

As palavras não são ordenadas, logo se perde a noção do contexto.

Bag of n-grams

Frase	Brasil acima	acima de	de todos	Deus acima	de tudo
Brasil acima de todos	1	1	1	0	0
Deus acima de tudo	0	1	0	1	1

Quando expandimos os n-grams obtemos matrizes esparsas de alta dimensionalidade. Ainda assim, é capaz apenas de observar pequenos contextos.

Word2Vec

Se tomar vacina [e] virar jacaré não
tenho nada a ver com isso...

Maximizar a seguinte função de
verossimilhança:

$$\frac{1}{T} \sum_{t=k}^{T-k} \log p(w_t | w_{t-k}, \dots, w_{t+k})$$

Ou seja, encontrar o melhor contexto
a partir da palavra central.

Paragraph2vec

Memória distribuída (PV-DM):

Se tomar vacina [e] virar jacaré não
tenho nada a ver com isso...

A partir do contexto predizemos o centro. O modelo matemático é semelhante ao Word2Vec, porém há a adição de um ente que se refere a memória do parágrafo.



ID Parágrafo (Contém informações faltantes no contexto) {Em verde}

Paragraph2vec

Bag of word distribuídas
(*PV-DBOW*):

Se tomar vacina [e] virar jacaré não
tenho nada a ver com isso...

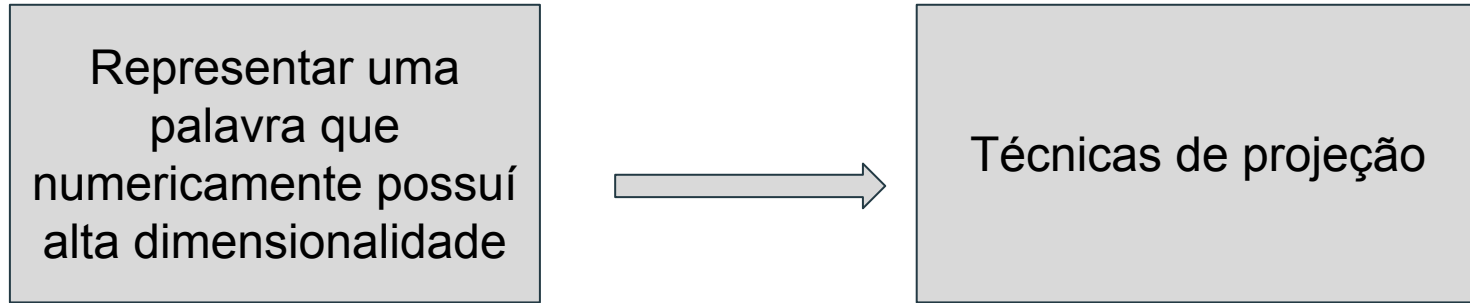
Treinamos um vetor que representa o parágrafo para
predizer todas as palavras no contexto em questão.



