

TEXT DATA

LARISSA SAYURI FUTINO CASTRO DOS SANTOS

14TH AUGUST 2018

Agenda



About/
Caracterís
ticas

IR

POS tagger

Tradução
Máquina

Wordnet

Conclusion/
Next steps

About

BACHELOR IN STATISTICS,
UNIVERSIDADE DE BRASÍLIA

MASTER IN STATISTICS,
UNIVERSIDADE FEDERAL DE MINAS GERAIS

DOCTORATE IN STATISTICS,
UNIVERSIDADE FEDERAL DE MINAS GERAIS

DOCTORAL SANDWICH,
UNIVERSITY OF WOLVERHAMPTON

2007

2011

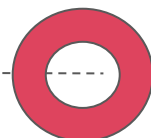
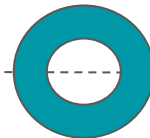
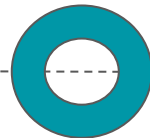
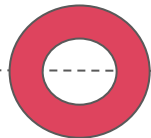
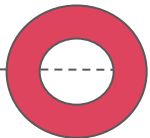
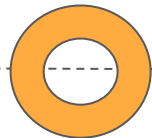
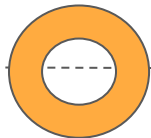
2013

2015

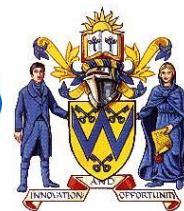
2017

2018

2019



CAPES



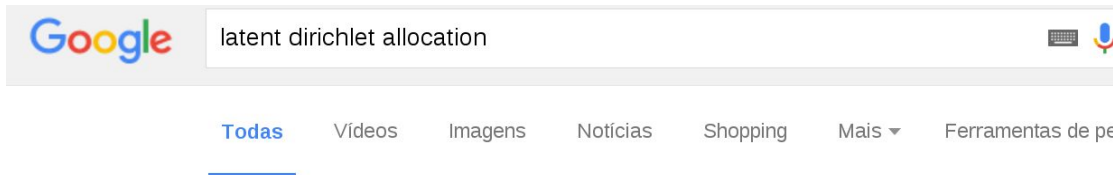
UFMG



UNIVERSIDADE FEDERAL
DE MINAS GERAIS

Characteristics

- ❖ Variabilidade de domínios
- ❖ Aspectos sócio econômicos
- ❖ Palavras são polissêmicas
- ❖ Formas abreviadas
- ❖ Gírias
- ❖ Linguagem informal
- ❖ Erros gramaticais
- ❖ Erros de ortografia
- ❖ Pontuação imprópria





Google  

[Todas](#) [Vídeos](#) [Imagens](#) [Notícias](#) [Shopping](#) [Mais ▾](#) [Ferramentas de pe](#)

Aproximadamente 215.000 resultados (0,27 segundos)

Artigos acadêmicos sobre **latent dirichlet allocation**

Latent dirichlet allocation - Blei - Citado por 14086

Latent dirichlet allocation - Blei - Citado por 262

Online learning for **latent dirichlet allocation** - Hoffman - Citado por 562

Latent Dirichlet allocation - Wikipedia, the free encyclopedia

https://en.wikipedia.org/.../Latent_Dirichlet_allocatio... [Traduzir esta página](#)

In natural language processing, Latent Dirichlet allocation (LDA) is a generative statistical model that allows sets of observations to be explained by unobserved ...

[Dirichlet distribution](#) - [Generative model](#) - [Topic model](#) - [Pachinko allocation](#)

[PDF] Latent Dirichlet Allocation - Computer Science

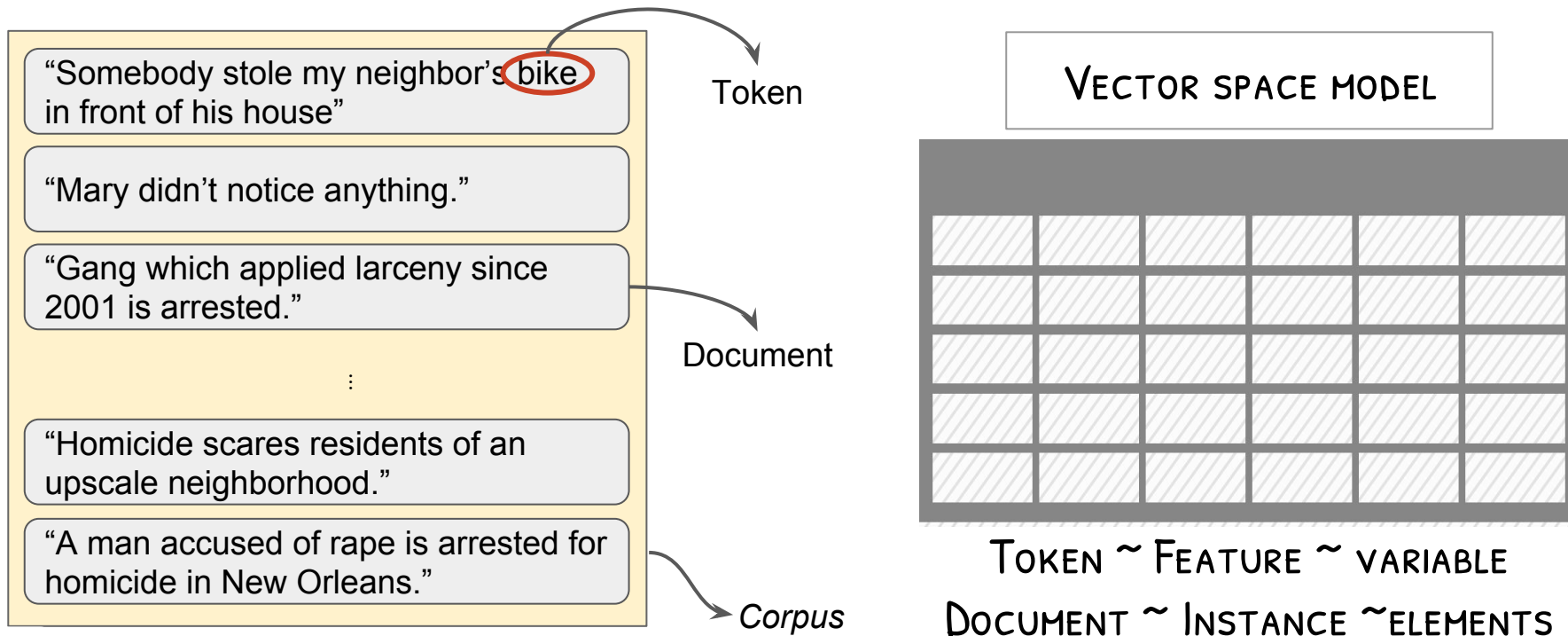
<https://www.cs.princeton.edu/.../BleiNgJordan2003.p...> [Traduzir esta página](#)

de DM Blei - 2003 - Citado por 14086 - [Artigos relacionados](#)

We describe latent Dirichlet allocation (LDA), a generative probabilistic model for collections ... LDA is a three-level hierarchical Bayesian model, in which each.

