

Introdução a Processamento de Linguagem Natural

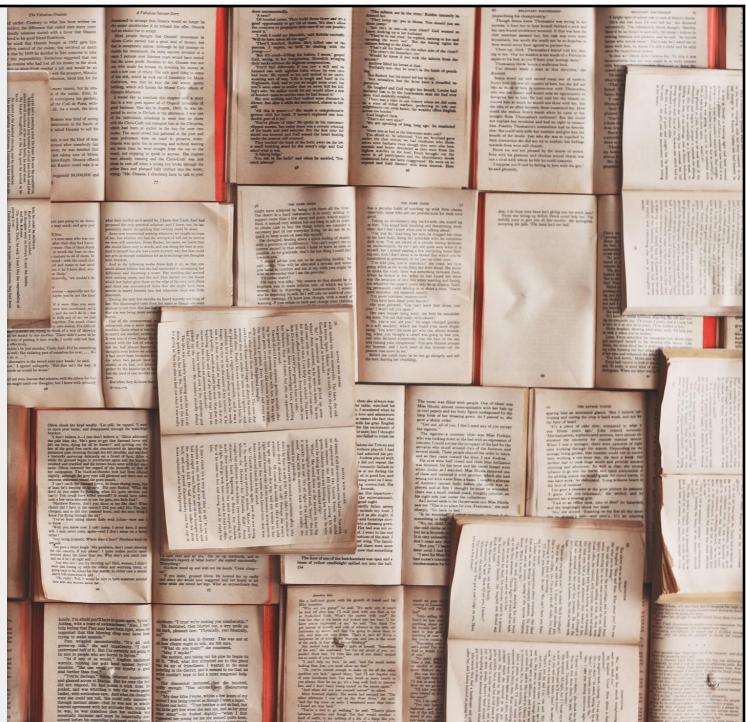
Jonatas Grosman

Luiz Schirmer

William Fernandes



DEPARTAMENTO
DE INFORMÁTICA
PUC-RIO



O que é Linguagem Natural?

Linguagem natural



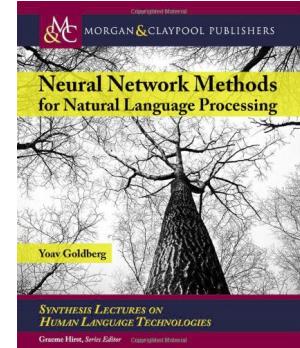
Linguagem artificial



O que é Processamento de Linguagem Natural?

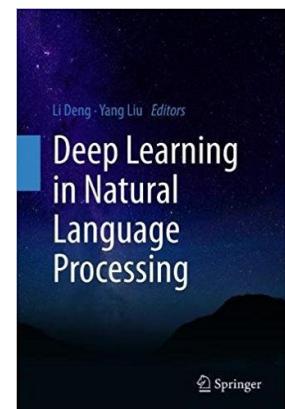
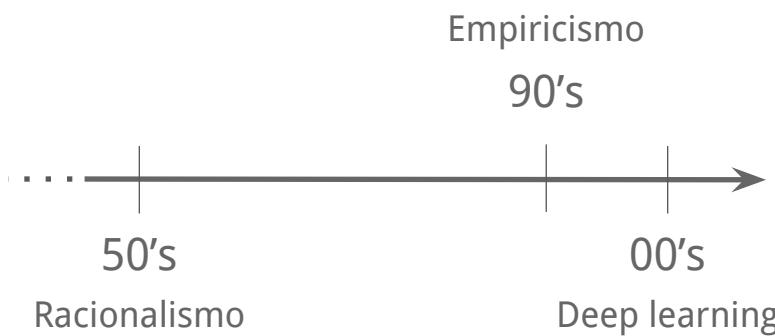
“Processamento de linguagem natural (PLN ou NLP do inglês *natural language processing*) é o campo que projeta métodos e algoritmos que recebem como entrada ou produzem como saída dados em linguagem natural não estruturados”

Goldberg, Y. (2017)



3

As ondas de NLP



Li, D., & Yang, L. (2014)

4

Onda racionalista

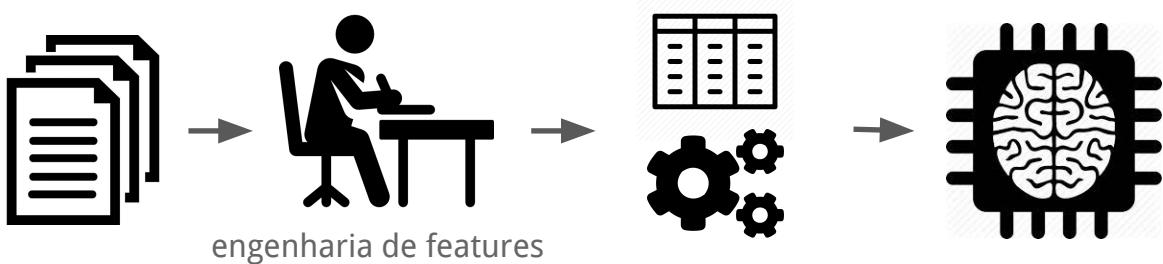
- Conhecimento da linguagem na mente humana é fixo
- Baseado em regras



5

Onda empiricista

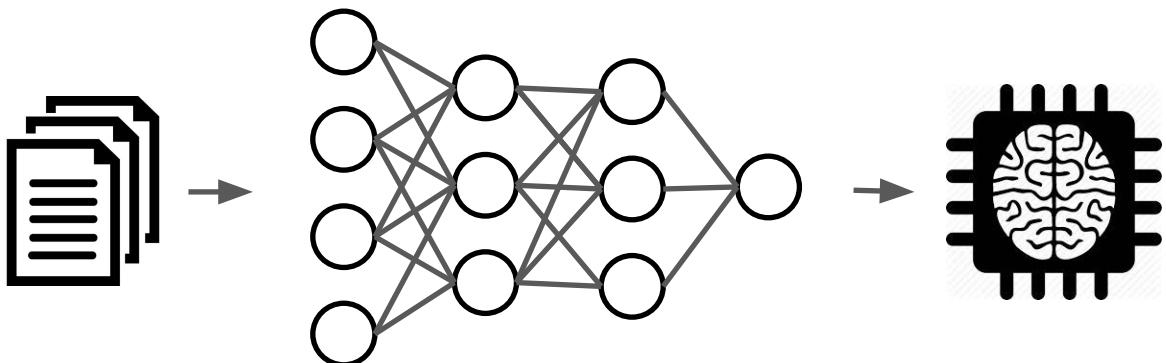
- A mente humana somente começa com operações gerais de associação, reconhecimento de padrões e generalização
- Baseada em corpora



6

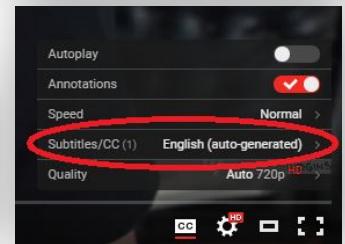
Onda de deep learning

- Independente de especialista de domínio
- Baseado em “grandes corpora”



7

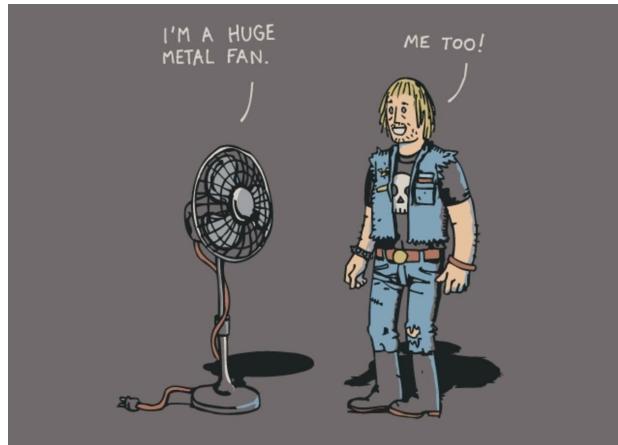
Exemplos de aplicação de NLP



8

Por que as tarefas de NLP são tão difíceis?

A linguagem natural é altamente ambígua



9

Por que as tarefas de NLP são tão difíceis?

A linguagem natural é dependente do contexto



Falkland Islands
@falklands_utd

Argentina doesn't do well against small islands, does it?

5:40 pm · 16 Jun 18



FIFA World Cup 2018
16 jun 18



10

Por que as tarefas de NLP são tão difíceis?

A linguagem natural está sempre evoluindo



11

Tarefas de NLP

Brasil

tarefa sintática

substantivo

6 letras

2 vogais

...

tarefa semântica



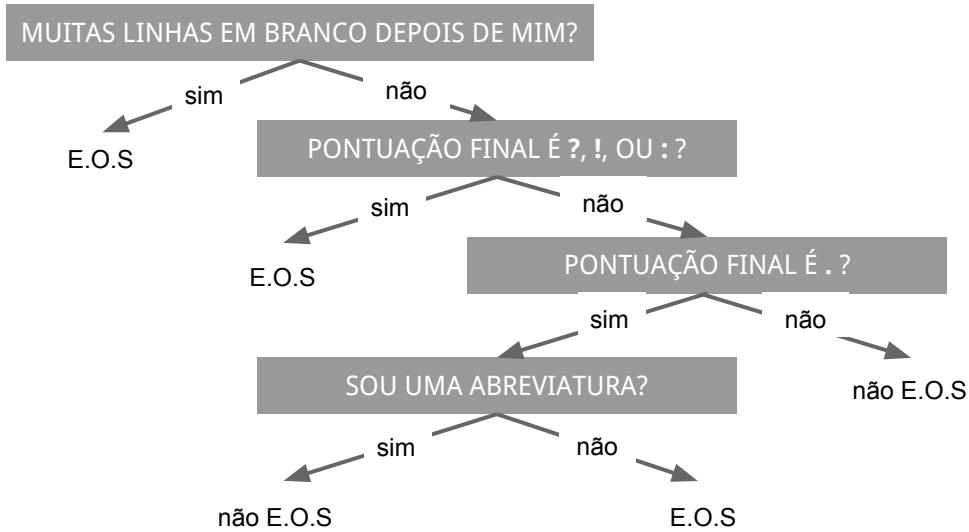
12



Segmentação de sentenças

- !,? não são muito ambíguos
- A pontuação “.” é muito ambígua
 - Números .04 or 5.0 ...
 - Limite da sentença
 - Abreviações (Inc. , Mr.)
- Construção de um classificador binário!
 - Olhe um “.”
 - Decide EndOfSentence/NotEndOfSentence
 - Classificadores: regras escritas manualmente, expressões regulares, machine learning

Segmentação de sentenças



15

Segmentação de palavras

- Separar um trecho de texto contínuo em palavras separadas

ogatonochapéu -> o | gato | no | chapéu

- Em Português é uma tarefa trivial -> espaços
- Em Chinês, Japonês etc as palavras não são limitadas desta maneira

16

Lematização

- Duas palavras têm a mesma raiz?

sou, são, é -> ser

carro, carros -> carro

Os carros de os meninos são de cores diferentes -> ???

17

Stemização

- Reduzir palavras flexionadas a sua base

voto, vota, votei -> vot

provimento, proveu, provido -> prov

Voto em o sentido de dar provimento a o recurso -> ???