ID Parágrafo (Contém
informações faltantes no
contexto) {Em verde}

Representação Gráfica

Problema - Solução



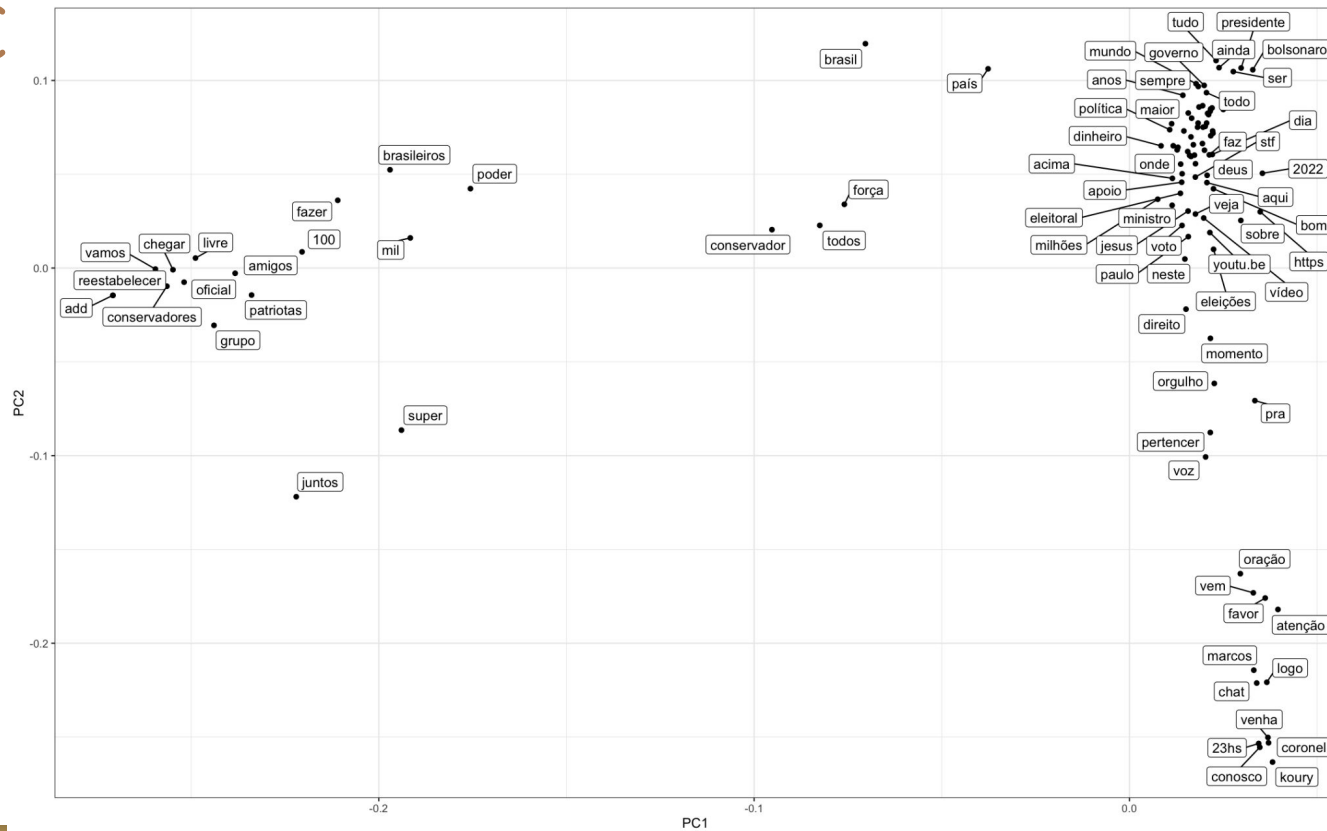
Soluções

- PCA;
- Umap;
- T-sne;

PCA

Objetivo: Encontrar um novo conjunto de variáveis menor que o conjunto original que preserve a maior parte da variabilidade presente nos dados (pautado na covariância).