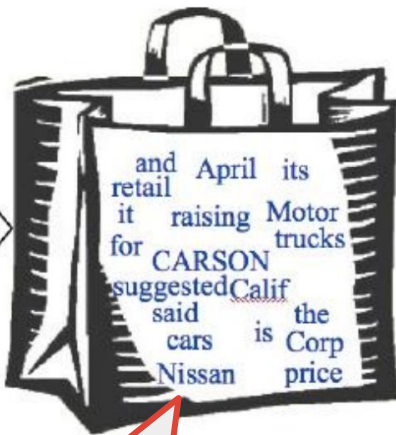


Representação



Modelo Bag of Words

CARSON, Calif., April 3 - Nissan Motor Corp said it is raising the suggested retail price for its cars and trucks.....



ESSE É UM MODELO
MUIIIITO SIMPLISTA...
MAS ÚTIL!

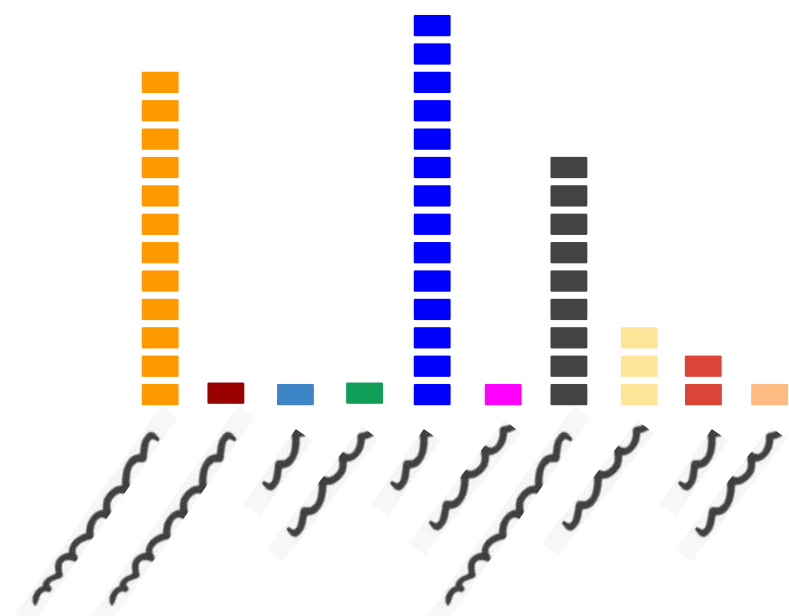
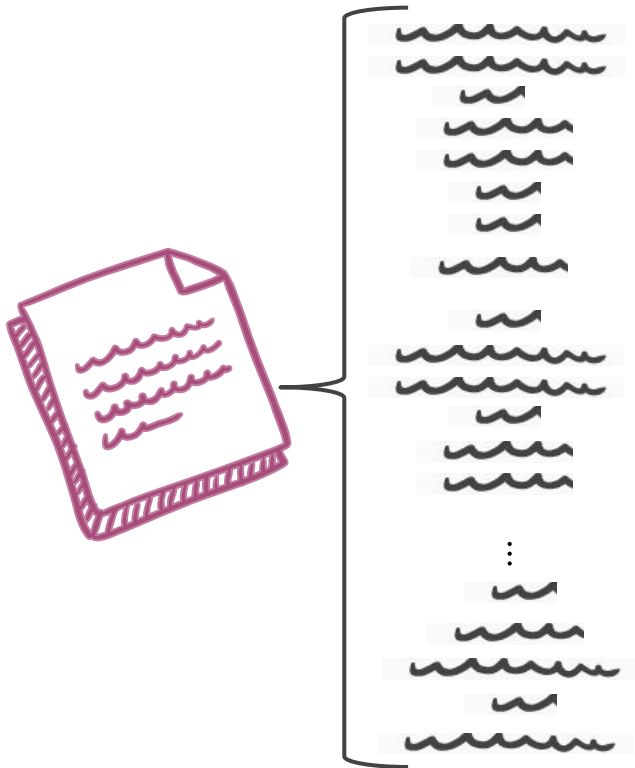
Cada palavra do texto, qualquer que seja a sua posição, é escolhida do dicionário com as mesmas probabilidades e independentemente das demais.

O texto resume-se ao vetor:

$$X = (N_1, N_2, N_3, \dots, N_M),$$

N_i = nº de vezes que a i -ésima palavra do dicionário ocorreu no texto

Modelo Bag of Words - Peso



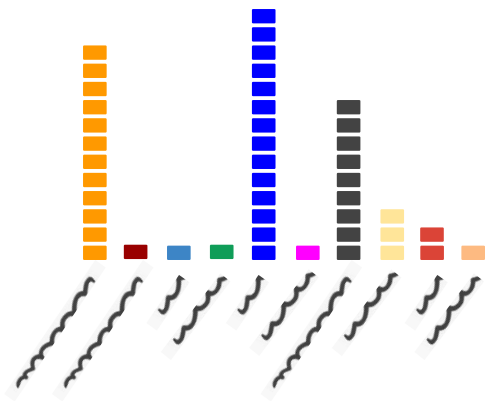
TF: FREQUÊNCIA DO TERMO NO DOCUMENTO

Modelo Bag of Words - Peso



IDF: INVERSO FREQUÊNCIA DO TERMO NO CORPUS

Modelo Bag of Words - Peso

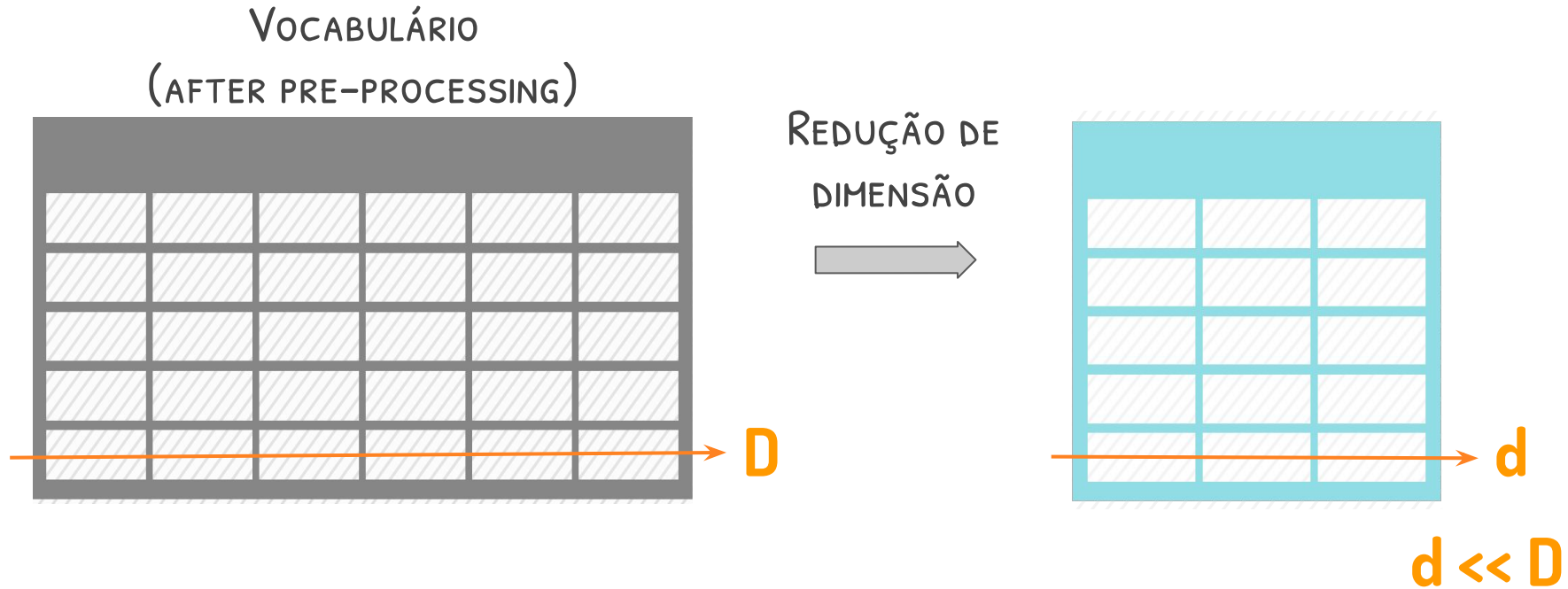


TF: FREQUÊNCIA DO TERMO NO DOCUMENTO

IDF: INVERSO FREQUÊNCIA DO TERMO NO CORPUS



Representação do texto



Vector Space Model

(FORTE) PRESSUPOSTO:

A ORDEM DOS TERMOS NÃO IMPORTA

DOC. REPRESENTADO PELA PRESENÇA DOS SEUS
TERMOS:

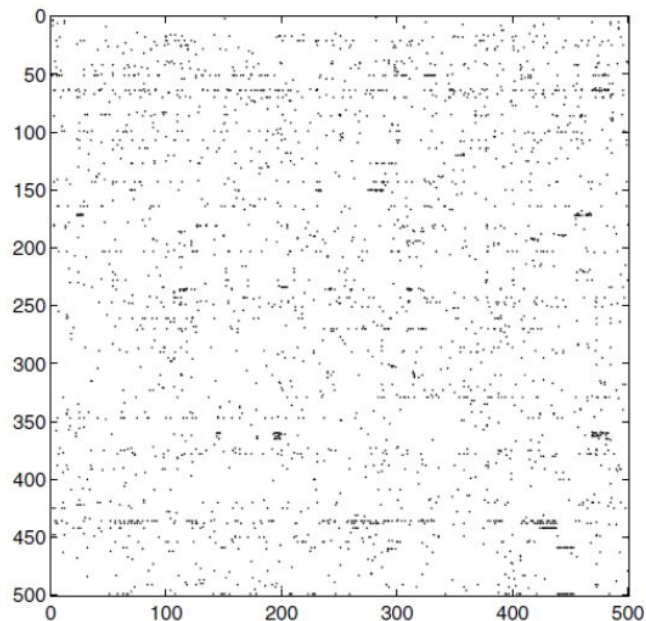
A LEBRE É MAIS VELOZ QUE A TARTARUGA

A TARTARUGA É MAIS VELOZ QUE A LEBRE

CARACTERÍSTICAS:

GRANDE DIMENSÃO

ESPARSIDADE



POS - tagger

POS TAGGING:

“O PROCESSO DE ATRIBUIR UM RÓTULO DE CLASSE GRAMATICAL PARA CADA TERMO EM UM CORPUS.”

Word	Lemma	Annotation	Code
ele	ele	Pron. Pes. 3aPes. Masc. Sing. Indef.	NCMS000
pára	parar	Verbo Princ. Indic. Pres. 3aPes. Sing.	VMIP3S0
todos os dias	todos os dias	Adv. Geral	RG00000
para	para	Prepos. Simples	SP00000
ir	ir	Verbo Principal Infinit.	VMN0000
para	para	Prepos. Simples	SP00000
a	o	Determ. Artigo Fem. Sing. Indefin.	DAIFS00
faculdade	faculdade	Nome Comum Fem. Sing.	NCFS00
.	.	Pontuação	Fd

POS - tagger

MODELO ESTATÍSTICO:

“UM MODELO PARA UMA SEQUÊNCIA DE PALAVRAS”



PREPOSIÇÕES
CONJUNÇÕES
ARTIGOS
PRONOMES
NUMERAIS



SUBSTANTIVOS
VERBOS
ADJETIVOS
ADVERBIOS

POS - tagger



THE



DOG

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO



ATE

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO



THE



BONE

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO

POS - tagger



THE



DOG

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO

$P(\text{Substantivo}|\text{Artigo}) = 0.64$

$P(\text{Verbo}|\text{Artigo}) = 0.07$

$P(\text{Adjetivo}|\text{Artigo}) = 0.29$

$P(\text{Advérbio}|\text{Artigo}) = 0$

POS - tagger



DOG

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO



ATE

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO

$P(\text{Substantivo}|\text{Substantivo}) = 0.08$

$P(\text{Verbo}|\text{Substantivo}) = 0.17$

$P(\text{Adjetivo}|\text{Substantivo}) = 0$

$P(\text{Advérbio}|\text{Substantivo}) = 0$

$P(\text{Substantivo}|\text{Adjetivo}) = 0.72$

$P(\text{Verbo}|\text{Adjetivo}) = 0$

$P(\text{Adjetivo}|\text{Adjetivo}) = 0.03$

$P(\text{Advérbio}|\text{Adjetivo}) = 0$

POS - tagger



ATE

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO



THE

$$P(\text{Artigo}|\text{Substantivo}) = 0.08$$

$$P(\text{Artigo}|\text{Verbo}) = 0.17$$

$$P(\text{Artigo}|\text{Adjetivo}) = 0$$

$$P(\text{Artigo}|\text{Advérbio}) = 0$$

POS - tagger



THE

BONE

SUBSTANTIVO

VERBO

ADJETIVO

ADVERBIO

$$P(\text{Substantivo}|\text{Artigo}) = 0.08$$

$$P(\text{Verbo}|\text{Artigo}) = 0.17$$

$$P(\text{Adjetivo}|\text{Artigo}) = 0$$

$$P(\text{Advérbio}|\text{Artigo}) = 0$$

THE DOG ATE THE BONE

P(SUBSTANTIVO|SUBSTANTIVO)

P(VERBO|SUBSTANTIVO)

P(ADJETIVO|SUBSTANTIVO)

P(ADVERBIO|SUBSTANTIVO)