18

Etiquetagem de classe gramatical (POS tagging)

- Determinar a classe gramatical de cada palavra em uma sentença

Palavra: O carro é preto

POS: DET N V ADJ

19

Etiquetagem de classe gramatical (POS tagging)

Exemplo de
conjunto de tags:

Nome Classe gramatical

N Substantivo

PROP Nome próprio

SPEC Especificador (pronome ou adjetivo)

DET Determinante (artigo, pronome ou adjetivo)

PERS Pronome pessoal

ADJ adjetivo

ADV advérbio

V verbo

NUM numeral

KS conjunção subordinativa

KC conjunção coordenativa

PRP preposição

IN interjeição

EC prefixo

20

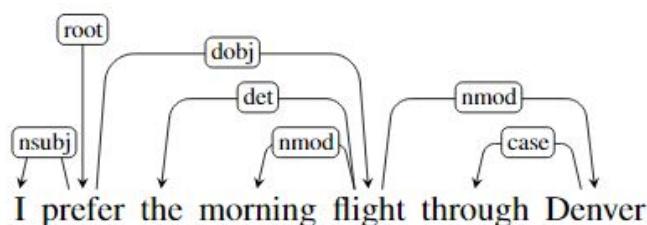
Etiquetagem de classe gramatical (POS tagging)

- Componente importante para análise sintática
- Predição de palavras em reconhecimento de fala
 - – Pronomes possessivos (meu, seu, dela) são geralmente seguidos por substantivos
 - – Pronomes pessoais (eu, você, ele) são geralmente seguidos por verbos
- Tradução automática

21

Dependency parsing

- Idéia básica:
 - A estrutura sintática consiste de relações binárias, assimétricas entre palavras de uma sentença



22

Dependency parsing

Exemplo de conjunto de tags:

Clausal Argument Relations	Description
NSUBJ	Nominal subject
DOBJ	Direct object
IOBJ	Indirect object
CCOMP	Clausal complement
XCOMP	Open clausal complement
Nominal Modifier Relations	Description
NMOD	Nominal modifier
AMOD	Adjectival modifier
NUMMOD	Numeric modifier
APPOS	Appositional modifier
DET	Determiner
CASE	Prepositions, postpositions and other case markers
Other Notable Relations	Description
CONJ	Conjunct
CC	Coordinating conjunction

Figure 14.2 Selected dependency relations from the Universal Dependency set. ([de Marneffe et al., 2014](#))

23

Tarefas semânticas



Predição de palavra

- Calcular a probabilidade de uma palavra vir a seguir



25

Geração de linguagem natural

bathrooms	bedrooms	year_built	num_floors	house_style
4	3.5	2007	1	Craftsman
3	2	1961	1	Bungalow
4	2	2004	1	Ranch
4	2	1930	2	Craftsman
3	2.5	1900	1	Ranch
2	2	1942	2	Bungalow
3	2	2008	2	Cape Cod
2	1.5	1946	2	Modern
3	1.5	1937	2	Craftsman
4	3	1900	1	Bungalow

Row 2
Located in Old Brighton, this lovely bungalow is ideal for buyers looking for the warmth and comfort of a cottage-inspired home.
Row 1
Located in Slope Heights, this lovely craftsman-style home is ideal for buyers looking for modern touches and mid-century character. This one-story home is perfect for large families or entertaining and features four bedrooms and three and a half bathrooms. Lovely hardwoods and luxurious carpeting are featured throughout and the stylish tin roof provides exterior flair.

26

Classificação de documentos

- SPAM?

Subject: Important notice!
From: Stanford University <newsforum@stanford.edu>
Date: October 28, 2011 12:34:16 PM PDT
To: undisclosed-recipients:;

Greats News!

You can now access the latest news by using the link below to login to Stanford University News Forum.

<http://www.123contactform.com/contact-form-StanfordNew1-236335.html>

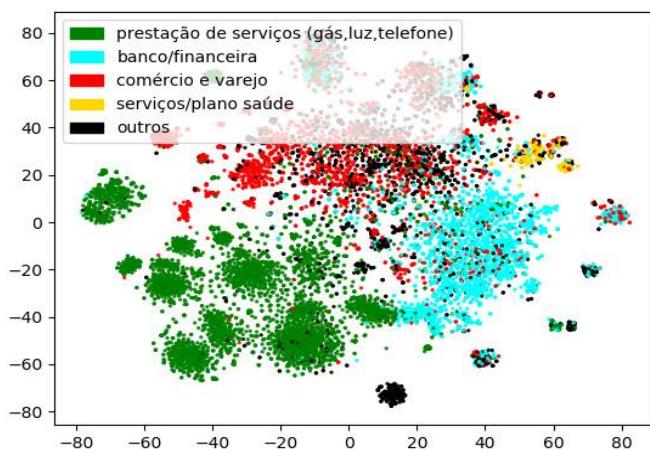
Click on the above link to login for more information about this new exciting forum. You can also copy the above link to your browser bar and login for more information about the new services.

© Stanford University. All Rights Reserved.

27

Classificação de documentos

Qual é o contexto destes processos judiciais?



28

Classificação de documentos

- Atribuir assunto, tópico ou gênero
- Detecção de spam
- Identificação de autoria
- Identificação de idade/gênero
- Identificação de língua
- Análise de sentimento
- ...

Entrada:

um documento d

um conjunto fixo de classes $C = \{c_1, c_2, \dots, c_j\}$

Saída: uma classe predita c de C

29

Classificação de documentos

Regras escritas manualmente

- Regras baseadas em combinações de palavras ou outros atributos
 - spam: endereço em lista negra, US\$
- Qualidade é alta -> especialista
- Construir e manter regras é caro

30

Classificação de documentos

Aprendizado de máquina

Entrada:

- um documento d
- um conjunto fixo de classes $C = \{c_1, c_2, \dots, c_j\}$
- um conjunto de documentos anotados

manualmente

Saída: um classificador treinado

Classificadores

- Naïve Bayes
- Regressão logística
- Support-vector machines
- k-Nearest Neighbors

31

Reconhecimento de entidades nomeadas

- Encontrar e classificar nomes em um texto

The decision by the independent MP Andrew Wilkie to withdraw his support for the minority Labor government sounded dramatic but it should not further threaten its stability. When, after the 2010 election, Wilkie, Rob Oakeshott, Tony Windsor and the Greens agreed to support Labor, they gave just two guarantees: confidence and supply.

32

Reconhecimento de entidades nomeadas

- Encontrar e classificar nomes em um texto

The decision by the independent MP Andrew Wilkie to withdraw his support for the minority Labor government sounded dramatic but it should not further threaten its stability. When, after the 2010 election, Wilkie, Rob Oakeshott, Tony Windsor and the Greens agreed to support Labor, they gave just two guarantees: confidence and supply.

33

Reconhecimento de entidades nomeadas

- Encontrar e classificar nomes em um texto

The decision by the independent MP Andrew Wilkie to withdraw his support for the minority Labor government sounded dramatic but it should not further threaten its stability. When, after the 2010 election, Wilkie, Rob Oakeshott, Tony Windsor and the Greens agreed to support Labor, they gave just two guarantees: confidence and supply.

34

Reconhecimento de entidades nomeadas

- Encontrar e classificar nomes em um texto

The decision by the independent MP Andrew Wilkie to withdraw his support for the minority Labor government sounded dramatic but it should not further threaten its stability. When, after the 2010 election, Wilkie, Rob Oakeshott, Tony Windsor and the Greens agreed to support Labor, they gave just two guarantees: confidence and supply.

Person
Date
Location
Organization

35

Extração de relações

- Texto:

"A professora Simone Barbosa trabalha para a PUC-Rio".
- Entidades nomeadas extraídas:
 - Pessoa Simone Barbosa
 - Organização PUC-Rio
- Relação extraída:
 - Trabalha-para Simone Barbosa -> PUC-Rio

36

Extração de relações

Centro de Operações @OperacoesRio · 20 h
AV ATLÂNTICA | há retenções no sentido Leme, desde a altura do Othon Palace até a Av Princesa Isabel.

(a)



(b)

37

Resolução de correferência

Resolução de correferência é a tarefa de encontrar todas as expressões que referem a mesma entidade em um texto

"I voted for Nader because he was most aligned with my values," she said.

38

Desambiguação

APPLE



39

Respostas a perguntas

AT&T 22:26 65 %

“ Do I need an umbrella tomorrow in San Francisco ”

Yes, San Francisco should get rain tomorrow:

Weekly Forecast

TUES	WED	THU	FRI	SAT	SUN
14° 9°	16° 7°	17° 6°	17° 6°	18° 7°	18° 8°



40

Sumariamento

cia documents reveal iot-specific televisions can be used to secretly record conversations .

cybercriminals who initiated the attack managed to commandeer a large number of internet-connected devices in current use .

cia documents revealed that microwave ovens can spy on you - maybe if you personally don't suffer the consequences of the sub-par security of the iot .

Internet of Things (IoT) security breaches have been dominating the headlines lately . WikiLeaks's trove of CIA documents revealed that internet-connected televisions can be used to secretly record conversations . Trump's advisor Kellyanne Conway believes that microwave ovens can spy on you - maybe she was referring to microwave cameras which indeed can be used for surveillance . And don't delude yourself that you are immune to IoT attacks , with 96 % of security professionals responding to a new survey expecting an increase in IoT breaches this year . Even if you personally don't suffer the consequences of the sub-par security of the IoT , your connected gadgets may well be unwittingly cooperating with criminals . Last October , Internet service provider Dyn came under an attack that disrupted access to popular websites . The cybercriminals who initiated the attack managed to commandeer a large number of internet-connected devices (mostly DVRs and cameras) to serve as their helpers . As a result , cybersecurity expert Bruce Schneier has called for government regulation of the IoT , concluding that both IoT manufacturers and their customers don't care about the security of the 8.4 billion internet-connected devices in current use . Whether because of government regulation or good old-fashioned self-interest , we can expect increased investment in IoT security technologies . In its recently-released TechRadar report for security and risk professionals , Forrester Research discusses the outlook for the 13 most relevant and important IoT security technologies , warning that "there is no single , magic security bullet that can easily fix all IoT security issues ." Based on Forrester's analysis , here's my list of the 6 hottest technologies for IoT security : IoT network security : Protecting and securing the network connecting IoT devices to back-end systems on the internet . IoT network security is a bit more challenging than traditional network security because there is a wider range of communication protocols , standards , and device capabilities , all of which pose significant issues and increased complexity . Key capabilities include traditional endpoint security features such as antivirus and antimalware as well as other features such as firewalls and intrusion prevention and detection systems . Sample vendors : Bayshore Networks , Cisco , Darktrace , and Sentri . IoT authentication : Providing the ability for users to authenticate an IoT device including managing multiple users of a single device (such as a connected car) ranging from simple static password/sins to more robust authentication mechanisms such as two-factor

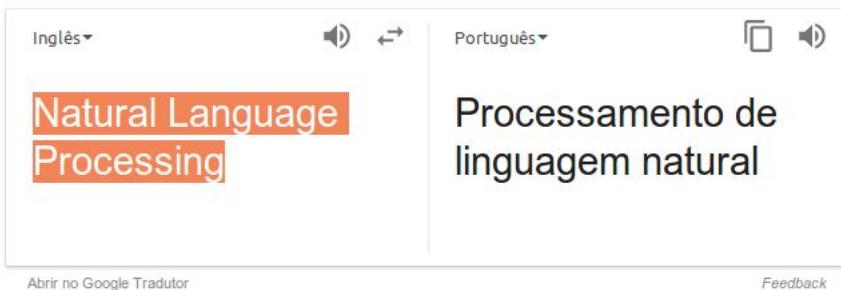
41

Sumariamento

- Aplicações
 - esboços ou resumos de qualquer documento
 - sumários de threads de email
 - itens de ação de uma reunião
 - simplificação de texto pela compressão de sentenças

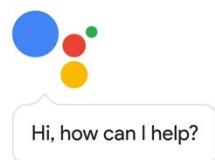
42

Tradução



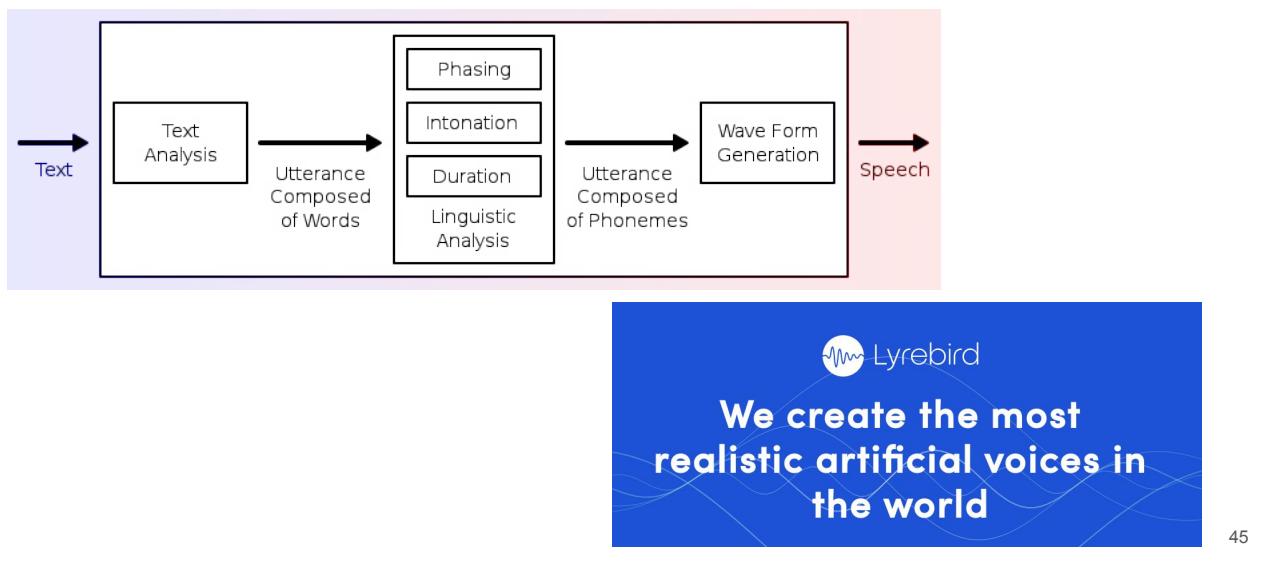
43

Reconhecimento de fala



44

Sintetização de fala



Regular expression

- A formal language for specifying text strings
- How can we search for any of these?
 - woodchuck
 - woodchucks
 - Woodchuck
 - Woodchucks



