

Análise de Sentimento em Textos Curtos Baseada em Processamento de Linguagem Natural e Aprendizado de Máquina

Breno Arosa

Natanael N. de Moura Junior
Luiz Pereira Calôba
Felipe Fink Grael

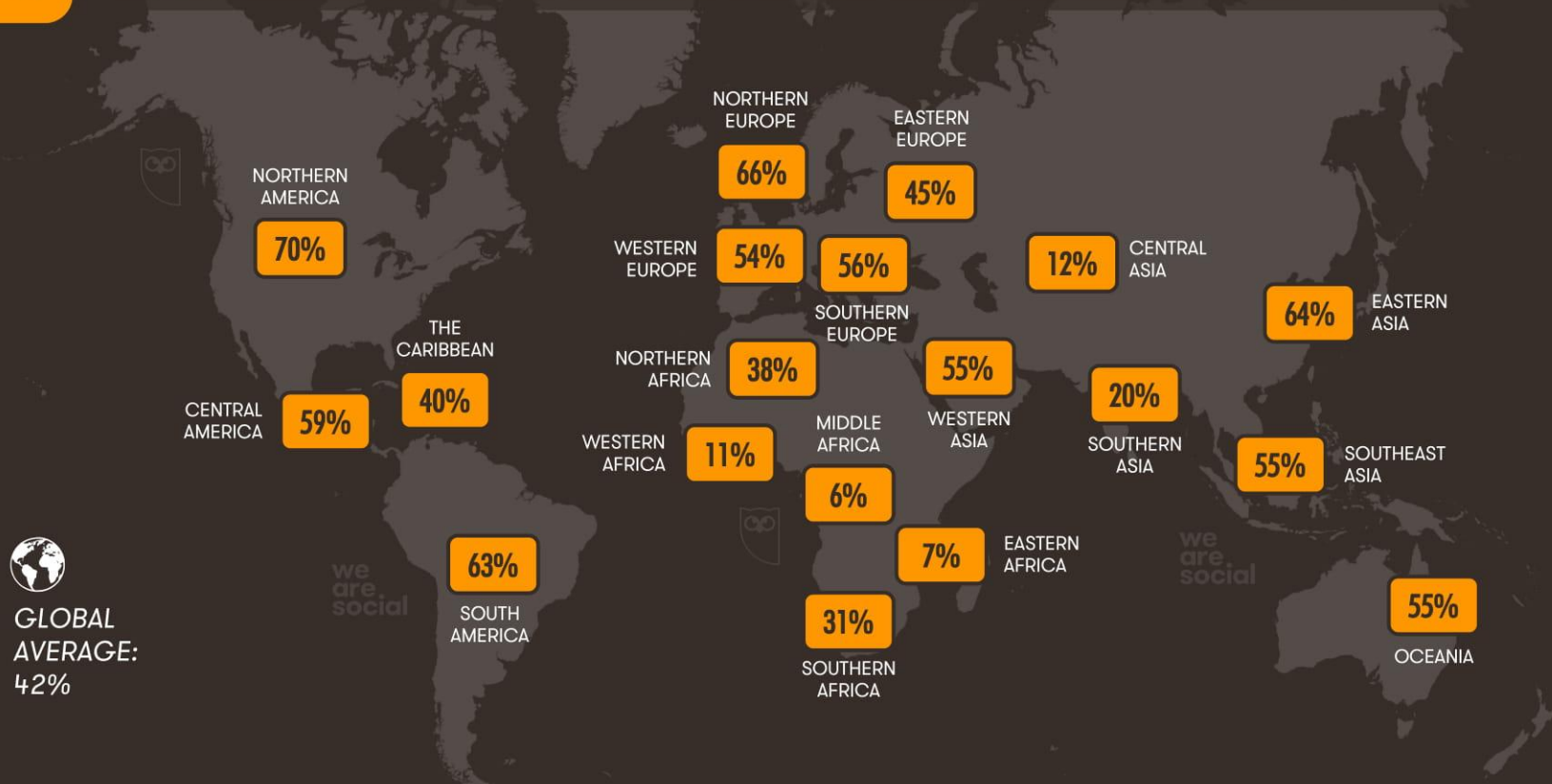
Agenda

- ▷ Motivação
- ▷ Processamento de Linguagem Natural
- ▷ Classificadores
- ▷ Método
- ▷ Resultados
- ▷ Conclusão

JAN
2018

SOCIAL MEDIA PENETRATION BY REGION

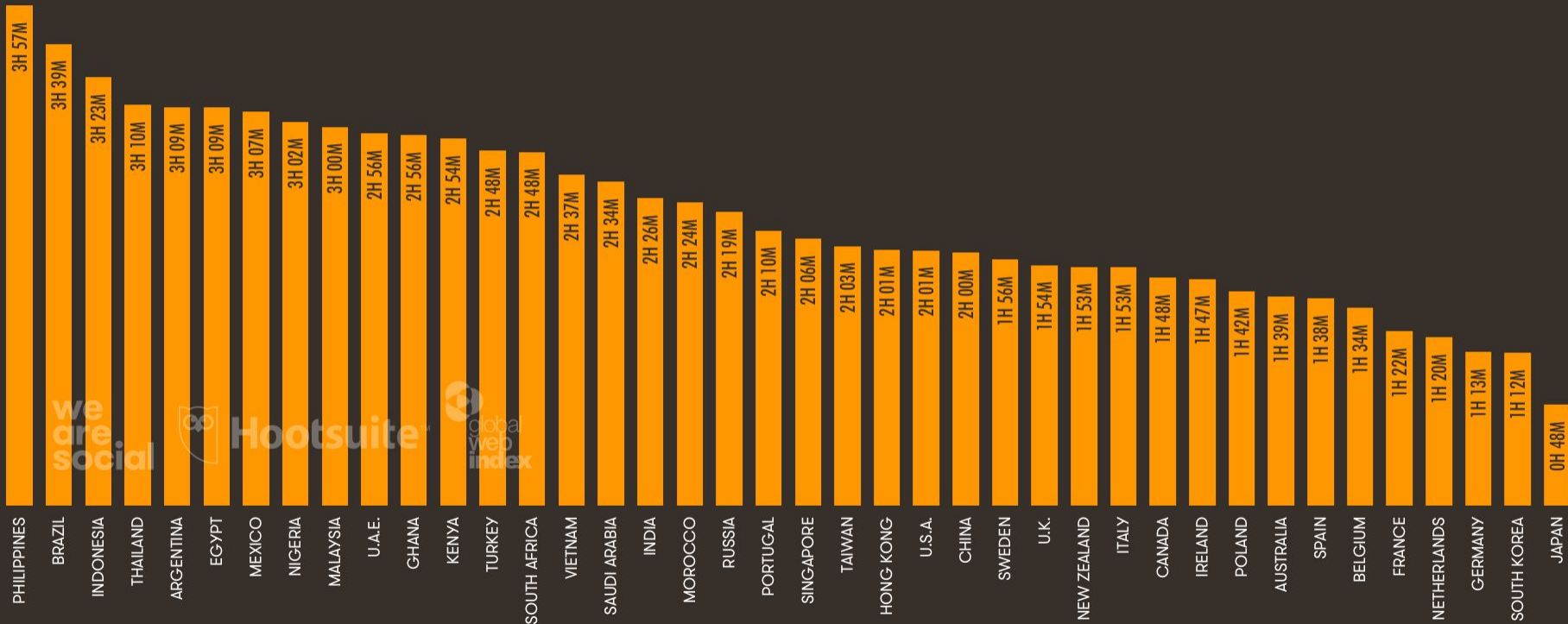
TOTAL ACTIVE ACCOUNTS ON THE MOST ACTIVE SOCIAL NETWORK IN EACH COUNTRY, COMPARED TO POPULATION



JAN
2018

TIME SPENT ON SOCIAL MEDIA

AVERAGE NUMBER OF HOURS THAT SOCIAL MEDIA USERS SPEND USING SOCIAL MEDIA EACH DAY VIA ANY DEVICE [SURVEY BASED]

