

# DEEP LEARNING

Anatomia de um projeto de Machine Learning  
aplicado à recomendação de notícias

Mar/2018

# Fabrício V Matos

[linkedin.com/in/fabriciovargasmatos](https://linkedin.com/in/fabriciovargasmatos)

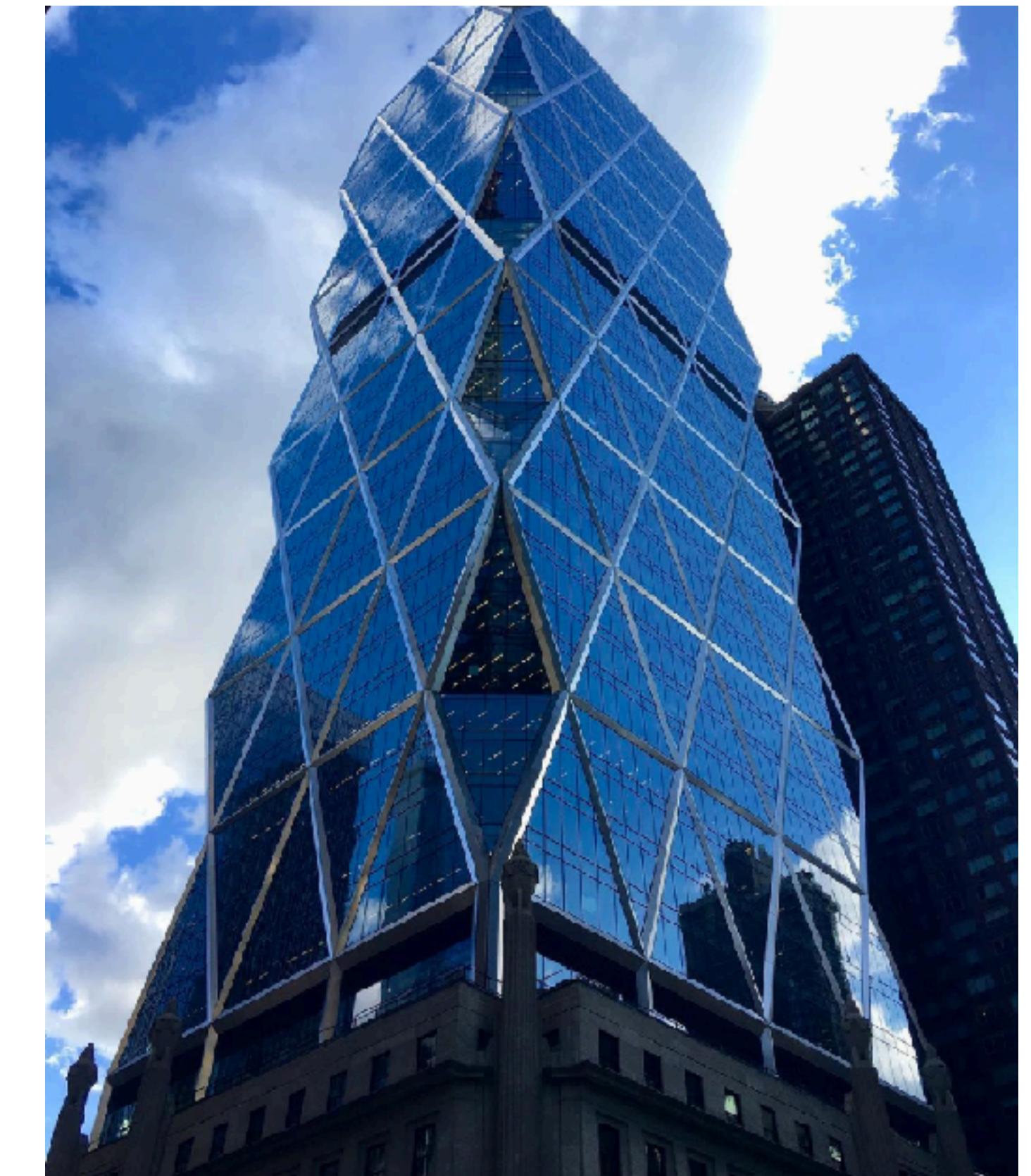
Vix Data *meetup* #7



# Recomendação Personalizada de Notícias



# 29 Estações de TV



# Sede em NY

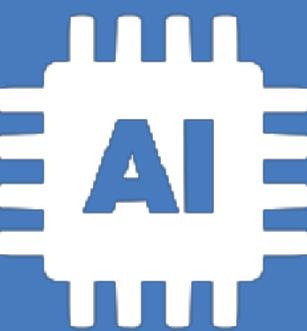


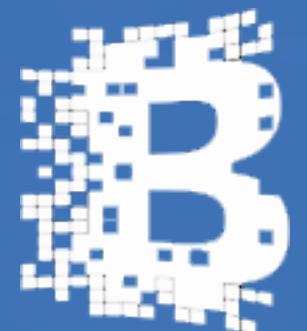
mutual.life

SOUNDS LIKE INSURANCE.  
BUT IT'S SHARED PROTECTION.

# Pilares

|  Mobile

|  Inteligência Artificial

|  Blockchain





# Amazon\_rainforest

## The Stanford Question Answering Dataset

The Amazon rainforest (Portuguese: Floresta Amazônica or Amazônia; Spanish: Selva Amazónica, Amazonía or usually Amazonia; French: Forêt amazonienne; Dutch: Amazoneregenwoud), also known in English as Amazonia or the Amazon Jungle, is a moist broadleaf forest that covers most of the Amazon basin of South America. This basin encompasses 7,000,000 square kilometres (2,700,000 sq mi), of which 5,500,000 square kilometres (2,100,000 sq mi) are covered by the rainforest. This region includes territory belonging to nine nations. The majority of the forest is contained within Brazil, with 60% of the rainforest, followed by Peru with 13%, Colombia with 10%, and with minor amounts in Venezuela, Ecuador, Bolivia, Guyana, Suriname and French Guiana. States or departments in four nations contain "Amazonas" in their names. The Amazon represents over half of the planet's remaining rainforests, and comprises the largest and most biodiverse tract of tropical rainforest in the world, with an estimated 390 billion individual trees divided into 16,000 species.

Human Performance: 82.304

AI from Microsoft Research (Fev/18): **82.849**

...

AI from Microsoft Research (Out/17): 76.461

Which name is also used to describe the Amazon rainforest in English?

Ground Truth Answers: also known in English as Amazonia or the Amazon Jungle, Amazonia or the Amazon Jungle Amazonia

Prediction: Amazonia

How many square kilometers of rainforest is covered in the basin?

Ground Truth Answers: 5,500,000 square kilometres (2,100,000 sq mi) are covered by the rainforest. 5,500,000 5,500,000

Prediction: 5,500,000

How many nations control this region in total?

Ground Truth Answers: This region includes territory belonging to nine nations. nine nine

Prediction: nine

How many nations contain "Amazonas" in their names?

Ground Truth Answers: States or departments in four nations contain "Amazonas" in their names. four four

Prediction: four

# ARTIFICIAL LAWYER

AI AND LEGAL AUTOMATION NEWS + VIEWS



SEARCH ...

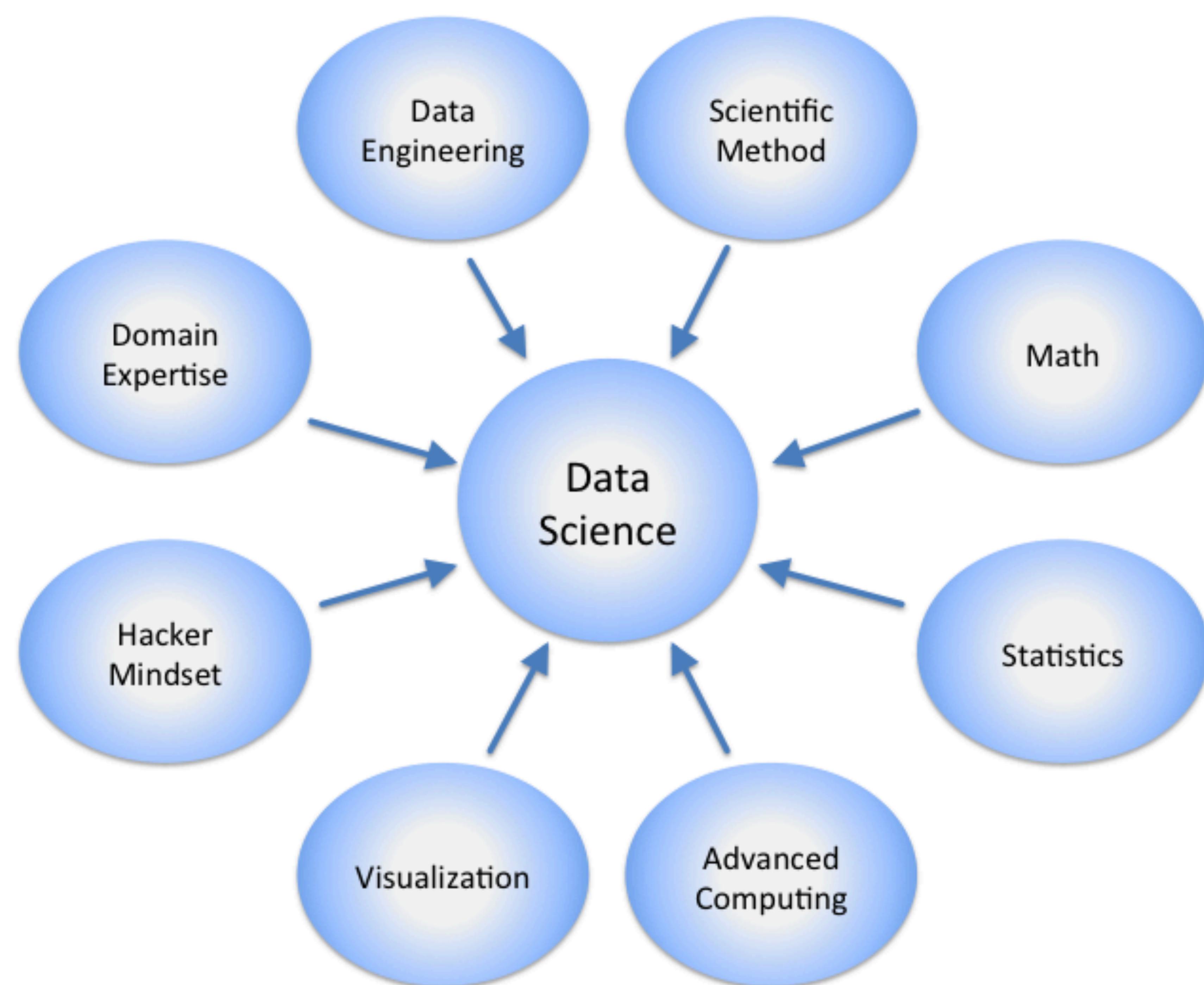
## AI Beats Human Lawyers in CaseCrunch Prediction Showdown + DATA UPDATES

⌚ 28th October 2017 🚑 artificiallawyer 📁 Legal AI Prediction 💬 9

**CaseCrunch** is proud to announce the results of the lawyer challenge. **CaseCruncher Alpha scored an accuracy of 86.6%. The lawyers scored an accuracy of 62.3%.**

Over 100 commercial London lawyers signed up for the competition and made over 750 predictions over the course of a week in an unsupervised environment.

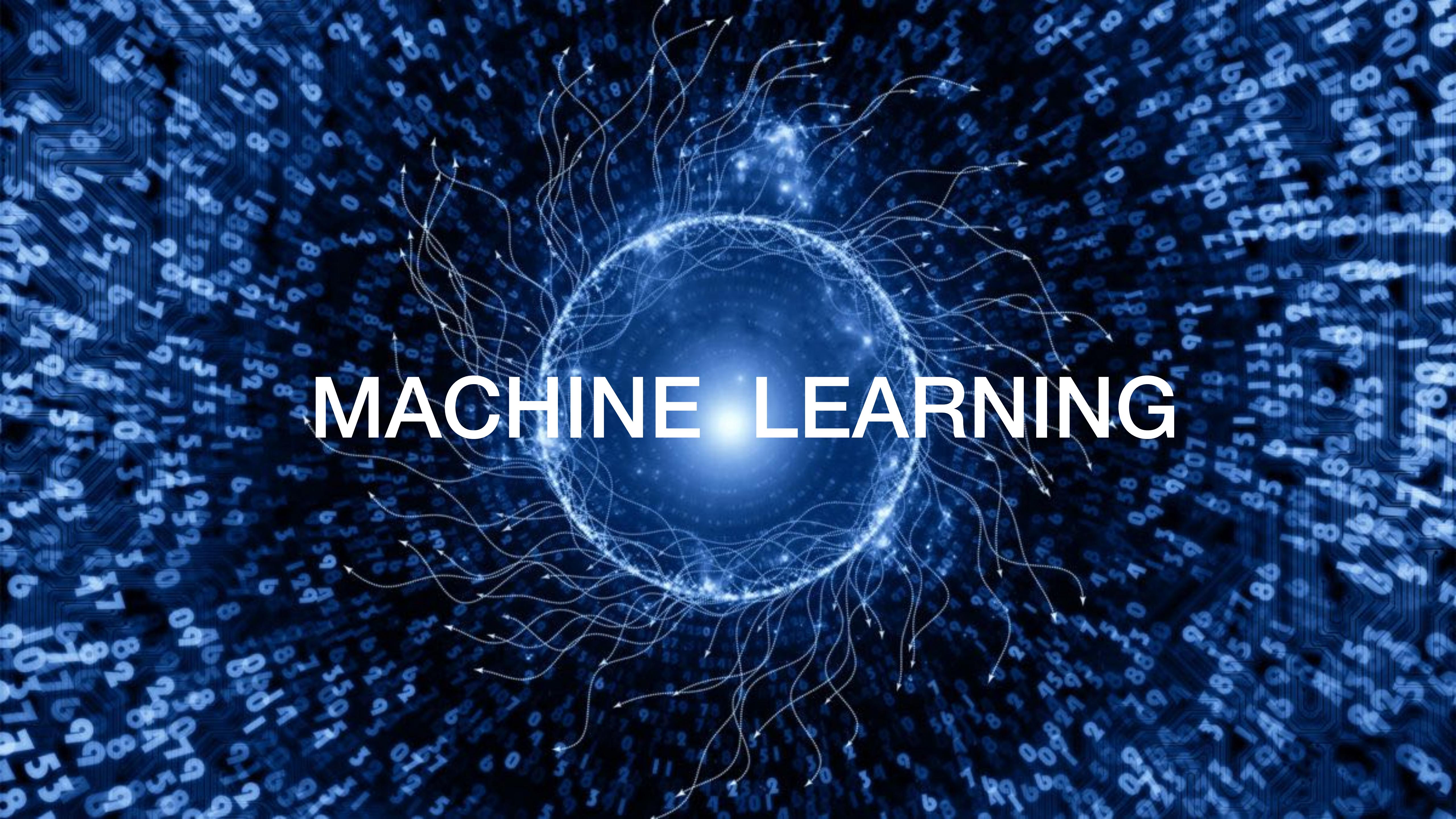
The problems were real complaints about PPI mis-selling decided and published by the Financial Ombudsman Service (FOS) under the Freedom Of Information Act. (*N.B. Payment Protection Insurance, or PPI, has been a massive issue in the UK, with banks having to pay customers billions of pounds in total refunds for making consumers take on insurance products they never required.*)





