



# WORD2VEC

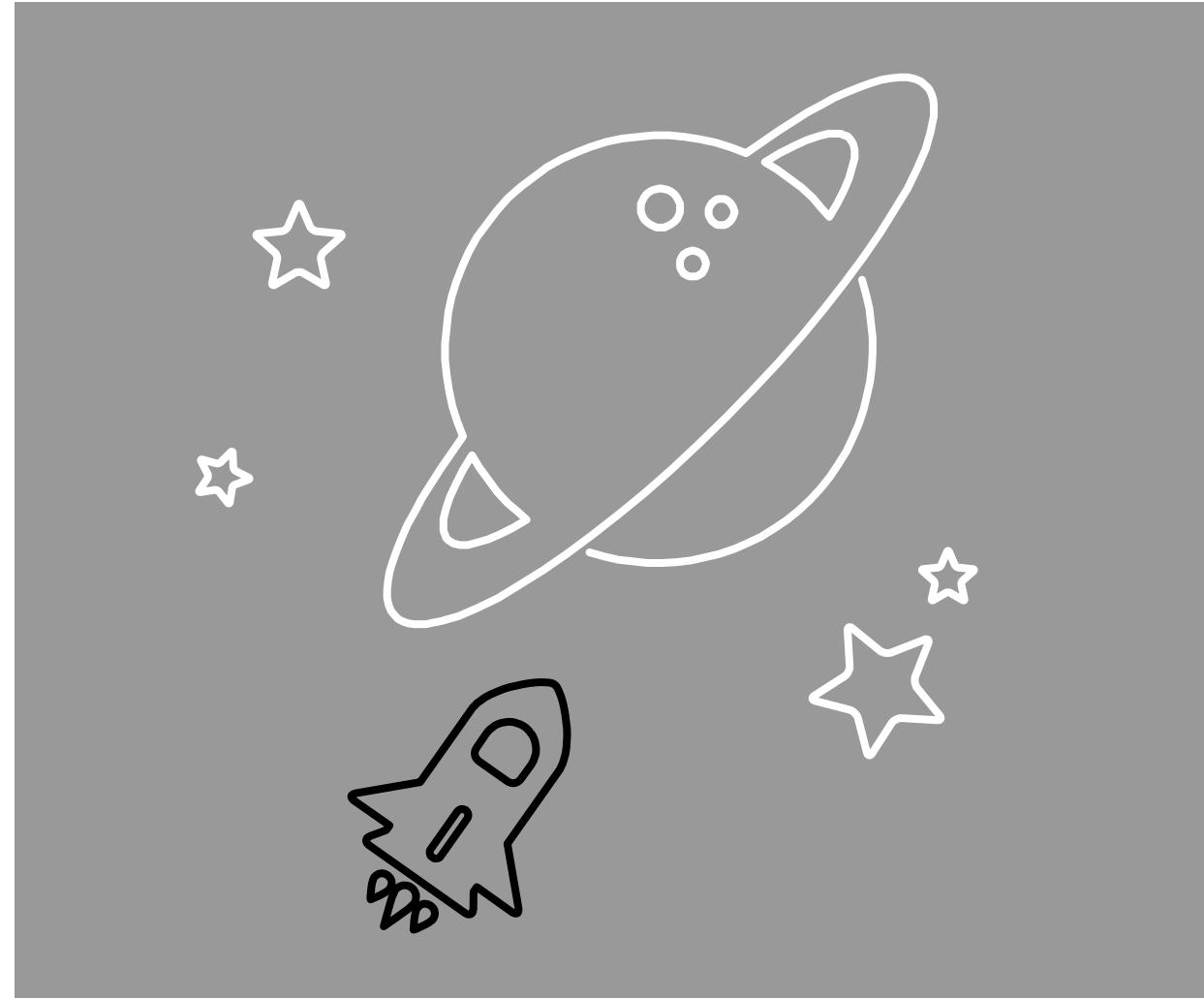
## O CASO MAIS FAMOSO DE WORD EMBEDDINGS

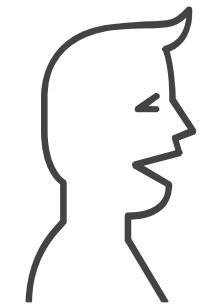
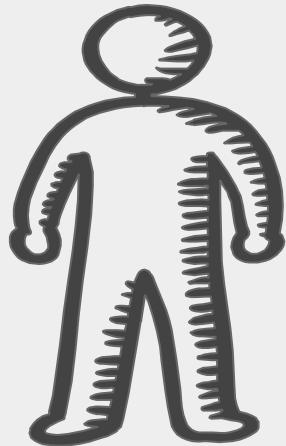
Larissa Sayuri Futino Castro dos Santos  
18 maio 2019

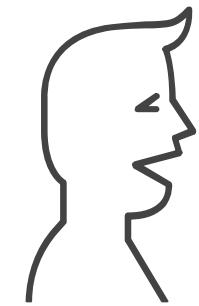
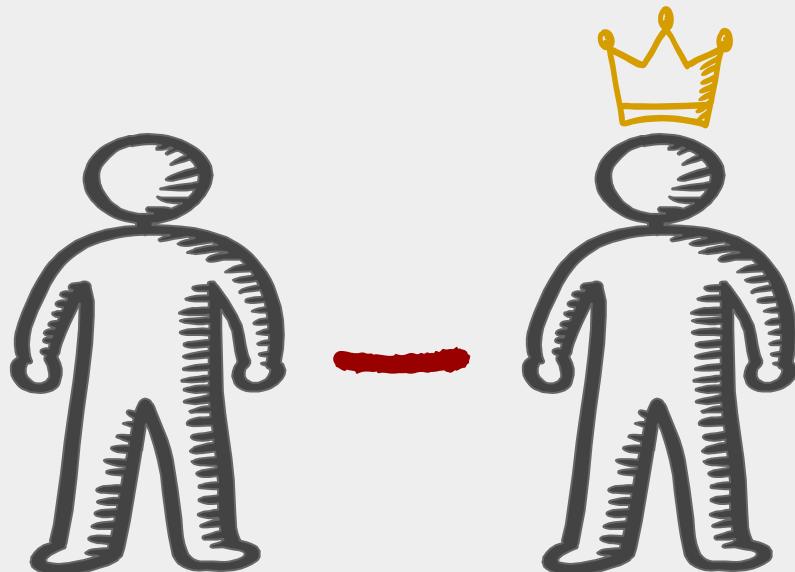


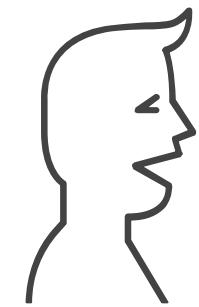
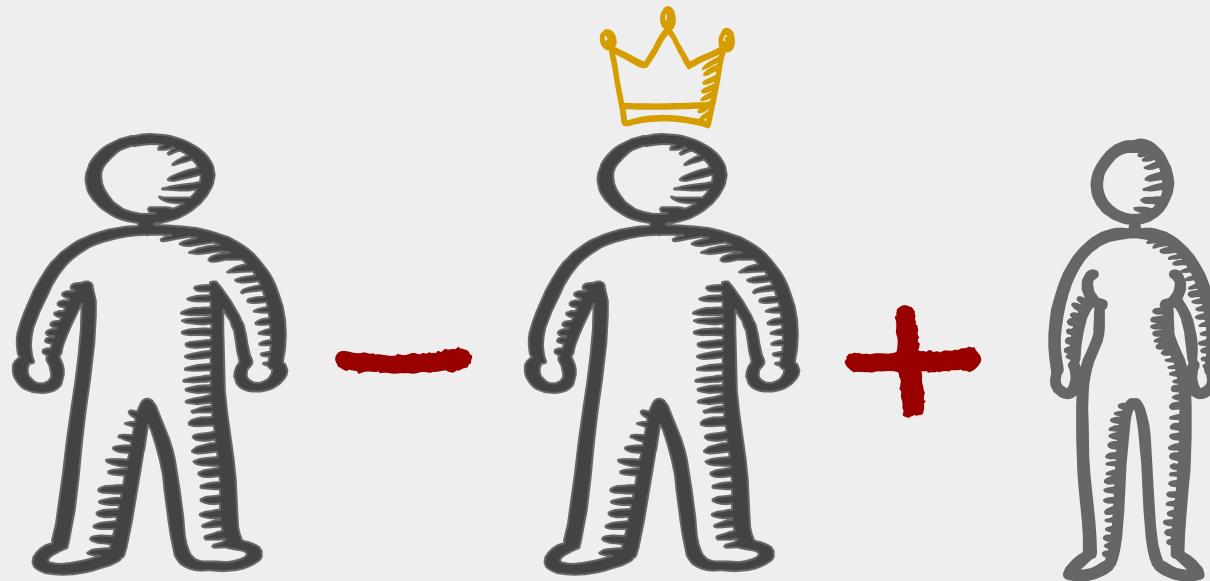
PARECE  
MÁGICA!

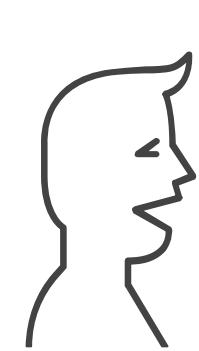
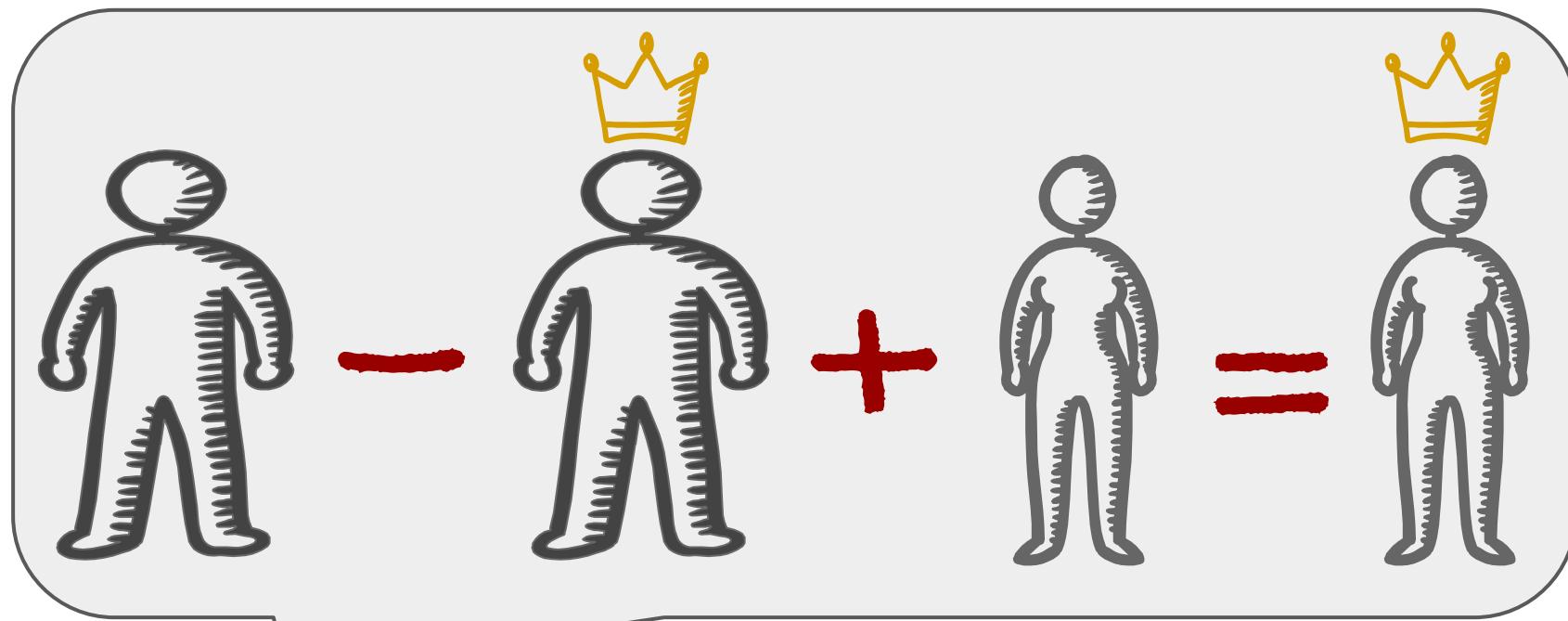
Usa matemática,  
Jeitão de linguística