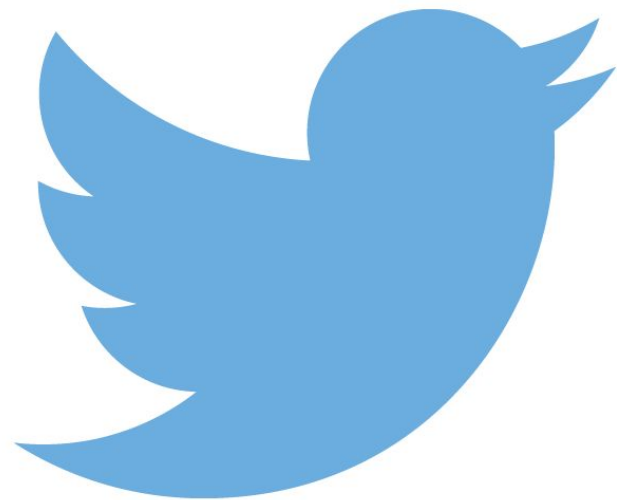


310M

Milhões de usuário ativos

500M

Tweets por dia



**Donald J. Trump** ✓

@realDonaldTrump

Follow



Just heard Foreign Minister of North Korea speak at U.N. If he echoes thoughts of Little Rocket Man, they won't be around much longer!

**Louis Tomlinson** ✓

@Louis_Tomlinson

Follow



Thoughts go out to everyone in Paris .
[#prayforparis](#)



Wil
@Wil_Joaquim



Follow

#ToNumaFaseQ não sei o q esta acontecenu



hashtag



abreviação



neologismo

Mídias

Data



||||

@secondplease

Follow



find yourself somebody who look at you the way Jennifer Lawrence looks at her glass of wine at the #oscars 🏆



2:03 AM - 5 Mar 2018

74 Retweets 182 Likes



Re-publicações e
Curtidas

Análise de sentimento

Análise de sentimento pode ser considerada como a extração de **polaridade** de mensagem. Esta consiste em classificar a mensagem entre **positiva**, neutra e **negativa**.

Polaridade de mensagem



Anne

@aeneenee



Following

Macarrão de arroz is the new miojo



luscas

@luscas



Follow

semana de provas eu queria ser dessas pessoa
q fica aaa mas eu so consigo ficar
AAAAAAAAAAAAAAAAAAAA o tempo todo



Camila Martins

@araujocami



Following

Tô igualmente fascinada e enojada



truetretha

@truetha



Following

Recomendo chegar pra dar aula e terem
mudado o horário sem te avisar



**Multiplicidade
de idiomas**



**Sentimento
definido pela
forma**



Ambiguidade



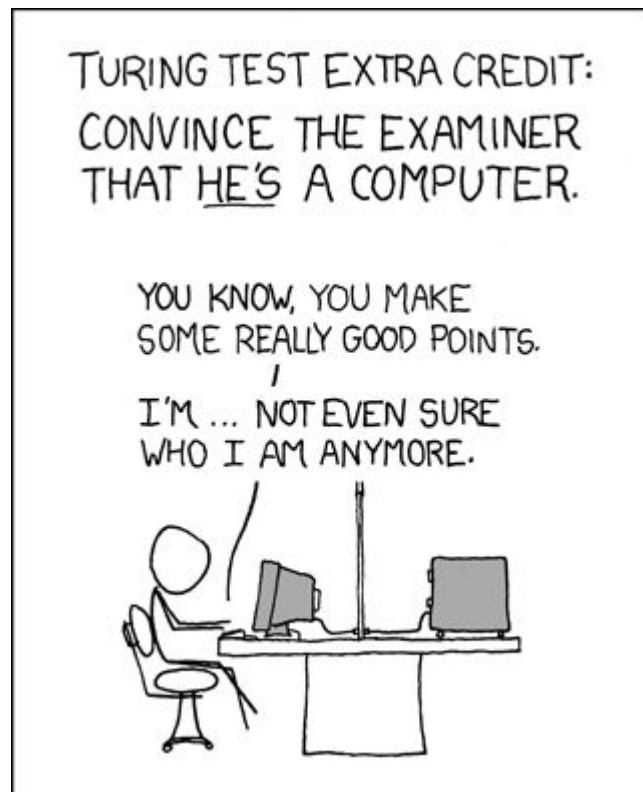
Ironia

Objetivos

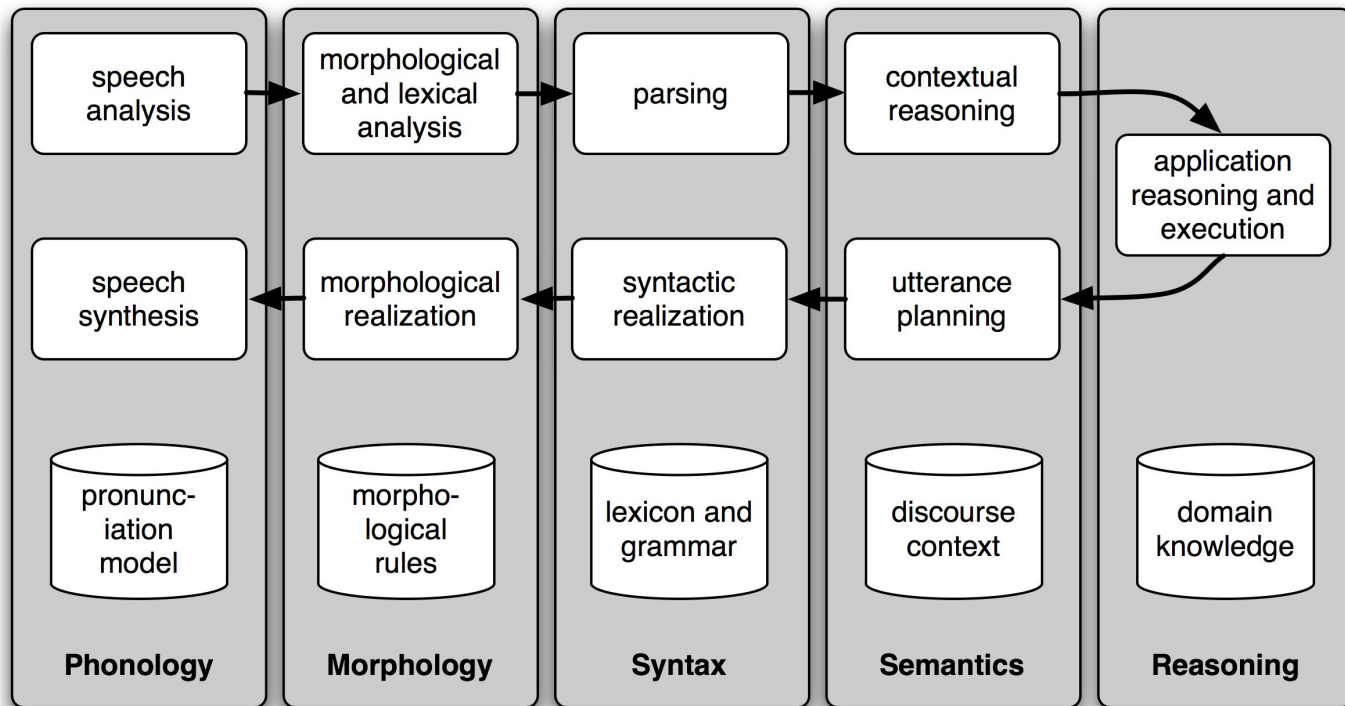
- ▶ Método de produção de classificadores de análise de sentimento.
- ▶ Avaliar técnicas de Deep Learning.