**Por onde começar?**

# MACHINE LEARNING



Artificial  
Intelligence

Machine  
Learning

Data Science

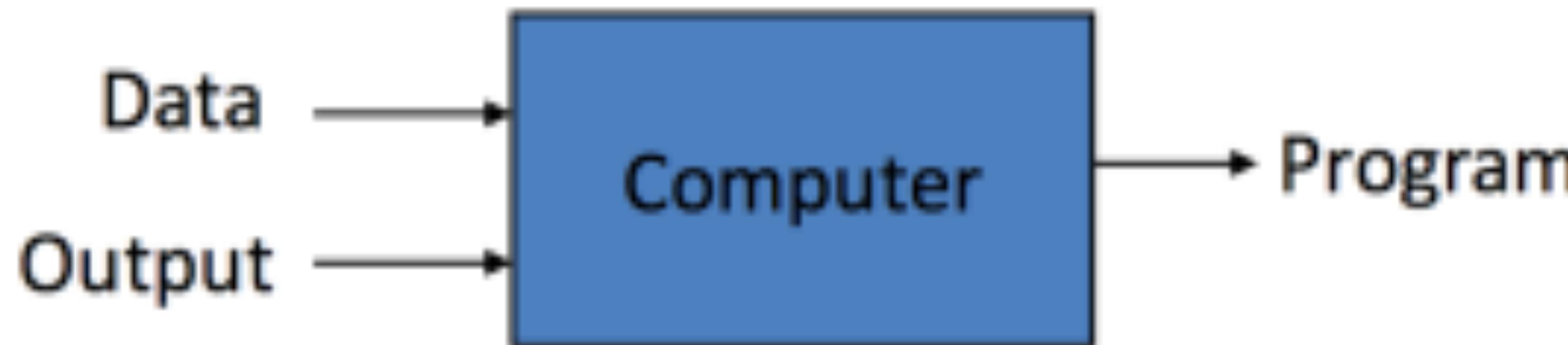
Deep  
Learning

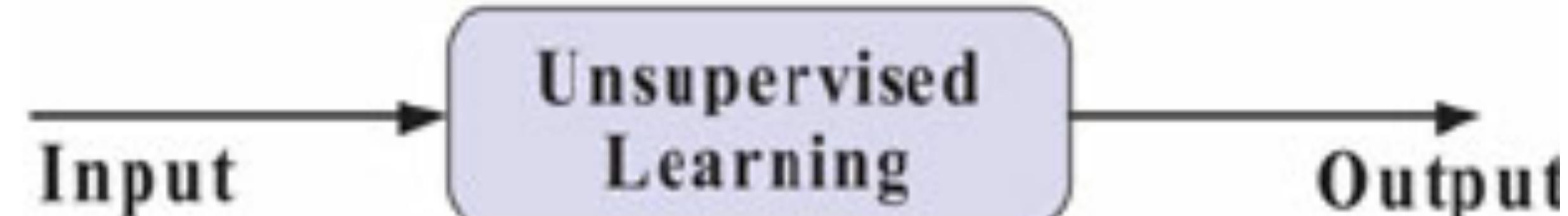
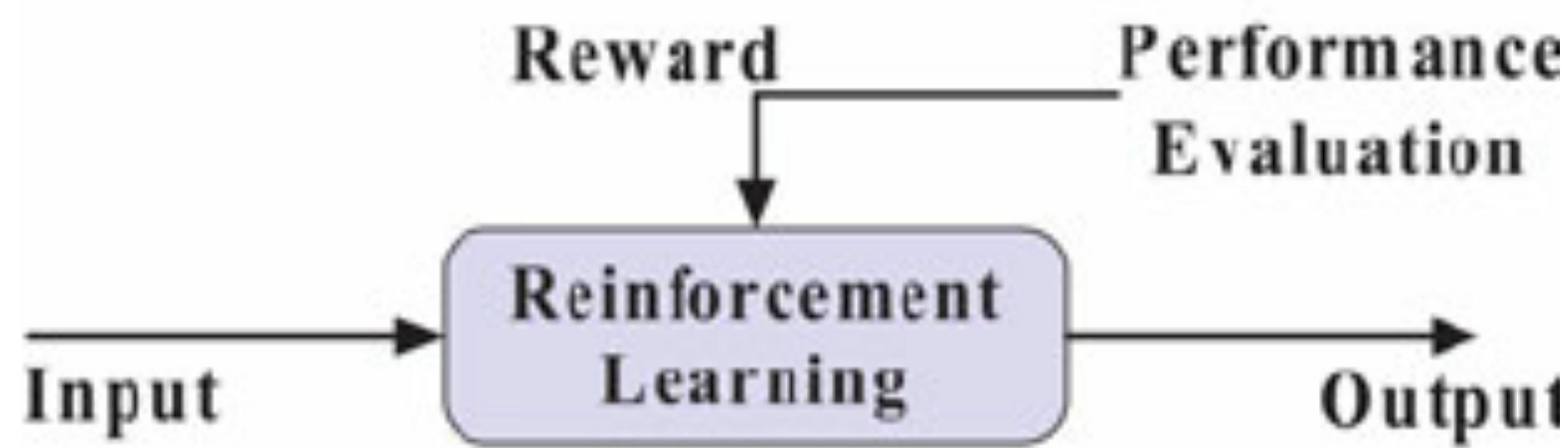
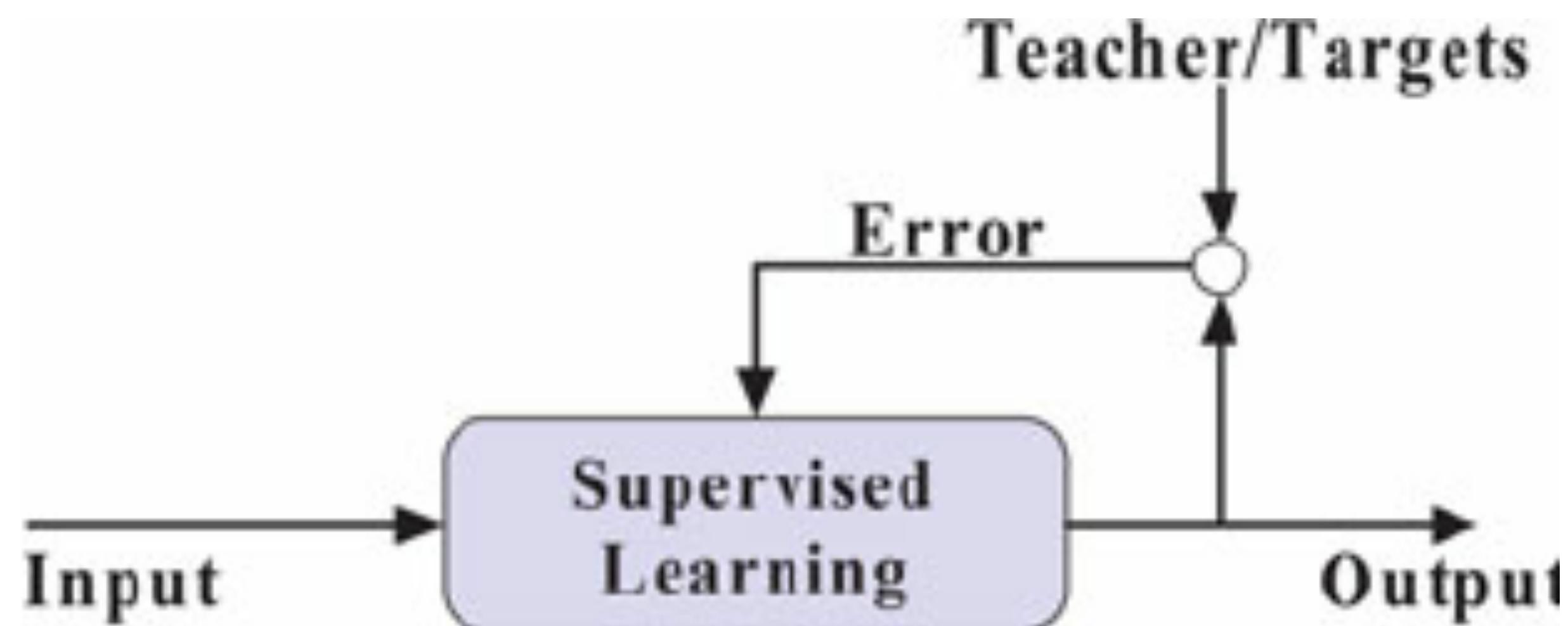
# ML: Novo Paradigma