47

Regular expression

Pattern	Matches
[wW] oodchuck	Woodchuck, woodchuck
[1234567890]	Any digit

Pattern	Matches	
[A-Z]	An upper case letter	Drenched Blossoms
[a-z]	A lower case letter	my beans were impatient
[0-9]	A single digit	Chapter 1: Down the Rabbit Hole

48

Regular expression

Pattern	Matches
[^A-Z]	Not an upper case letter
[^Ss]	Neither 'S' nor 's'
[^e^]	Neither e nor ^
a^b	The pattern a carat b

49

Regular expression

Pattern	Matches
groundhog woodchuck	
yours mine	yours mine
a b c	= [abc]
[gG] roundhog [Ww]oodchuck	

50

Regular expression

Pattern	Matches	
colou?r	Optional previous char	<u>color</u> <u>colour</u>
oo*h!	0 or more of previous char	<u>oh!</u> <u>ooh!</u> <u>oooh!</u> <u>ooooh!</u>
o+h!	1 or more of previous char	<u>oh!</u> <u>ooh!</u> <u>oooh!</u> <u>ooooh!</u>
baa+		<u>baa</u> <u>baaa</u> <u>baaaa</u> <u>baaaaa</u>
beg.n		<u>begin</u> <u>begun</u> <u>begun</u> <u>beg3n</u>

51

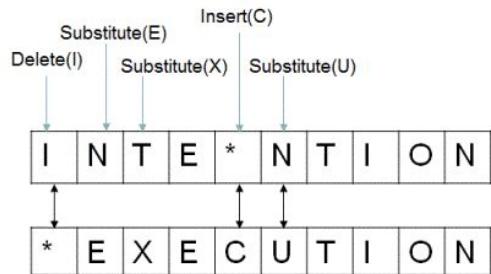
Regular expression

Pattern	Matches
^ [A-Z]	<u>Palo</u> <u>Alto</u>
^ [^A-Za-z]	<u>1</u> <u>"Hello"</u>
\. \$	The end <u>.</u>
. \$	The end <u>?</u> The end <u>!</u>

52

Edit distance

- How similar are two strings?
- Spell correction:
 - The user typed “graffe”
- What is the closest?
 - graffite
 - graf
 - grail
 - giraffe
- Machine translation, Speech Recognition..



53

Bag of words

Also known as vector space models (VSMs)

Sentence 1 (S1): The car is blue.

Sentence 2 (S2): The color of her eyes is blue.

	the	of	is	her	blue	color	eyes	car
S1:	1	0	1	0	1	0	0	1
S2:	1	1	1	1	1	1	1	0

54

Bag of words

$$\mathbf{D} = \begin{pmatrix} \text{tf}(t_1, d_1) & \cdots & \text{tf}(t_N, d_1) \\ \vdots & \ddots & \vdots \\ \text{tf}(t_1, d_\ell) & \cdots & \text{tf}(t_N, d_\ell) \end{pmatrix}$$

Where:

I is the number of documents

N is the number of terms

$\text{tf}(t_i, d_j)$ is the frequency of term t_i in document d_j

55

Bag of words > input matrix

		the	of	is	her	blue	color	eyes	car
D	S1	1	0	1	0	1	0	0	1
	S2	1	1	1	1	1	1	1	0
		S1	S2						
D'	the	1	1						
	of	0	1						
	is	1	1						
	her	0	1						
	blue	1	1						
	color	0	1						
	eyes	0	1						
	car	1	0						

56

Bag of words > input matrix

Term-by-term matrix $D'D$

Document-by-document matrix DD'

Where

$$D' = D^T$$

57

Bag of words > semantic issues

VSM representation ignores semantic relations

Word position, context, synonyms, homonyms, etc.

58

Bag of words > improving the quality

Associate different weights (w_i) to each term t_i

Remove uninformative terms (*and, or, the*, etc.)

Stemming (structural and structured to structural)

Remove the effect of the length of the document

Operations can be performed in sequence

59

Bag of words > kernel

$K = DD'$

Uses the document-based representation

$$\kappa(d_1, d_2) = \langle \phi(d_1), \phi(d_2) \rangle = \sum_{j=1}^N \text{tf}(t_j, d_1) \text{tf}(t_j, d_2)$$

60

Bag of words > term weighting

Not all terms have the same importance

Stop words (*and, or, the*, etc.) are removed before the analysis

The frequency of a word across the documents → amount of information

61

Bag of words > term weighting

We consider *inverse document frequency (idf)*

$$w(t) = \ln \left(\frac{\ell}{\text{df}(t)} \right)$$

Where

ℓ : number of documents in the corpus

$\text{df}(t)$: number of documents containing term t

Stop words can be treated here ($w(t) = 0$), but for efficiency it is better to remove them earlier

62

Bag of words > term weighting

Given a term weighing, we define a new VSM

We choose the diagonal matrix \mathbf{R} as

$$\mathbf{R}_{tt} = w(t)$$

The kernel computes

$$\tilde{\kappa}(d_1, d_2) = \phi(d_1) \mathbf{R} \mathbf{R}' \phi(d_2)' = \sum_t w(t)^2 \text{tf}(t, d_1) \text{tf}(t, d_2)$$

The evaluation of this kernel involves

Term frequencies and inverse document frequencies (tf-idf)

63

Bag of words > term proximity

tf-idf is not capable of recognizing two terms that are semantically related

Two documents that share no terms, but share synonyms are not connected

Therefore we need to establish semantic similarities between terms

We need a non-zero off-diagonal matrix \mathbf{P}

Where

$$P_{ij} > 0, \text{ if term } i \text{ is semantically related to term } j$$

64

Bag of words > term proximity

Given \mathbf{P} , the vector space kernel

$$\tilde{\kappa}(d_1, d_2) = \phi(d_1) \mathbf{P} \mathbf{P}' \phi(d_2)'$$

Related to an IR technique known as *query expansion*

We can view $\mathbf{P} \mathbf{P}'$ as encoding semantic strength between terms

We will present some methods for obtaining semantic relationships

65

Bag of words > term proximity > explicit construction

Use a semantic network such as Wordnet

- S: (n) **wife**, married woman
- direct hyponym / full hyponym
 - has instance
 - direct hypernym / inherited hypernym / sister term
 - S: (n) woman, adult female
 - S: (n) spouse, partner, married person, mate, better half
 - antonym
 - derivationally related form

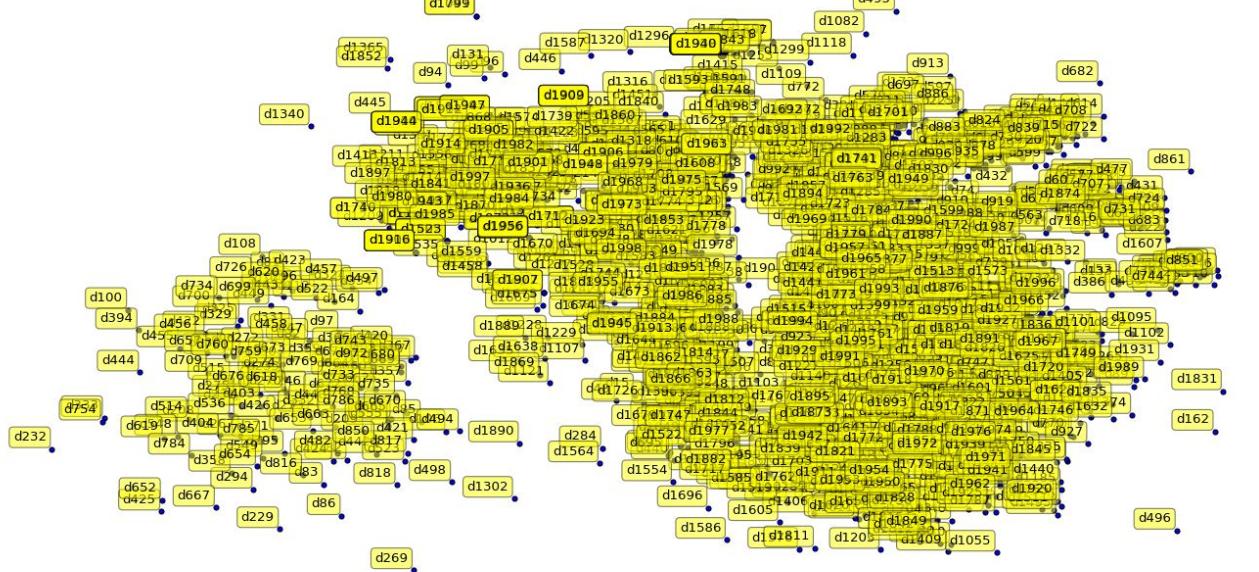
We handcraft matrix \mathbf{P}

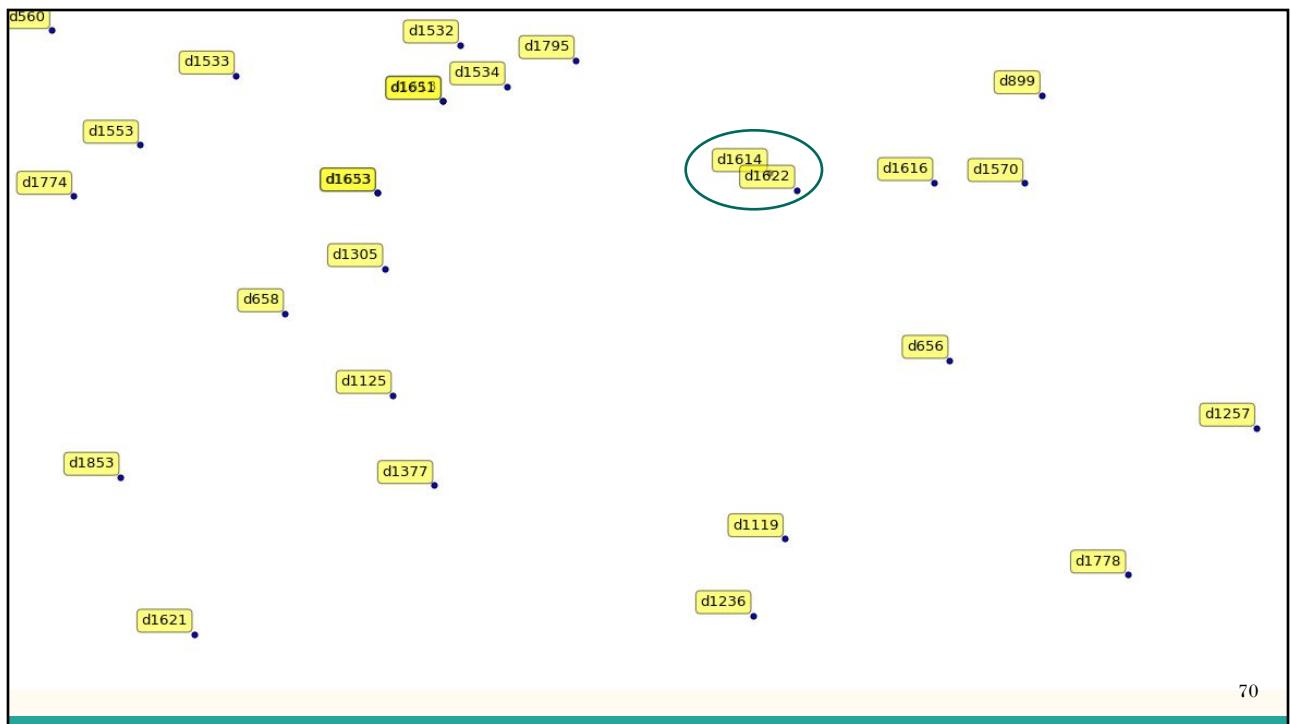
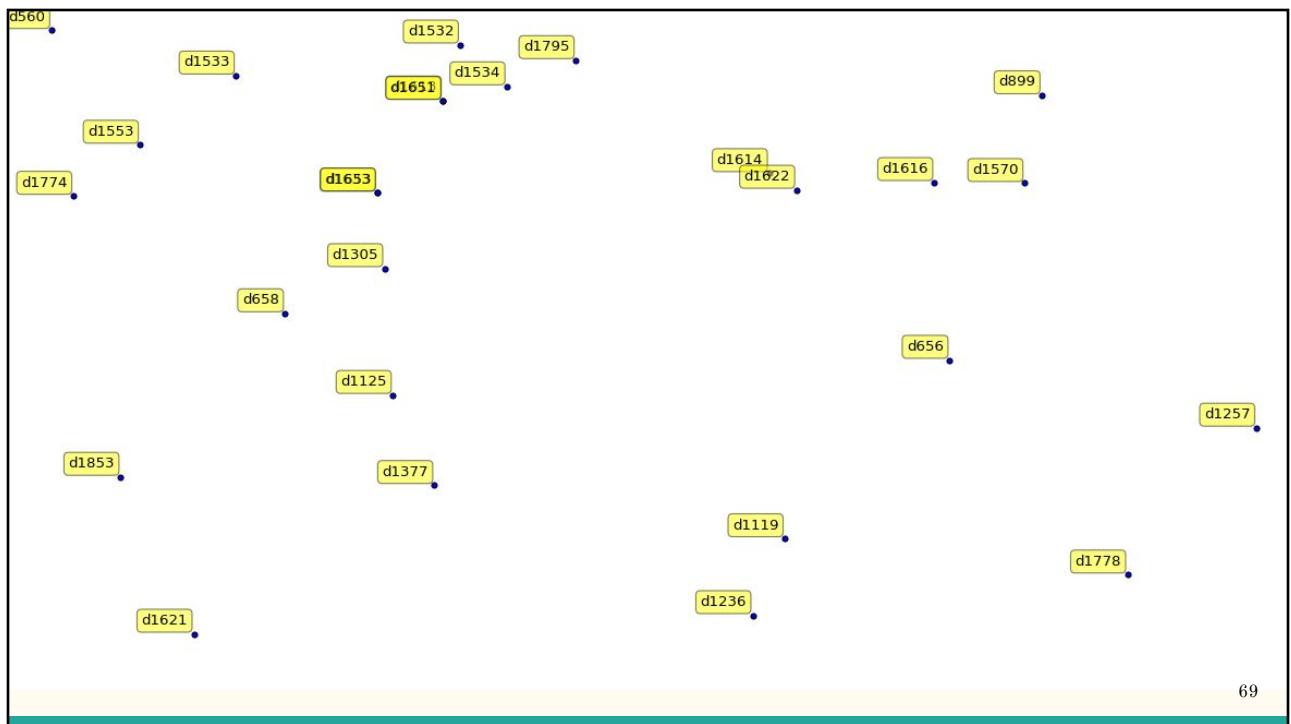
66