Processamento de linguagem natural

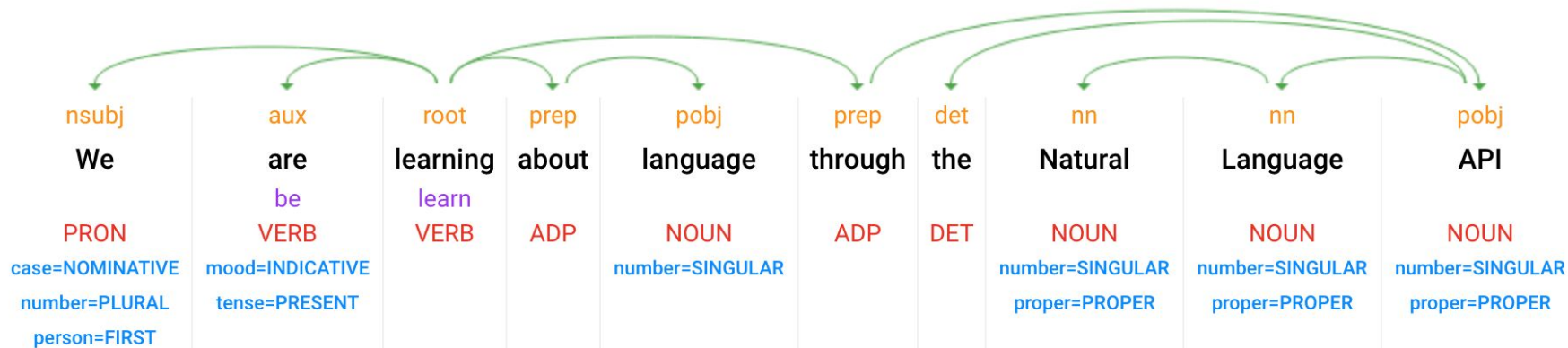
NLP é o ramo da inteligência artificial que trata da **análise**, **compreensão** e **geração** de linguagem que os seres humanos usam para comunicação em contextos escritos e/ou falados.



Níveis de Processamento



☒ Dependency
 ☒ Parse Label
 ☒ Part of Speech
 ☒ Lemma
 ☒ Morphology



Reconhecimento de Entidade

CODEPEN

But **Google** **ORG** is starting from behind. The company made a late push into hardware, and **Apple** **ORG** 's **Siri** **PRODUCT**, available on **iPhones** **PRODUCT**, and **Amazon** **ORG** 's **Alexa** **PRODUCT** software, which runs on its **Echo** **PRODUCT** and **Dot** **PRODUCT** devices, have clear leads in consumer adoption.

The screenshot shows a Gmail inbox interface. At the top is the Google logo, a search bar, and the user's name 'Jess Hardy' with a notification badge '4'. Below the header is a 'Gmail' label and navigation buttons for 'Compose', 'Refresh', and 'More'. The inbox is categorized into 'Primary', 'Social', 'Promotions', and 'Updates'. The 'Social' tab is selected, showing three new notifications. The left sidebar includes a 'COMPOSE' button, 'Inbox (7)', 'Starred', 'Drafts', 'Sent Mail', and a contact card for 'Jenny Kang'.

Google [Search Bar] Jess Hardy 4

Gmail [Compose] [Refresh] [More] 1-100 of 2,067

COMPOSE

Inbox (7)

Starred

Drafts

Sent Mail

[Contact Card: Jenny Kang]

Primary **Social** 3 new Google+, YouTube, Emi... **Promotions** 2 new Google Offers, Zagat **Updates** 2 new Shoehop, Blitz Air +

<input type="checkbox"/>	★	Google+	You were tagged in 3 photos on Google+ - Google+ You were tagged in three photos in an album title
<input type="checkbox"/>	★	YouTube	LauraBlack just uploaded a video. - Jess, have you seen the video LauraBlack uploaded...
<input type="checkbox"/>	★	Emily Million (Google+)	[Knitting Club] Are we knitting tonight? - [Knitting Club] Are we knitting tonight?
<input type="checkbox"/>	★	Sean Smith (Google+)	Photos of the new pup - Sean Smith shared an album with you. View album be thoughtful about who
<input type="checkbox"/>	★	Google+	Kate Baynham shared a post with you - Follow and share with Kate by adding her to a circle. Don't know
<input type="checkbox"/>	★	Google+	Danielle Hoodhood added you on Google+ - Follow and share with Danielle by adding her to a circle.

Julian Assange

UK prosecutors admit destroying key emails in Julian Assange case

Correspondence between CPS and its Swedish counterparts about WikiLeaks founder deleted after lawyer retired in 2014

↑ [-] [autotldr](#) [score hidden] 11 hours ago

↓ This is the best tl;dr I could make, [original](#) reduced by 91%. (I'm a bot)

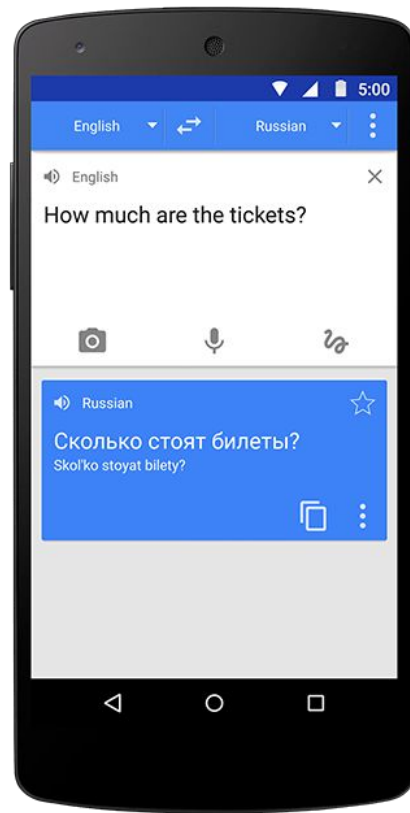
“ The Crown Prosecution Service is facing embarrassment after admitting it destroyed key emails relating to the WikiLeaks founder Julian Assange, who is holed up in Ecuador's London embassy fighting extradition.

Email exchanges between the CPS and its Swedish counterparts over the high-profile case were deleted after the lawyer at the UK end retired in 2014.

The CPS, responding to questions from the Guardian, denied there were any legal implications of the data loss for an Assange case if it were to come to court in the future.

[Extended Summary](#) | [FAQ](#) | [Feedback](#) | Top keywords: [CPS](#)^{#1} [Assange](#)^{#2} [case](#)^{#3} [Swedish](#)^{#4} [Email](#)^{#5}

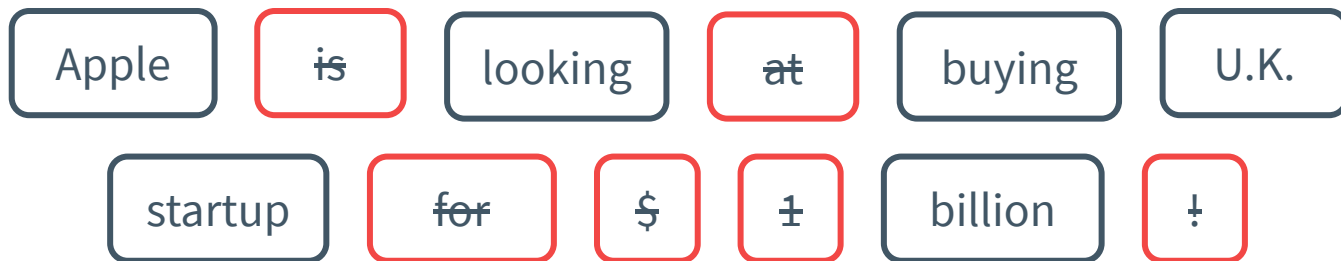
[permalink](#) [embed](#) [save](#) [report](#) [give gold](#) [REPLY](#)



Conversação Chatbot



Apple is looking at buying U.K. startup for \$1 billion!

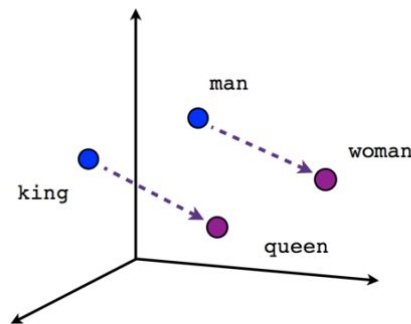


Bag of words

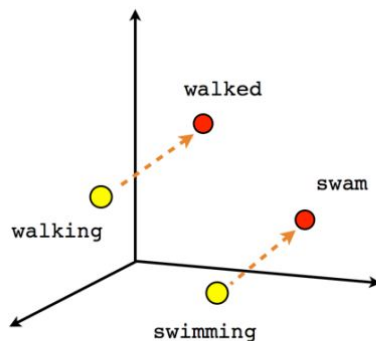
Representa-se uma mensagem por um vetor binário referente a todas suas palavras contidas.