## Traditional Programming

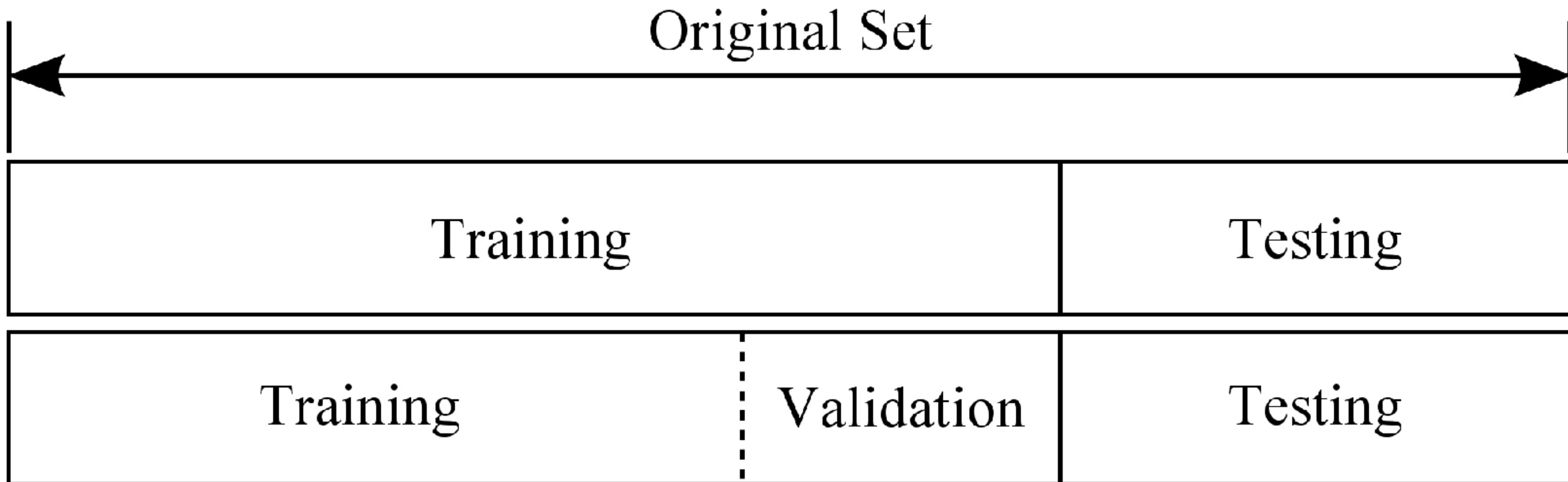


## Machine Learning

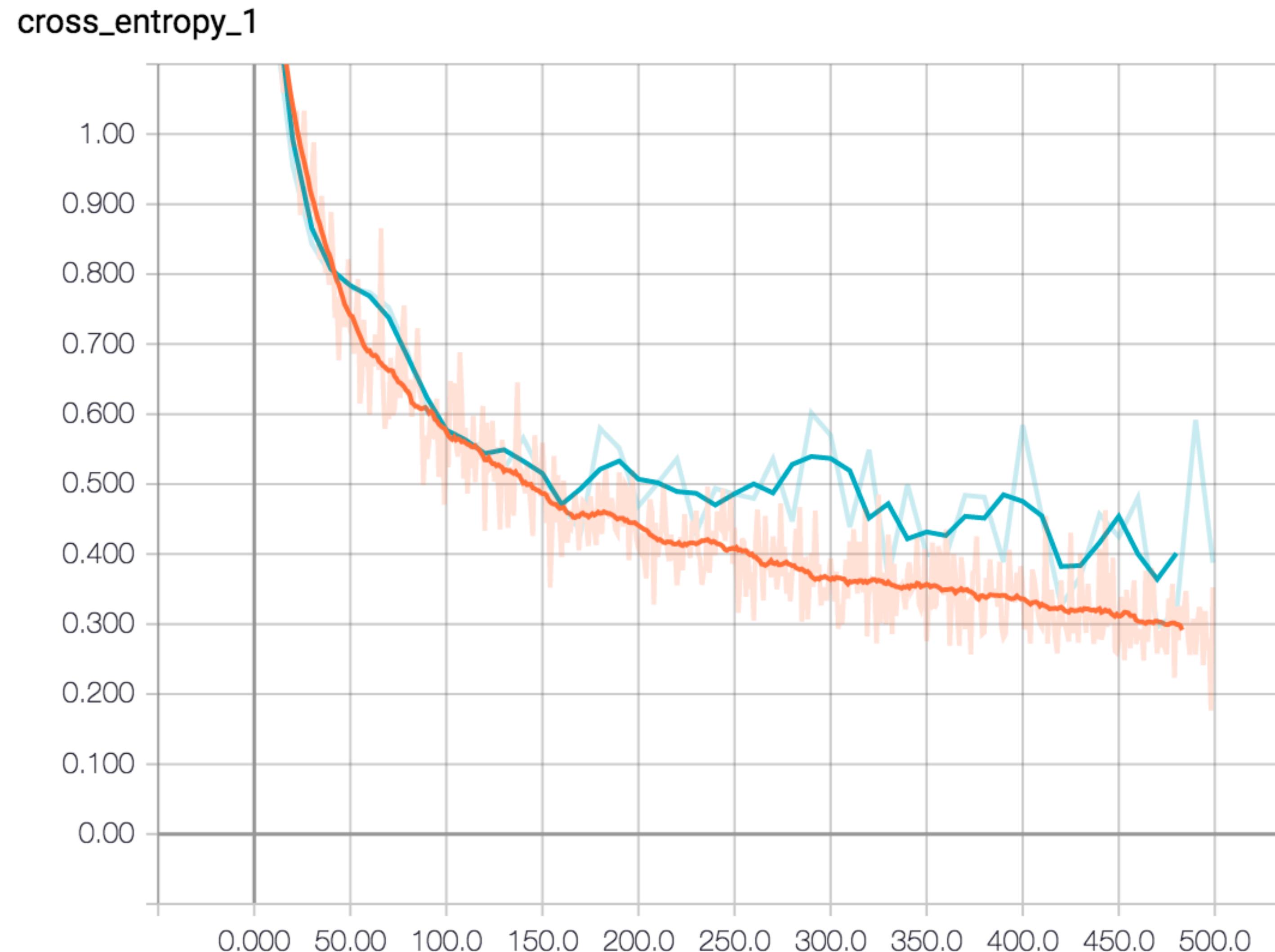




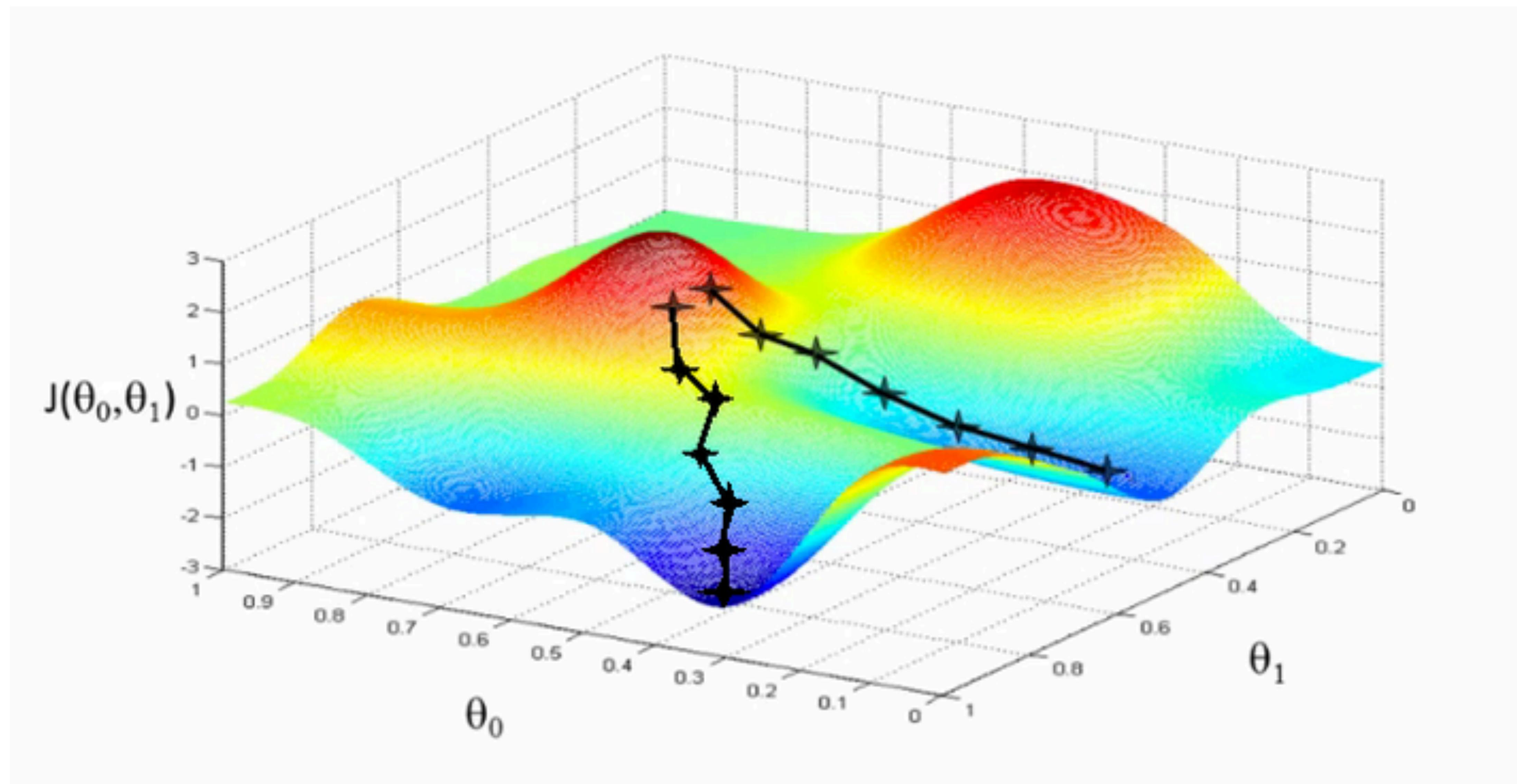
# Treino supervisionado: Datasets



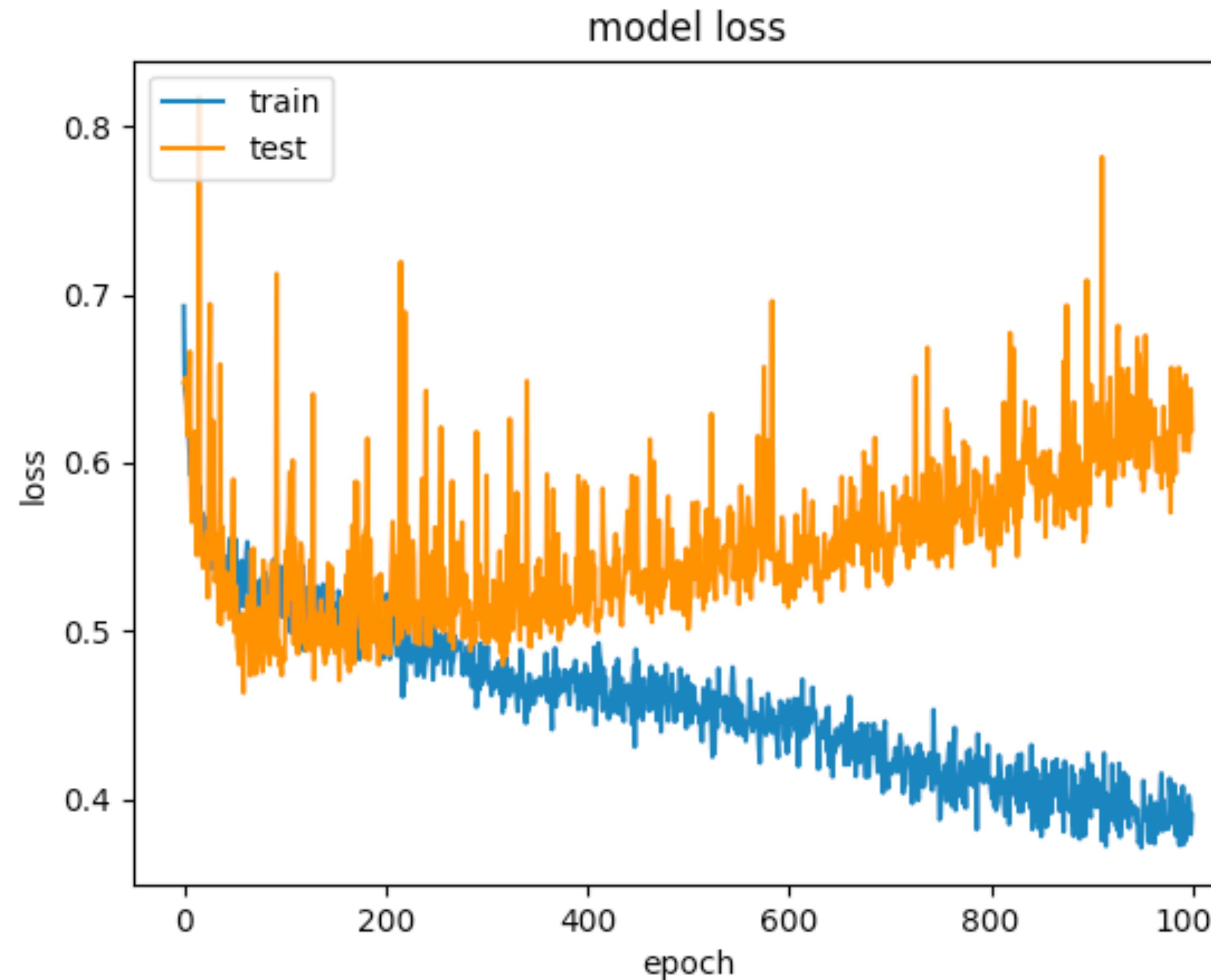
# Treino supervisionado: Cost function



# Otimizador: Stochastic Gradient Descent



# Overfit (“mente fechada”)



# **Estudo de caso:**

## **Recomendação personalizada de notícias**

# Aplicação

The screenshot shows a news article from WCVB.com. At the top, there's a banner for "NOWCAST WCVB On Demand" and a "Watch on Demand" button. The main header reads "SEVERE WEATHER: There are currently 3 active weather alerts". Below this, the main headline is "Trump says Comey's 'leaks' are more prevalent". A subtext below the headline states: "Trump is again challenging Comey after the ousted FBI director's testimony before the Senate Intelligence committee last week". An AP logo indicates the source was updated at 9:41 AM EDT on Jun 11, 2017. At the bottom, there's a video thumbnail showing two men at a podium with flags, and social sharing icons for Facebook and Twitter.

The screenshot shows a news article from WCVB.com. The headline is "Comey's 'leaks will be far more prevalent than anyone ever thought possible.' Trump says on Twitter, 'Totally illegal? Very 'cowardly!'" Below the headline, there's a summary: "Trump is again challenging Comey after the ousted FBI director's testimony before the Senate Intelligence committee last week. While many of Trump's Republican allies have found Comey's testimony credible, the president has called the man he fired a liar and a 'leaker.' Comey said during his testimony that he asked a friend to release contents of the memos he'd written about his conversations with the president to a reporter. He contended that information was not classified or otherwise protected." At the bottom, there's a red arrow pointing to a section titled "Notícias recomendadas" which includes a link to an article about a truck explosion.