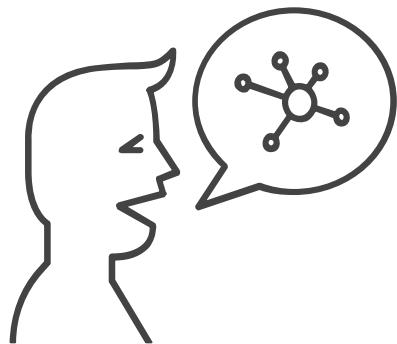






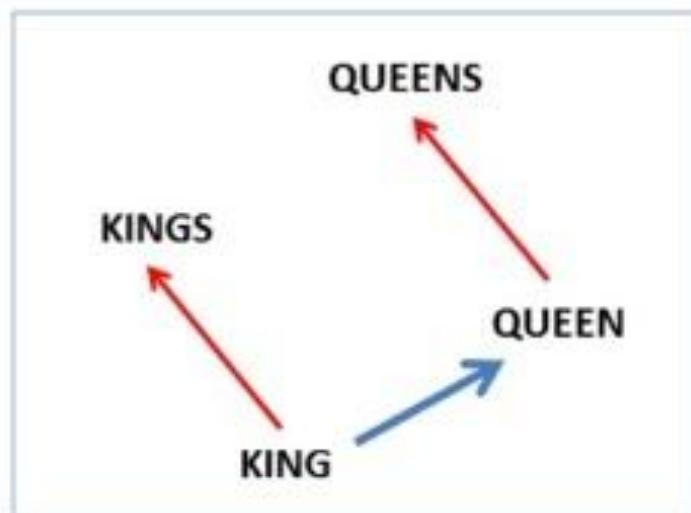
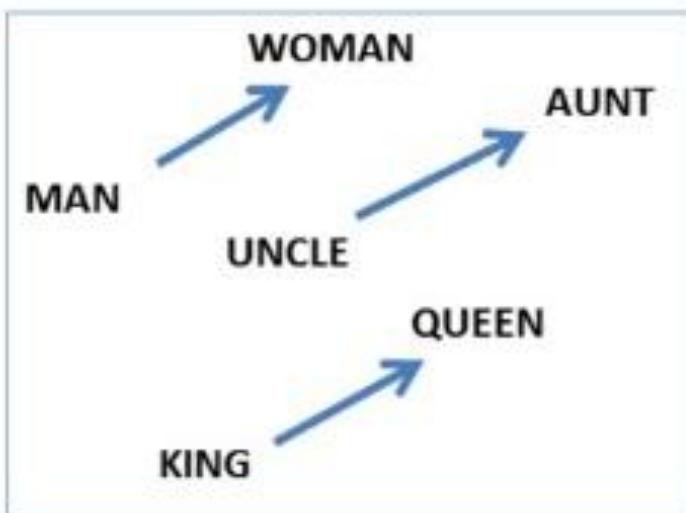
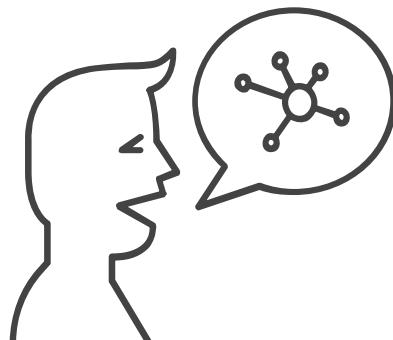


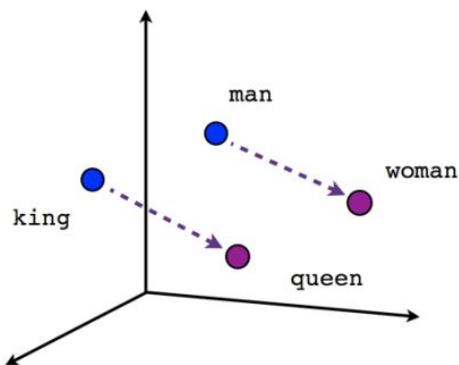



$$\text{vec("man")} - \text{vec("king")} + \text{vec("woman")} = \text{vec("queen")}$$


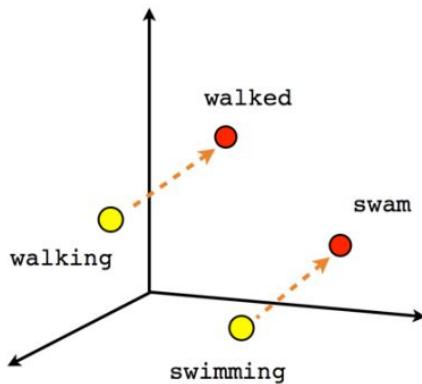


$$\text{vec}(\text{"man"}) - \text{vec}(\text{"king"}) + \text{vec}(\text{"woman"}) = \text{vec}(\text{"queen"})$$

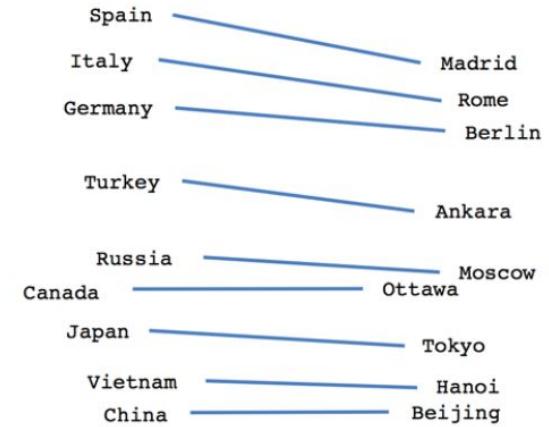




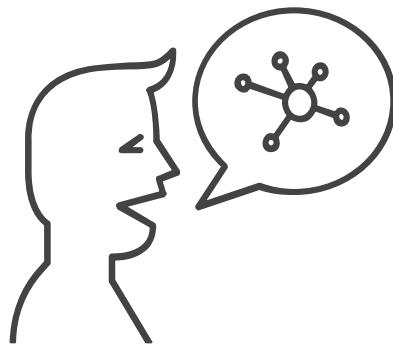
Male-Female

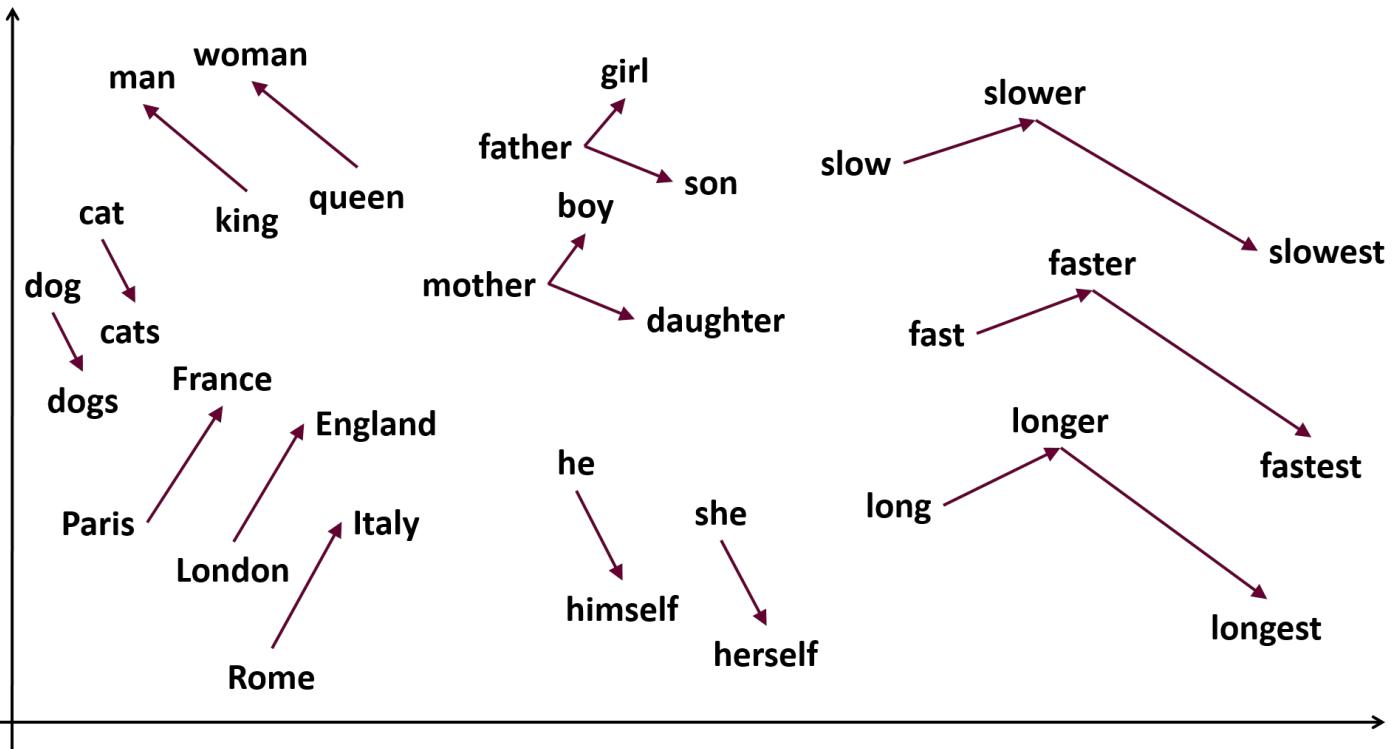
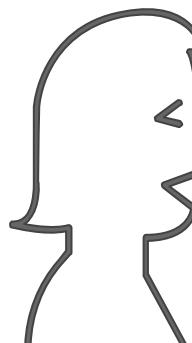


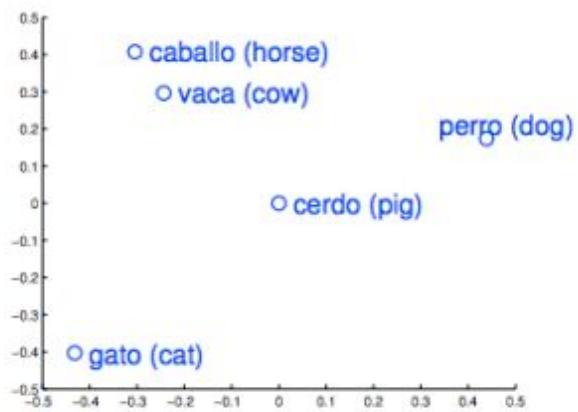
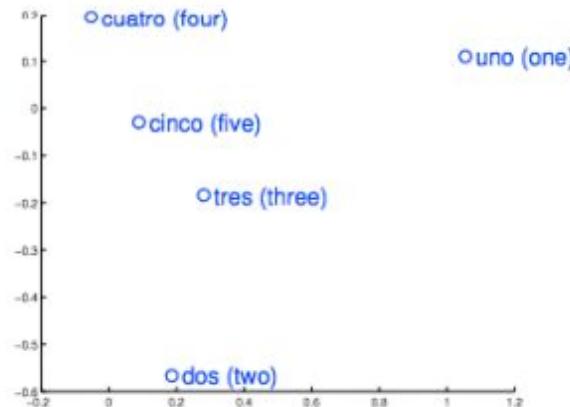
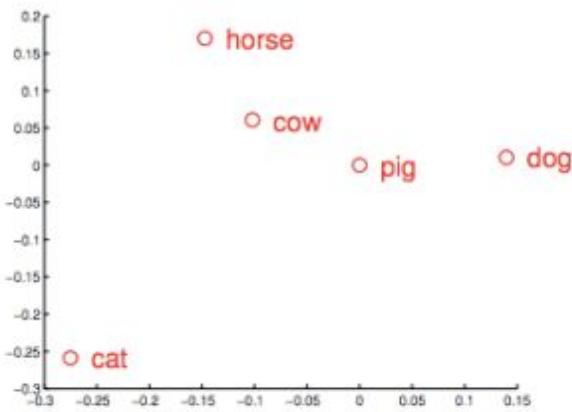
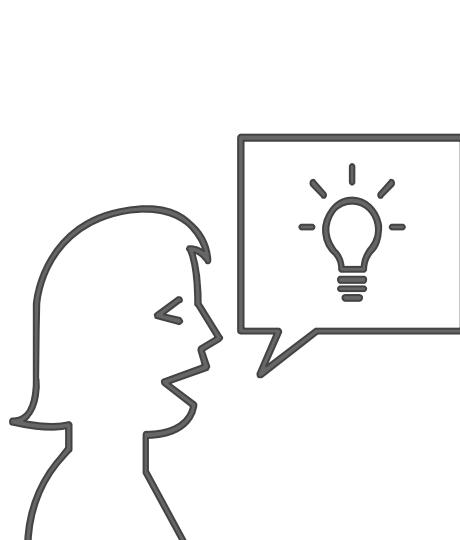
Verb tense



Country-Capital







# Trabalhos

MIKOLOV, Tomas et al. Efficient estimation of word representations in vector space. **arXiv preprint arXiv:1301.3781**, 2013.