Apple will buy.

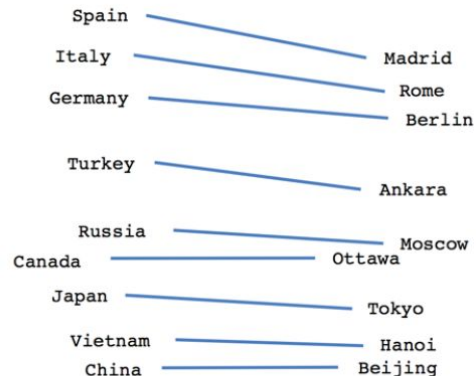
google	0
buy	1
apple	0
walk	0
shoes	0
apple	1
startup	0
will	1



Male-Female



Verb tense



Country-Capital

Czech + currency	Vietnam + capital	German + airlines	Russian + river	French + actress
koruna	Hanoi	airline Lufthansa	Moscow	Juliette Binoche
Check crown	Ho Chi Minh City	carrier Lufthansa	Volga River	Vanessa Paradis
Polish zolty	Viet Nam	flag carrier Lufthansa	upriver	Charlotte Gainsbourg
CTK	Vietnamese	Lufthansa	Russia	Cecile De

GloVe: Global Vectors for Word Representation

Jeffrey Pennington, Richard Socher, Christopher D. Manning

The logo for fastText, featuring the word "fast" in a red, italicized sans-serif font and the word "Text" in a blue, bold sans-serif font.

Library for efficient text classification and representation learning

Assume-se independência entre as probabilidades dos pares de palavras.

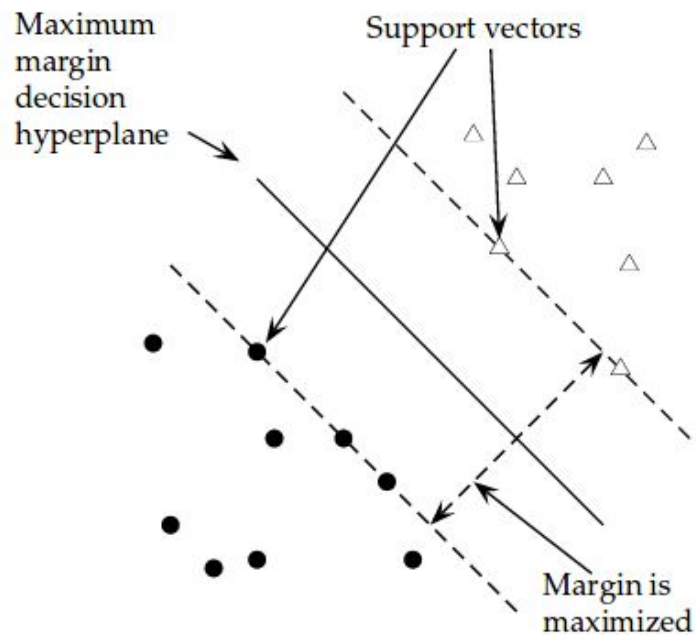
$$P(c|d) \propto P(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$$

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} \hat{P}(c|d) = \arg \max_{c \in \mathbb{C}} \hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(t_k|c)$$

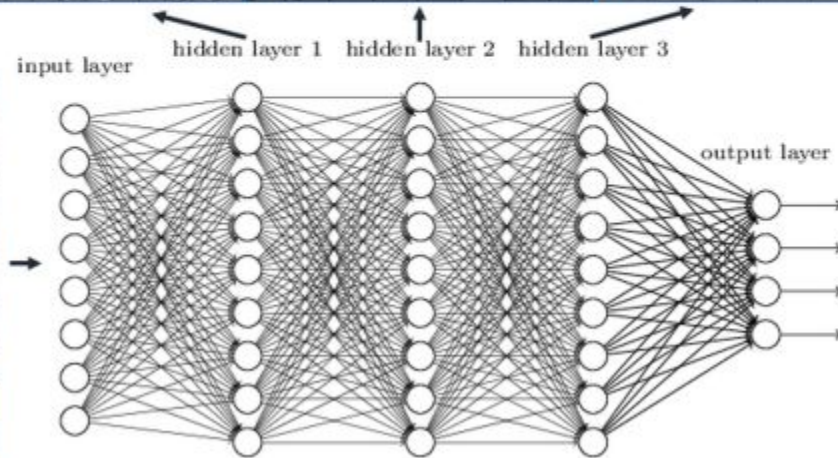
$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