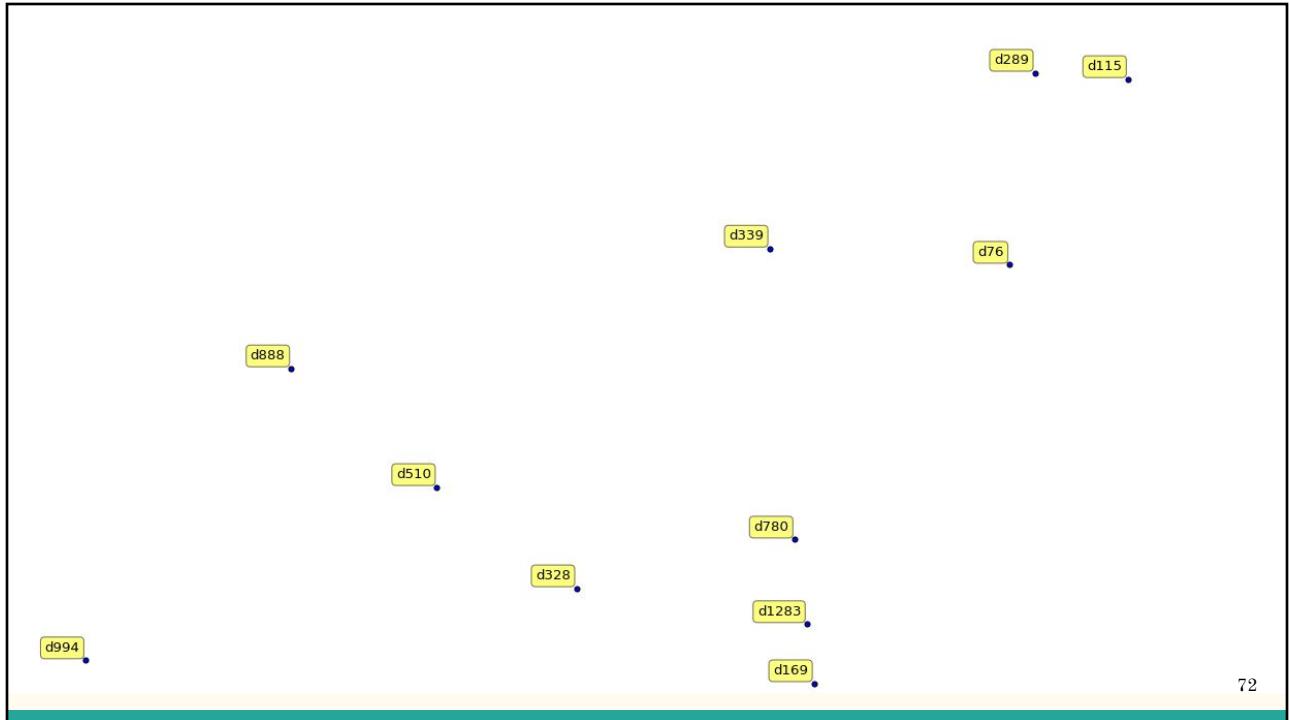
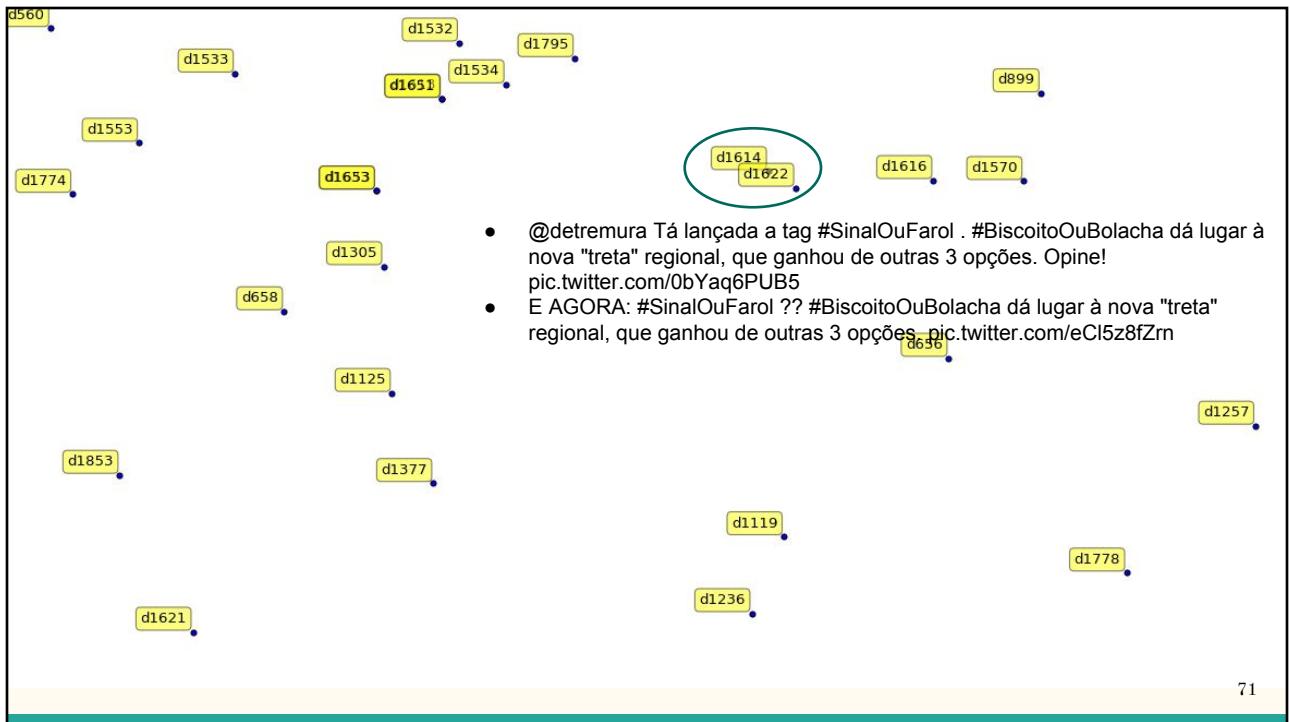
Bag of words > application

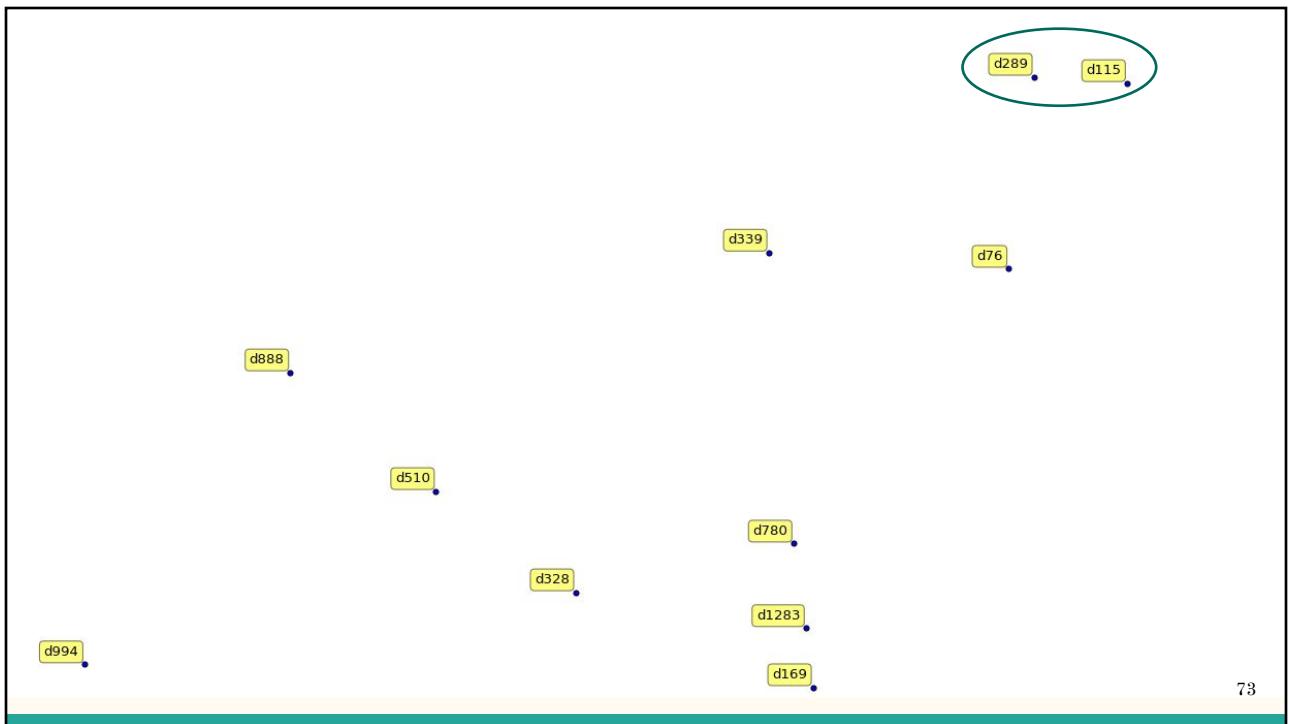
Document-based kernel

$$\mathbf{K} = \mathbf{D}\mathbf{D}'$$

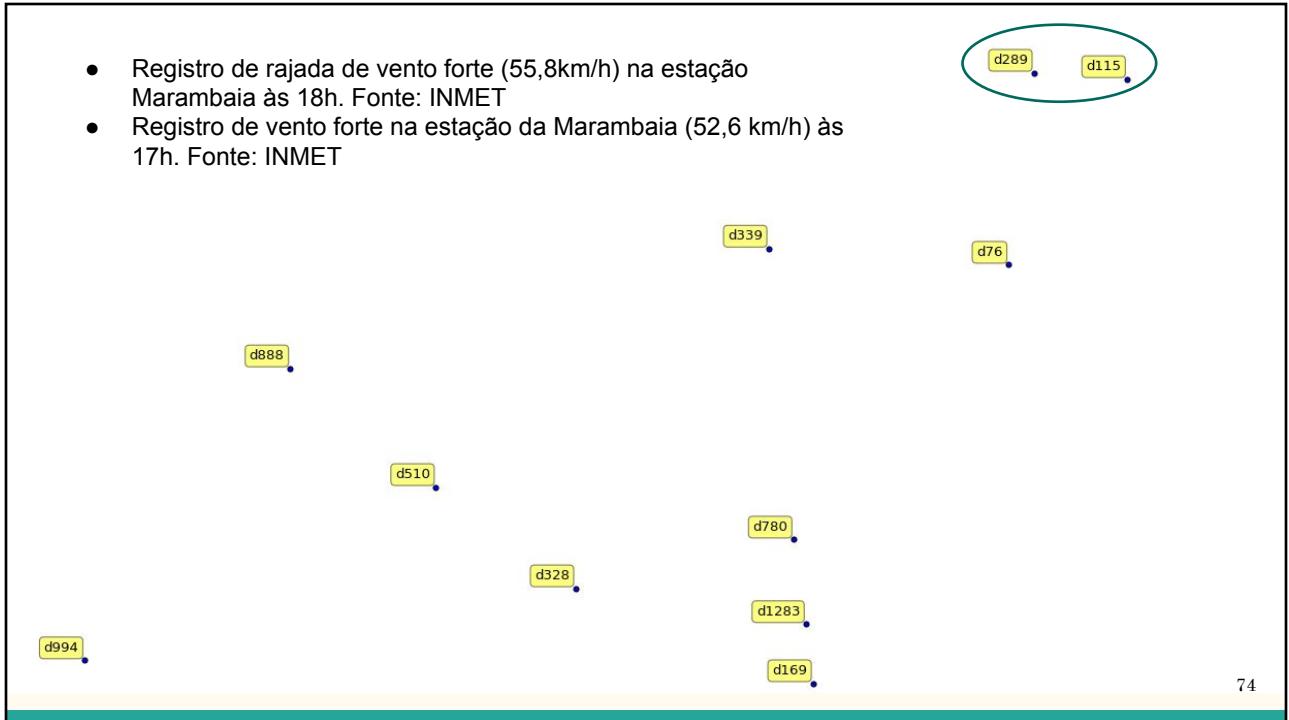








- Registro de rajada de vento forte (55,8km/h) na estação Marambaia às 18h. Fonte: INMET
- Registro de vento forte na estação da Marambaia (52,6 km/h) às 17h. Fonte: INMET