P(SUBSTANTIVO|ARTIGO)

P(VERBO|ARTIGO)

P(ADJETIVO|ARTIGO)

P(ADVERBIO|ARTIGO)

P(SUBSTANTIVO|VERBO)

P(VERBO|VERBO)

P(ADJETIVO|VERBO)

P(ADVERBIO|VERBO)

P(SUBSTANTIVO|ADJETIVO)

P(VERBO|ADJETIVO)

P(ADJETIVO|ADJETIVO)

P(ADVERBIO|ADJETIVO)

P(SUBSTANTIVO|ADVERBIO)

P(VERBO|ADVERBIO)

P(ADJETIVO|ADVERBIO)

P(ADVERBIO|ADVERBIO)

P(ARTIGO|SUBSTANTIVO)

P(ARTIGO|VERBO)

P(ARTIGO|ADJETIVO)

P(ARTIGO|ADVERBIO)

P(SUBSTANTIVO|ARTIGO)

P(VERBO|ARTIGO)

P(ADJETIVO|ARTIGO)

P(ADVERBIO|ARTIGO)

P(ARTIGO)

THE DOG ATE THE BONE

P(ARTIGO)	P(SUBSTANTIVO ARTIGO)	P(SUBSTANTIVO SUBSTANTIVO)	P(ARTIGO SUBSTANTIVO)	P(SUBSTANTIVO ARTIGO)
	P(VERBO ARTIGO)	P(VERBO SUBSTANTIVO)	P(ARTIGO VERBO)	P(VERBO ARTIGO)
	P(ADJETIVO ARTIGO)	P(ADJETIVO SUBSTANTIVO)	P(ARTIGO ADJETIVO)	P(ADJETIVO ARTIGO)
	P(ADVERBIO ARTIGO)	P(ADVERBIO SUBSTANTIVO)	P(ARTIGO ADVERBIO)	P(ADVERBIO ARTIGO)
		P(SUBSTANTIVO VERBO)		
		P(VERBO VERBO)		
		P(ADJETIVO VERBO)		
		P(ADVERBIO VERBO)		
		P(SUBSTANTIVO ADJETIVO)		
		P(VERBO ADJETIVO)		
		P(ADJETIVO ADJETIVO)		
		P(ADVERBIO ADJETIVO)		
		P(SUBSTANTIVO ADVERBIO)		
		P(VERBO ADVERBIO)		
		P(ADJETIVO ADVERBIO)		
		P(ADVERBIO ADVERBIO)		

A
SEQUÊNCIA DE
ANOTAÇÕES MAIS
PLAUSÍVEL É A DE
MAIOR
PROBABILIDADE

POS - tagger

PROBABILIDADES ESTIMADAS COMO FREQUÊNCIAS
RELATIVAS

EXIGE O EMPREGO DE UM CORPUS REPRESENTATIVO DA
LÍNGUA

A ORDEM DE DEPENDÊNCIA DA CADEIA (N-GRAM) É UM
PARÂMETRO A SER CONSIDERADO.

GRANDES DIMENSÕES SÃO RECONHECIDAMENTE POUCO
REPRESENTATIVAS DA LÍNGUA



Aplicação - Tradução de Máquina

português inglês espanhol Detectar idioma ▾



inglês português espanhol ▾

Traduzir

Natural language processing (NLP) is a field of computer science, artificial intelligence, and computational linguistics concerned with the interactions between computers and human (natural) languages. As such, NLP is related to the area of human-computer interaction. Many challenges in NLP involve: natural language understanding, enabling computers to derive meaning from human or natural language input; and others involve natural language generation.



processamento de linguagem natural (NLP) é um campo da ciência da computação, inteligência artificial e lingüística computacional preocupados com as interações entre computadores e linguagens humanas (naturais). Como tal, PNL está relacionada com a área de interação humano-computador. Muitos desafios em PNL envolvem: compreensão da linguagem natural, permitindo computadores para extrair significado de entrada linguagem humana ou natural; e outros envolvem geração de linguagem natural.



Sugerir uma edição



Aplicação - Tradução de Máquina

FONTE:

COMPLETO ENTENDIMENTO



ALVO:

CONHECIMENTO SOFISTICADO,
POÉTICO E CRIATIVO

Aplicação - Tradução de Máquina

“DIFERENÇAS SINTÁTICAS ESTÃO RELACIONADAS COM DIFERENÇAS SEMÂNTICAS EM COMO AS LÍNGUAS MAPEIAM CONCEITOS EM PALAVRAS.” (JURAFSKY)

English	<i>brother</i>	Japanese	<i>otooto</i> (younger)
		Japanese	<i>oniisan</i> (older)
		Mandarin	<i>gege</i> (older)
		Mandarin	<i>didi</i> (older)
English	<i>wall</i>	German	<i>Wand</i> (inside)
		German	<i>Mauer</i> (outside)
English	<i>know</i>	French	<i>connaître</i> (be acquainted with)
		French	<i>savoir</i> (know a proposition)
English	<i>they</i>	French	<i>ils</i> (masculine)
		French	<i>elles</i> (feminine)
German	<i>berg</i>	English	<i>hill</i>
		English	<i>mountain</i>
Mandarin	<i>tā</i>	English	<i>he, she, or it</i>

Figure 21.1 Differences in specificity.

“ALGUMAS LÍNGUAS DIVIDEM UM DADO CONCEITO EM MAIS DETALHES DO QUE OUTRAS.”
(JURAFSKY)

Tradução de Máquina

SVO:

ALEMÃO

FRANCÊS

INGLÊS

SOV:

HINDU

JAPONÊS

VSO:

IRLANDÊS

ÁRABE CLÁSSICO

HEBRAICO BÍBLICO

IDEIA GERAL:

ALTERAR A ESTRUTURA DA LÍNGUA FONTE PARA QUE ESTEJA EM CONFORMIDADE COM AS REGRAS DA LÍNGUA ALVO.

Tradução de Máquina

IDEIA GERAL:

ALTERAR A ESTRUTURA DA LÍNGUA FONTE PARA QUE ESTEJA EM CONFORMIDADE COM AS REGRAS DA LÍNGUA ALVO.