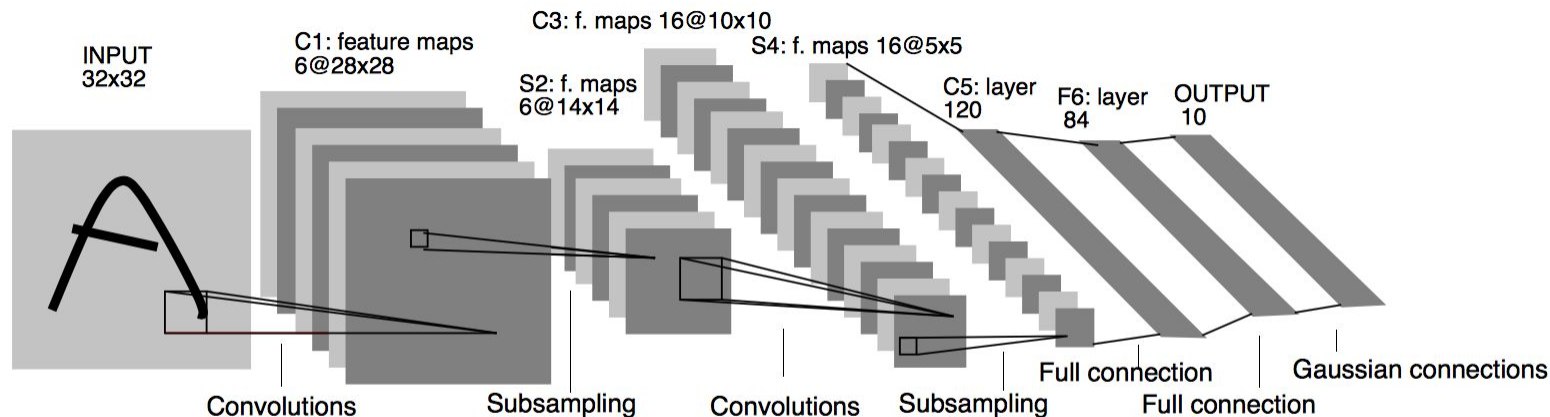
Support Vector Machine



Deep neural networks learn hierarchical feature representations



Redes Neurais Convolucionais



CNN

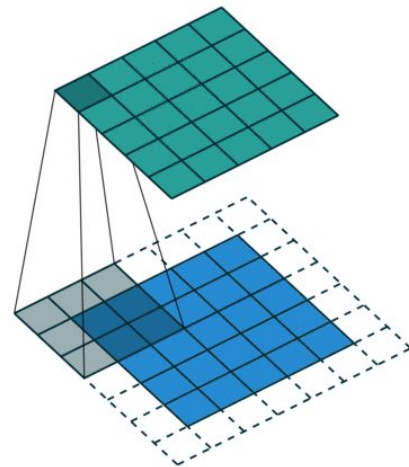
Convolução

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature



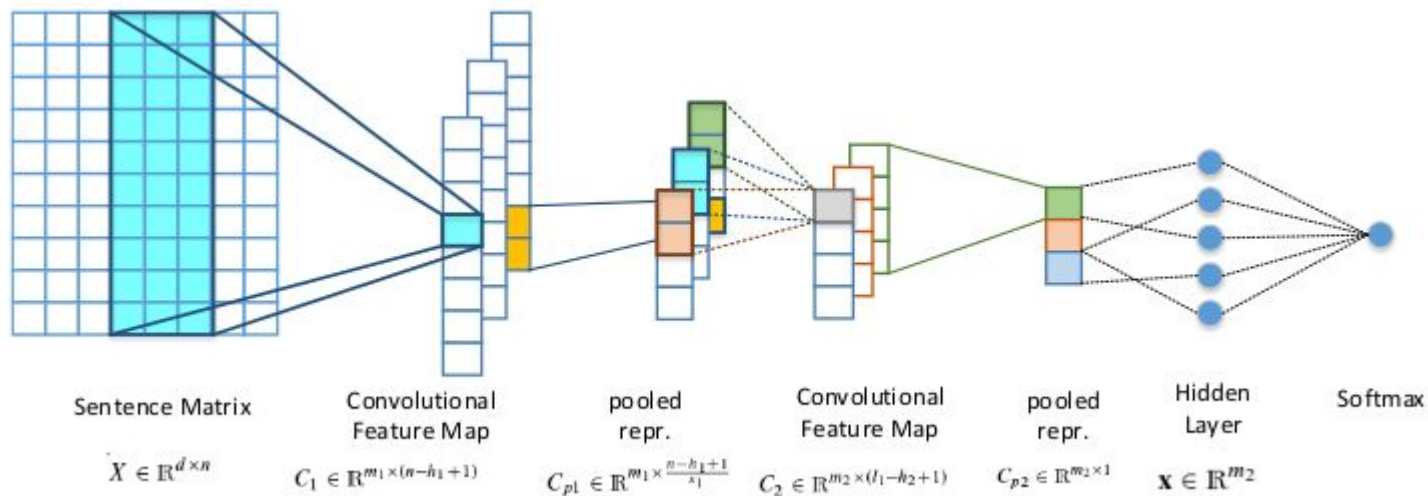
CNN

Max-Pooling

1	3	2	9
7	4	1	5
8	5	2	3
4	2	1	4

CNN

Aplicadas em Texto



Emoticons positivos	Emoticons negativos
:)	:(
:-)	:-)
:D	: (

**Arnold Varona**

@ArnyVee



Follow



TheWDB.com - Very cool to hear old Walt interviews! :D 🎵 <http://blip.fm/~8bmta>

**Karoli**

@Karoli



Follow



@nationwideclass no, it's not behaving at all. i'm mad. why am i here? because I can't see you all over there. :(

Acurácia dos classificadores do Sentiment140.

	Máxima Entropia	Naive Bayes	SVM
Unigrama	80,5%	81,3%	82,2%
Bigrama	79,1%	81,6%	78,8%
Unigrama + Bigrama	83,0%	82,7%	81,6%
Unigrama + POS	79,9%	79,9%	81,9%

- ▶ Conferência anual de análise semântica computacional.
- ▶ Realiza competições em análise de sentimento de tweets desde 2013.
- ▶ Participantes disponibilizam artigos com metodologia, técnicas e resultados obtidos.

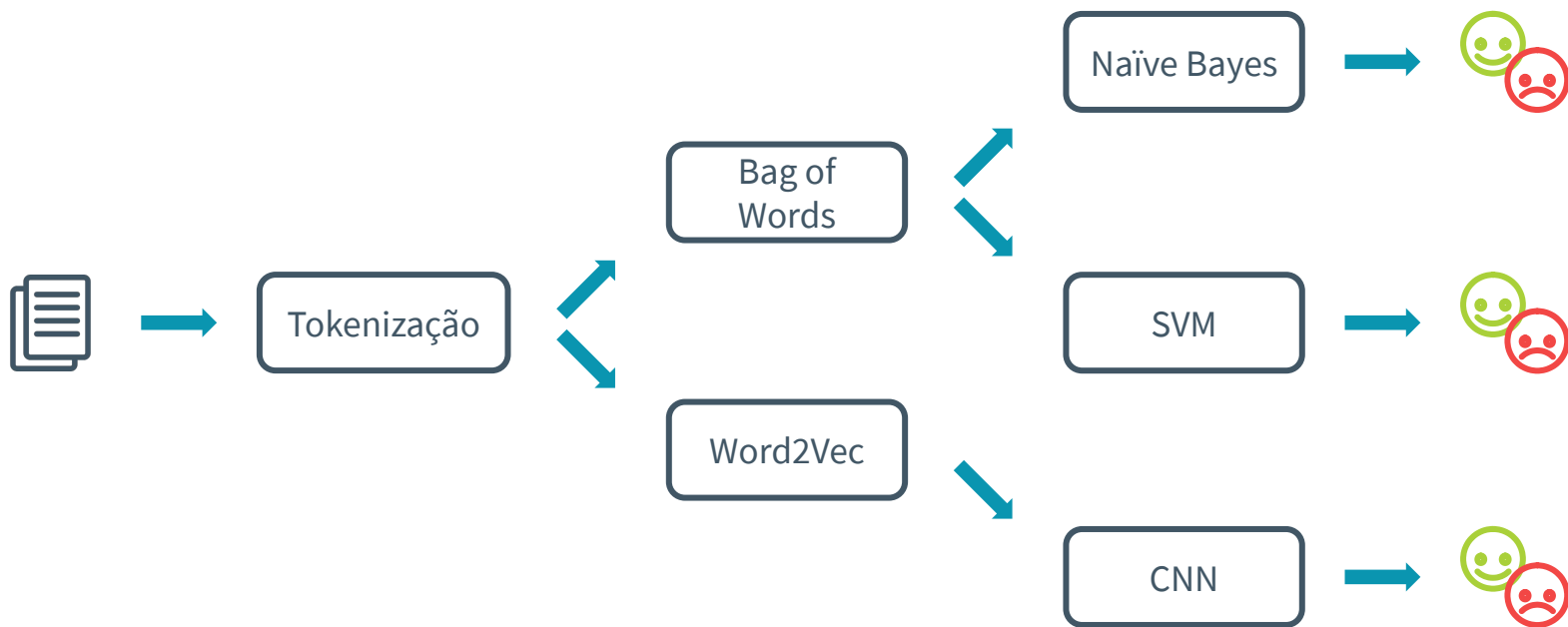
Método

	Tamanho	Anotação	Ano
Sentiment140 Treino	1,6M	supervisão distante	2009
Sentiment140 Teste	360	manual	2009
Base Própria	40M	supervisão distante	2015 - 2017
SemEval	60K	manual	2013 - 2017

	Objetivo	Treinamento	Teste
Etapa 1	Replicar classificadores do Sentiment140.	Sentiment140 Treino	Sentiment140 Teste
Etapa 2	Treinar os classificadores a partir de base de dados própria.	Base Própria	SemEval
Etapa 3	Aplicar técnicas de Deep Learning para classificação de sentimento.	Base Própria	SemEval

Classificadores

Visão Geral



Etapa 3

Word2Vec

- ▷ Vocabulário de **3M termos**
- ▷ Embedding de **300 dimensões**
- ▷ Treinado em 100B palavras de notícias.
- ▷ Contém **89k** dos **1,1M** termos únicos dos datasets de treino e teste.