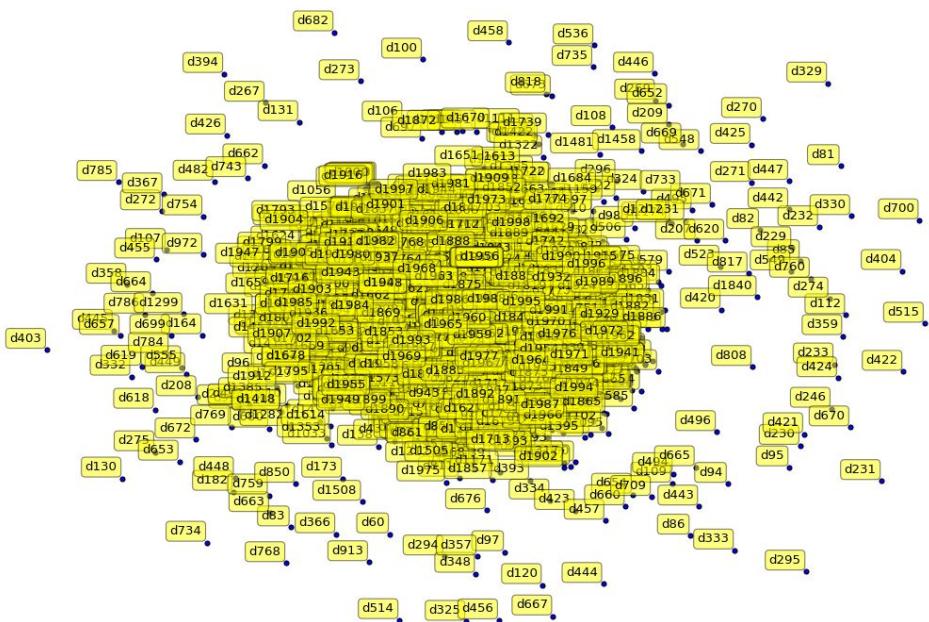


Bag of words > application

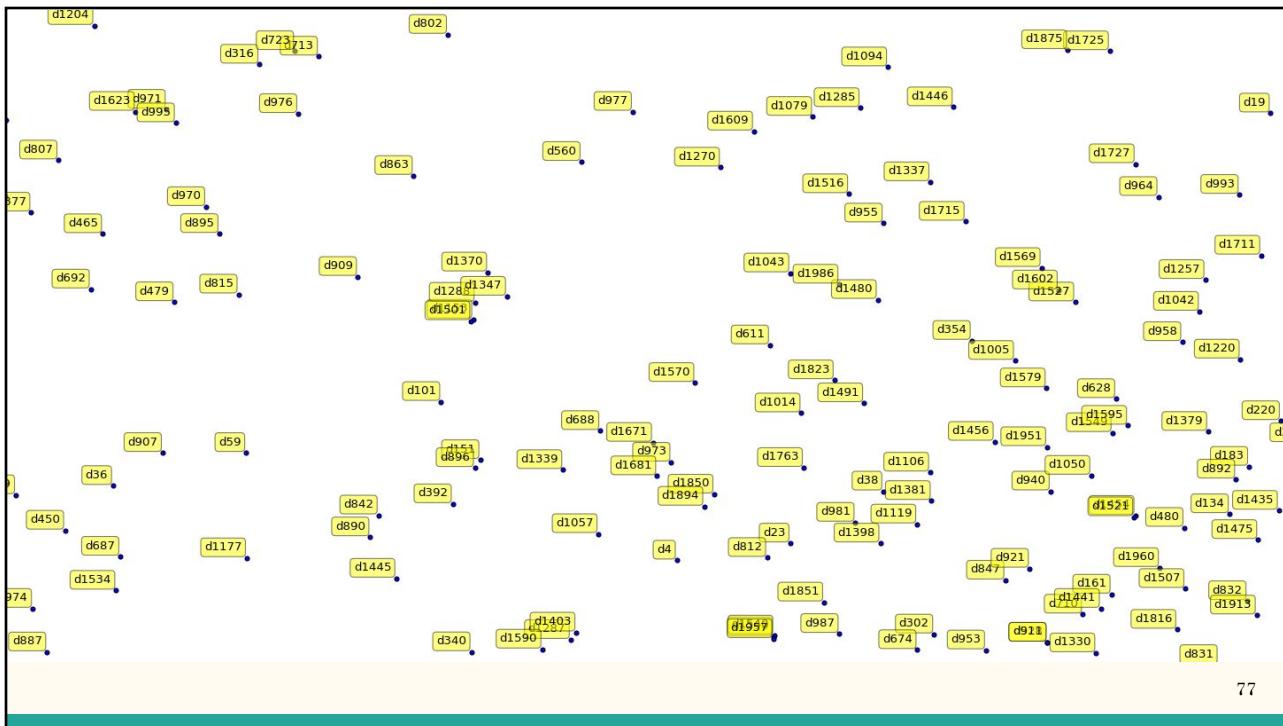
Document-based tf-idf kernel

$$\mathbf{K} = \mathbf{D}\mathbf{R}\mathbf{R}'\mathbf{D}'$$

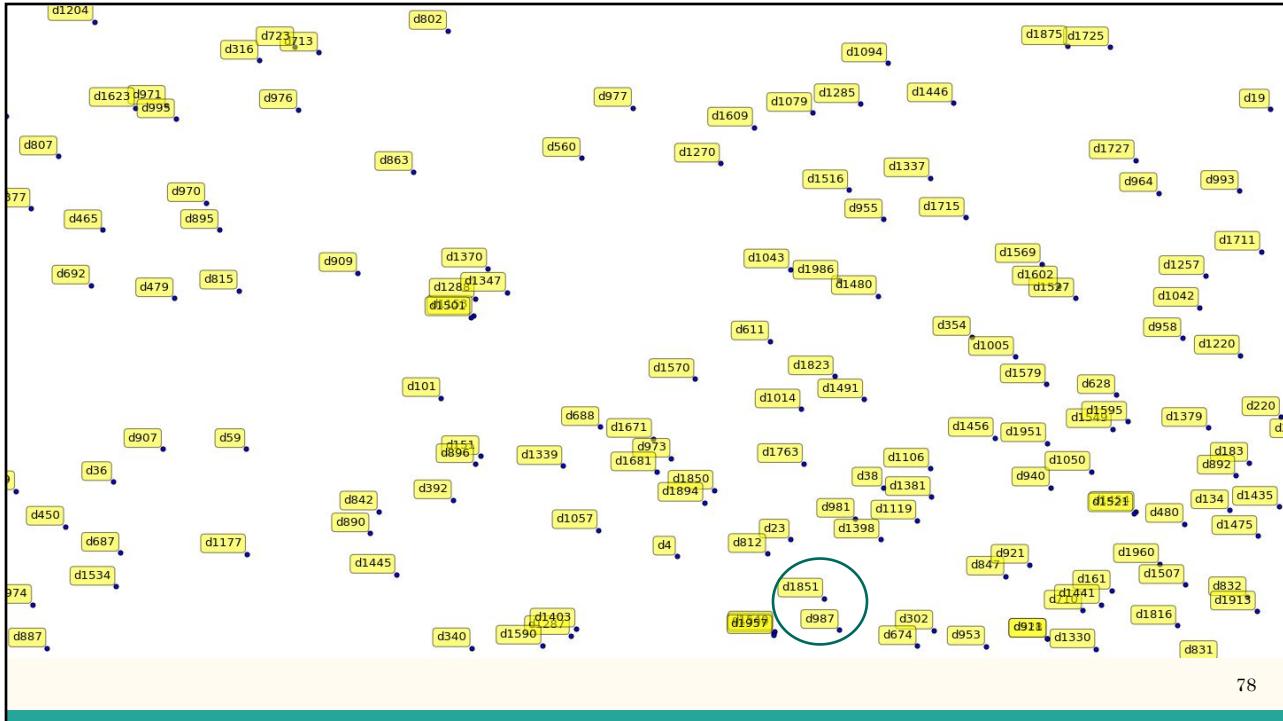
75



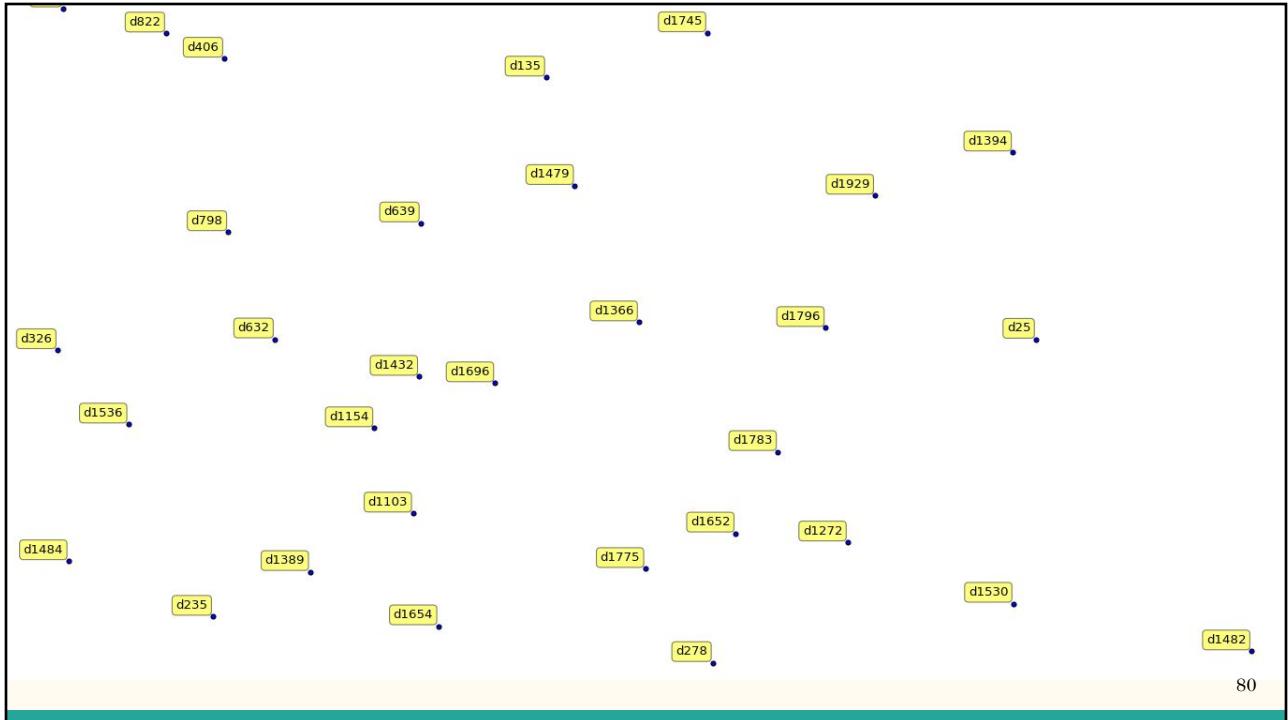
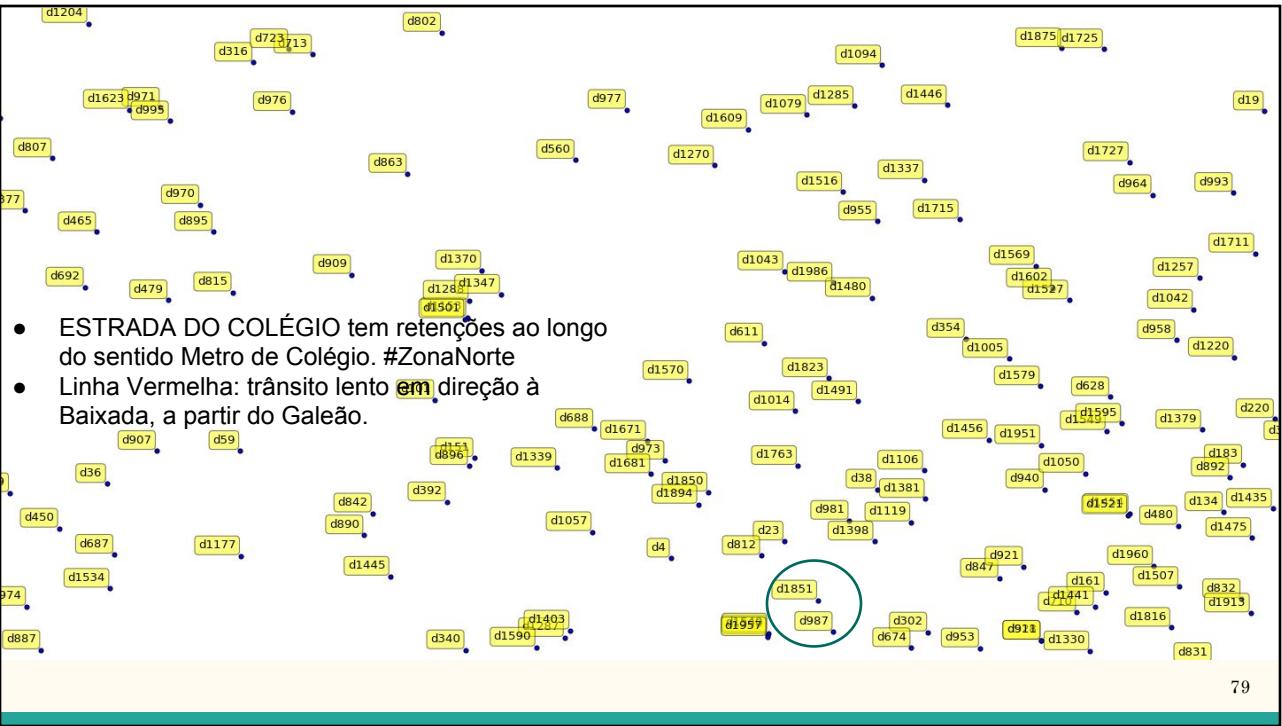
76