# Desafios específicos

- Grande volume de novos items todos os dias
- Notícias perdem relevância muito rapidamente
- Usuários não se autenticam (escassez de dados históricos)
- Conteúdos devem ser atualizados pelo menos 50 vezes/dia
- API de recomendação escalável
- Implementar v1.0 em 3 meses

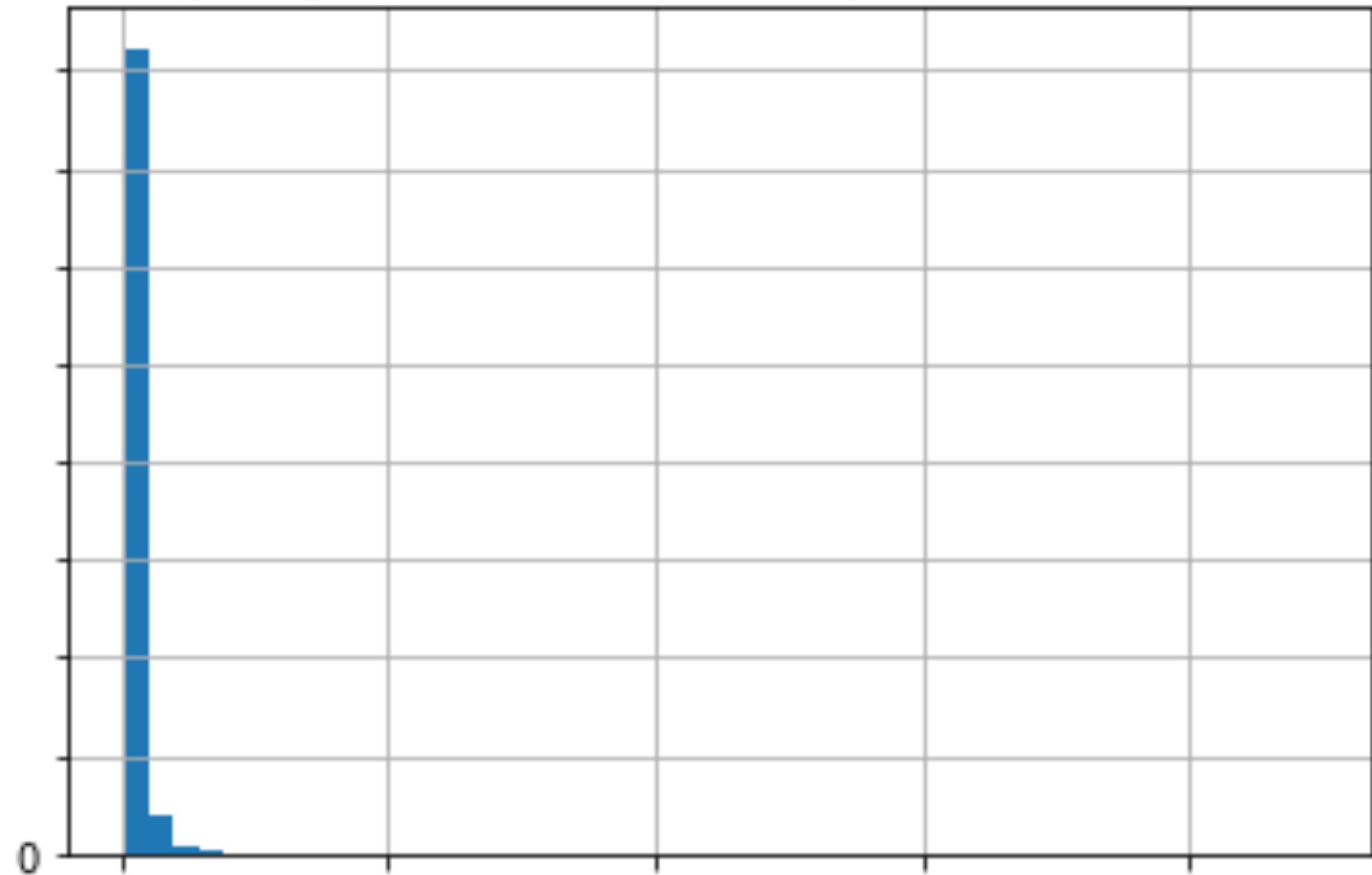
# Resultados

- Grupo de teste - um pequeno percentual de usuários de 2 estações de TV: 800 mil usuários (teste + controle)
- Recomendações no rodapé da home: CTR aumentou em 15%
- Recomendações nas páginas de artigo (infinite scroll): Tempo médio de leitura aumentou em 168%
- Sinal verde para versão 2 mais sofisticada (deep learning)

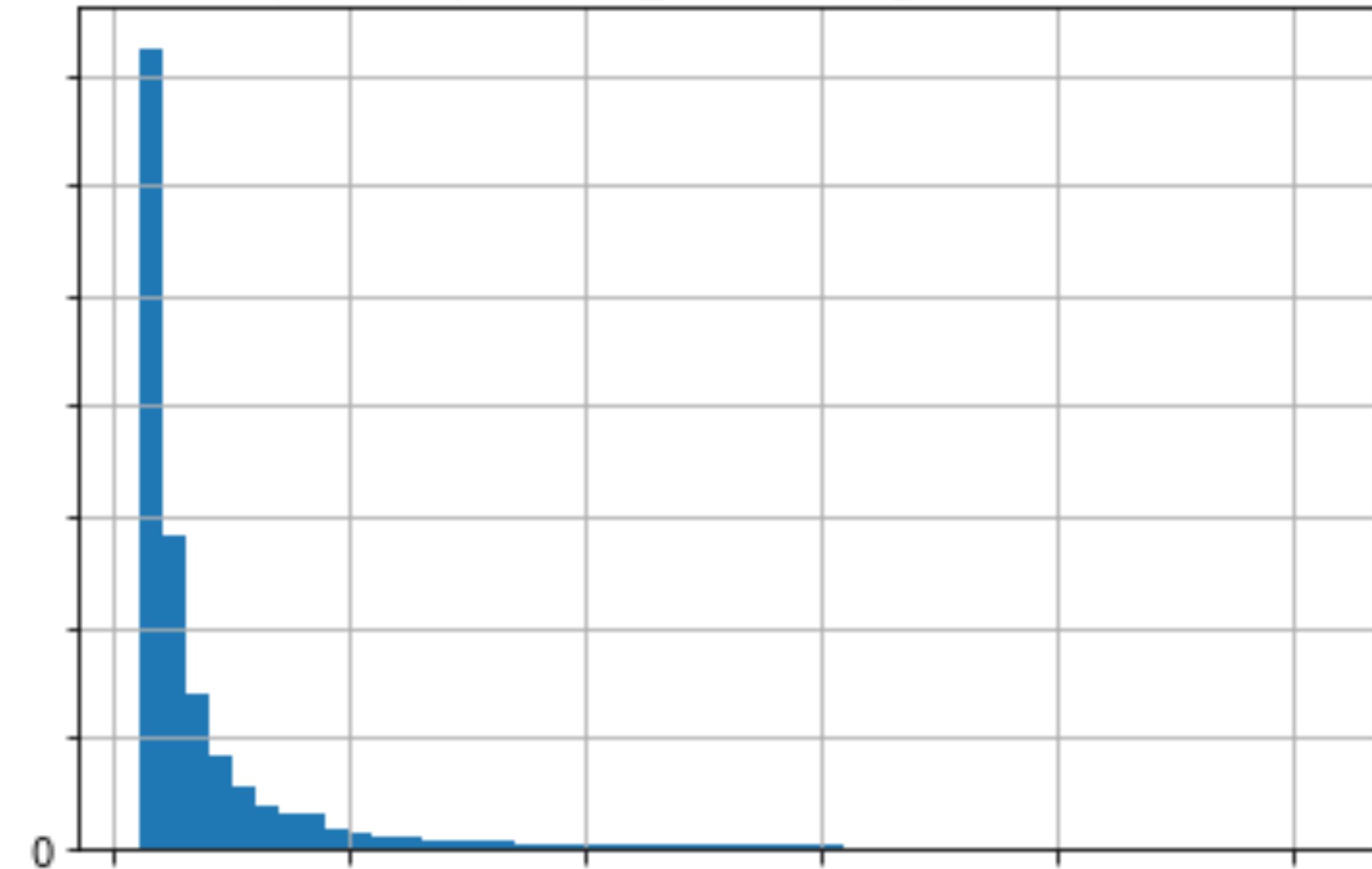
# Análise Exploratória dos Dados

# Page-views e sessões

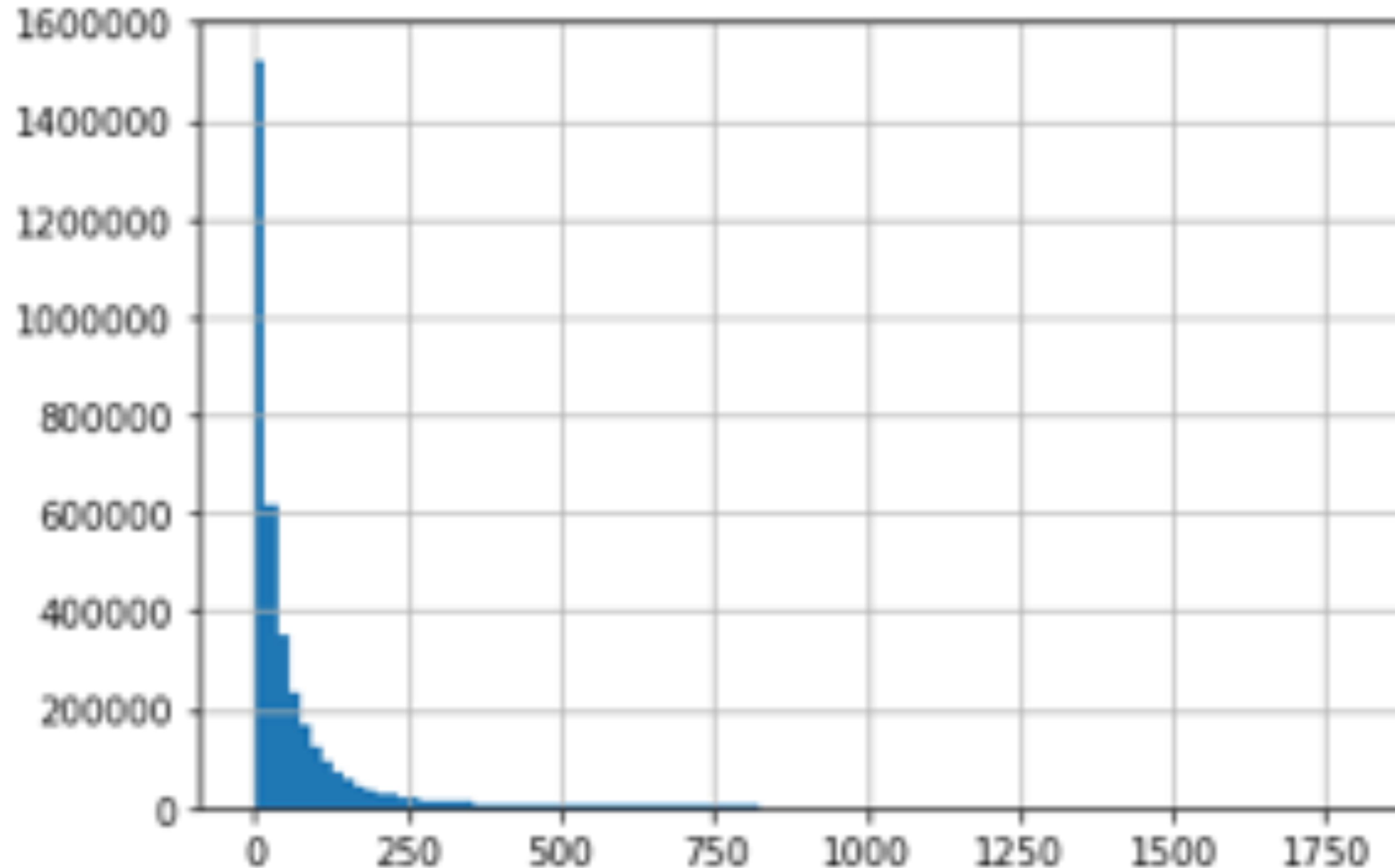
PV history length of user sessions (input for session behaviour)



#sessions per unique user\_id (for long-term user behaviour)