Etapa 3

Hiperparâmetros

- ▷ ReLU
- ▷ Entropia Cruzada
- ▷ Adam
- ▷ L2
- ▷ Dropout
- ▷ Early-Stopping
- ▷ N° Camadas: 1 e 2
- ▷ N° Filtros Conv.: 100 e 200
- ▷ Tam. Filtros Conv.: 2 e 3
- ▷ Tam. Pooling: 2, 3 e 5

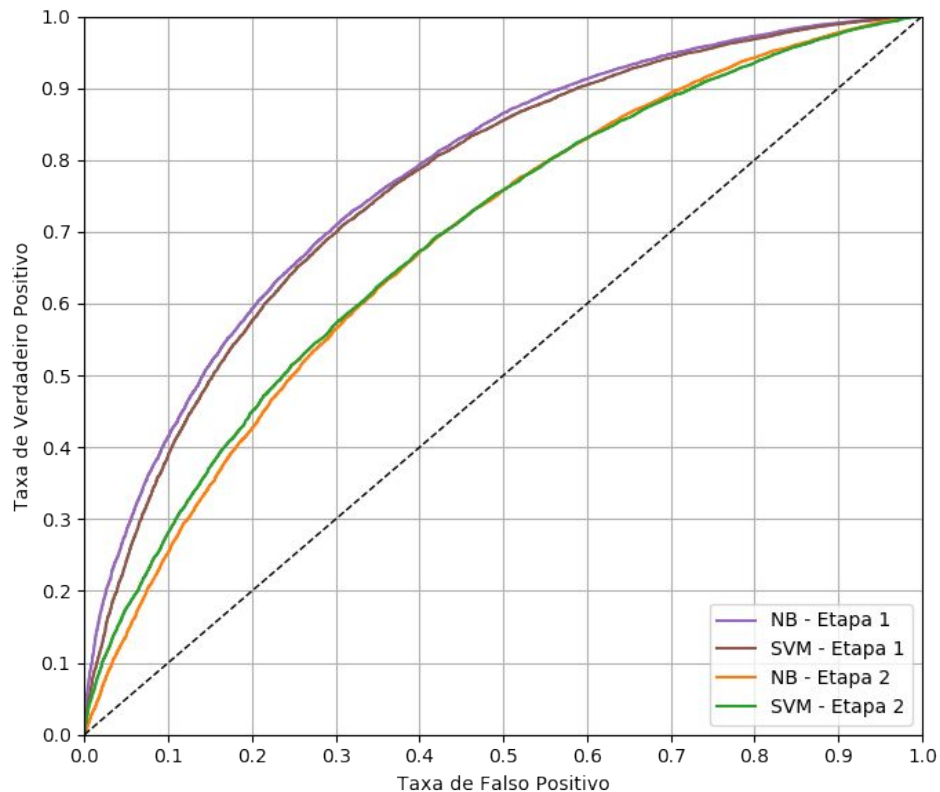
Resultados

Etapa 1 Resultados

Acurácia

	Referência	Reprodução
Naive Bayes	81,3%	83,3%
SVM	82,2%	83,0%

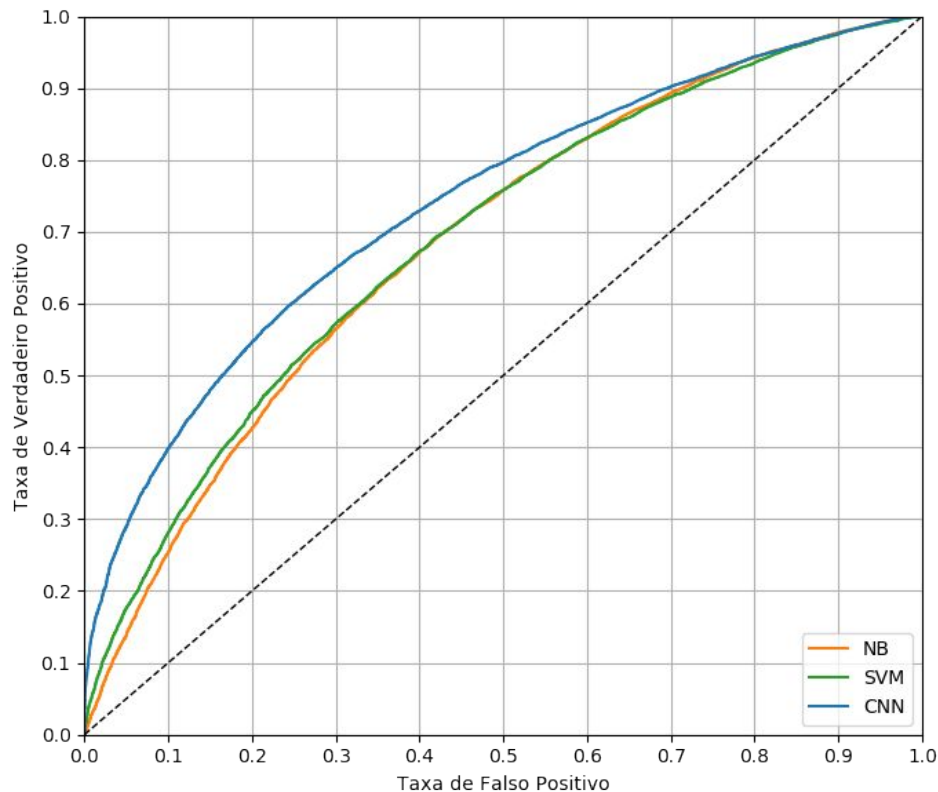
Etapa 2 Resultados



	AUC	SP
Naïve Bayes	0,688	0,639
SVM	0,693	0,640

Etapa 3

Resultados



	AUC	SP
Naïve Bayes	0,688	0,639
SVM	0,693	0,640
CNN	0,738	0,675

1. Reproduzimos com sucesso classificadores de sentimento de redes sociais.
2. Se treinaram classificadores eficientes a partir de bases anotadas por supervisão distante.
3. Deep Learning obteve desempenho superior a técnicas clássicas.

- ▶ Aplicar o processo para português.
- ▶ Adicionar sentimento neutro.
- ▶ Avaliar outras técnicas de Deep Learning.

Obrigado!

Perguntas?

X.

Extras

Polaridade de mensagem



João Silva
@tchallaz



 Follow

Grande notícia, João Sousa e Gastão Elias representam Portugal nos Jogos Olímpicos do Rio de Janeiro [#Rio2016](#)

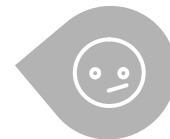


Esportes - Olímpicos
@brenobarros



 Follow

entrevistei hoje a Priscilla Carnaval do bmx. Classificada para o [#Rio2016](#), ela tem novidades na preparação olímpica.

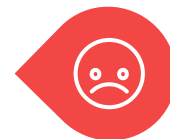


Rogério Faria
@Carabutim



 Follow

Estão transformando Olimpíadas que é algo sério num espetáculo triste e de mau gosto [#Rio2016](#) [#TourDaTocha](#)



Codificação One-Hot

É criado um espaço vetorial no qual cada palavra do vocabulário corresponde a uma dimensão.

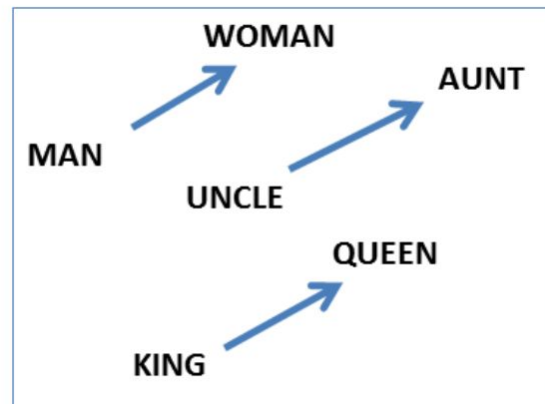
startup

buy

[illegible]