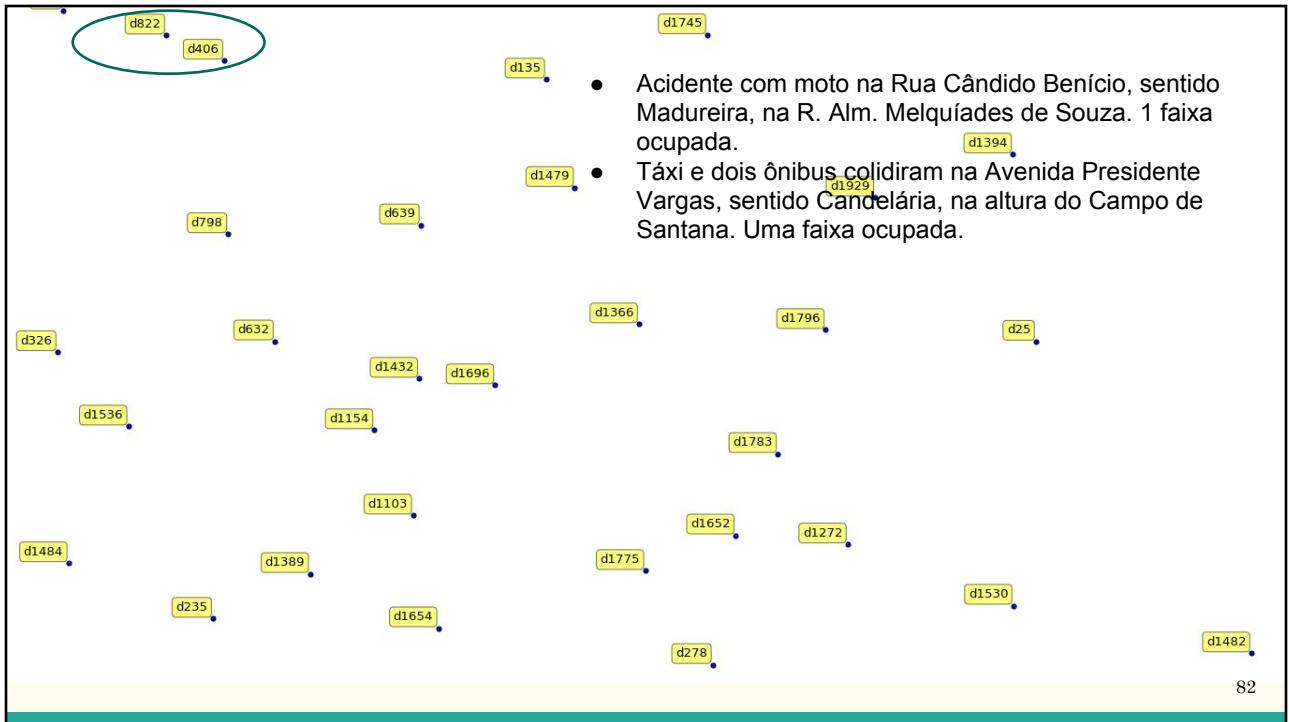
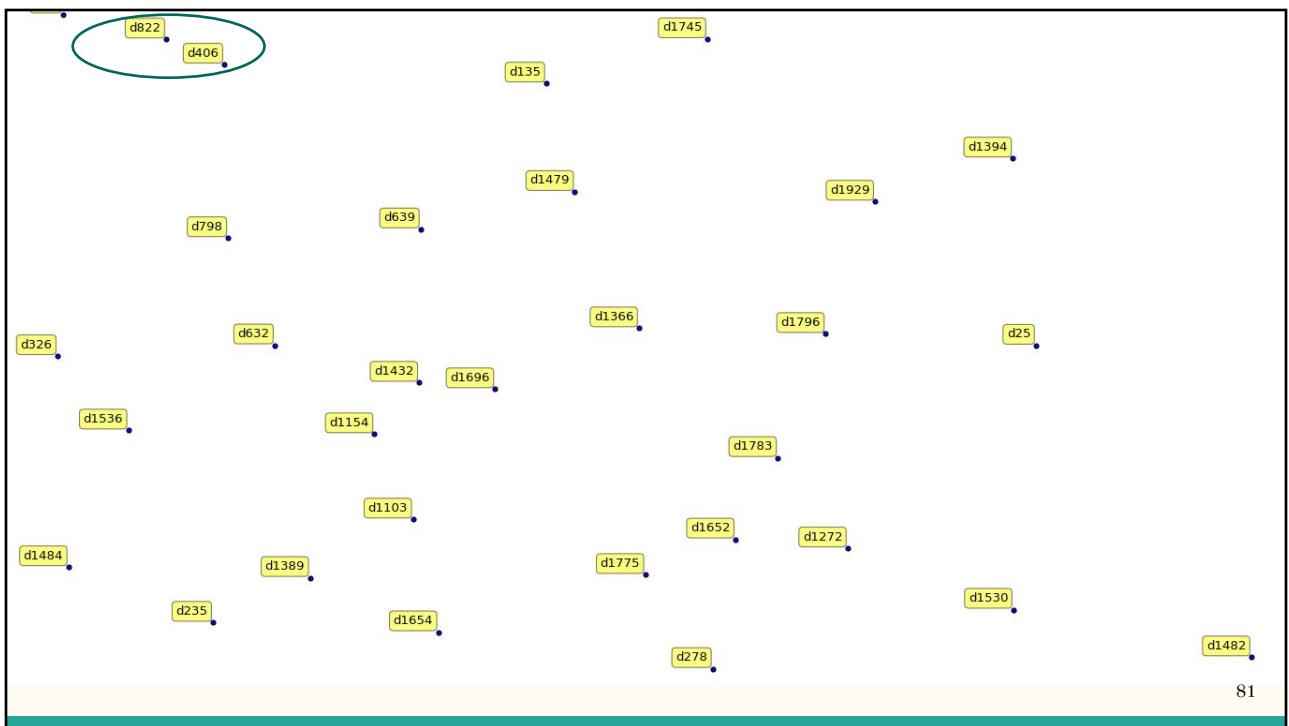