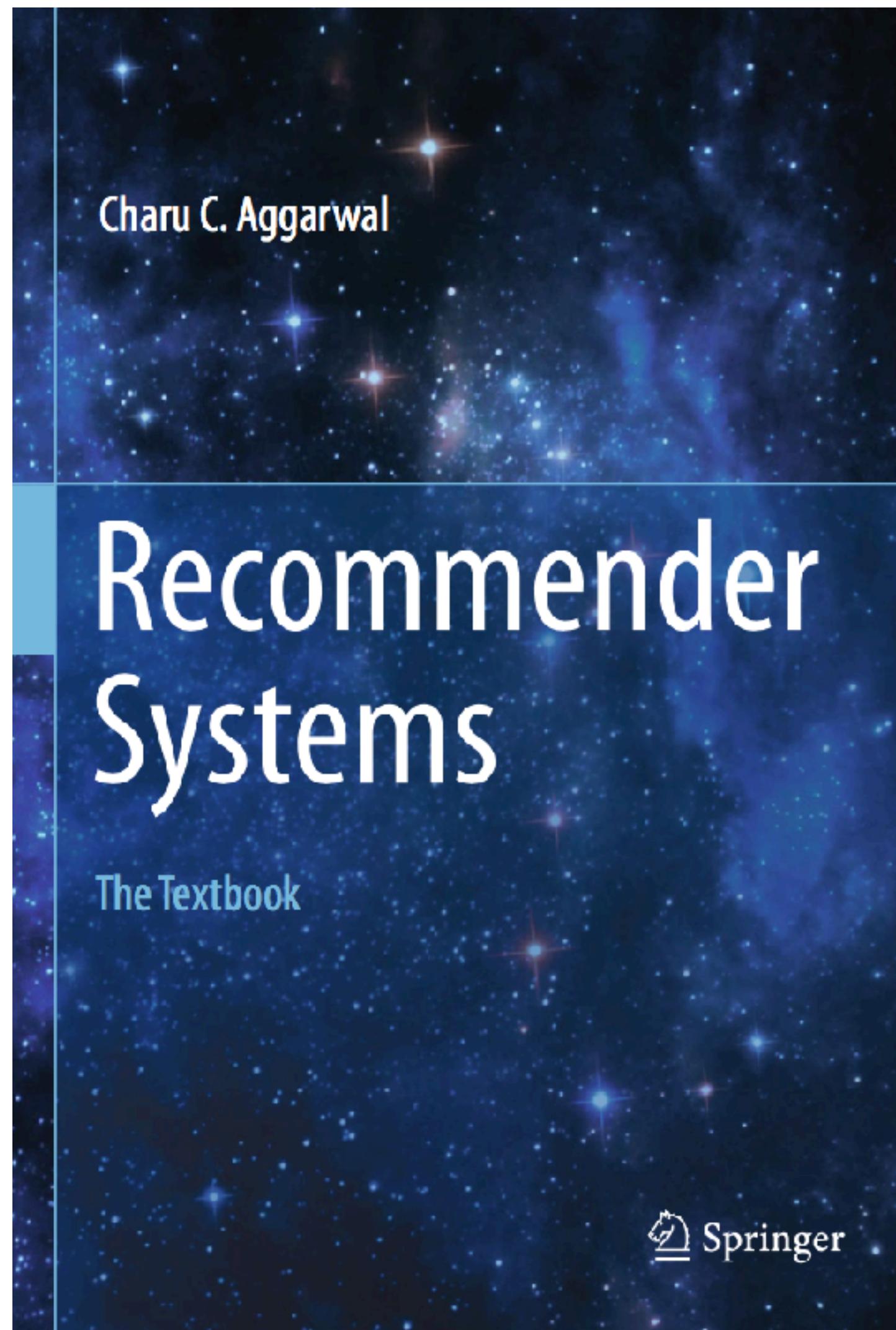
# Duração (seg) de cada page-view



# Sparse Matrix: Users x Articles

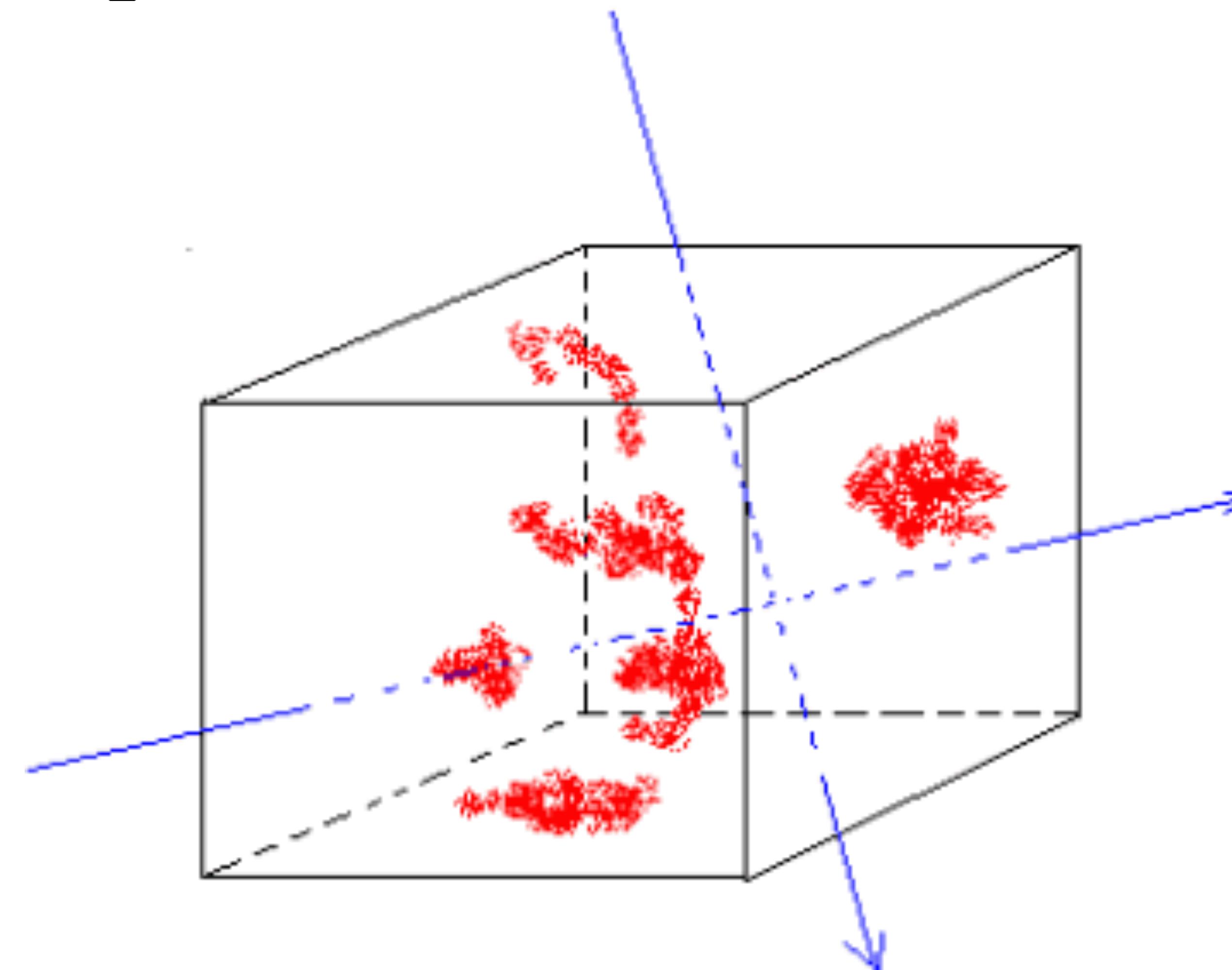
Dataset	Preenchimento
MovieLens (movies)	1,39%
Netflix (movies)	1,18%
TV Stations (news) 20x menos dados que Netflix	0,06%

# Referências

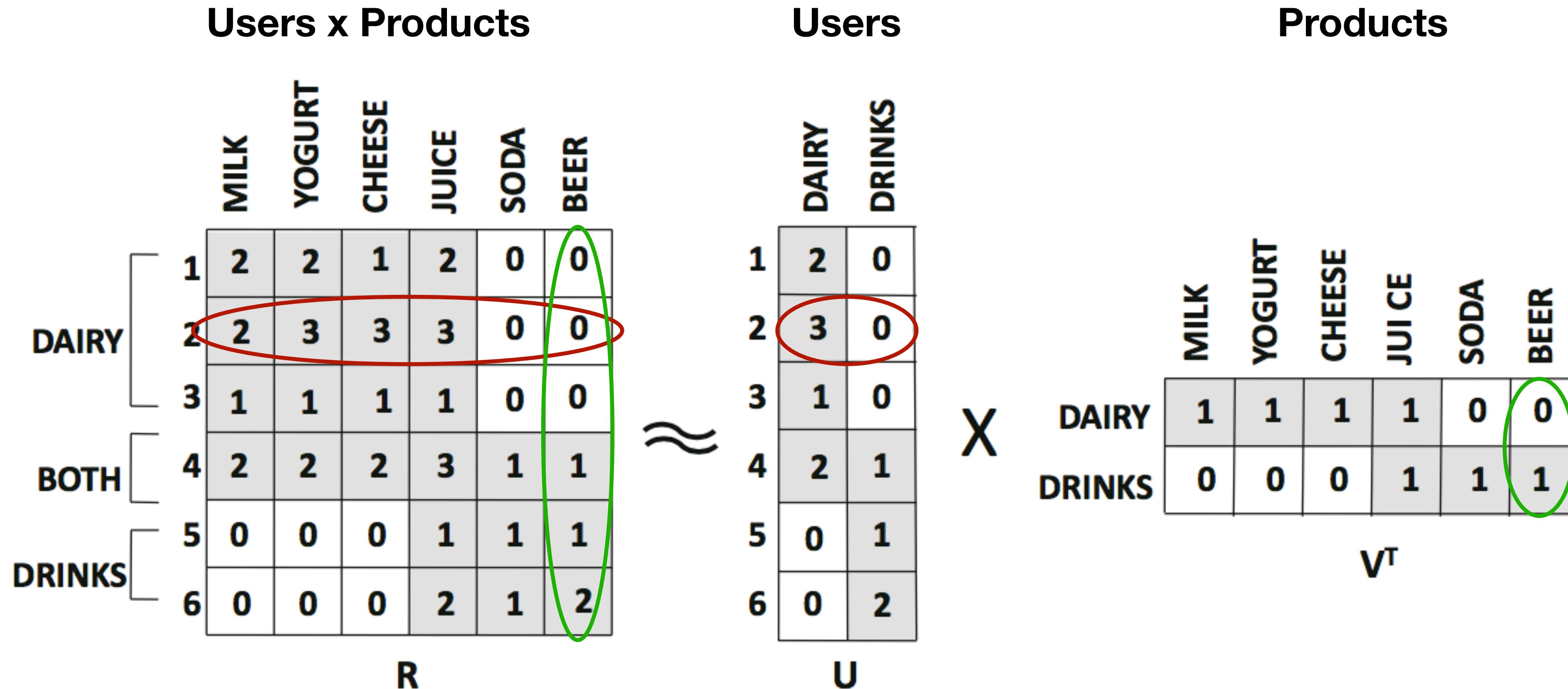


3.6	Latent Factor Models . . . . .	90
3.6.1	Geometric Intuition for Latent Factor Models . . . . .	91
3.6.2	Low-Rank Intuition for Latent Factor Models . . . . .	93
3.6.3	Basic Matrix Factorization Principles . . . . .	94
3.6.4	Unconstrained Matrix Factorization . . . . .	96
3.6.4.1	Stochastic Gradient Descent . . . . .	99
3.6.4.2	Regularization . . . . .	100
3.6.4.3	Incremental Latent Component Training . . . . .	103
3.6.4.4	Alternating Least Squares and Coordinate Descent . . . . .	105
3.6.4.5	Incorporating User and Item Biases . . . . .	106
3.6.4.6	Incorporating Implicit Feedback . . . . .	109

# Espaço de Fatores Latentes



# Model-based Collaborative Filtering



# Frameworks e linguagens

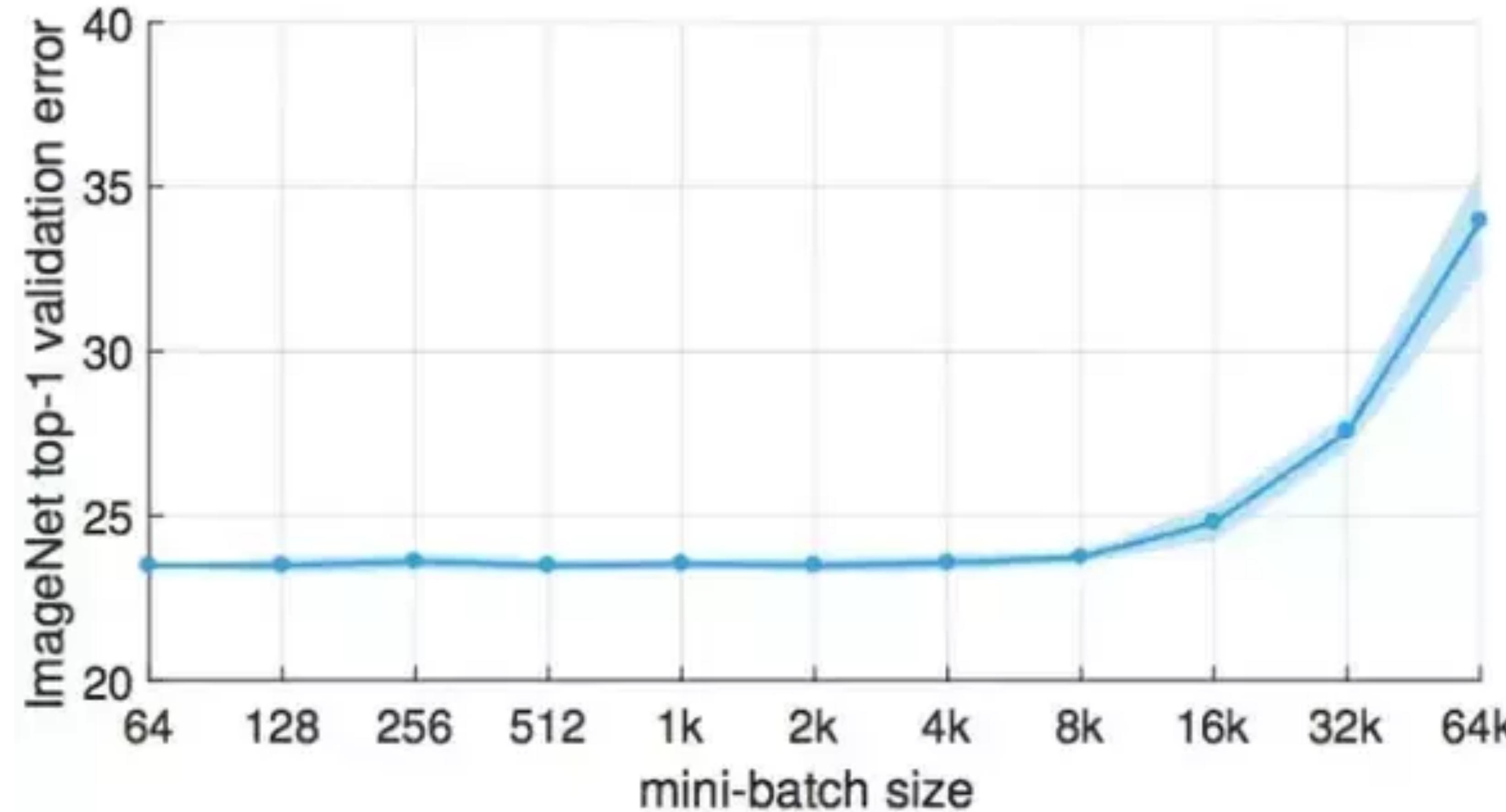
Top  
Frameworks



Programming  
languages



# Mini-batch size x Validation Error



ImageNet dataset: 1h de treino com mini-batches de 8192 imagens em 256 GPUs.