MIKOLOV, Tomas; YIH, Wen-tau; ZWEIG, Geoffrey. Linguistic regularities in continuous space word representations. In: **Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies**. 2013. p. 746-751.

MIKOLOV, Tomáš et al. Recurrent neural network based language model. In: **Eleventh Annual Conference of the International Speech Communication Association**. 2010.

BENGIO, Yoshua et al. A neural probabilistic language model. **Journal of machine learning research**, v. 3, n. Feb, p. 1137-1155, 2003.

**EU FAÇO PARTE**

de uma das 100 empresas  
mais inovadoras do país.



2019





Doutorado Sanduíche  
UFMG/CAPES/  
University of  
Wolverhampton

2017

2018



2013

2015

Bacharelado  
UNB

2007



2010

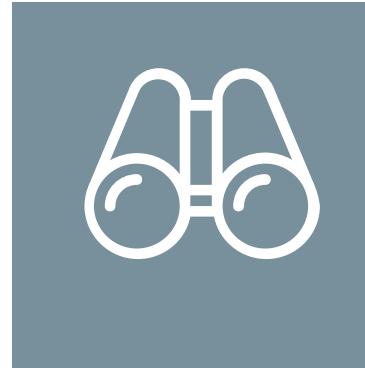


# About: MIKOLOV, Tomas



Tomas Mikolov

A grande referência



42.756  
citações  
(17/05/19)

Doutorado  
Brno  
University of  
Technology  
(Czech  
Republic)





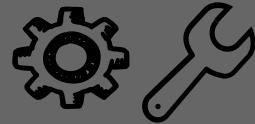
## BONS EXEMPLOS

Exemplos com  
diferentes formas  
de similaridade  
(semântica e  
sintática)



## POPULAR

Google  
desenvolveu uma  
ferramenta,  
tensorflow, e a  
disponibilizou



## GRANDES DATASETS

A técnica lida bem  
com bases grandes,  
típicas da realidade  
atual

(6B tokens, 1M  
termos mais freq)



# RÁPIDO

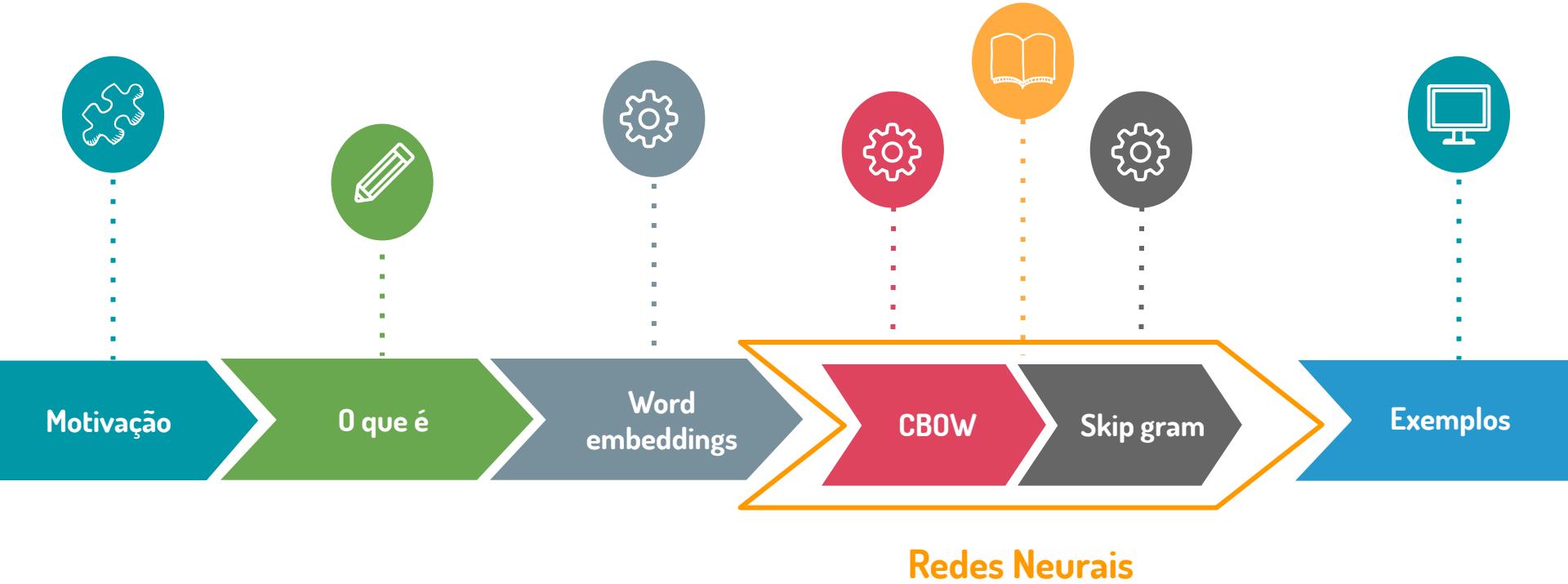
De: 08 semanas  
Para: 01-03 dias



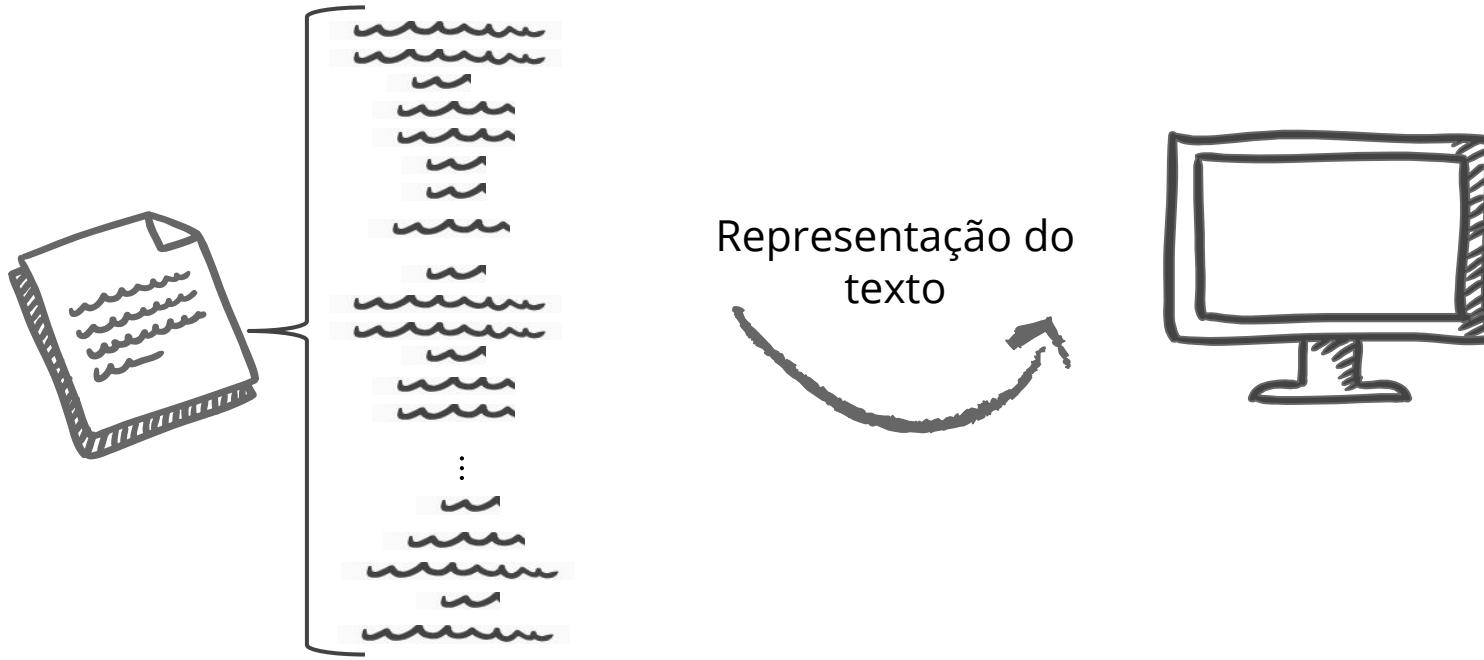
# RELACIONES PRESERVADAS ENTRE LÍNGUAS

Permite tradução  
termo-a-termo

# Agenda

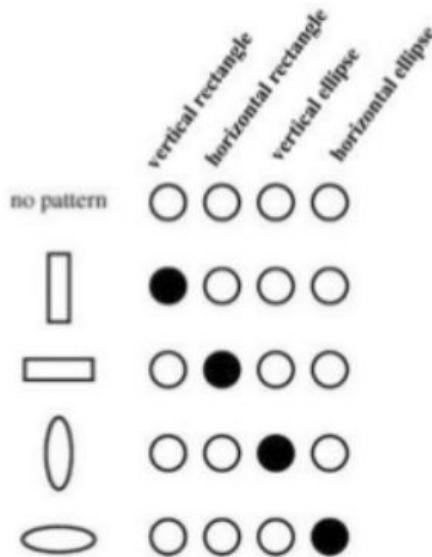


# Linguagem de Máquina



# Linguagem de Máquina

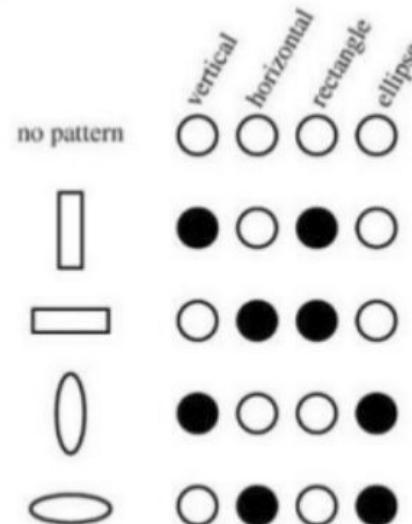
sparse representations



"one hot" vector

A distributed representation is dense

1. One concept is represented by more than one neuron firing
2. One neuron represents more than one concept



Fonte: [slideshare\\_UnivGothenburg](#)

# O QUE FAZ?

Cria vetores que são representações numéricas dos atributos (termos/palavras) sem intervenção humana

# OBJETIVO

Criar representação de termos (palavras) que capte seus significados, relações semânticas e os diferentes tipos de contexto em que são empregadas

# OBJETIVO (GEOMETRIA)

Posiciona termos (palavras) em um espaço dimensional tal que as localizações dos termos são determinadas pelos seus significados e os vazios e distâncias no espaço também tem significado

# COMO FUNCIONA?

Rede neural de duas camadas que processa texto

# PARA QUÊ SERVE?

Cria representação de strings para algoritmos de Machine Learning ou Deep Learning

# RESULTADO

Agrupar vetores de termos similares no espaço de atributos

# Word Embeddings

## Definição:

Representações numéricas para textos

Textos convertidos em números

D1: He is a lazy boy. She is also lazy.