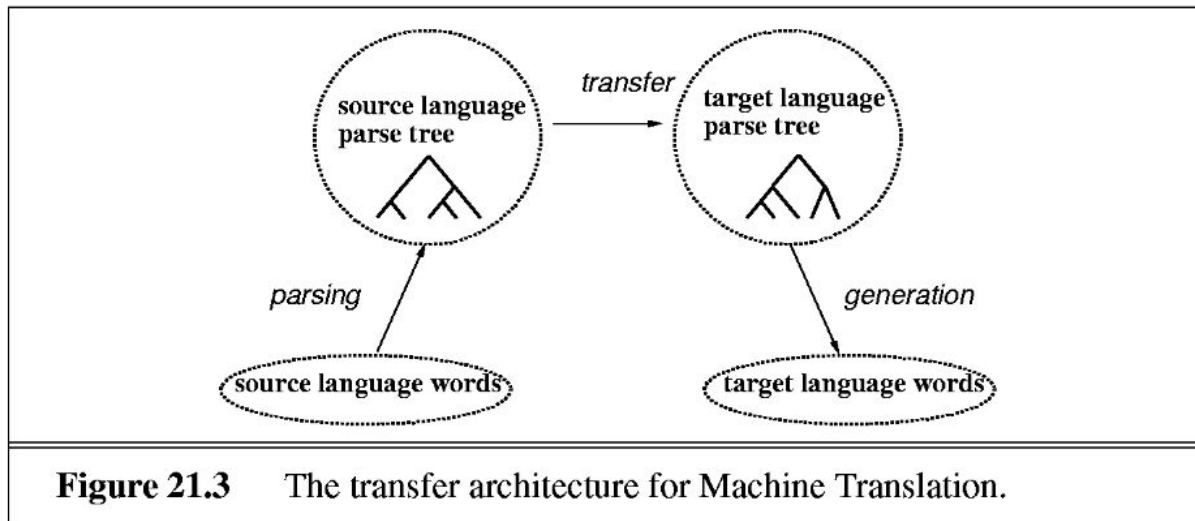
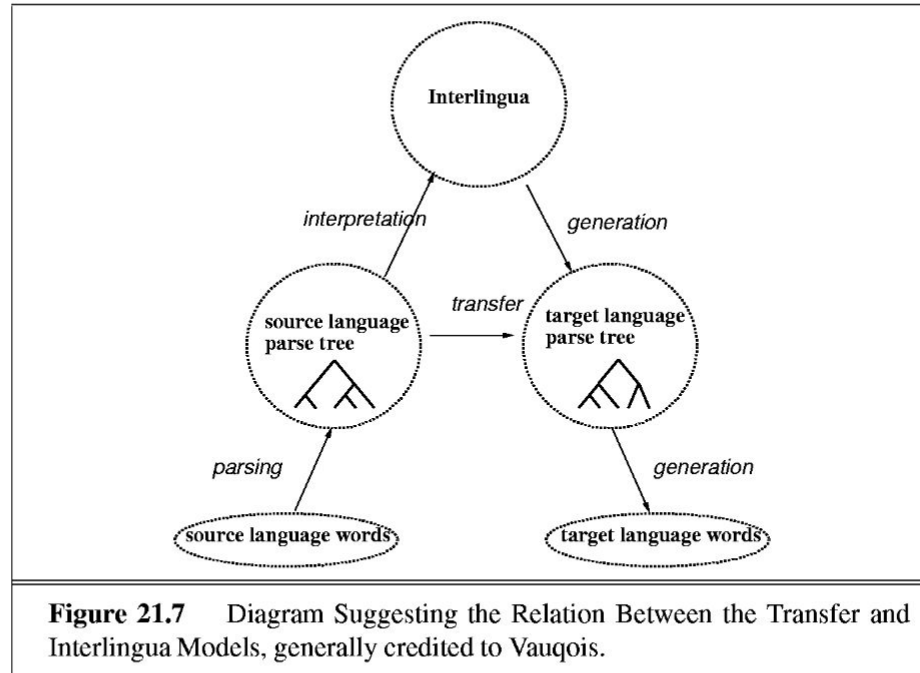




Figure 21.4 A simple transformation that reorders adjectives and nouns

CONJUNTO DISTINTO DE REGRAS DE REFERÊNCIA POR PAR DE LÍNGUA

Tradução de Máquina - Interlingua



NECESSIDADE DE UMA ONTOLOGIA: INVENTÁRIO DE CONCEITOS E RELAÇÕES

Tradução de Máquina - Tradução Direta

English to French:	
1.	NP \rightarrow Adjective ₁ Noun ₂ \Rightarrow NP \rightarrow Noun ₂ Adjective ₁
Japanese to English:	
2.	Existential-There-Sentence \rightarrow There ₁ Verb ₂ NP ₃ Postnominal ₄ \Rightarrow Sentence \rightarrow (NP \rightarrow NP ₃ Relative-Clause ₄) Verb ₂
3.	NP \rightarrow NP ₁ Relative Clause ₂ \Rightarrow NP \rightarrow Relative-Clause ₂ NP ₁
Figure 21.5 An informal description of some transformations.	

TENDÊNCIA A SEREM CONSERVADORES

Input:	watashihatsukuenouenopenwojonniageta.
After stage 1:	watashi ha tsukue no ue no pen wo jon ni ageru PAST.
After stage 2:	I ha desk no ue no pen wo John ni give PAST.
After stage 3:	I ha pen on desk wo John to give PAST.
After stage 4:	I give PAST pen on desk John to.
After stage 5:	I give PAST the pen on the desk to John.
After stage 6:	I gave the pen on the desk to John.

Figure 21.9 An Example of Processing in a Direct System

Tradução de Máquina - Tradução Estatística

TRADUZIR: ADONA ROI

THE LORD IS MY SHEPHERD

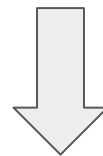
THE LORD WILL LOOK AFTER ME

THE LORD IS FOR ME LIKE SOMEBODY WHO
LOOKS AFTER ANIMALS WITH COTTON-LIKE
HAIR

CONCLUSÃO: SEMPRE HÁ UMA METÁFORA,
CONSTRUÇÃO, PALAVRA OU TEMPO VERBAL
SEM UM PARALELO NA LÍNGUA ALVO

TRANSFER, INTERLINGUA, DIRECT MODELS:
QUAL REPRESENTAÇÃO EMPREGAR
QUAIS PASSOS SEGUIR PARA TRADUZIR

FOCAR NO RESULTADO E NÃO NO PROCESSO



TRADUÇÃO ESTATÍSTICA

Tradução de Máquina

O QUE DEFINE UMA BOA TRADUÇÃO?

FIDELIDADE À LÍNGUA
FONTE



NATURALIDADE/FLUÊNCIA
PARA A LÍNGUA ALVO

Tradução de Máquina

O QUE DEFINE UMA BOA TRADUÇÃO?

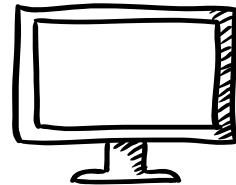
FIDELIDADE À LÍNGUA
FONTE



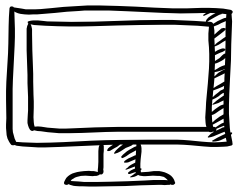
NATURALIDADE/FLUÊNCIA
PARA A LÍNGUA ALVO



ENCONTRAM O TOLERÁVEL PARA OS
DOIS CRITÉRIOS



?