0
1
0
0
0
0
0
0

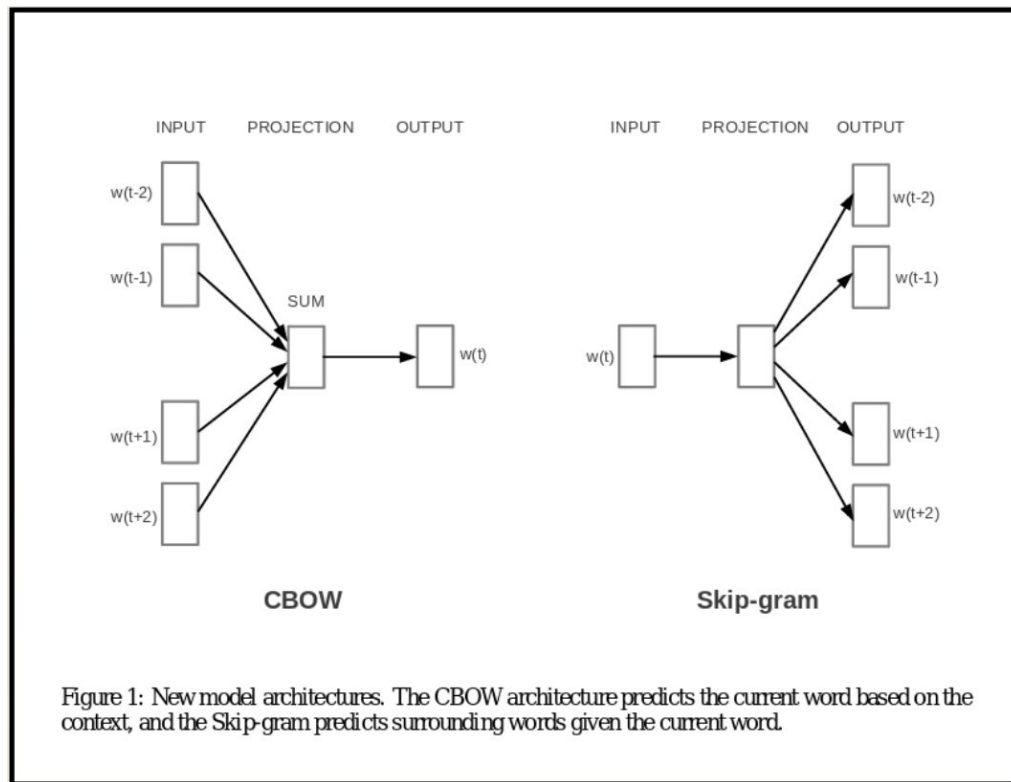
- ▶ Embedding que representa uma palavra em um vetor real de tamanho arbitrário.
- ▶ É produzido a partir de janelas de contexto ao redor das palavras.



In this paper we present several extensions that improve both the quality of the vectors and the training speed.