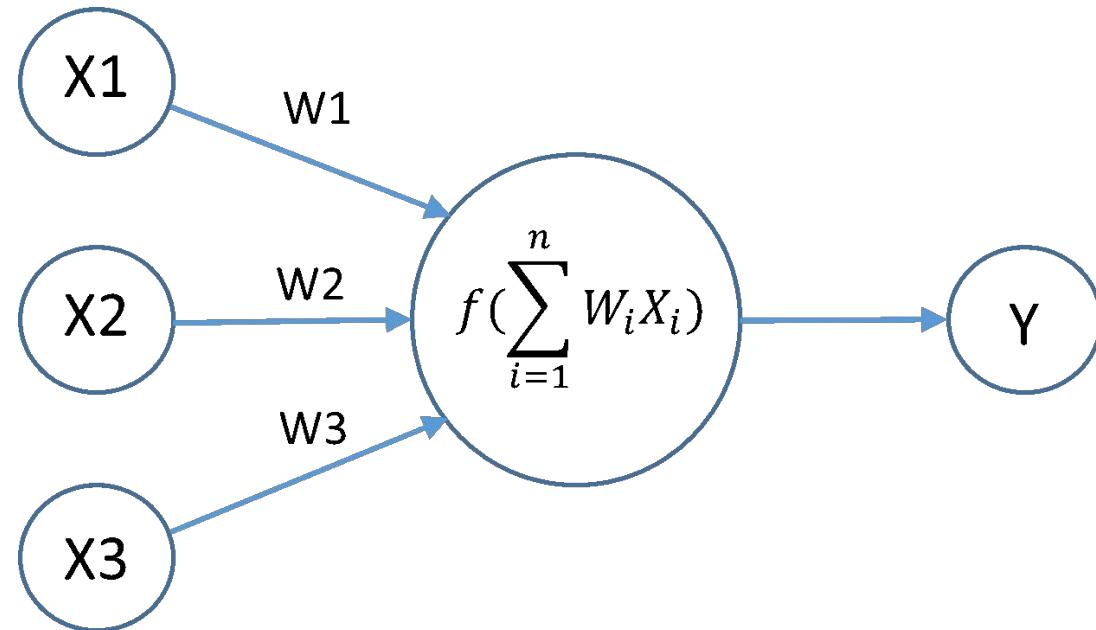
D2: Neeraj is a lazy person.

	He	She	lazy	boy	Neeraj	person
D1	1	1	2	1	0	0
D2	0	0	1	0	1	1

# Redes Neurais

## Definição:

Técnicas  
**computacionais**  
inspiradas na estrutura  
neural de organismos  
inteligentes e que  
**adquirem**  
**conhecimento através**  
**da experiência**

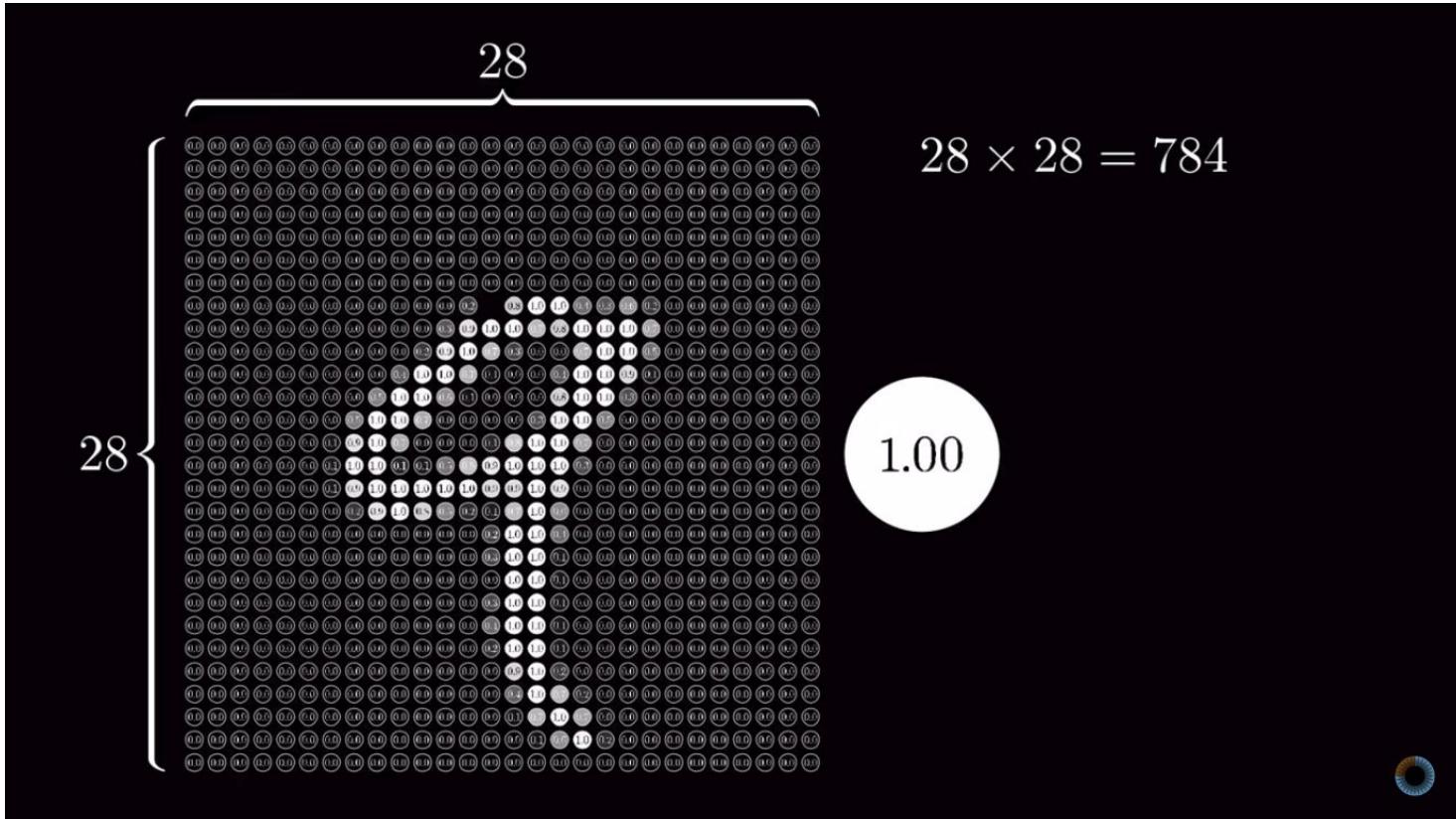


# Redes Neurais

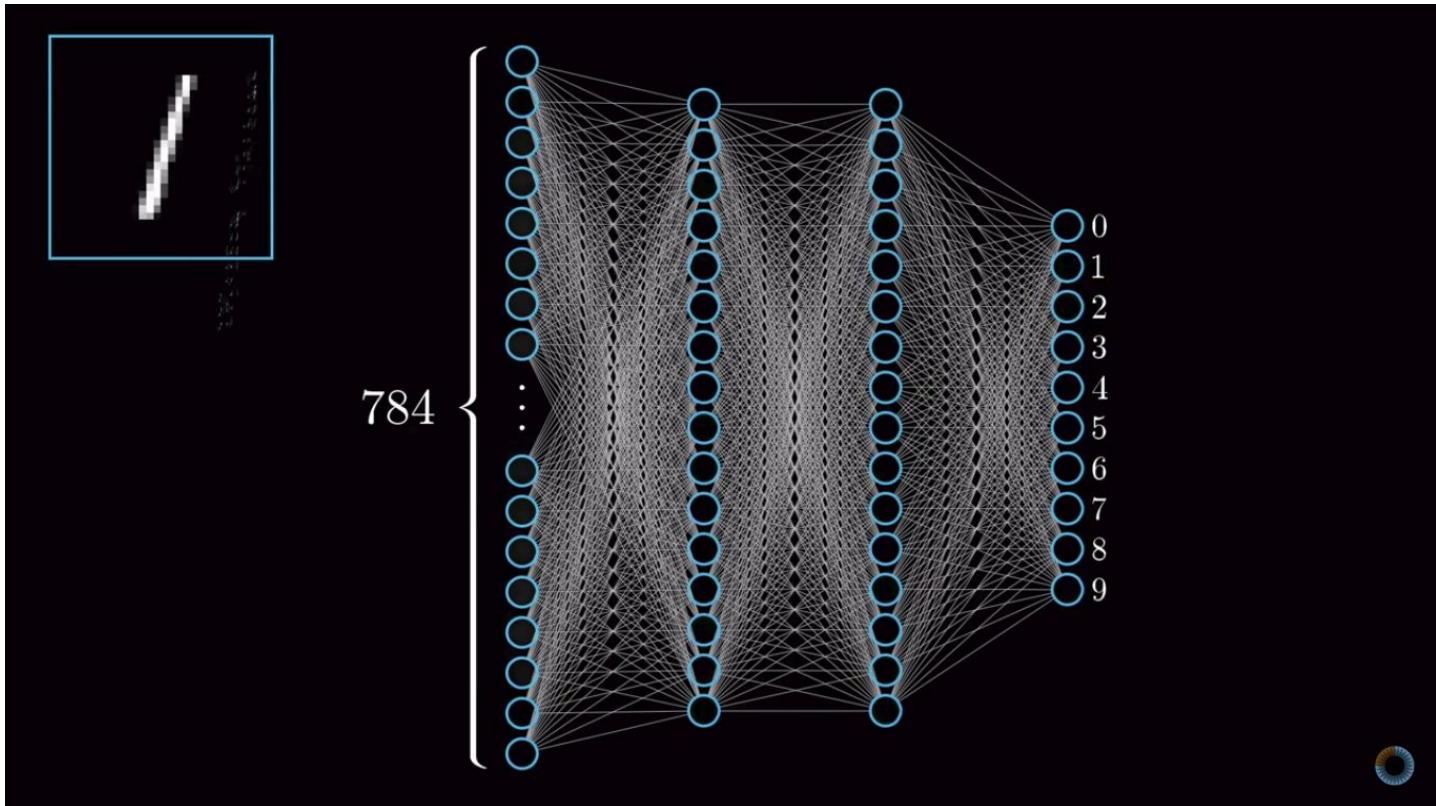
Motivação: Um modelo para reconhecer dígitos numéricos escritos à mão.



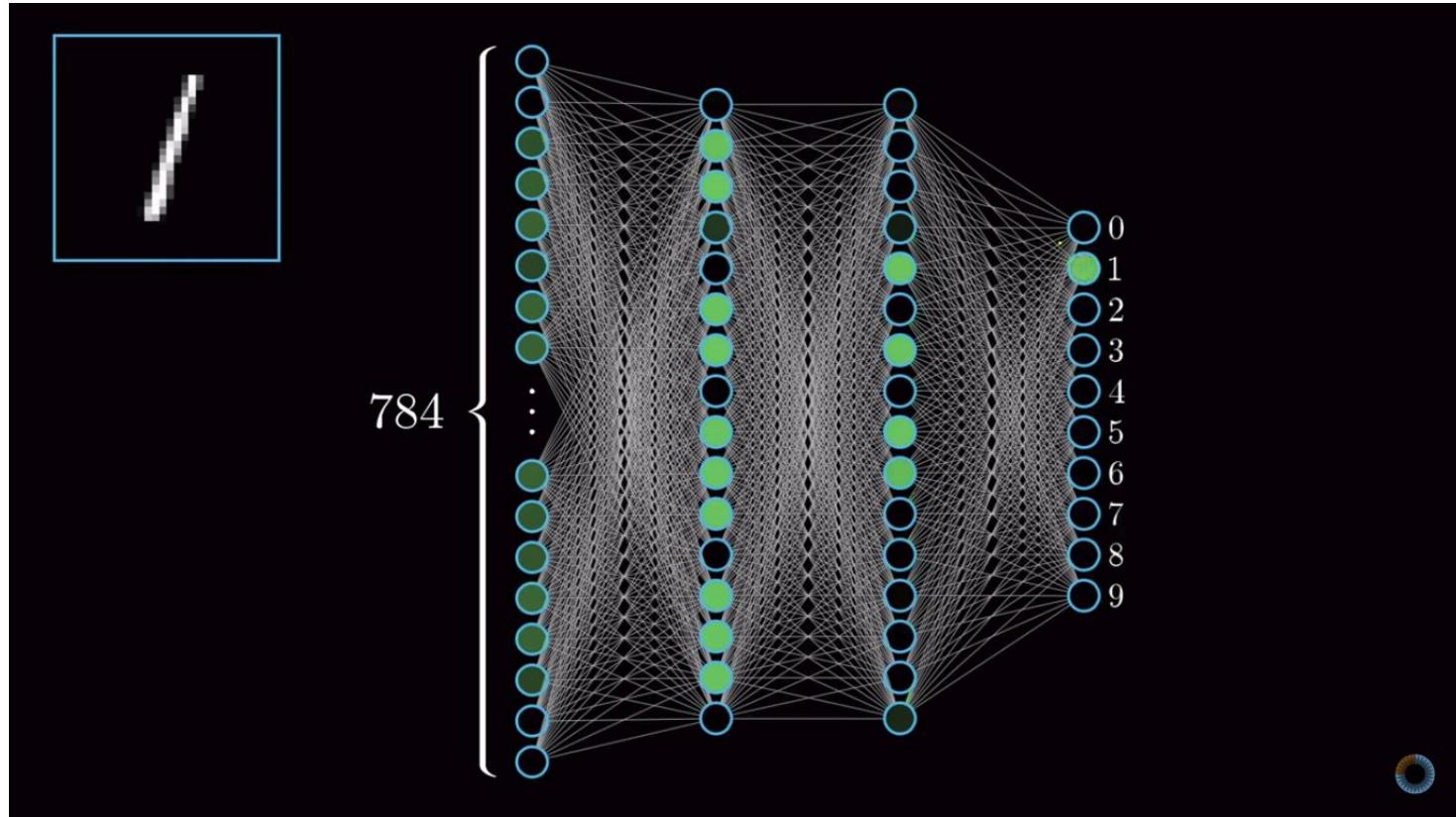
# Redes Neurais - Dígitos escritos à mão