FUNÇÃO QUE REPRESENTA A IMPORTÂNCIA DA FIDELIDADE
E DA FLUÊNCIA DA TRADUÇÃO

$$\text{melhorTraducao}(\hat{T}) = \operatorname{argmax}_T \text{fluencia}(T) * \text{fidelidade}(T, S)$$

$$\text{melhorTraducao}(\hat{T}) = \operatorname{argmax}_T P(T) * P(T, S)$$

Tradução de Máquina

$P(T) \sim$ FLUÊNCIA:

PODE SER ESTIMADA COM MODELOS DO TIPO N-GRAM

QUAL TRADUÇÃO É MAIS
FLUENTE?

1. THE CAR WAS FAST
2. THE FAST CAR WAS
3. FAST THE WAS CAR

2-GRAM (RELIABLE CORPUS):

$$P(\text{CAR}|\text{THE}) = A$$

$$P(\text{WAS}|\text{CAR}) = B$$

$$P(\text{FAST}|\text{WAS}) = C$$

$$P(\text{FAST}|\text{THE}) = D$$

$$P(\text{CAR}|\text{FAST}) = E$$

$$P(\text{THE}|\text{FAST}) = F$$

$$P(\text{WAS}|\text{THE}) = G$$

$$P(\text{CAR}|\text{WAS}) = H$$

1. THE CAR WAS FAST 

2. THE FAST CAR WAS

3. FAST THE WAS CAR

Tradução de Máquina - P(T) ~ Fluência

DAS AUTO WAS SCHNELL

1. THE CAR WAS FAST
2. THE CAR WAS RED
3. THE BICYCLE WAS BLUE

BILINGUAL PROBABILISTIC DICTIONARY

Alemão	Inglês	Prob
Das	the	1
auto	car	0.9
auto	bicycle	0.1
war	was	1
schnell	fast	0.9
schnell	red	0.05
schnell	blue	0.05

ALINHAMENTO A NÍVEL DE PALAVRA

1. THE(1) CAR(.9) WAS(1) FAST(.9)
2. THE(1) CAR(.9) WAS(1) RED(.05)
3. THE(1) BICYCLE(.1) WAS(1) BLUE(.05)

Semantic tagger

SEMANTIC TAGS:

“SETS OF WORDS RELATED
AT SOME LEVEL OF
GENERALITY WITH
THE SAME MENTAL
CONCEPT.”

A general and abstract terms	B the body and the individual	C arts and crafts	E emotion
F food and farming	G government and public	H architecture, housing and the home	I money and commerce in industry
K entertainment, sports and games	L life and living things	M movement, location, travel and transport	N numbers and measurement
O substances, materials, objects and equipment	P education	Q language and communication	S social actions, states and processes
T Time	W world and environment	X psychological actions, states and processes	Y science and technology
Z names and grammar			

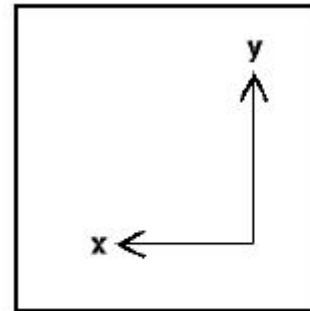
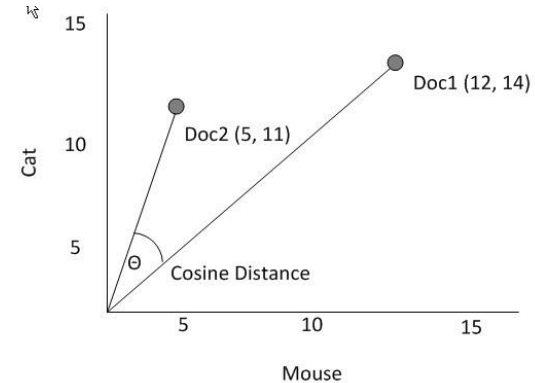
Semantic tagger

QUAL A SIMILARIDADE ENTRE:

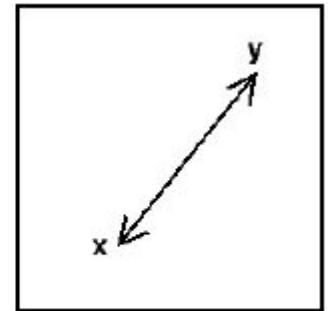
MENINA E MENINO?

ABELHA E COLMÉIA?

BOLA E FUTEBOL?