# Tensorflow Framework

High-Level  
TensorFlow APIs

Mid-Level  
TensorFlow APIs

Low-level  
TensorFlow APIs

TensorFlow  
Kernel

Estimators

Layers

Datasets

Metrics

Python

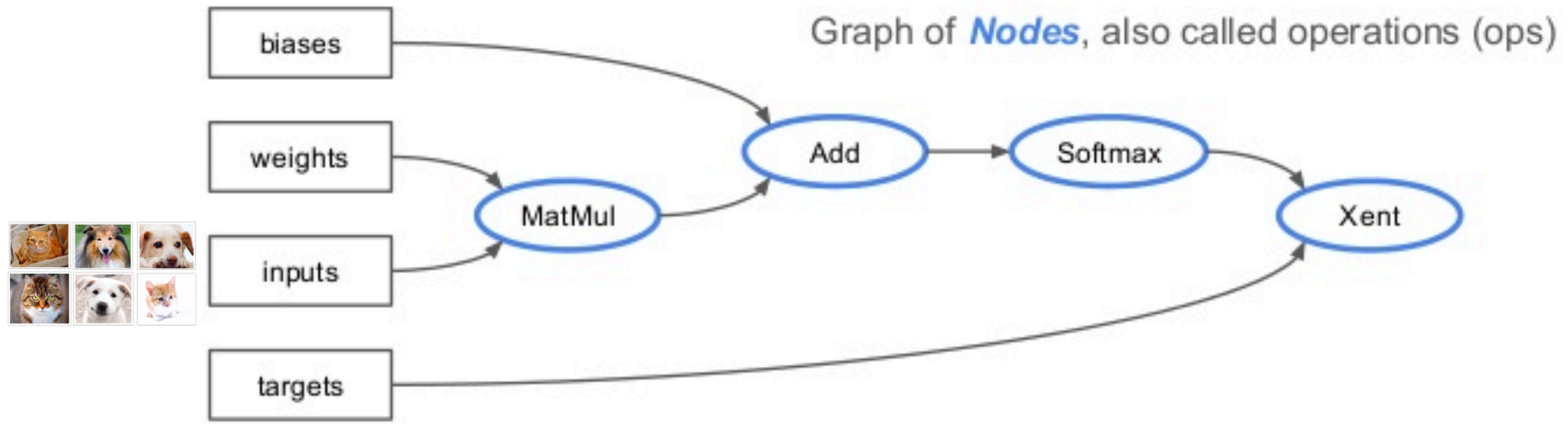
C++

Java

Go

TensorFlow Distributed Execution Engine

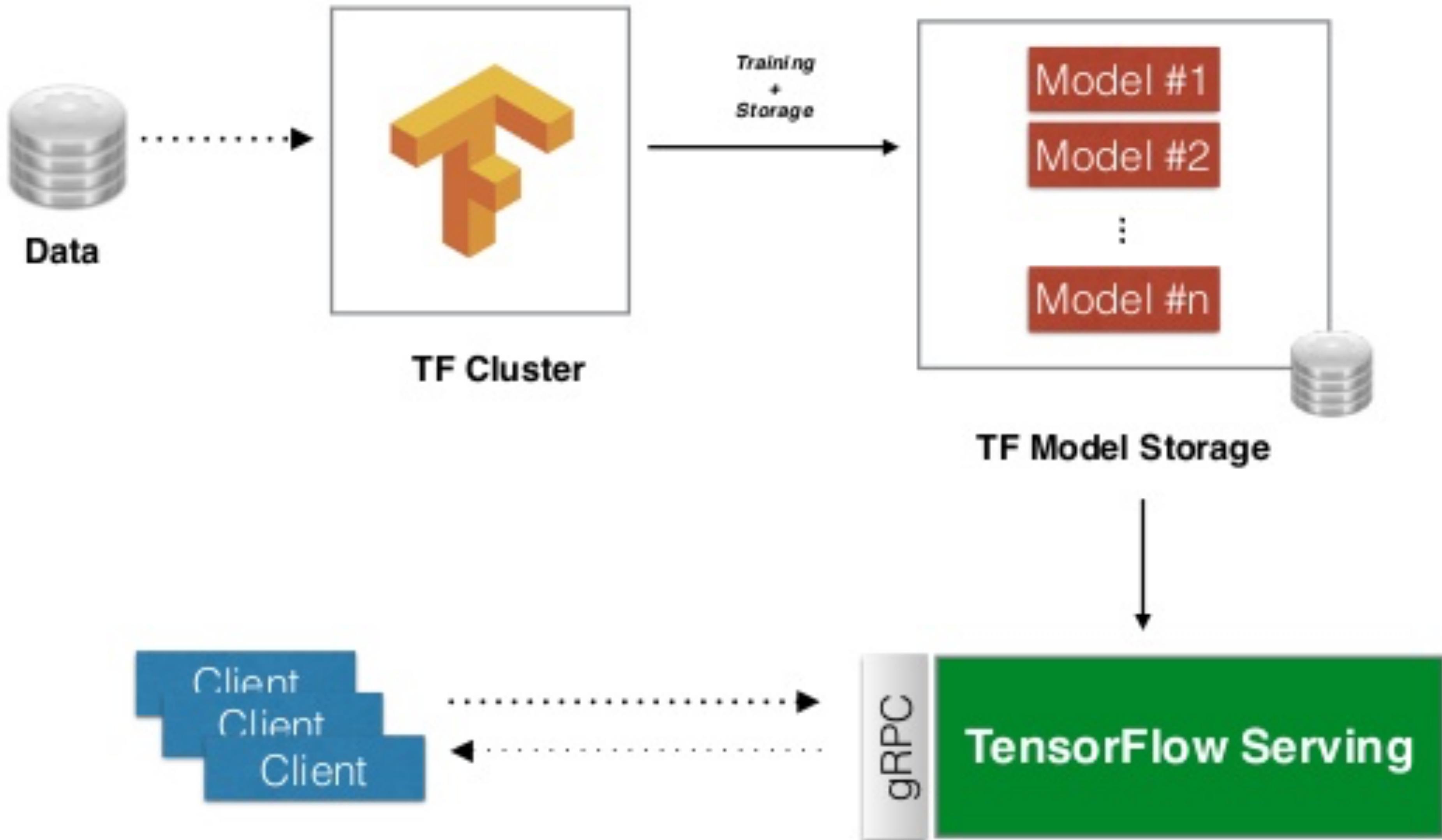
# Computation as a Dataflow Graph



classifier = **softmax**(inputs \* **weights** + **biases**)

cost = **cross\_entropy**(classifier, targets)

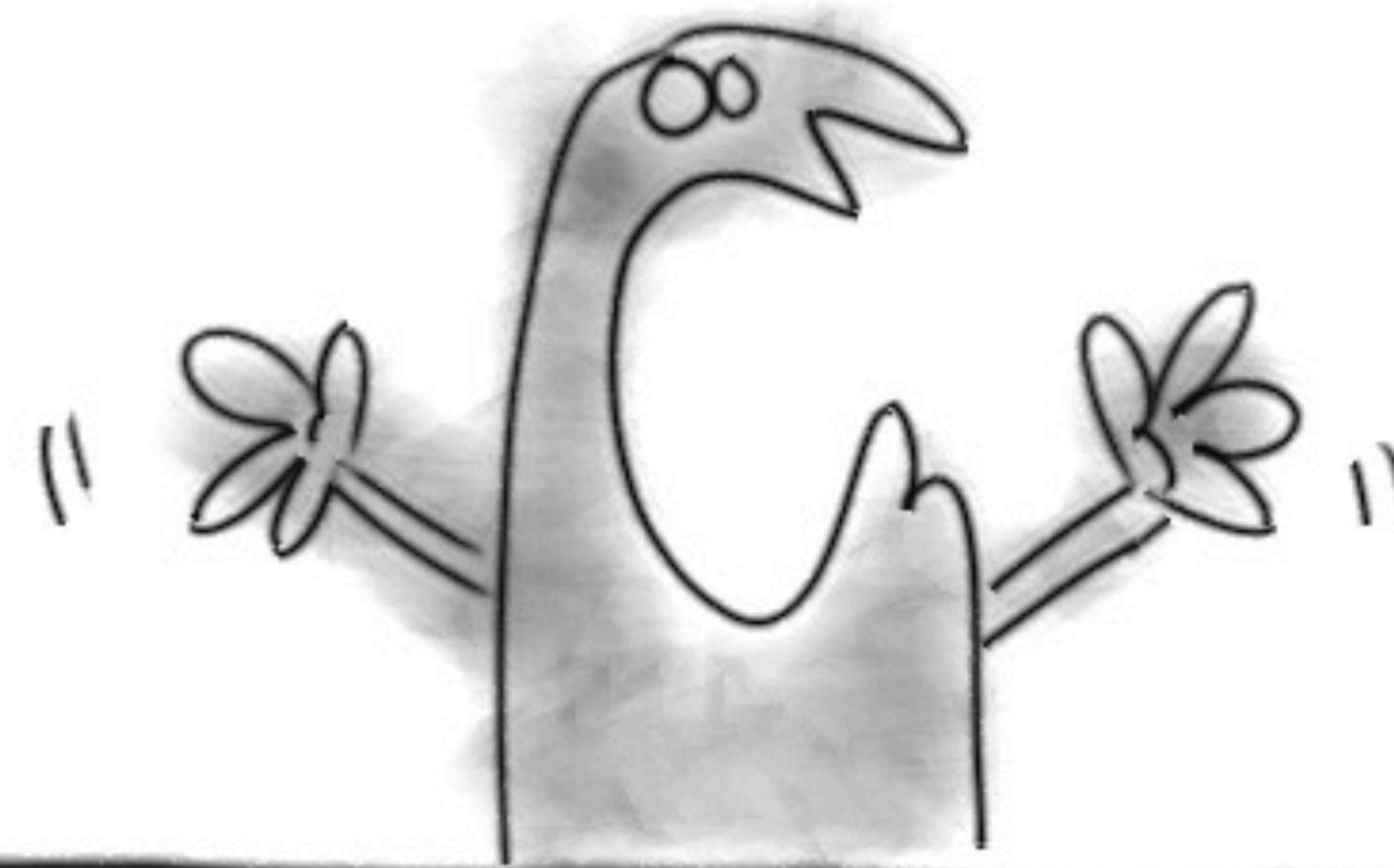
# Tensorflow Serving



# Como testar o modelo?

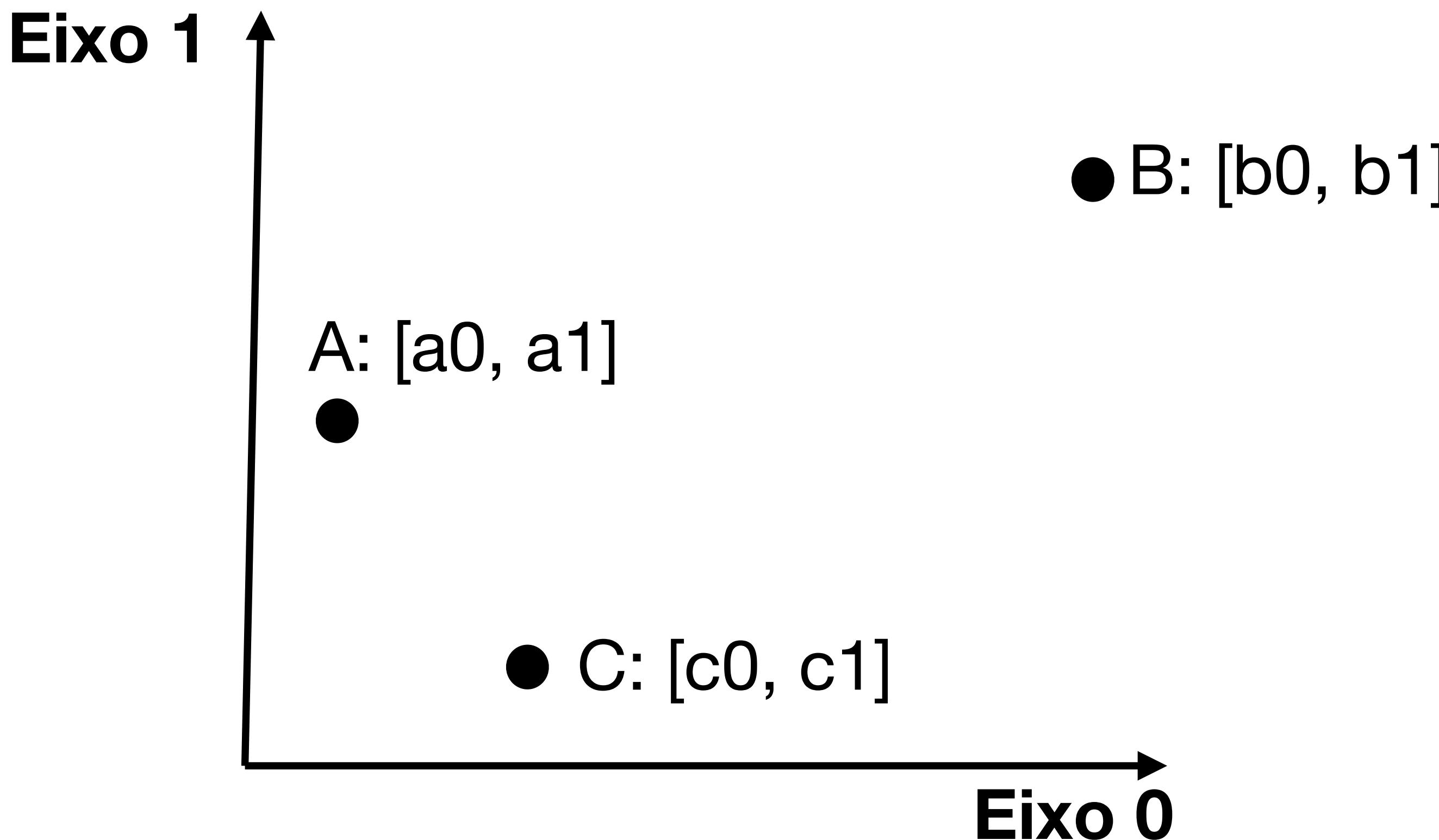
- Benchmark com dados do Movie Lens - Ok
- Apoio Google: Code review e comparação com algoritmo deles - Ok
- "Surrogate problem"
- Personas com histórico fake para testes offline
- Plugin Chrome para teste online
- **Resultado:** Recomendações pouco relevantes/correlacionadas ao histórico das personas