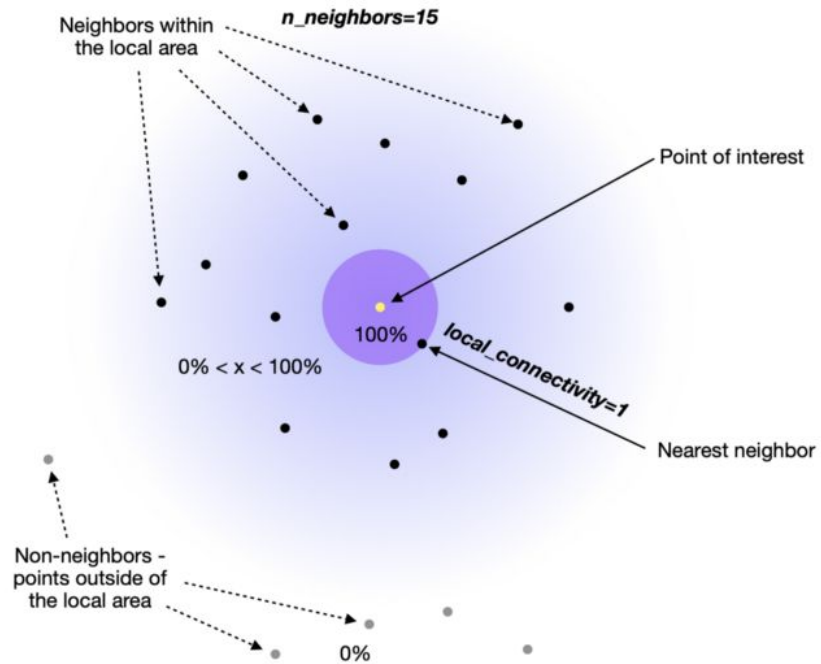
PC



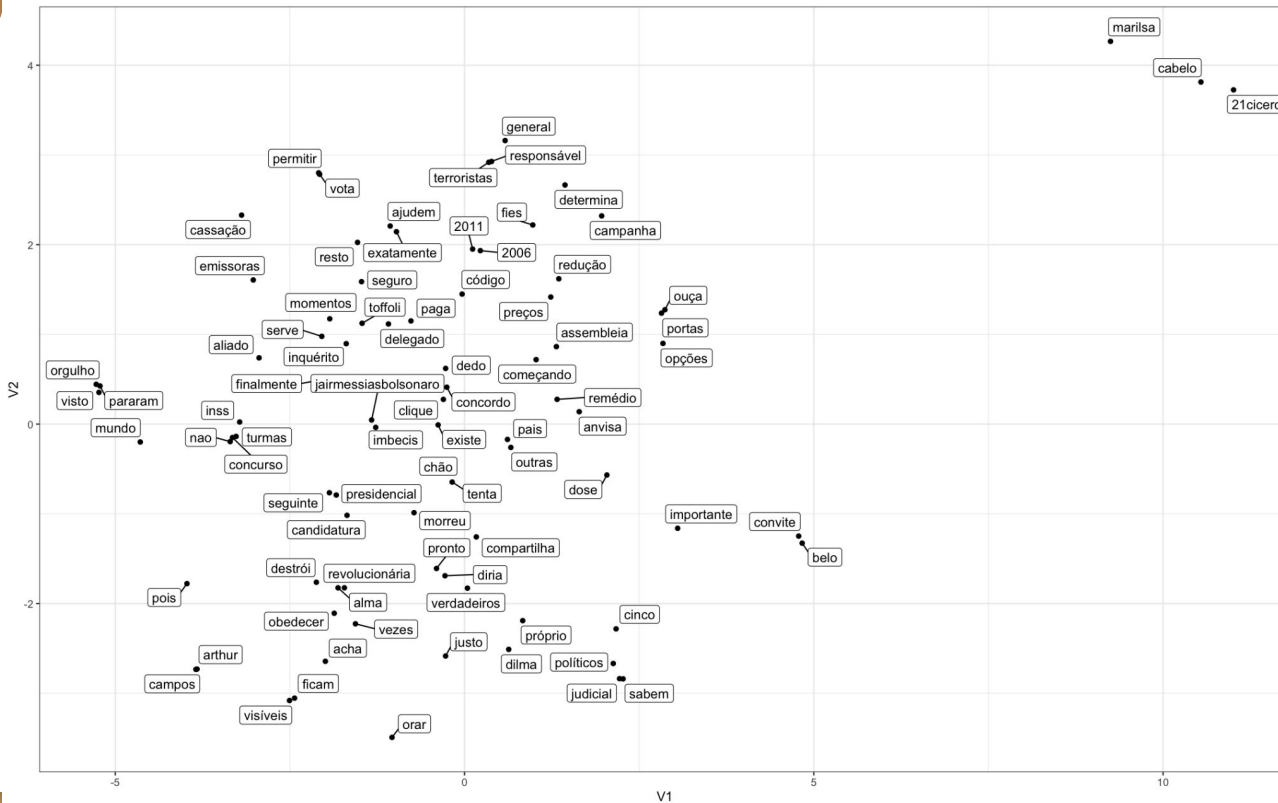
UMAP

UMAP

UMAP



UMAD



TSNE

T-SNE

Para transformarmos distâncias em probabilidades, imaginamos cada ponto como o centro de uma distribuição e usamos a distância entre este ponto e os outros como a probabilidade de ocorrência daquela distância dentro da distribuição escolhida.

b)

TSNE