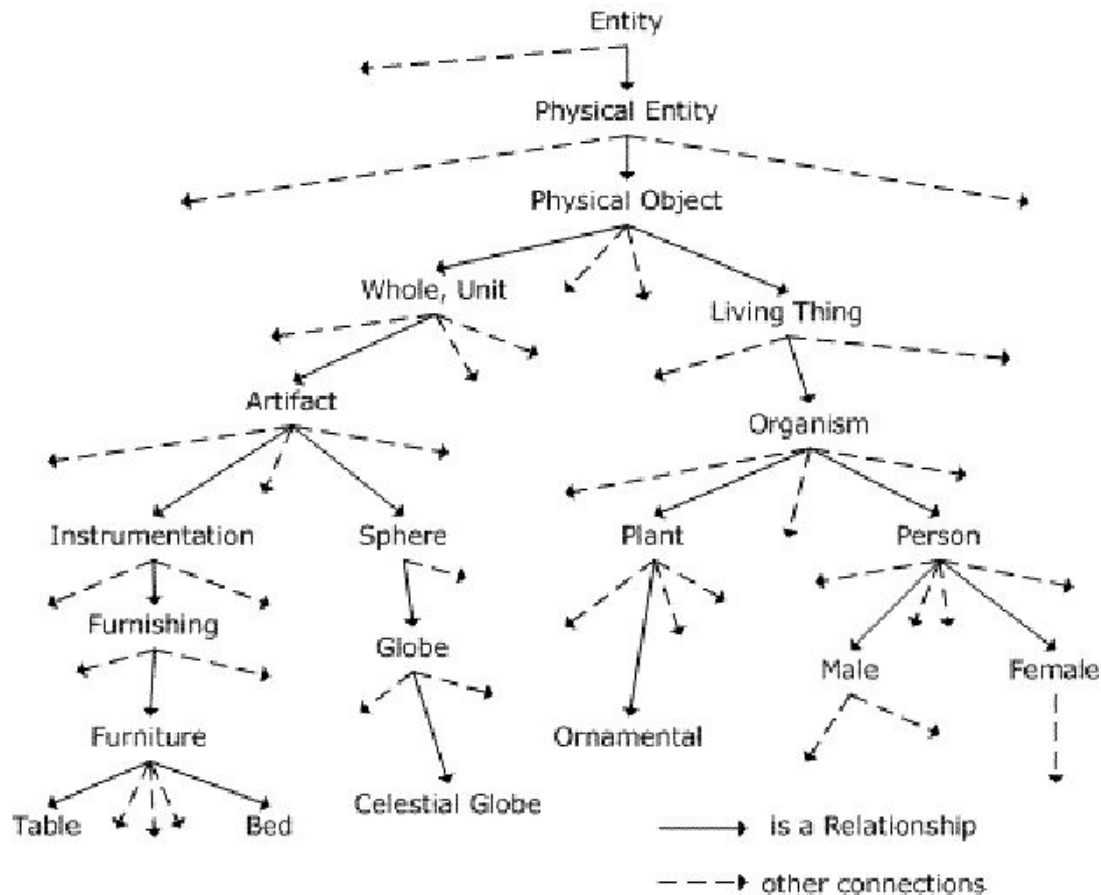


Manhattan

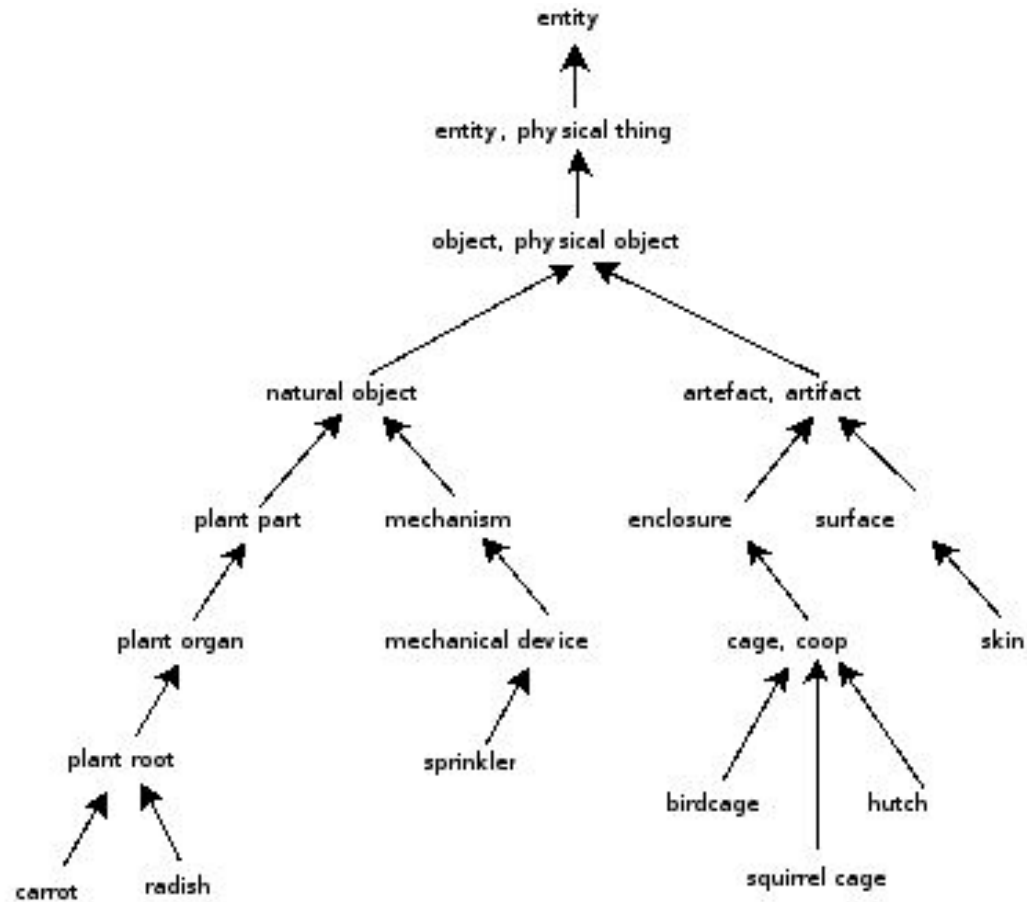


Euclidean

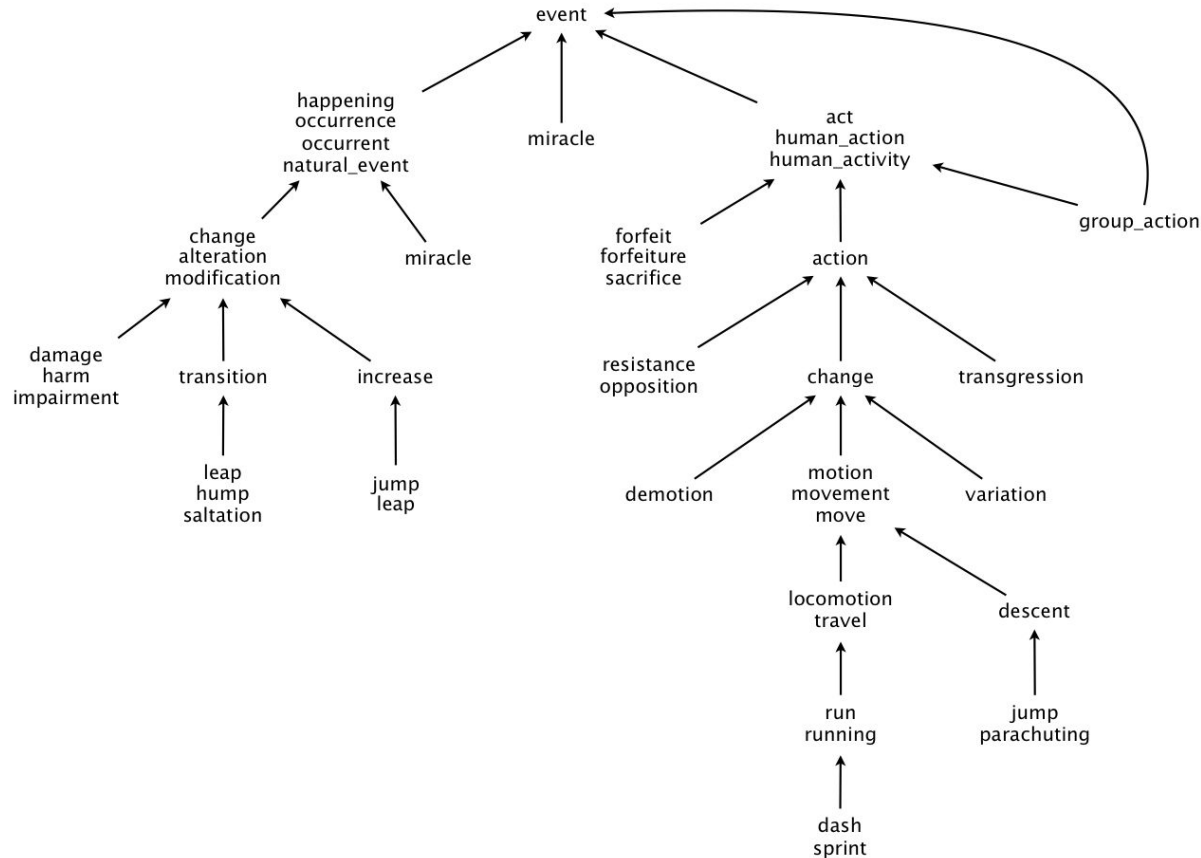
Wordnet/Semantics



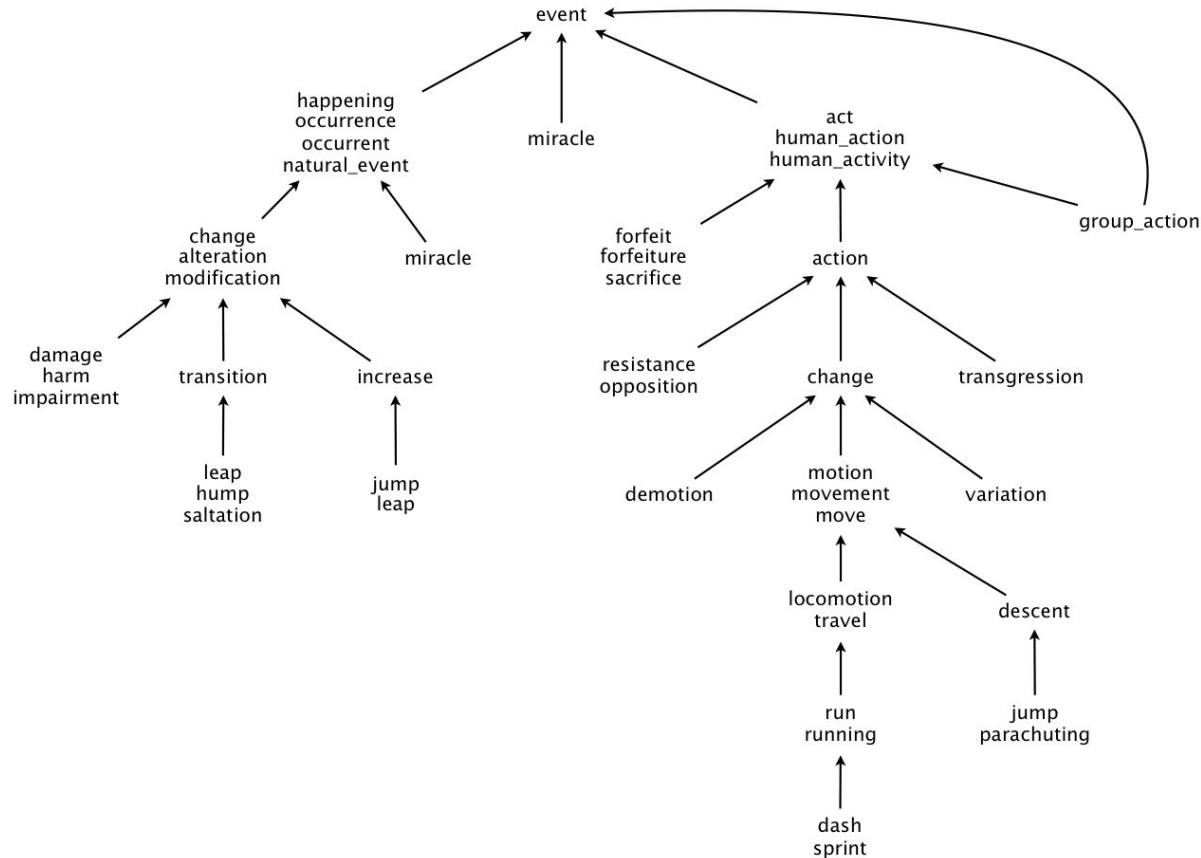
Wordnet/Semantics



Wordnet/Semantics

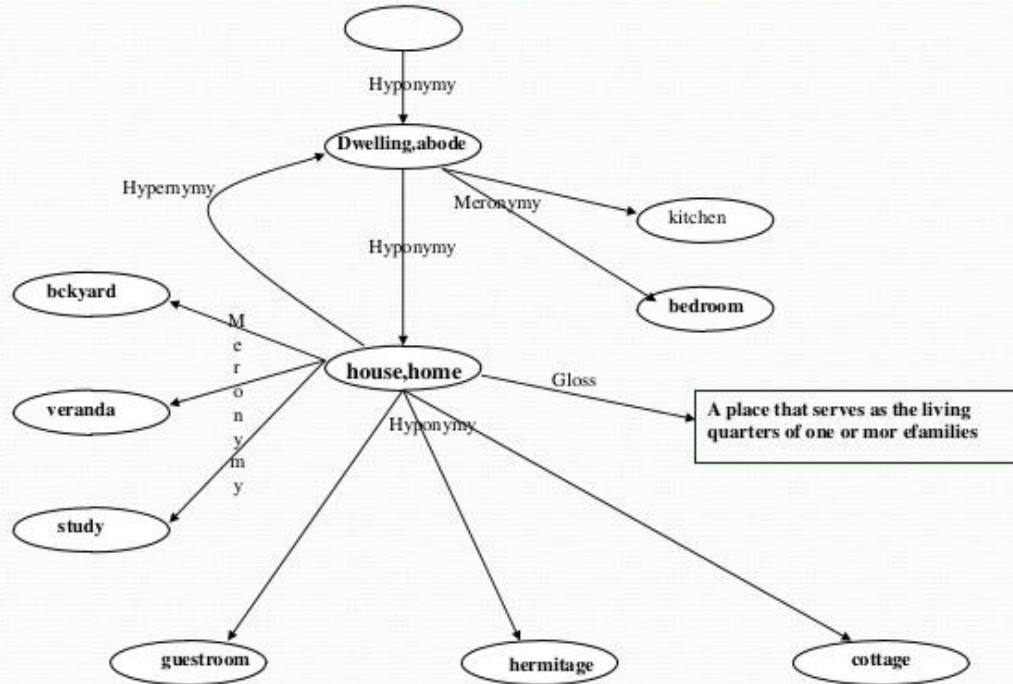


Wordnet/Semantics



Wordnet/Semantics

WordNet Sub-Graph (English)



Wordnet/Semantics

```
In [32]: # we can calculate the similarity by looking at the shortest path  
right.path_similarity(minke)
```

```
Out[32]: 0.25
```

```
In [33]: right.path_similarity(orca)
```

```
Out[33]: 0.16666666666666666
```


```
In [34]: # Leacock-Chodorow Similarity takes into consideration the distance and  
# the depth of the taxonomy  
right.lch_similarity(minke)
```

```
Out[34]: 2.2512917986064953
```

```
In [35]: right.lch_similarity(orca)
```

```
Out[35]: 1.845826690498331
```

Wordnet/Semantics

- DICIONÁRIO DIGITAL (MUITO COMPLETO)
- FEITO POR LEXICÓGRAFOS
- DISPONÍVEL PARA UM NÚMERO LIMITADO DE LÍNGUAS 

Conclusions and Future Work

Uma ferramenta de NLP é input para outra ferramenta de NLP. Os erros só aumentam.

Ferramentas de NLP crescentemente incorporam redes Neurais

Ferramentas de NLP são desenvolvidas com especificidade para domínios.

Ferramentas de NLP dependem (demais) da base de dados. Muita relevância para quem constrói os Corpus.

Ferramentas de NLP são desenvolvidas prioritariamente no mercado, não na academia.