# Now What?!!

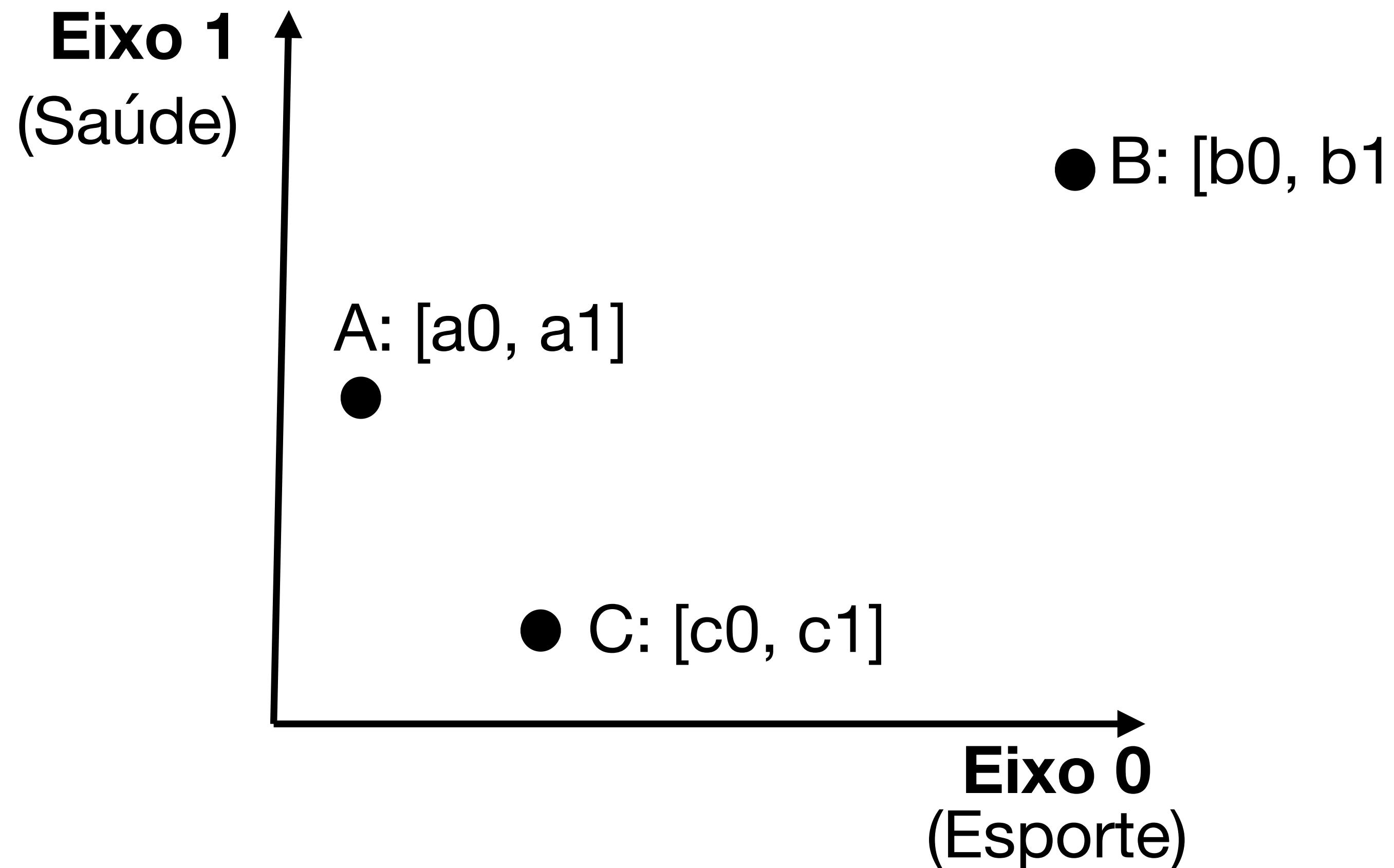


**Embeddings: Notícias devem virar  
vetores em um espaço  
n-dimensional de assuntos**

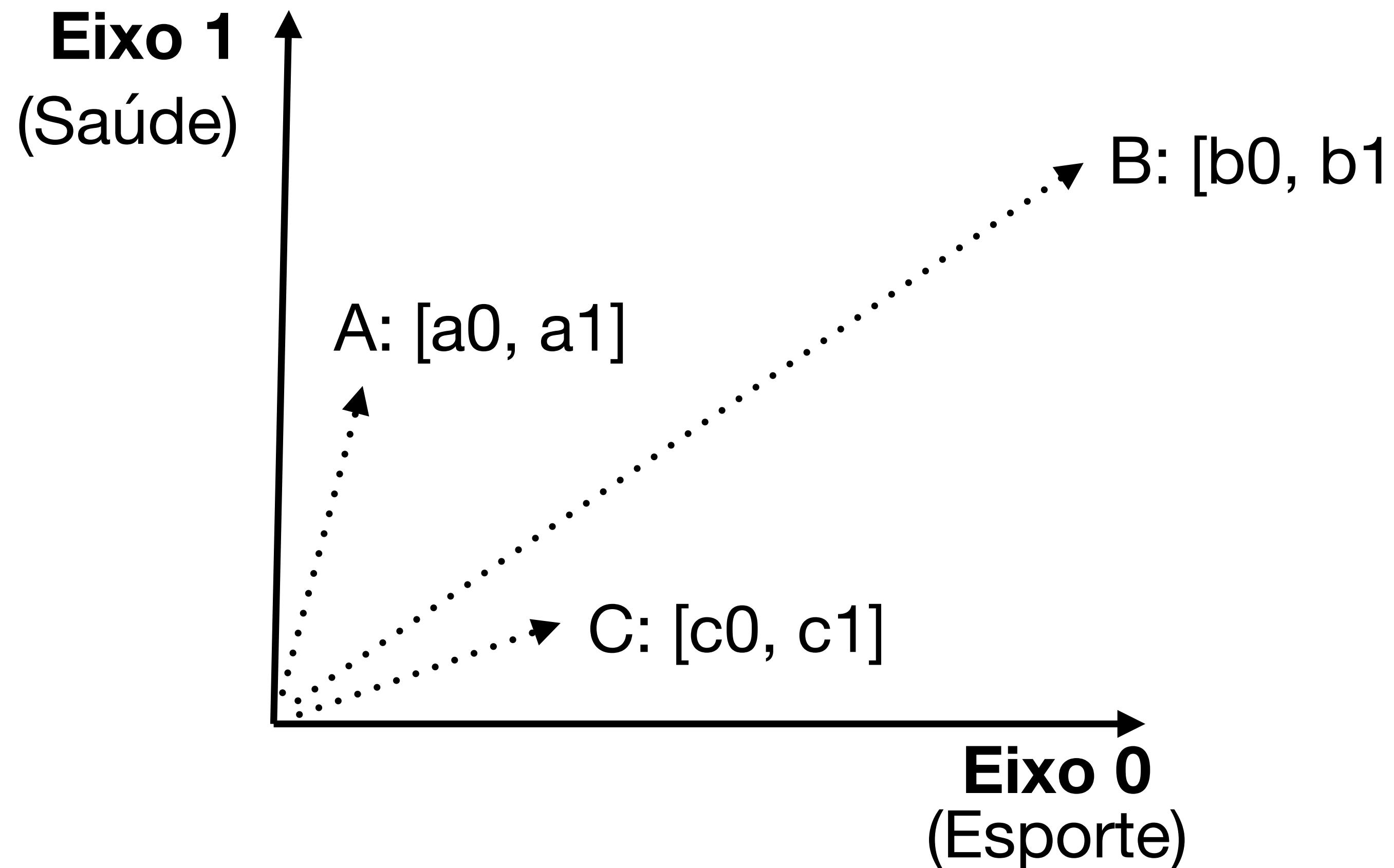
# Embeddings (Intuição)