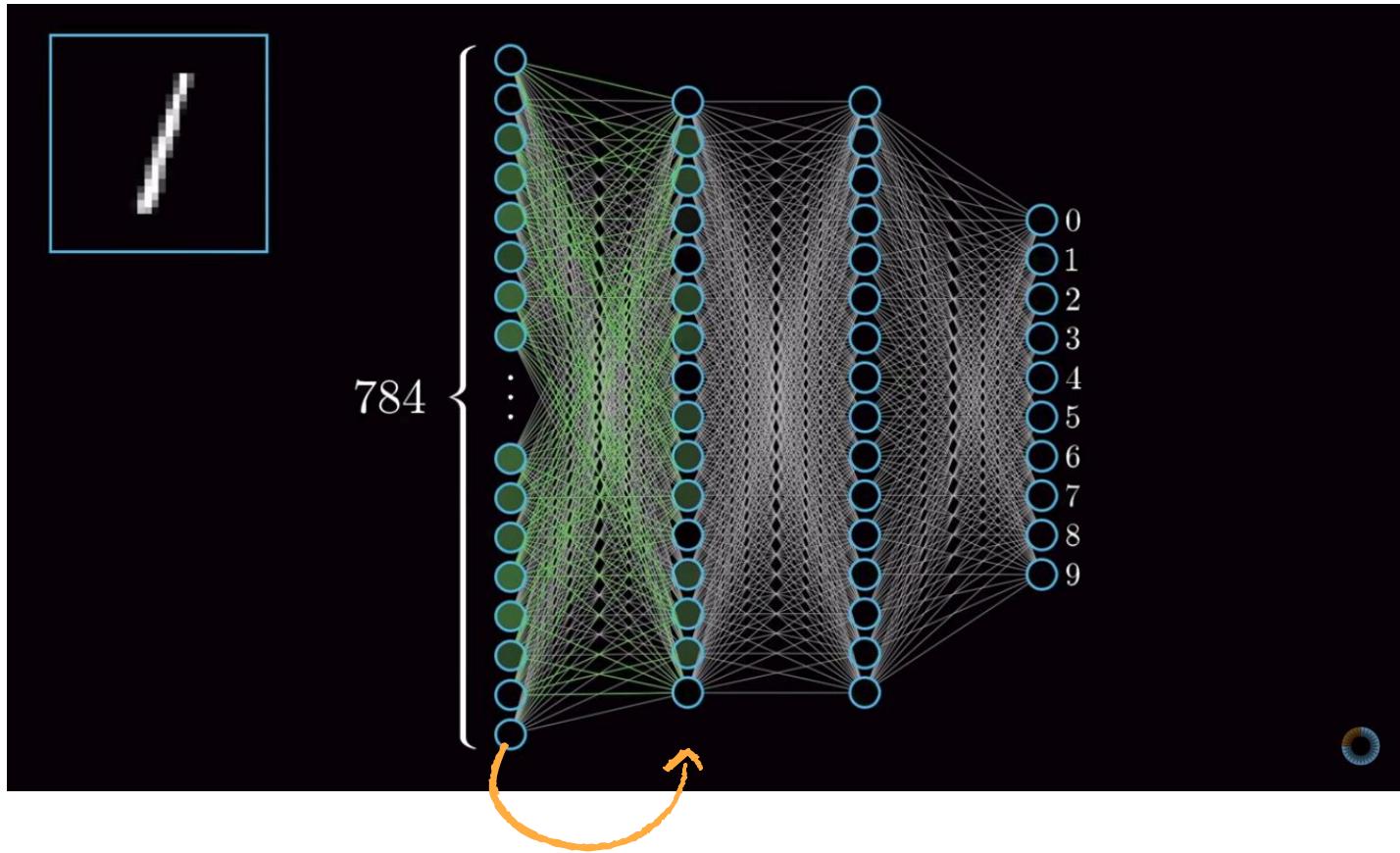
# Redes Neurais - Dígitos escritos à mão



# Redes Neurais - Dígitos escritos à mão

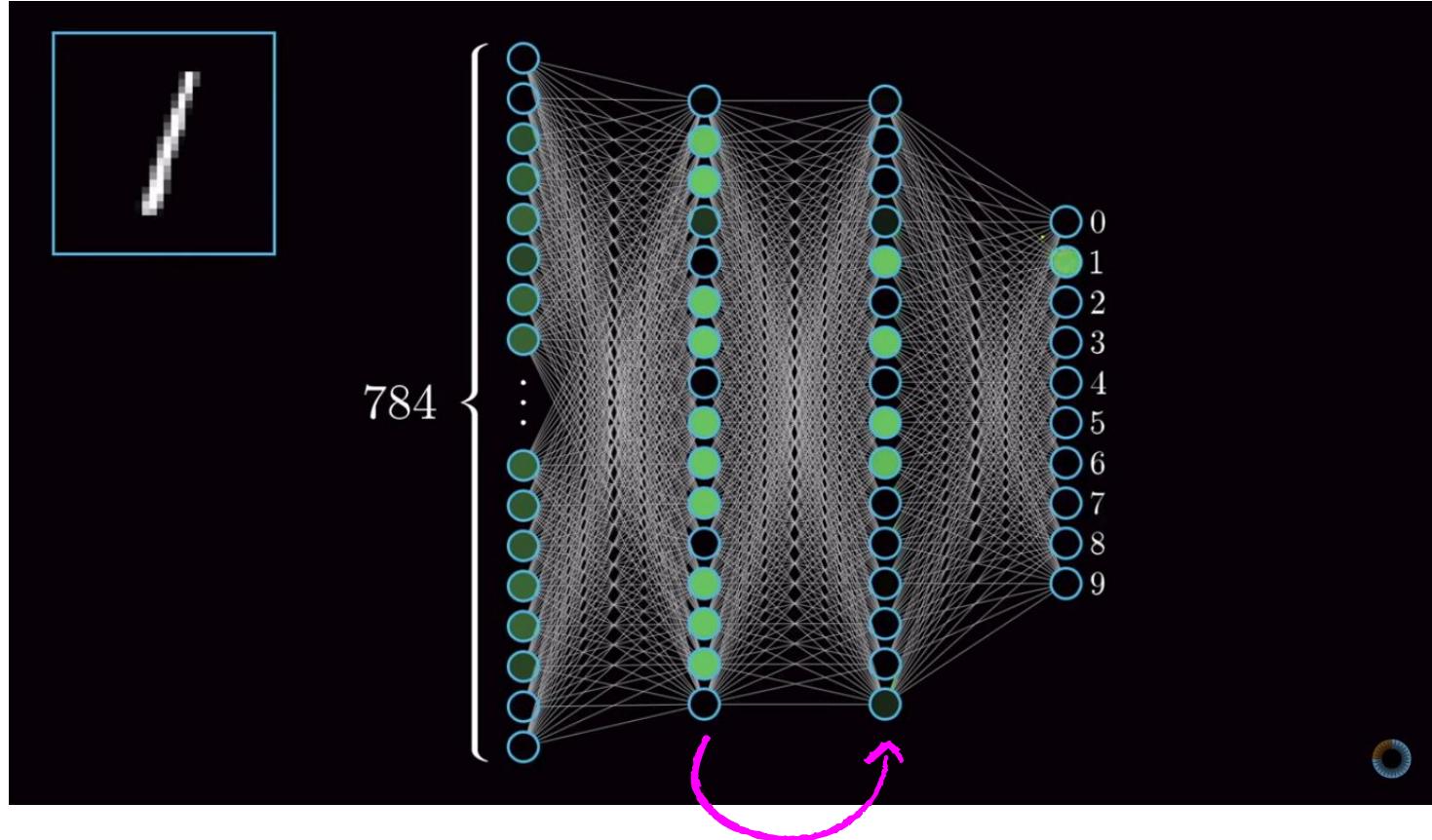


Fonte: [3blue1brown](#)



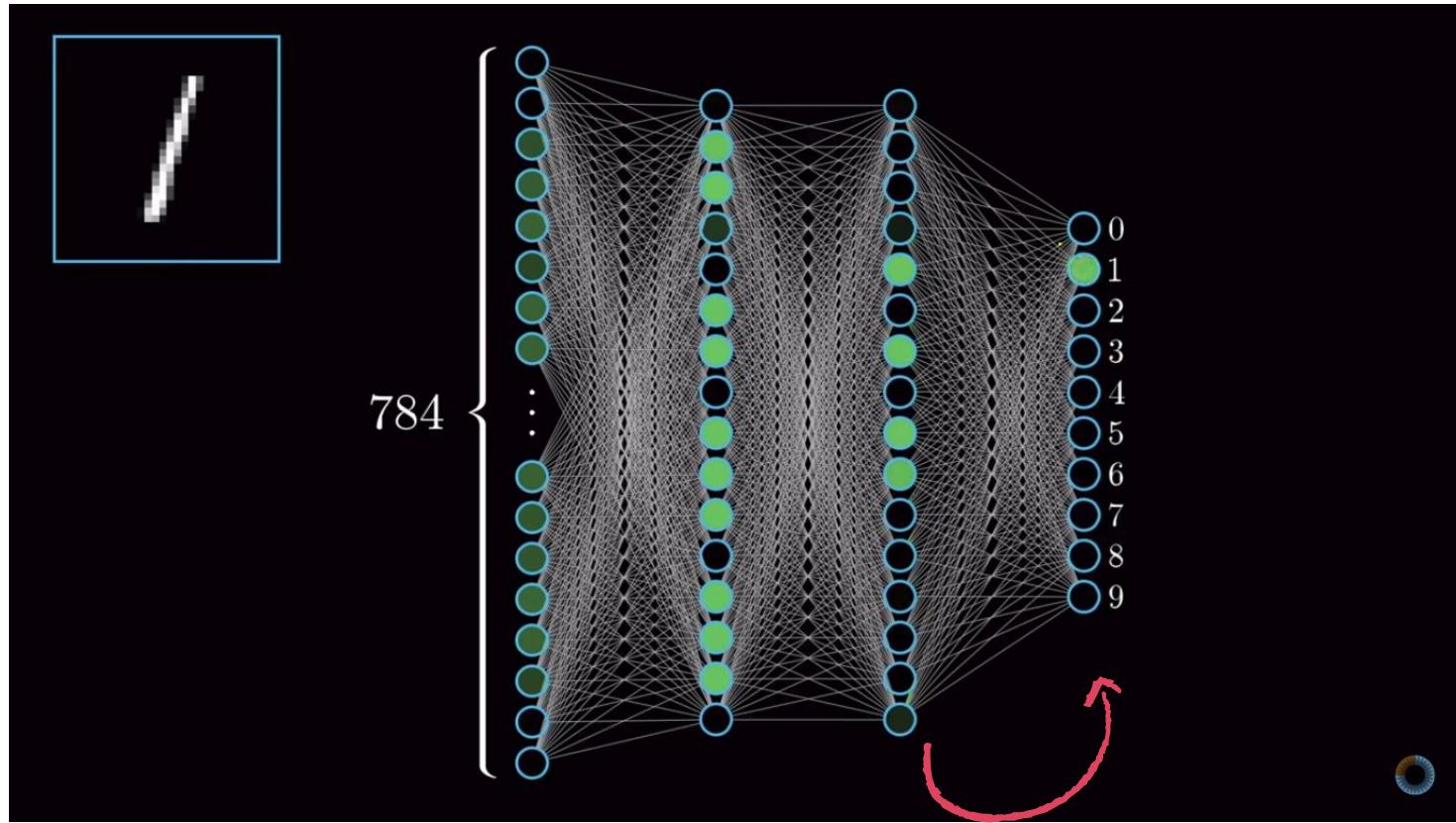
Um padrão de neurônios na camada inicial ativa um  
padrão de neurônios na 1<sup>a</sup> camada oculta