77



78



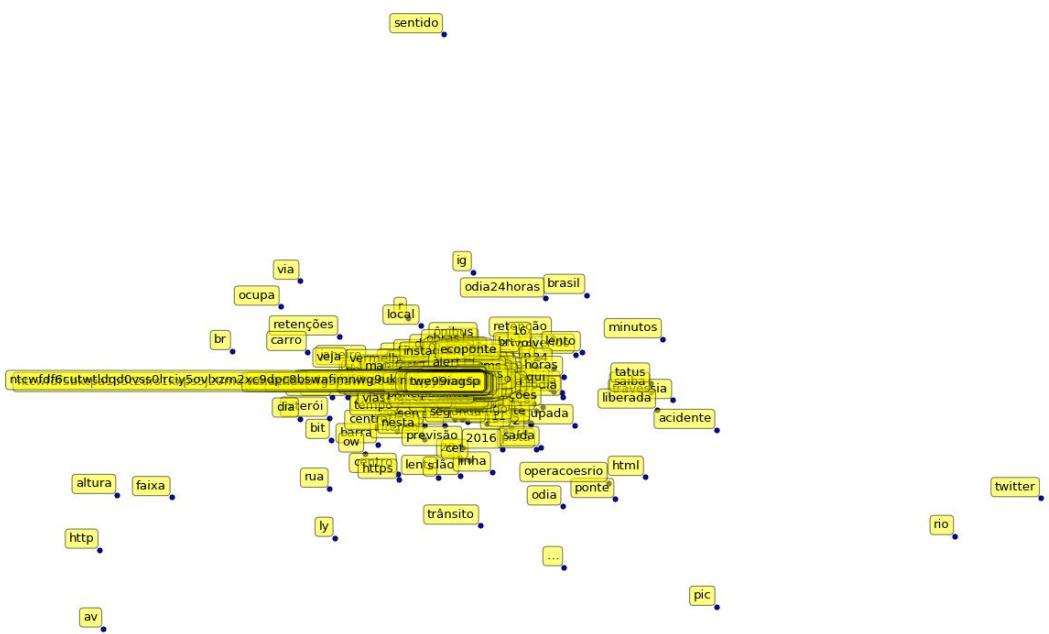


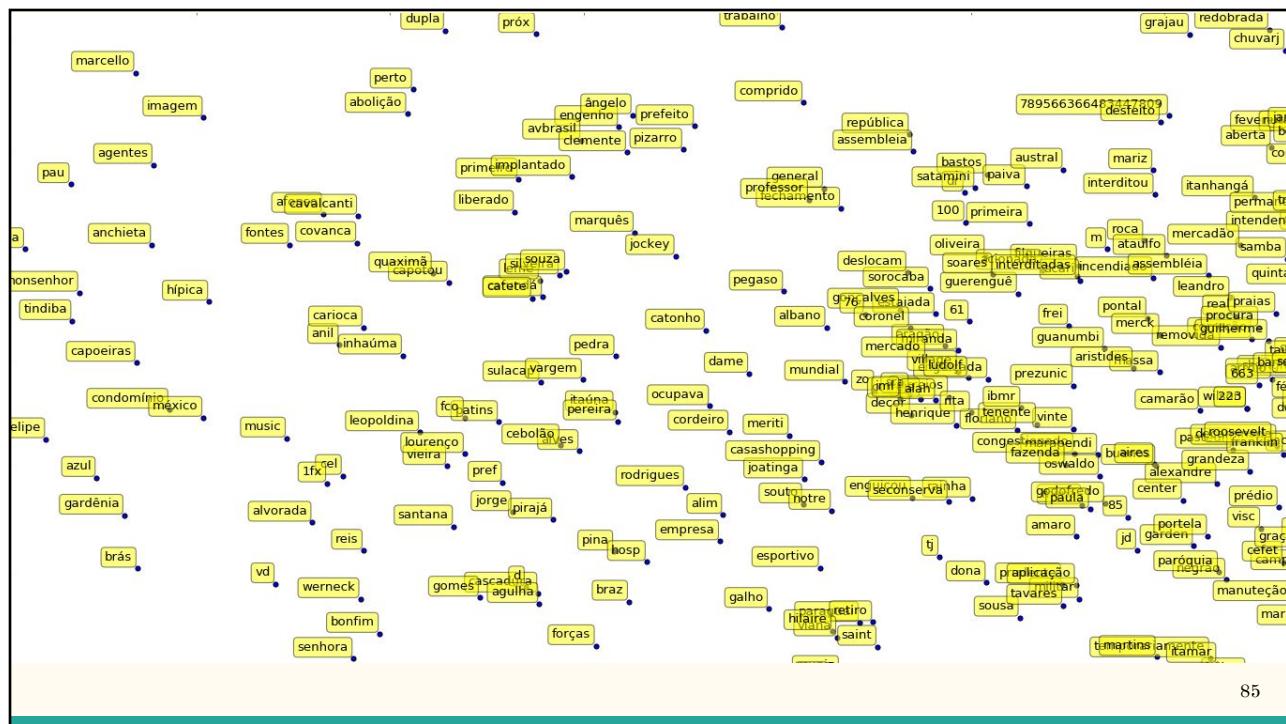
Bag of words > application

Term-based kernel

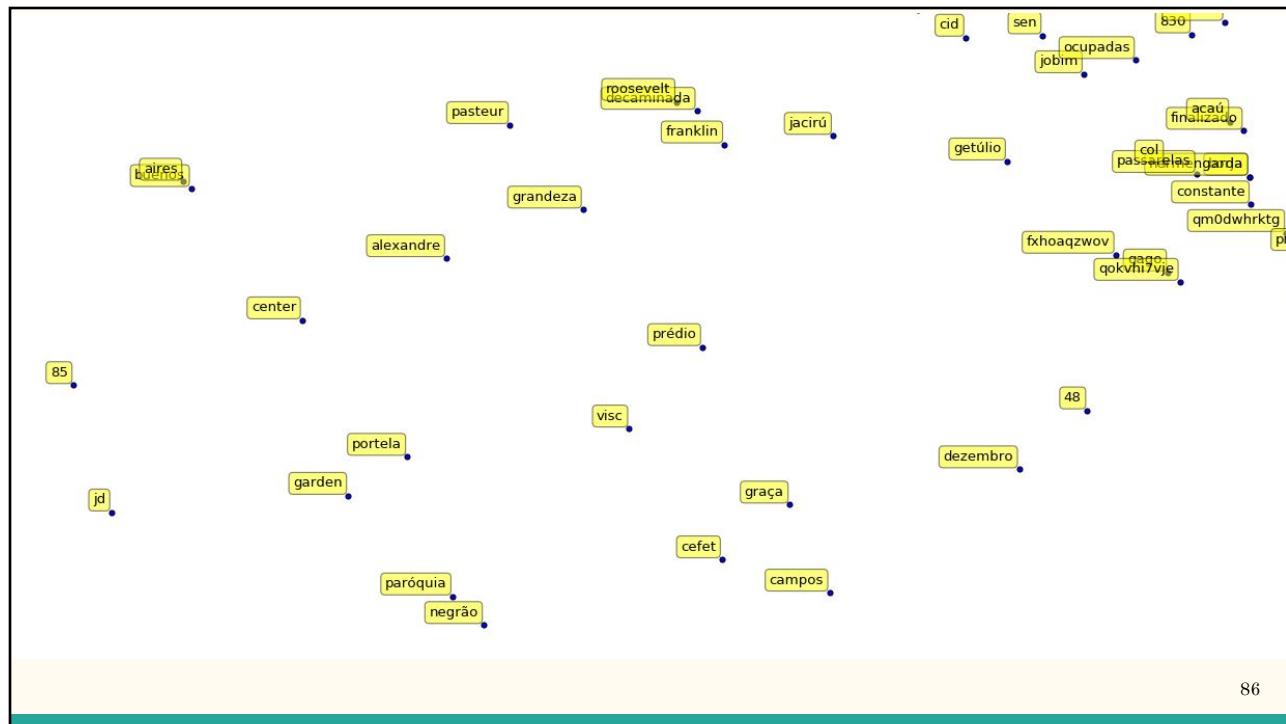
$$\mathbf{K} = \mathbf{D}'\mathbf{D}$$

83





85



86

Bag of words > application

Document-based kernel with proximity

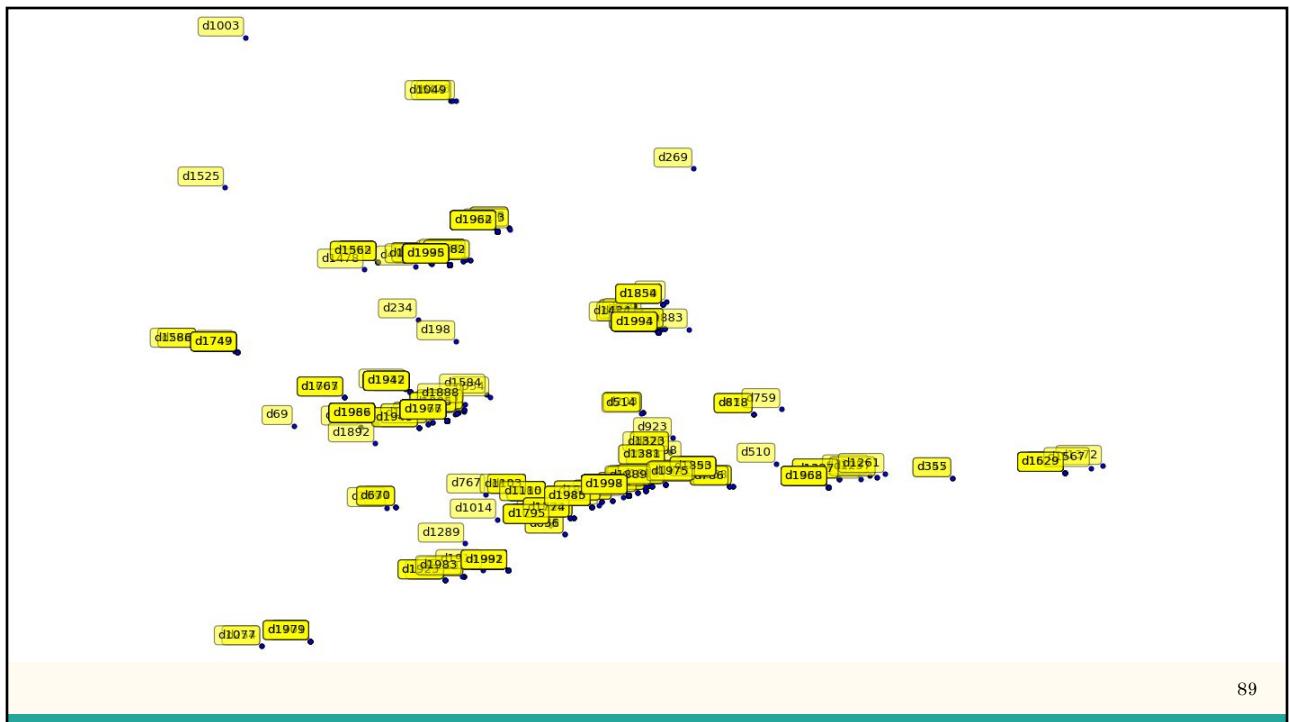
$$K = D P P' D'$$

87

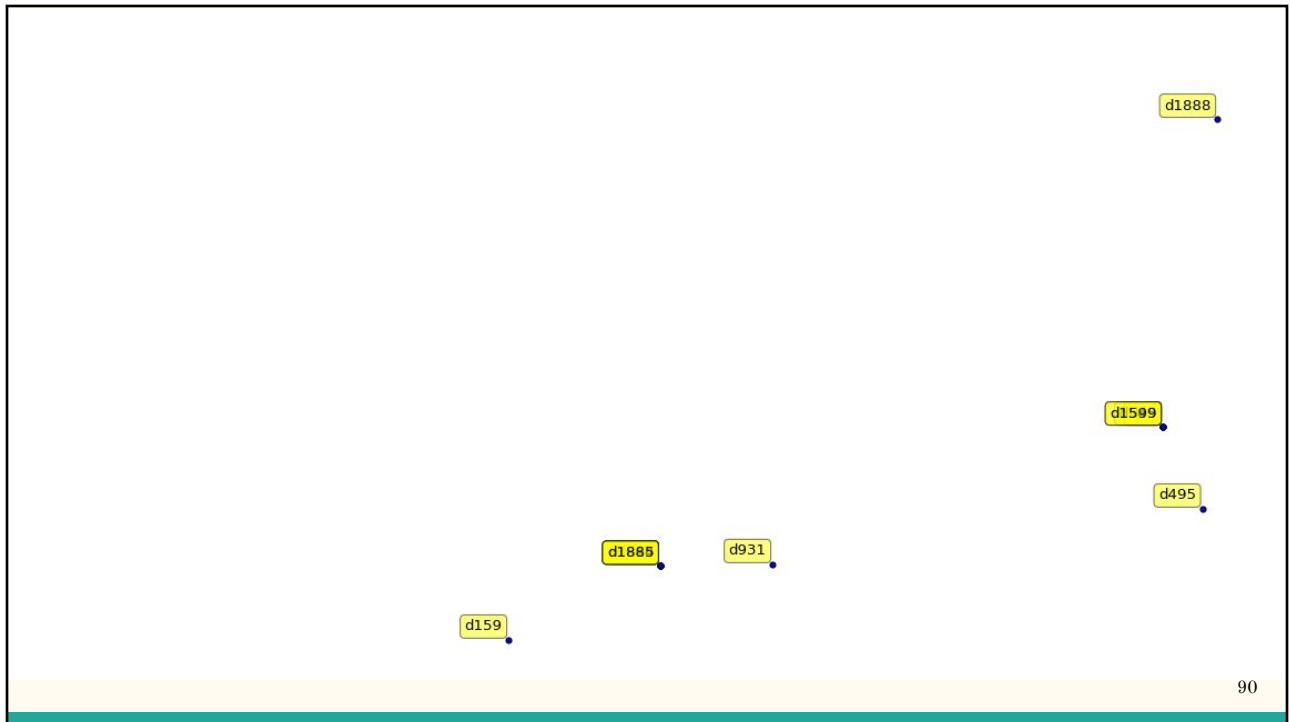
Bag of words > application > synonyms

- interdicao - fechada
- colisao - engavetamento -
choque - capotamento
- queda - atropelamento
- enguicado - pane
- trafeago - transito
- chuva - alagamento - bolsao
d agua
- nas duas direcoes - nos dois
sentidos
- direcao - sentido
- altura - perto de - proximo de
- neste momento - agora - em 1
hora
- carro - onibus - motocicleta
- bom - livre
- intenso - comretencoes
- lento - lentidao - parado -
congestionamento - retencao
- causou - gerou - occasionou -
complica

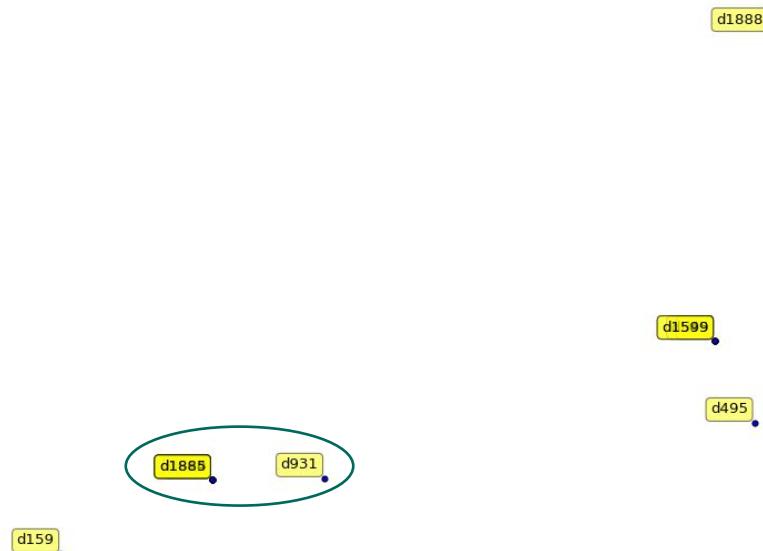
88



89



90



- Atropelamento envolvendo moto na Estrada das Capoeiras, próximo ao McDonald's, sentido Av. Brasil. 1 faixa ocupada.
- PRAÇA SIBÉLIUS: há congestionamento, no sentido Humaitá, devido a um acidente mais a frente, na R. Jardim Botânico. pic.twitter.com/fU5n7R2RBr



Pareamento de palavras

Similar to VSMs

Sentence 1: New York is very crowded.