Word2vec

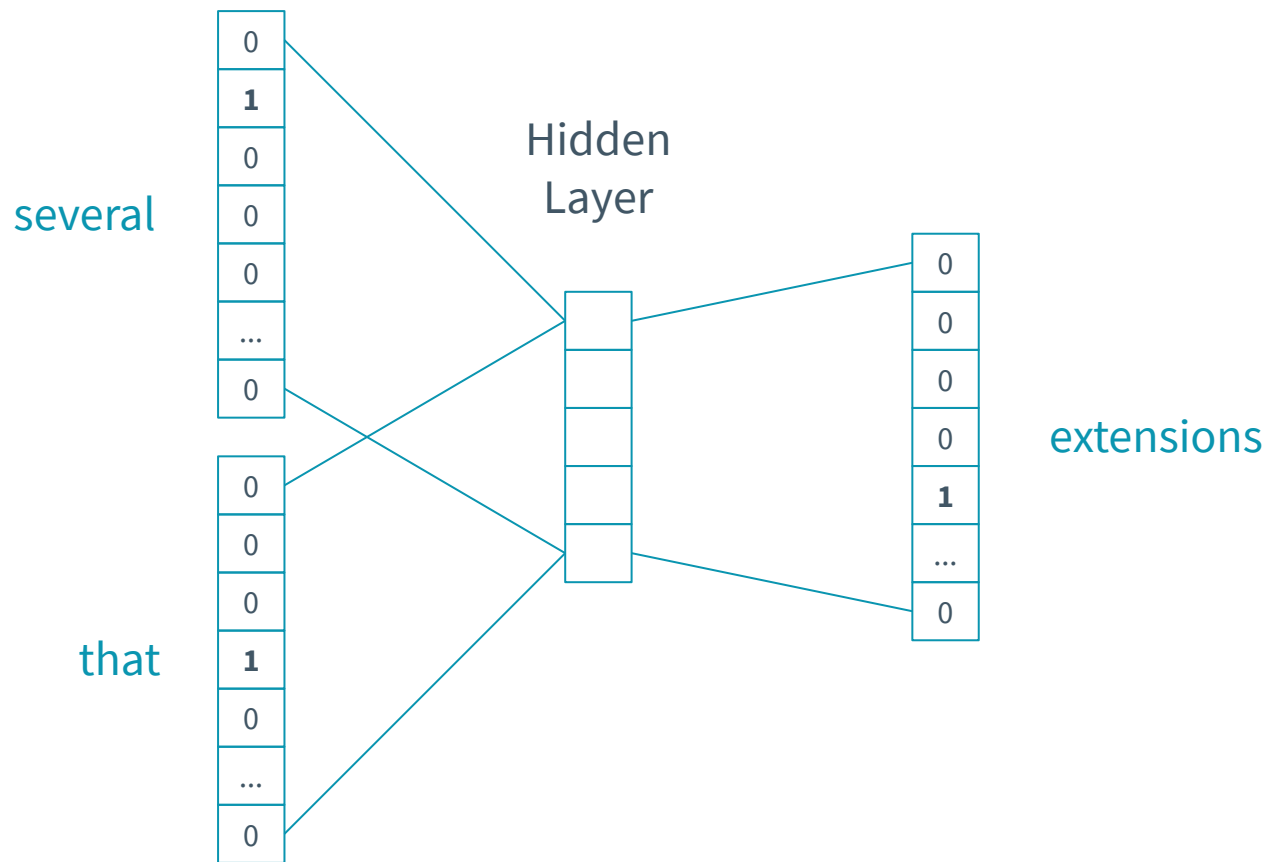
Arquiteturas

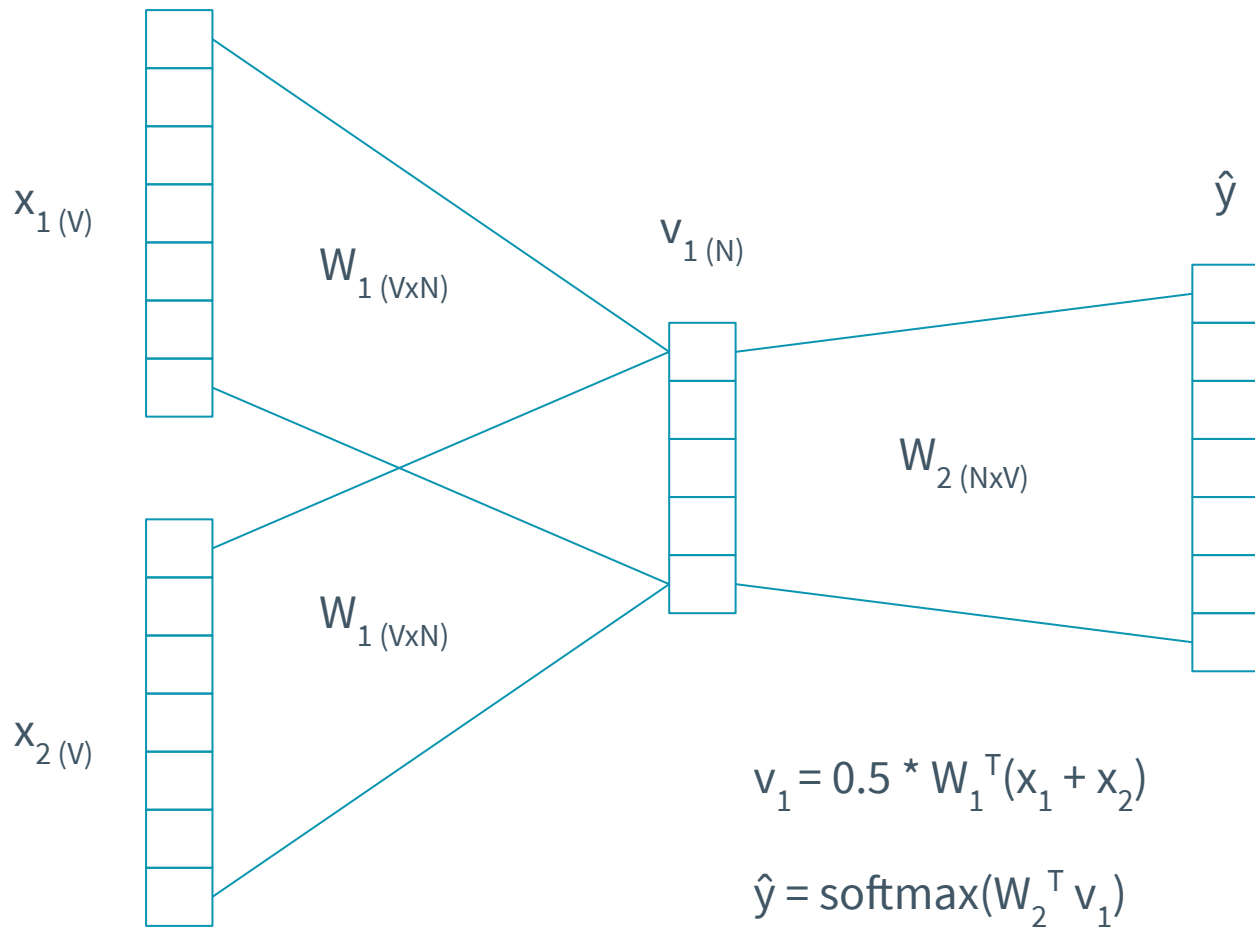


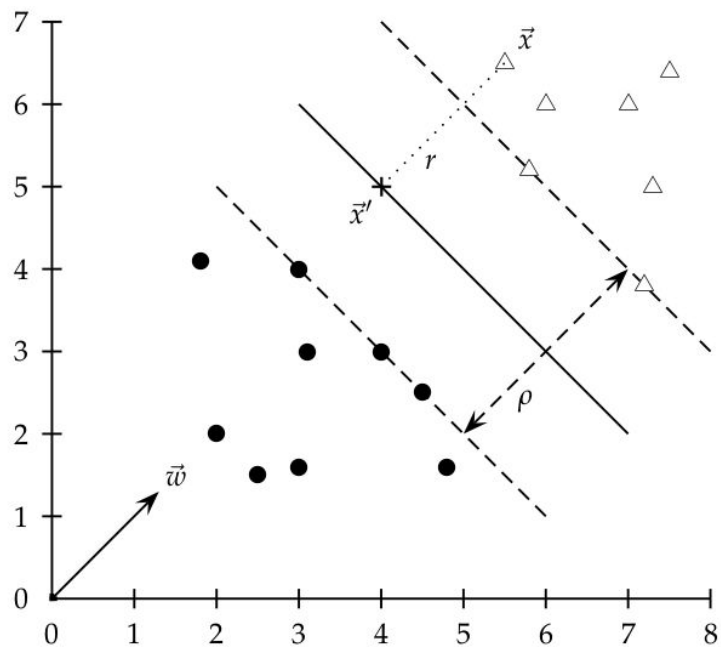
54

Treinamento
Word2vec

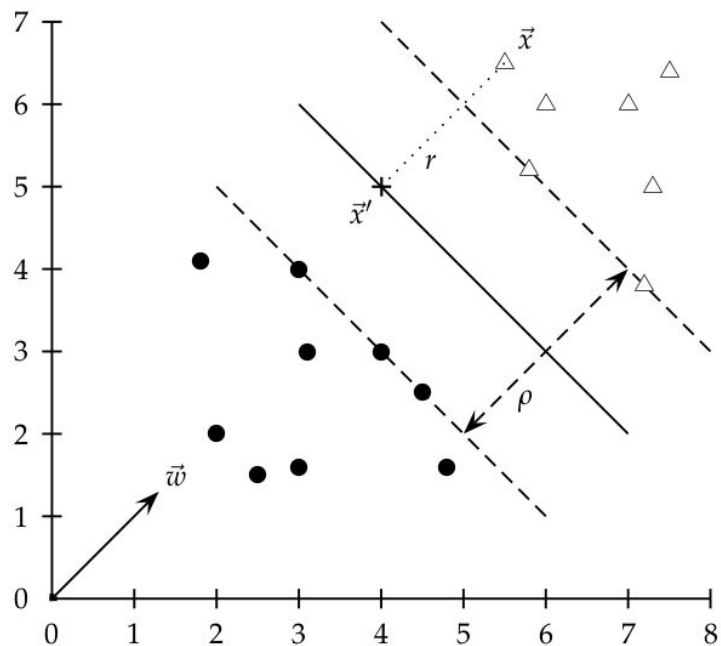
CBOW







$$f(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b)$$



$$\vec{x}' = \vec{x} - yr \frac{\vec{w}}{|\vec{w}|}$$

$$\vec{w}^T \left(\vec{x} - yr \frac{\vec{w}}{|\vec{w}|} \right) + b = 0$$

$$r = y \frac{\vec{w}^T \vec{x} + b}{|\vec{w}|}$$

$\rho = 2/|\vec{w}|$ is maximized

For all $(\vec{x}_i, y_i) \in \mathbb{D}$, $y_i(\vec{w}^T \vec{x}_i + b) \geq 1$

Neurônios: inicialização Glorot uniforme¹

Bias: inicialização em 0

$$^1 \textit{limites} = \pm \textit{ganho} \cdot \sqrt{\frac{2}{\textit{entradas} + \textit{saidas}}} \cdot \sqrt{3}$$

$$\textit{ganho} = \begin{cases} \sqrt{2} & \text{se } \varphi = \text{ReLU} \\ 1 & \text{se } \varphi = \text{Sigmoid} \end{cases}$$

Regra de atualização

Otimização Adam

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)

$v_0 \leftarrow 0$ (Initialize 2nd moment vector)

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

end while

return θ_t (Resulting parameters)

- ▶ Função custo: Entropia cruzada
Permite interpretar \hat{y} como probabilidade

$$\frac{- \sum [(\mathbf{y} \cdot \log(\hat{\mathbf{y}}) + (1 - \mathbf{y}) \cdot \log(1 - \hat{\mathbf{y}})) + \frac{\lambda}{2} \cdot \mathbf{W}^2]}{n}$$

Hiperparâmetros

Otimização

- ▷ $\alpha = 0,001$
- ▷ $\beta_1 = 0,9$
- ▷ $\beta_2 = 0,999$

Regularização

- ▷ L2: $\lambda = 0,001$
- ▷ Dropout: 0,5

Tweets podem ser facilmente coletados através de API. Podem ser coletados tweets de um contexto específico, a partir de filtro por hashtag ou palavras chave, ou não.

9.

Aplicações

Repercussion

3.13



37 16 6

Pesquisar

Data de publicação	Título	Positiva	Neutra	Negativa
10/10/2015 04:51:58	Corey Taylor acredita que o Slipknot pode entrar para a história como Deep Purple e Black Sabbath			
09/10/2015 09:12:00	Fernando, do 'BBB15', se declara para Aline em aniversário: 'Mudou minha vida'			
08/10/2015 10:00:00	Seal lança clipe de 'Every Time I'm With You'; assista ao vídeo			
08/10/2015 08:19:00	Baixo Manhattan: Rock in Rio USA é eleito um dos piores festivais do ano			
08/10/2015 06:26:00	Rihanna revela capa do CD 'Anti': acompanhe as mudanças da artista através dos álbuns			
08/10/2015 06:07:00	Bruna Marquezine comenta cenas de briga em 'I Love Paraisópolis': 'Me achando'			
08/10/2015 06:05:00	Rihanna mostra capa de novo disco na rede. Confira!			
08/10/2015 03:01:00	SpokFrevo Orquestra lança terceiro disco em apresentação gratuita em Porto Alegre			
08/10/2015 01:40:23	Alice Cooper toca com a formação original de sua banda em show surpresa. Veja!			
07/10/2015 12:06:00	Rock in Rio e feriado puxam alta de 23% das passagens aéreas, diz IBGE			
07/10/2015 11:52:14	Lzzy Hale (Halestorm) se veste de freira em show do Ghost			
07/10/2015 11:45:00	Inflação acelera em 4 dos 9 grupos do IPCA em setembro			
07/10/2015 11:01:00	Inflação acelera em 4 dos 9 grupos do IPCA em setembro			
07/10/2015 10:54:00	Inflação acelera em 4 dos 9 grupos do IPCA em setembro			
07/10/2015 10:42:00	Inflação acelera em 4 dos 9 grupos do IPCA em setembro			

Mostrando de 1 até 15 de 3.283 registros

[← Anterior](#) [1](#) [2](#) [3](#) [4](#) [5](#) [Próximo →](#)