Fonte: [3blue1brown](#)



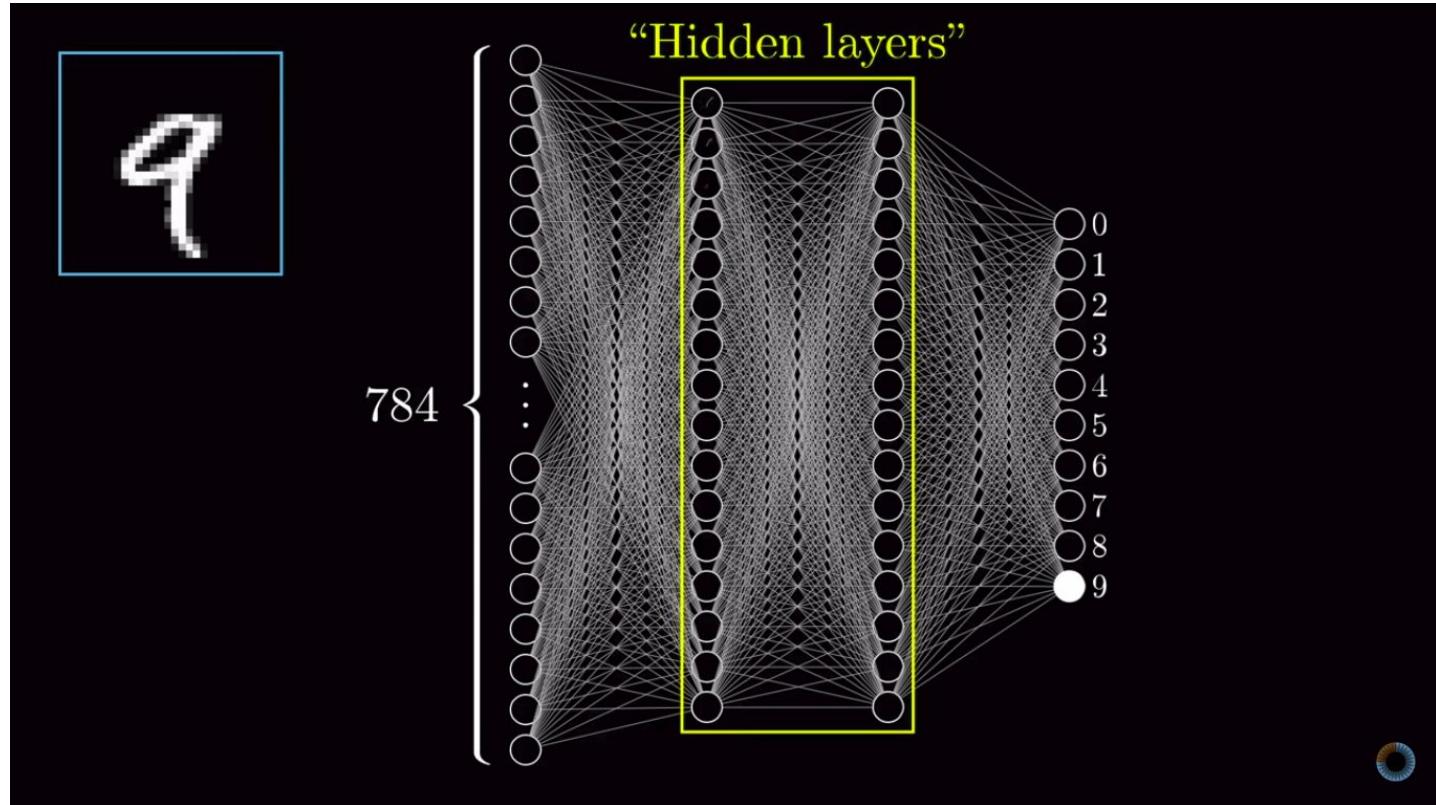
Um padrão de neurônios na 1<sup>a</sup> camada oculta ativa  
um padrão de neurônios na 2<sup>a</sup> camada oculta

Fonte: [3blue1brown](#)

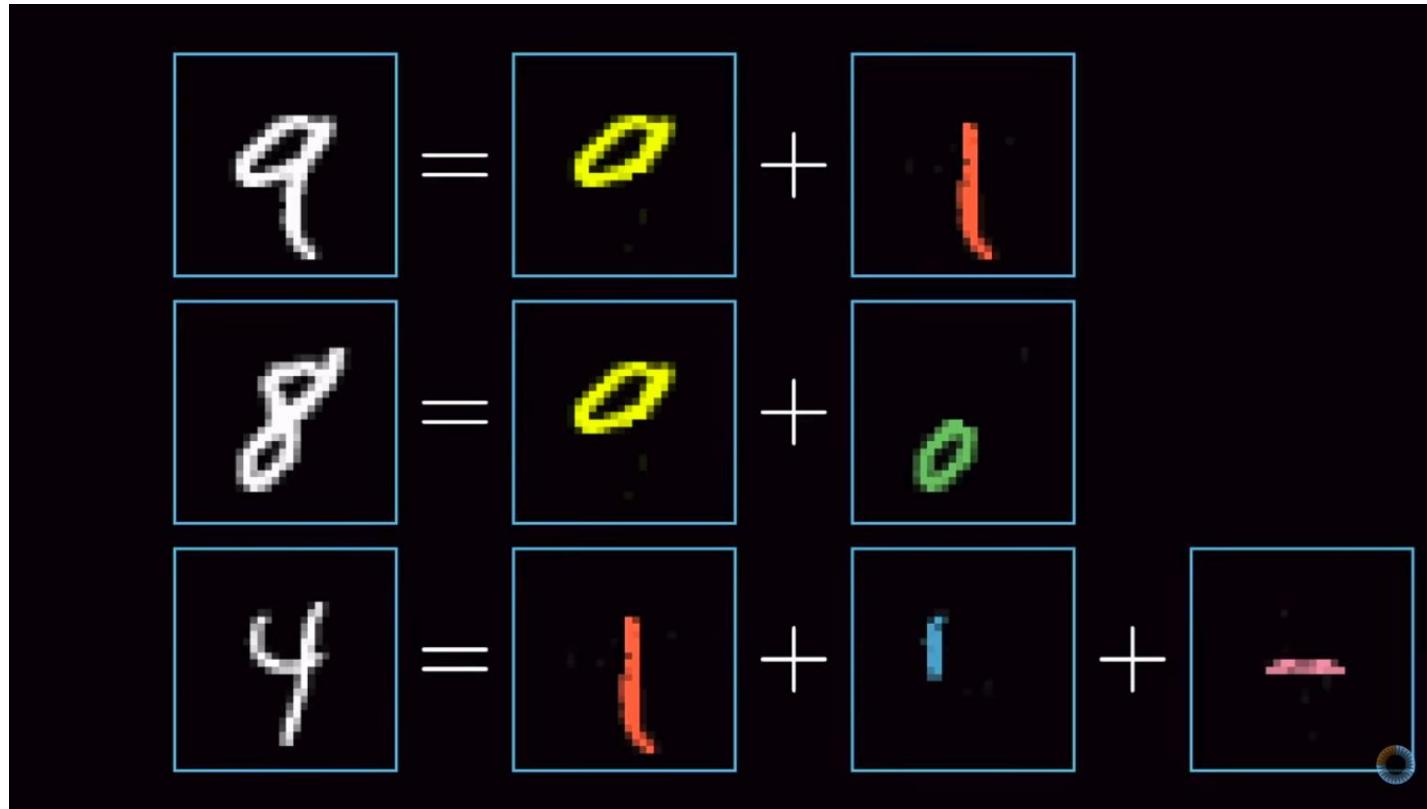


Um padrão de neurônios na 2<sup>a</sup> camada oculta ativa  
um label final

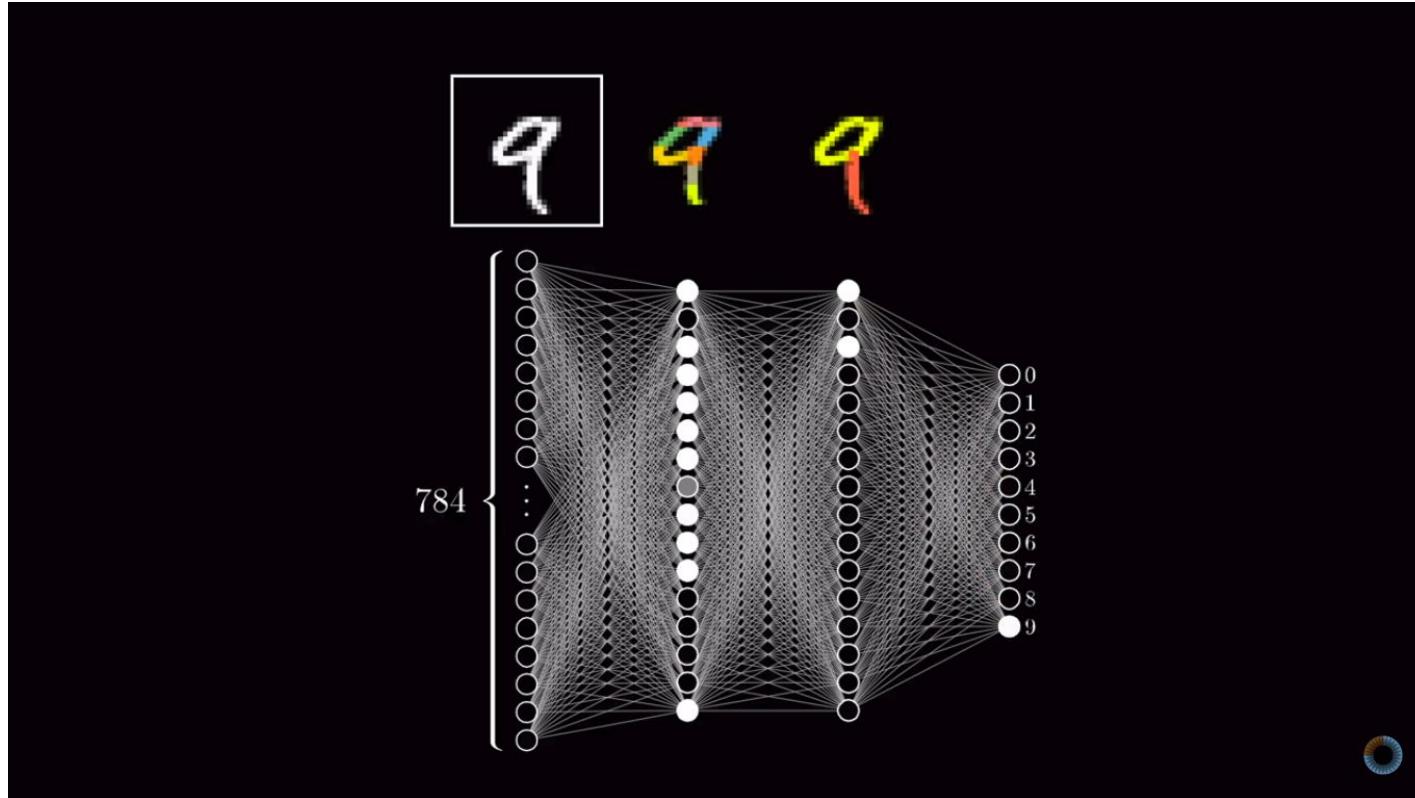
# Redes Neurais - Camadas ocultas



# Redes Neurais - Camadas ocultas



# Redes Neurais - Camadas ocultas



(forte) Suposição: a 1<sup>a</sup> camada reconhece vértices,  
a 2<sup>a</sup> camada reconhece padrões de formas

Fonte: [3blue1brown](#)

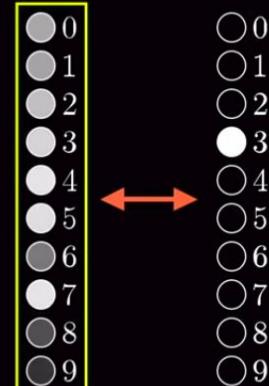
# Redes Neurais - Erro/Custo

Average cost of  
all training data...