Sentence 2: Washington Post is a leading newspaper.

	new	york	is	very	crowded	washington	post	...
S1	1	1	1	1	1			
S2						1	1	

93

Pareamento de palavras

Similar to VSMs

Sentence 1: New York is very crowded.

Sentence 2: Washington Post is a leading newspaper.

	new	...	new	york	york	is	is	very	very	crowded	washington	post	post	is	is	a	a	leading	leading	newspaper
S1	1		1	1	1	1	1	1	1											
S2											1	1	1	1	1	1	1	1	1	1

94

Pareamento de palavras

Similar to VSMs

Sentence 1: New York is very crowded.

Sentence 2: Washington Post is a leading newspaper.

	new	...	<u>new</u>	york	is	is very	very crowded	<u>washington</u>	<u>post</u>	post	is a	a leading	<u>leading</u>	<u>newspaper</u>
S1	1		1	1	1		1							
S2									1	1	1	1	1	1

95

Pareamento de palavras

Bi-gram: Santa Maria, Best Buy, Universidade Federal

Tri-gram: Rio de Janeiro, Juiz de Fora, New York Times, Pontifícia Universidade Católica

4-gram: United States of America, Folha de São Paulo

96

Machine learning

Naïve Bayes

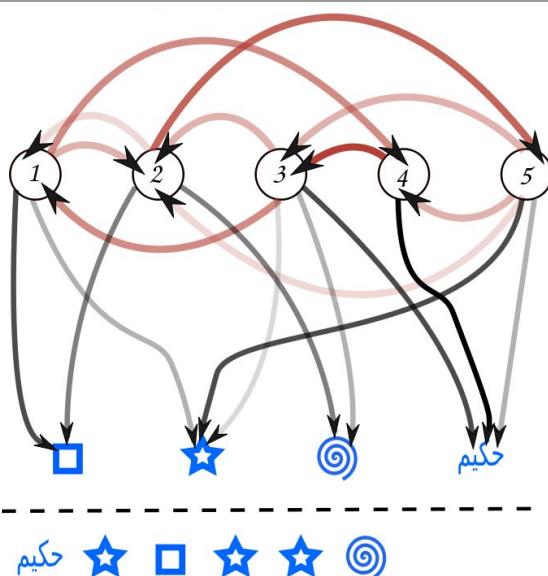
Logistic Regression

Hidden Markov Models (HMM's)

Long Short Term Memory Networks (LSTM's)

97

Machine learning > HMM's

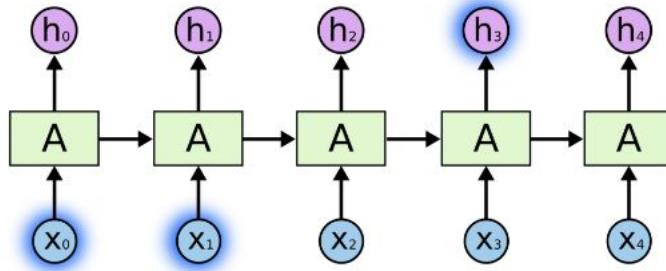


98

Machine learning > LSTM's

Long Short Term Memory Networks are a special kind of Recurrent Neural Networks (RNN's)

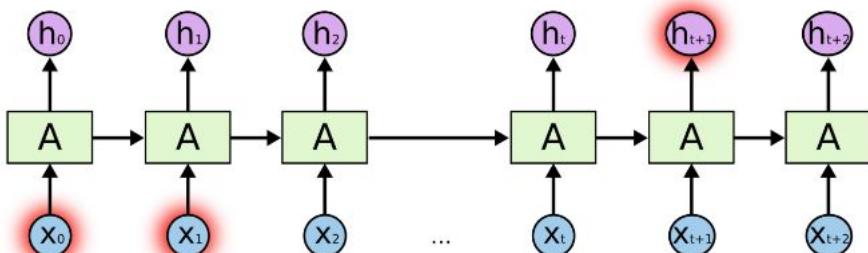
The idea of RNN's is to use past information to the present task



99

Machine learning > LSTM's

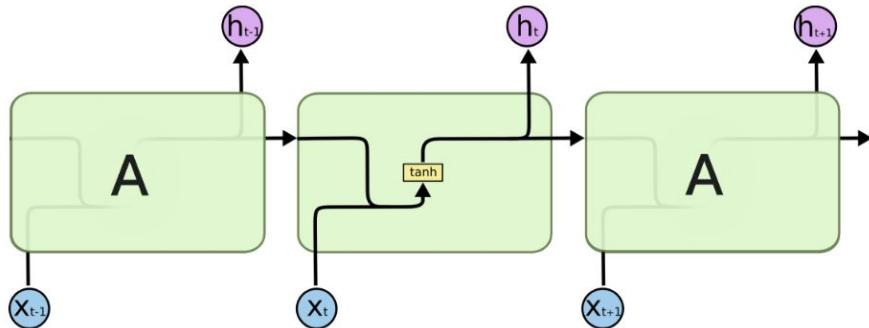
However RNN's do not deal with long distances



100

Machine learning > LSTM's

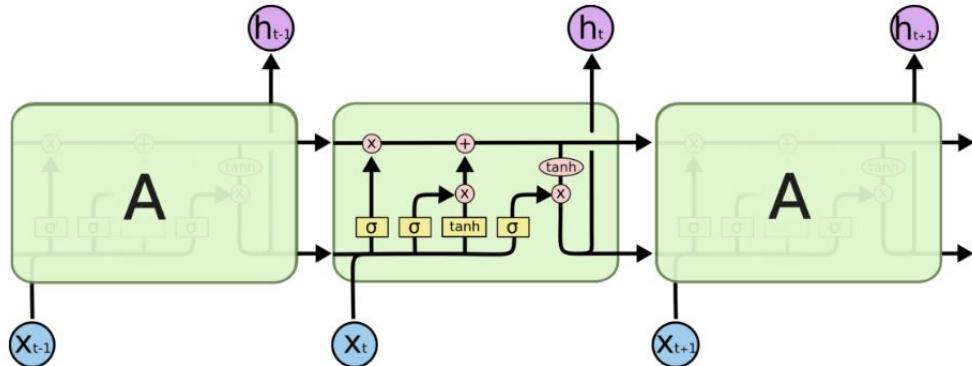
The repeating module in a RNN has one layer



101

Machine learning > LSTM's

The repeating module in a LSTM has four layers



102

word2vec

word2vec is a group of related models used to produce word embeddings

These models are used to reconstruct linguistic contexts of words

Takes as input a large corpus of text

Produces a vector space

For each unique word is assigned a vector

It was created at Google by a team of researchers led by Tomas Mikolov

103

word2vec

Simple neural network

One hidden layer

Train the neural network to execute a task

Take only the hidden layer

104

word2vec

Given a specific word in the middle of a sentence

Pick up a word at random

The network will tell the probability for every word in the vocabulary of being the “nearby word”

Input word *Soviet* → output probability of *Union* and *Russia* is higher than *watermelon* and *bottle*

105

word2vec > dataset

Source Text

The quick brown fox jumps over the lazy dog. →

Training Samples

(the, quick)
(the, brown)

The quick brown fox jumps over the lazy dog. →

(quick, the)
(quick, brown)
(quick, fox)

The quick brown fox jumps over the lazy dog. →

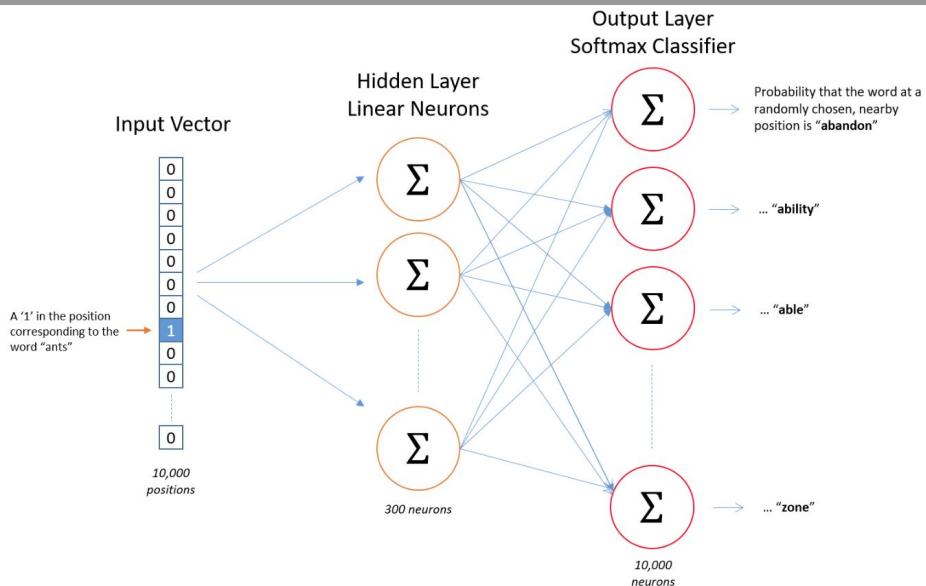
(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

The quick brown fox jumps over the lazy dog. →

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

106

word2vec > model



word2vec > training

Word pairs

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

Input is a one-hot vector

True output is a one-hot vector

Network output is a probability distribution

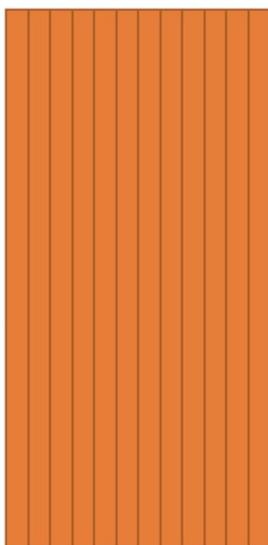
Hidden Layer
Weight Matrix



Word Vector
Lookup Table!

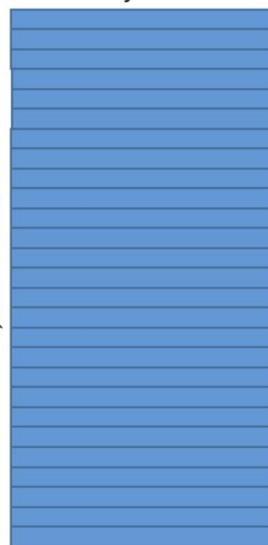
300 neurons

10,000 words



300 features

10,000 words



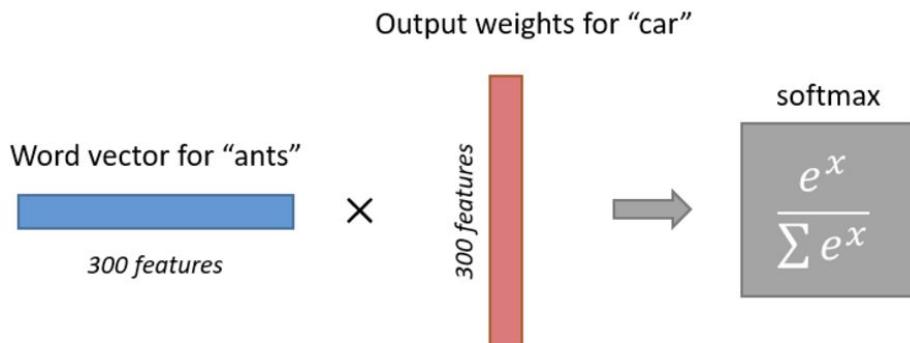
word2vec > hidden layer

The one-hot vector is almost all zeros

What is the effect of that?

$$\begin{bmatrix} 0 & 0 & 0 & \textcolor{green}{1} & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ \textcolor{green}{10} & \textcolor{green}{12} & \textcolor{green}{19} \\ 11 & 18 & 25 \end{bmatrix} = [10 \quad 12 \quad 19]$$

Word embedding > word2vec > output layer



111

Word embedding > word2vec > intuition

Two different words with similar contexts

The model needs to output very similar results

Two word vectors are very similar

Synonyms like intelligent and smart will have similar contexts

Related words like engine and transmission will have similar contexts

Ant and ants will have similar contexts

112

Word embedding > word2vec > modifications

Skip-gram model for word2vec is a huge neural network

With 300 components and 10,000 words

$300 \times 10,000 = 3$ million weights in each layer

Running gradient descent will be slow

Then the authors of word2vec proposed three innovations in their second paper

113

Word embedding > word2vec > innovations

Treating common word pairs as single words

Subsampling frequent words to decrease the number of training samples

Using a technique called *Negative Sampling* to update only a small percentage of model weights

114

Thanks for your attention :)

Jonatas Grosman jgrosman@inf.puc-rio.br
Luiz Schirmer lschirmer@inf.puc-rio.br
William Fernandes wfernandes@inf.puc-rio.br