Cost of **3**

$$\left\{ \begin{array}{l} (0.67 - 0.00)^2 + \\ (0.66 - 0.00)^2 + \\ (0.77 - 0.00)^2 + \\ (0.85 - 1.00)^2 + \\ (0.90 - 0.00)^2 + \\ (0.88 - 0.00)^2 + \\ (0.50 - 0.00)^2 + \\ (0.91 - 0.00)^2 + \\ (0.32 - 0.00)^2 + \\ (0.19 - 0.00)^2 \end{array} \right.$$

What's the “cost”  
of this difference?



Utter trash

Existe um rótulo  
para comparar.

Trata-se de uma  
abordagem  
supervisionada



Erro: diferença de ativação na camada final e como deveria ser a  
ativação para classificação correta

Fonte: [3blue1brown](#)

# Word2Vec



Estamos interessados na noção de **significado**.

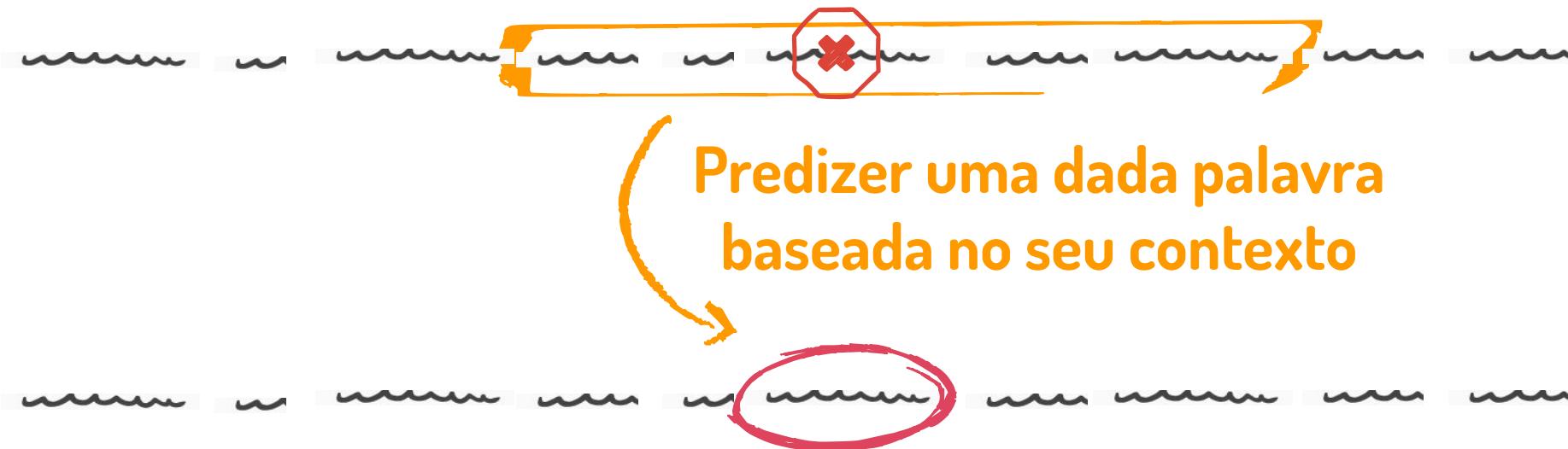
Vamos aprender os significados das palavras com uma ideia já usual:

A partir do **contexto**

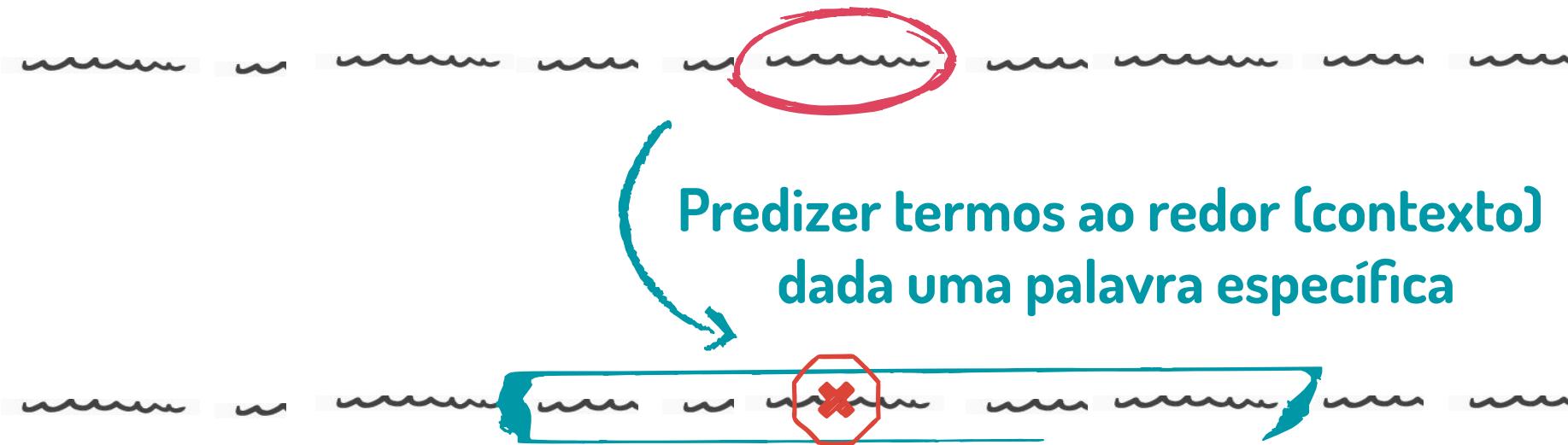
# Word2Vec



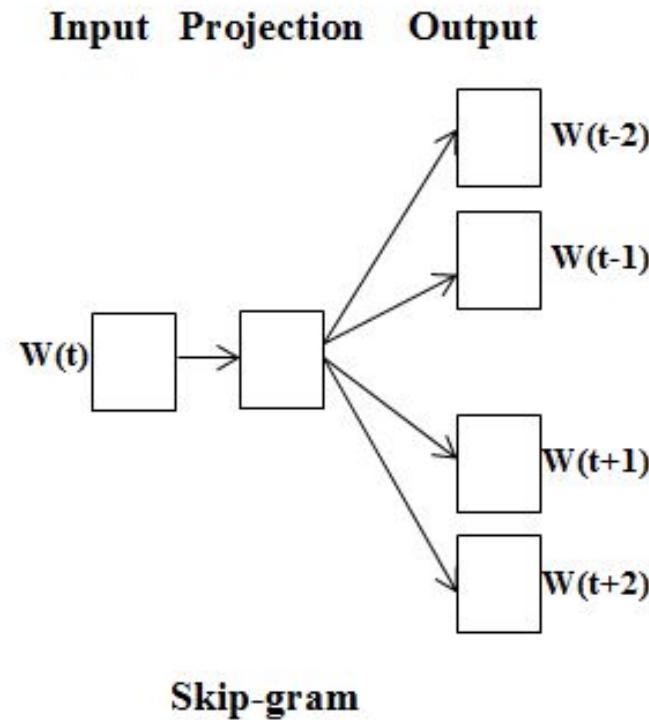
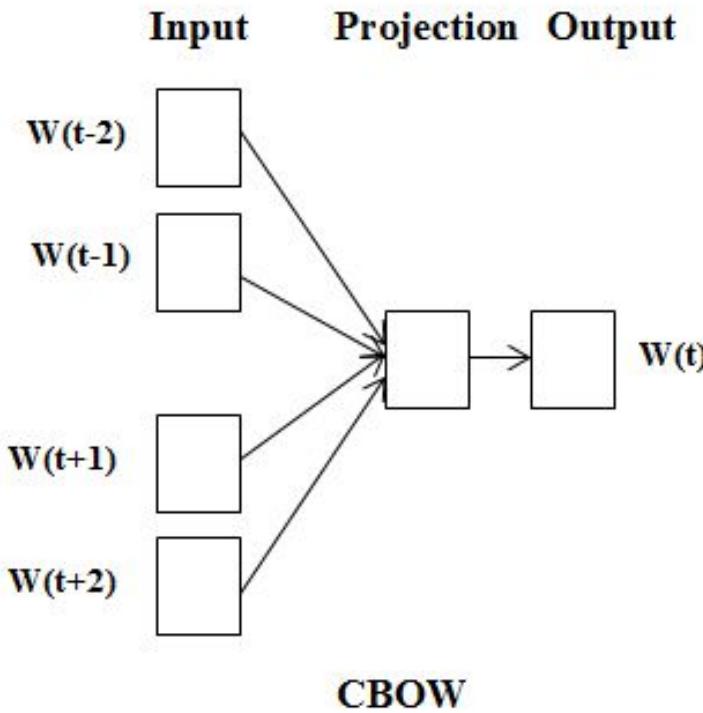
# Word2Vec



# Word2Vec



# Arquiteturas para modelar palavras

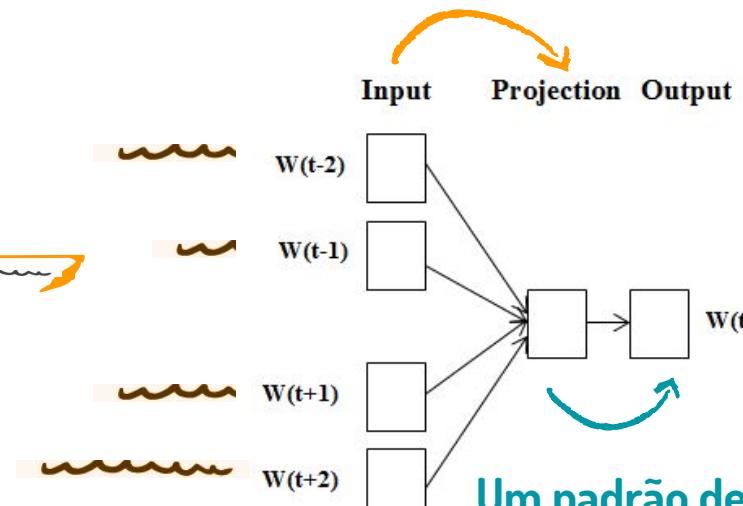


# CBOW

Um padrão de neurônios na camada inicial ativo a partir dos termos vizinhos ativa

um padrão de neurônios na camada oculta

Erro/Custo do emprego dos parâmetros utilizados

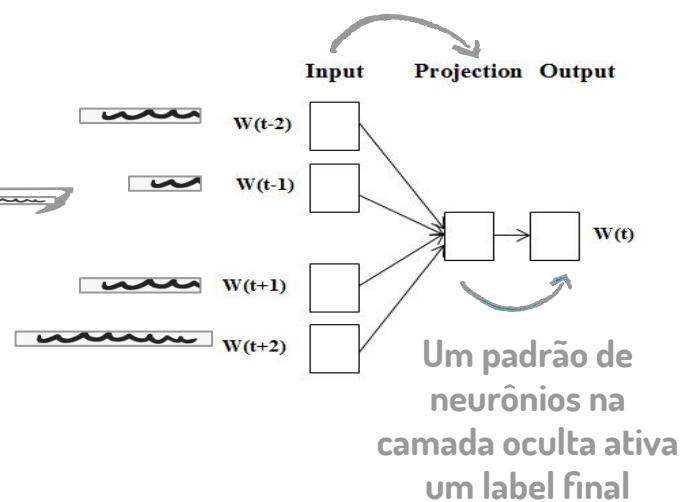


Um padrão de neurônios na camada oculta ativa um label final

# CBOW

Um padrão de neurônios na camada inicial ativo a partir dos termos vizinhos  
ativa

um padrão de neurônios na camada oculta

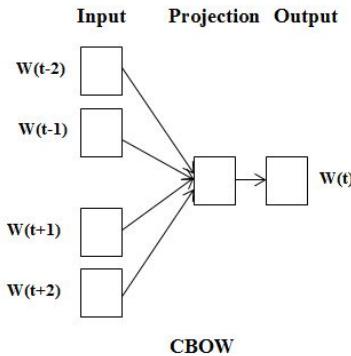


Erro/Custo do emprego dos parâmetros utilizados



**BACKPROPAGATION**

# CBOW



Melhor para associações  
Sintáticas

CBOW

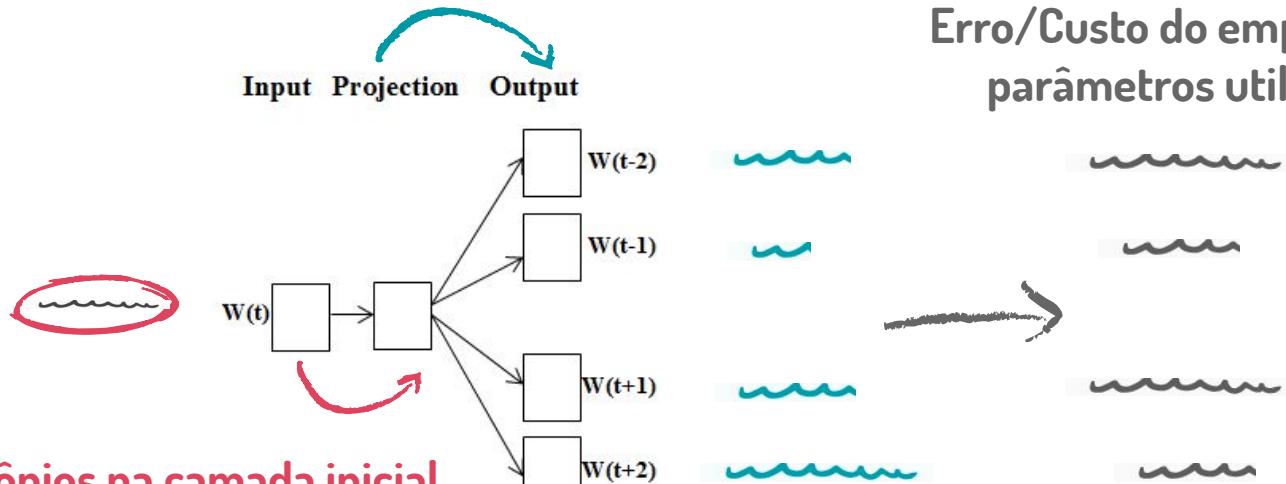
Skip  
gram

Type of relationship	Word Pair 1		Word Pair 2	
Adjective to adverb	apparent	apparently	rapid	rapidly
Opposite	possibly	impossibly	ethical	unethical
Comparative	great	greater	tough	tougher
Superlative	easy	easiest	lucky	luckiest
Present Participle	think	thinking	read	reading
Nationality adjective	Switzerland	Swiss	Cambodia	Cambodian
Past tense	walking	walked	swimming	swam
Plural nouns	mouse	mice	dollar	dollars
Plural verbs	work	works	speak	speaks

Fonte: Artigo 2013

# Skip-gram

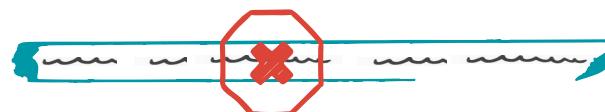
Um padrão de neurônios na camada oculta ativa N vizinhos padrões de label final



Erro/Custo do emprego dos parâmetros utilizados

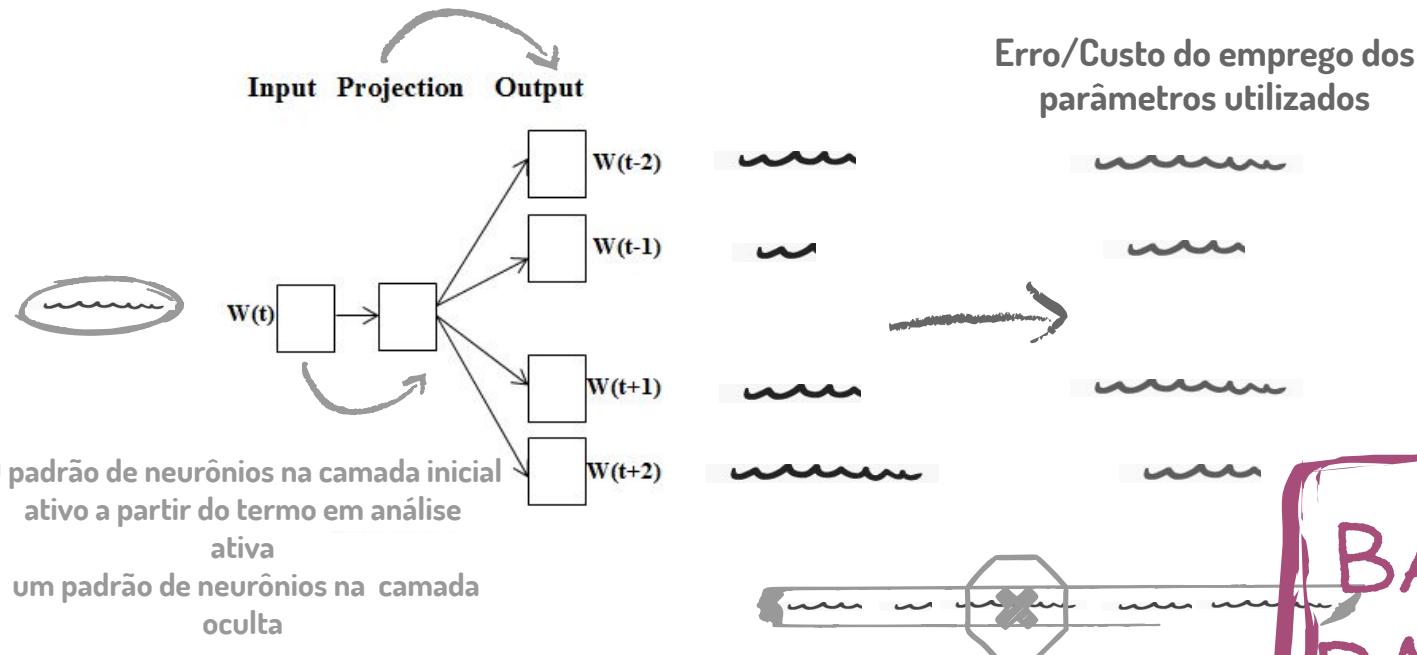
O padrão de neurônios na camada inicial ativo a partir do termo em análise ativa

um padrão de neurônios na camada oculta



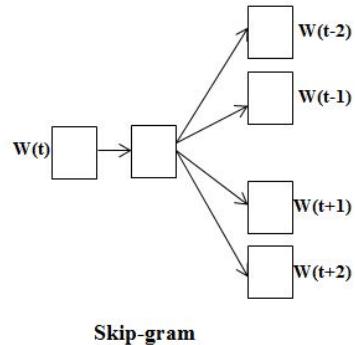
# Skip-gram

Um padrão de neurônios na camada oculta ativa N vizinhos padrões de label final



# Skip-gram

Input   Projection   Output



Melhor para associações  
Semânticas

Type of relationship	Word Pair 1		Word Pair 2	
Common capital city	Athens	Greece	Oslo	Norway
All capital cities	Astana	Kazakhstan	Harare	Zimbabwe
Currency	Angola	kwanza	Iran	rial
City-in-state	Chicago	Illinois	Stockton	California
Man-Woman	brother	sister	grandson	granddaughter

# Redes Neurais - Erro/Custo e Word2Vec

56

Average cost of all training data...  
Cost of 3

$$\left\{ \begin{array}{l} (0.67 - 0.00)^2 + \\ (0.66 - 0.00)^2 + \\ (0.77 - 0.00)^2 + \\ (0.85 - 1.00)^2 + \\ (0.90 - 0.00)^2 + \\ (0.88 - 0.00)^2 + \\ (0.50 - 0.00)^2 + \\ (0.91 - 0.00)^2 + \\ (0.32 - 0.00)^2 + \\ (0.19 - 0.00)^2 \end{array} \right.$$

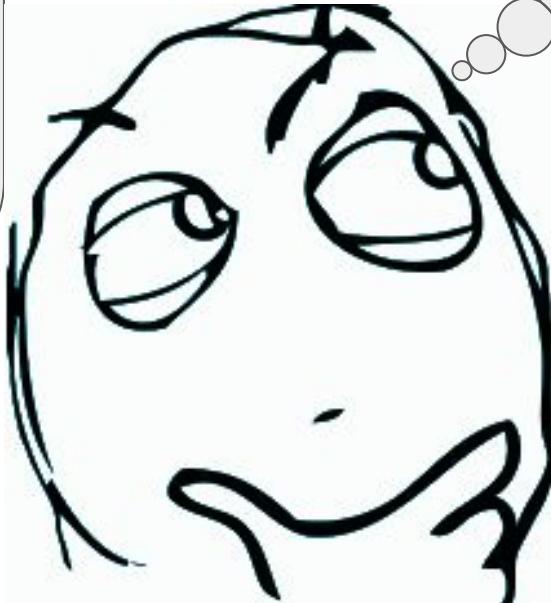
What's the "cost" of this difference?  
Utter trash

Existe um rótulo para comparar.

Trata-se de uma abordagem supervisionada



É impossível rotular contexto, ainda mais com relações sintáticas e semânticas



# Comentários

## GloVe- Global Vectors

Representação  
de **Stanford**

Muito utilizado

[Text2vec - Glove](#)

## Referência teórica

[IllustratedWord2  
Vec](#)

## Embeddings treinados em Português

[linkUSP](#)

[Parpinelli](#)

## Tensorflow

Ferramenta da  
**Google** para Mac  
Learn  
[detalhes](#)

[Rword2vec  
\(Google\)](#)

[Train\\_word2vec](#)

[wordVectors/  
primaryobjects](#)

[Tensorflow.rstudio](#)

# Word2Vec - Representação CBOW 50

```

Activities Text Editor *
Loading cbow_50.txt from ~/Downloads
sex 23:14 ● *cbow_50.txt
Save Cancel

929606 50
</s> -0.0188088 -0.088949 -0.088157 0.087802 -0.084775 0.088168 0.0880275 0.0883503 0.082693 -0.000599 0.082116 0.0807509 0.006414 0.0803036 0.080344 0.0805029 -0.005577 -0.007318 0.001901 0.007887 -0.006728
0.000236 -0.000243 -0.000462 0.000253 -0.000403 -0.000337 0.000473 0.0006524 0.006939 -0.000160 0.002124 0.002150 -0.003164 -0.0003439 -0.0009705 -0.0009343 -0.009820 0.009972 -0.0009068 -0.006134 -0.0008558
, 0.094044 0.195893 -0.155710 -0.1445906 -0.160736 0.000748 -0.130807 -0.004585 -0.121695 -0.087481 0.018703 0.081959 0.106759 -0.016001 -0.009896 -0.162624 0.014549 0.010382 0.208517 -0.029570 0.180953
, 0.123465 0.019036 -0.049734 0.019022 0.131606 -0.082430 0.0006308 0.158126 -0.151579 0.035982 0.015035 0.064862 0.000235 0.076560 0.160227 0.067308 0.141098 0.022217 -0.103135 -0.029701 0.006607 0.287858
, 0.100301 -0.087588 0.054029 0.039782 0.035716 0.059295 0.183569
...
The matrix continues with many more rows and columns of numerical values representing word embeddings. The terminal also shows standard file navigation commands like Open, Save, and Cancel.

```

PlainText Tab Width: 8 Ln 10658, Col 22 INS

# Recap - Definições

0 que faz	Definição	Como funciona	Para quê serve	Resultado
Cria vetores que são representações numéricas dos atributos (termos/palavras) sem intervenção humana	Posiciona termos (palavras) em um espaço dimensional tal que as localizações dos termos são determinadas pelos seus significados e os vazios e distâncias no espaço também tem significado	Rede neural de duas camadas que processa texto	Criar representação de strings para algoritmos de Machine Learning ou Deep Learning	Agrupar vetores de termos similares no espaço de atributos

# Recap - Vantagens



## Bons exemplos

Exemplos com diferentes formas de similaridade (semântica e sintática)



## Popular

Google desenvolveu uma ferramenta, tensorflow, e a disponibilizou



## Grandes Datasets

A técnica lida bem com bases grandes, típicas da realidade atual (6B tokens, 1M termos mais freq)



## Rápido

De: 08 semanas  
Para: 01-03 dias



## Relações preservadas entre línguas

Permite tradução termo-a-termo



# Obrigada! Dúvidas?

[larissa.sayuri.fcs@gmail.com](mailto:larissa.sayuri.fcs@gmail.com)

[larissa@maiemarketing.com.br](mailto:larissa@maiemarketing.com.br)