# Embeddings (Intuição)

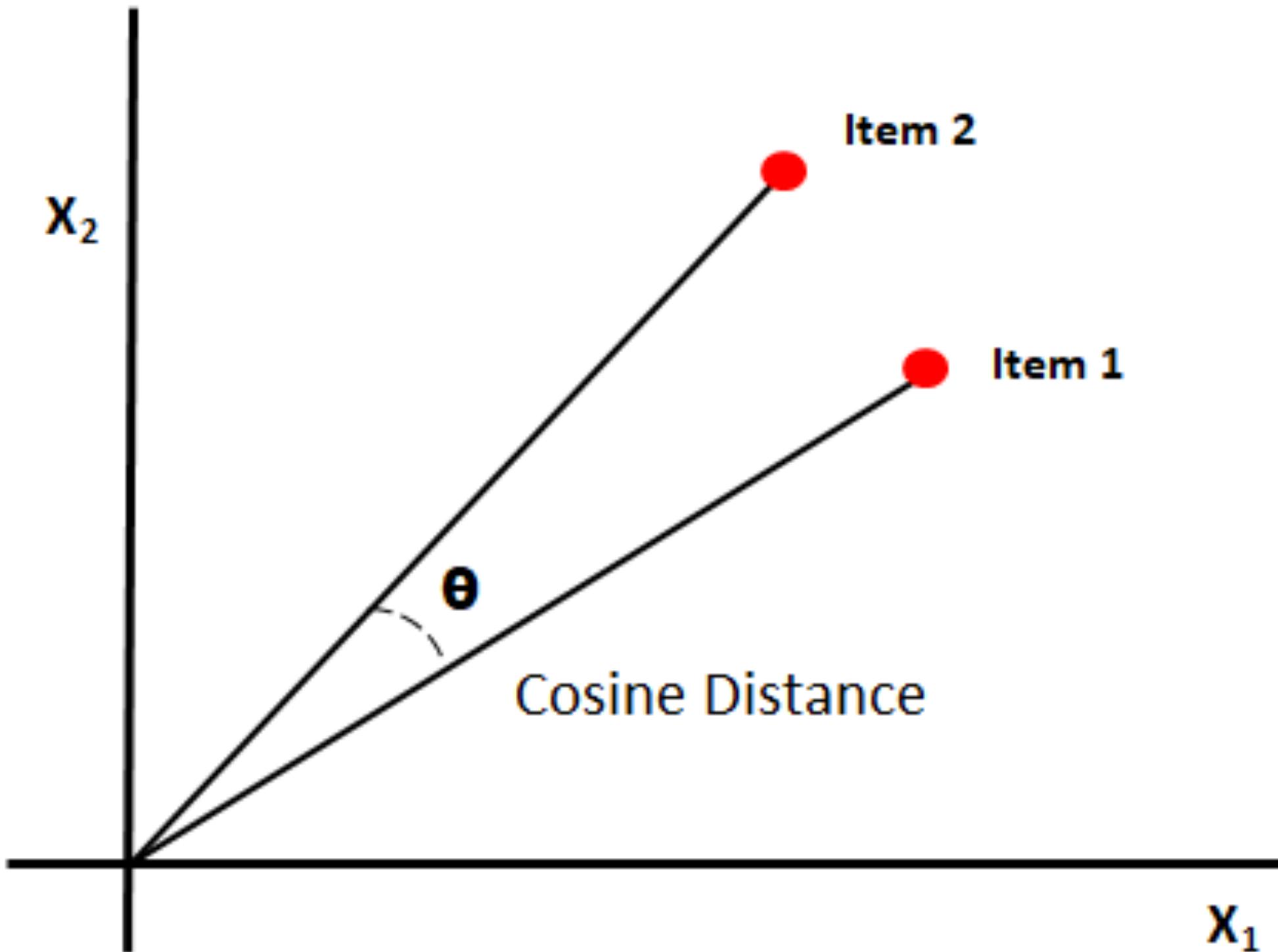


# Embeddings (Intuição)



# Semelhança entre notícias

*Cosine Distance/Similarity*

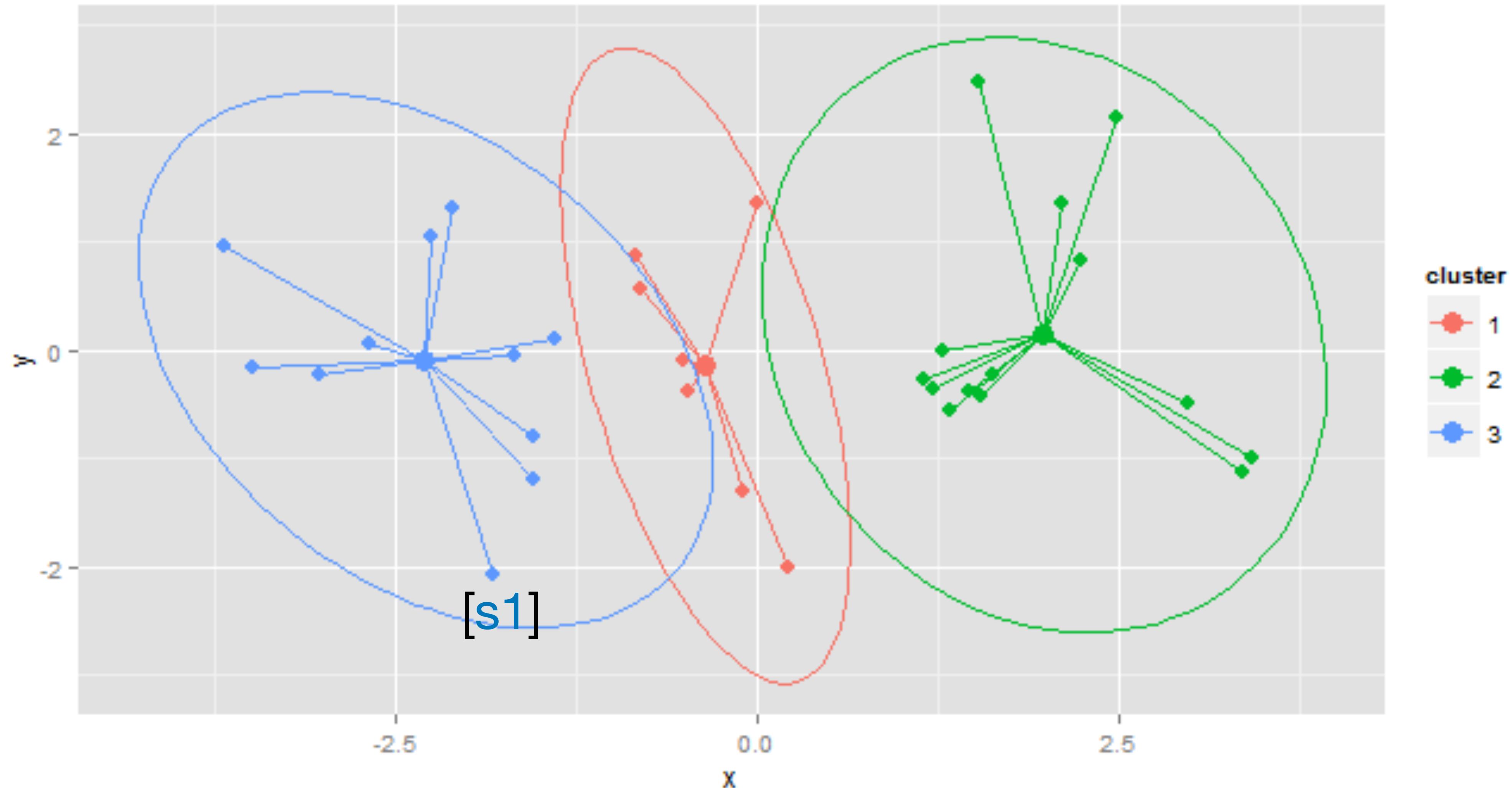


$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

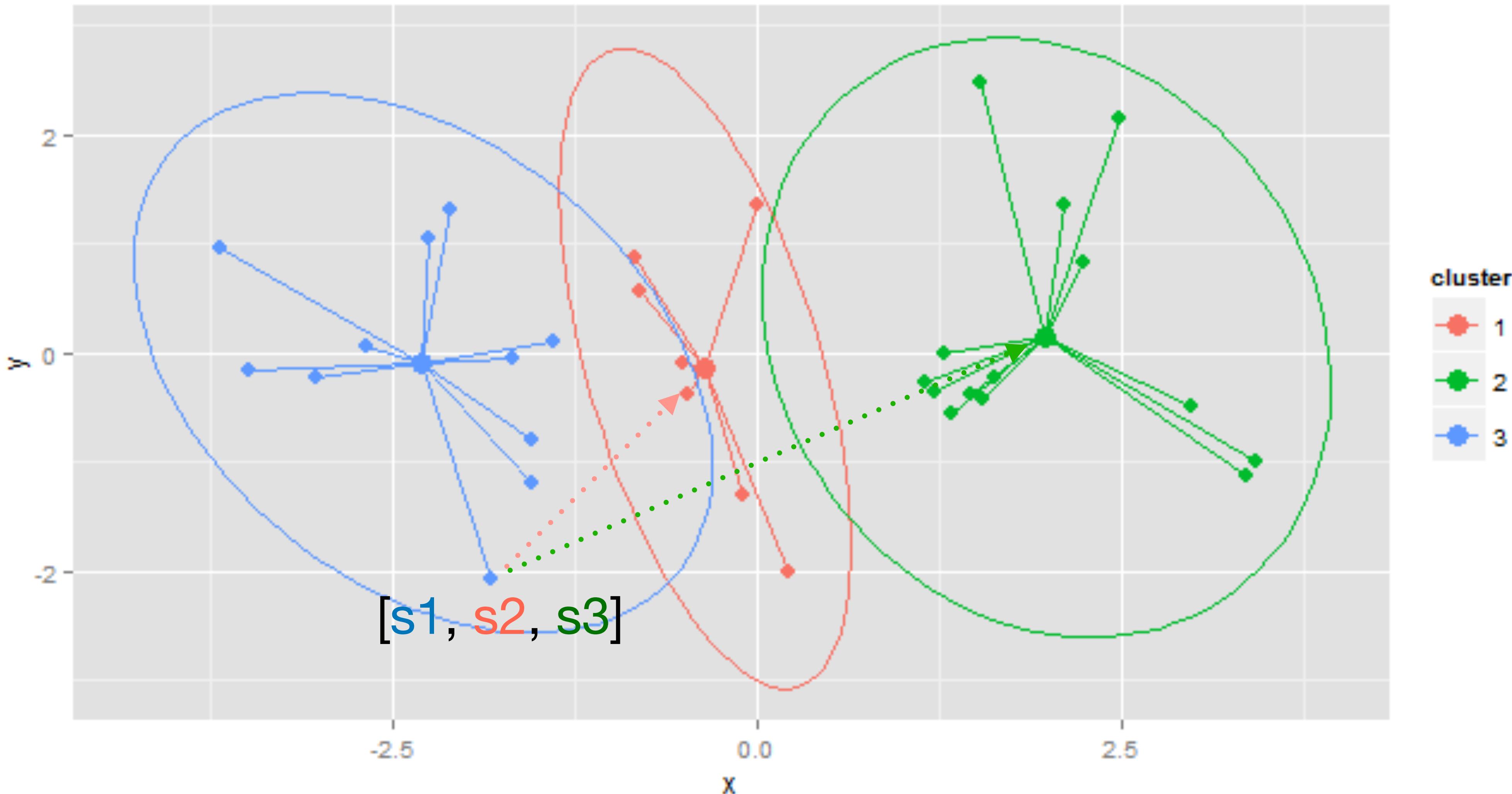
# Content (text) Embeddings

- Objetivo: Transformar texto da notícia em um array de tamanho fixo
- Tags
- BOW (bag of words)
- TF-IDF (term-frequency x inverse document frequency)
- Word Embeddings (word2vec, FastText) + PV-DM (paragraph vectors distributed memory)
- Seq2Seq Autoencoder
- HANs (Hierarchical Attention Network)

# Tópicos via clustering (KMeans)

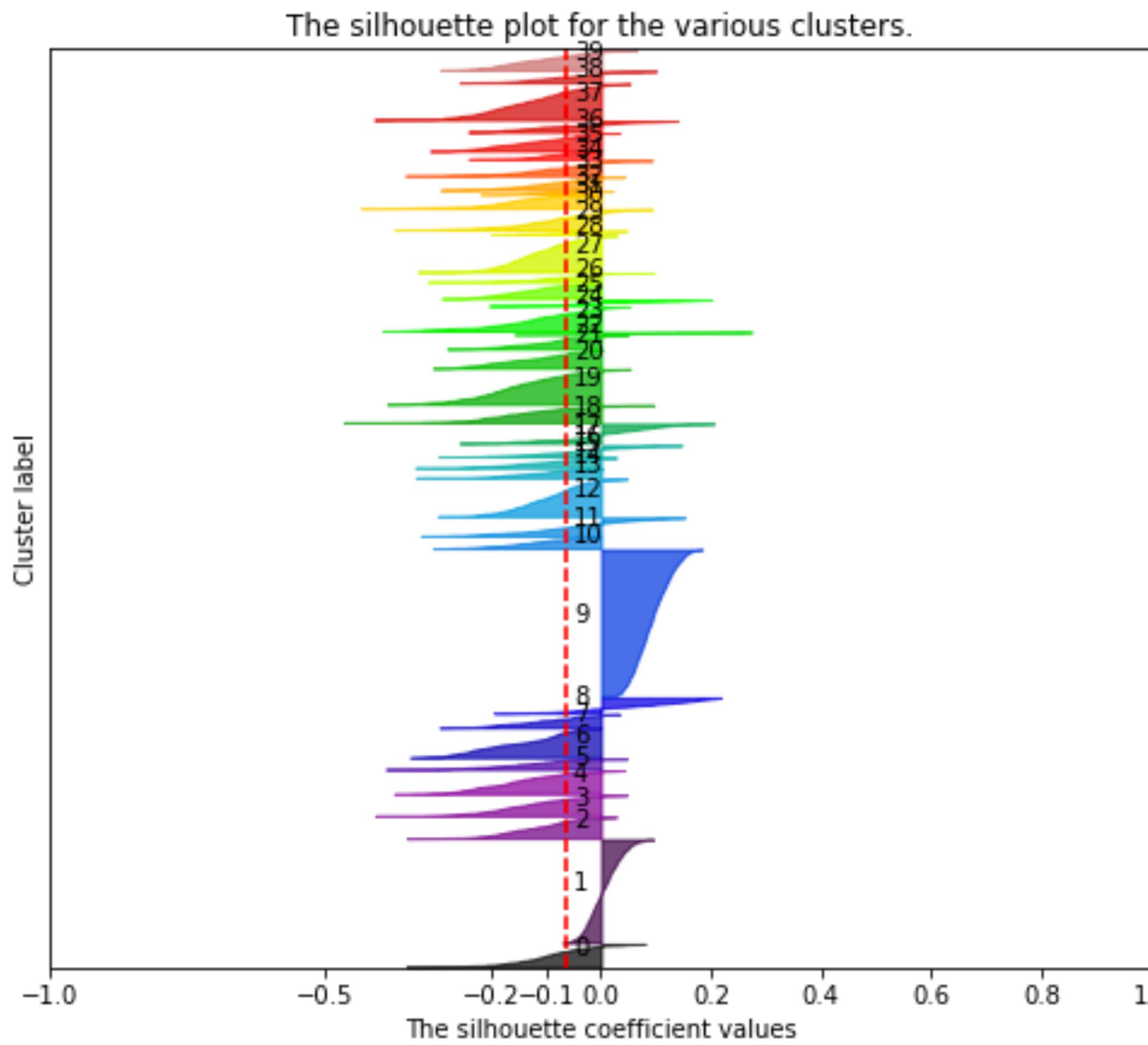


# Tópicos via clustering (KMeans)

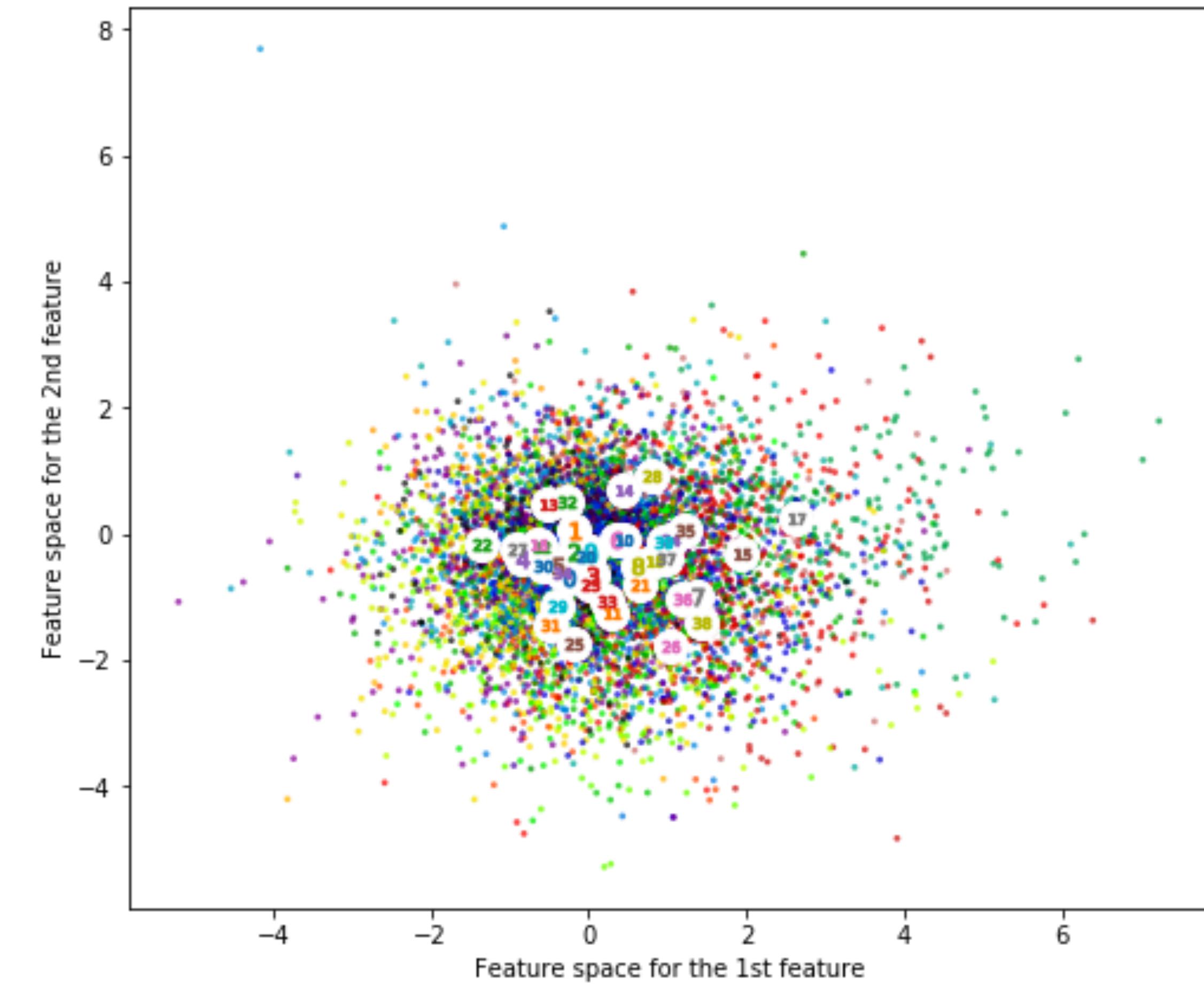


# Qtd de clusters (assuntos latentes)

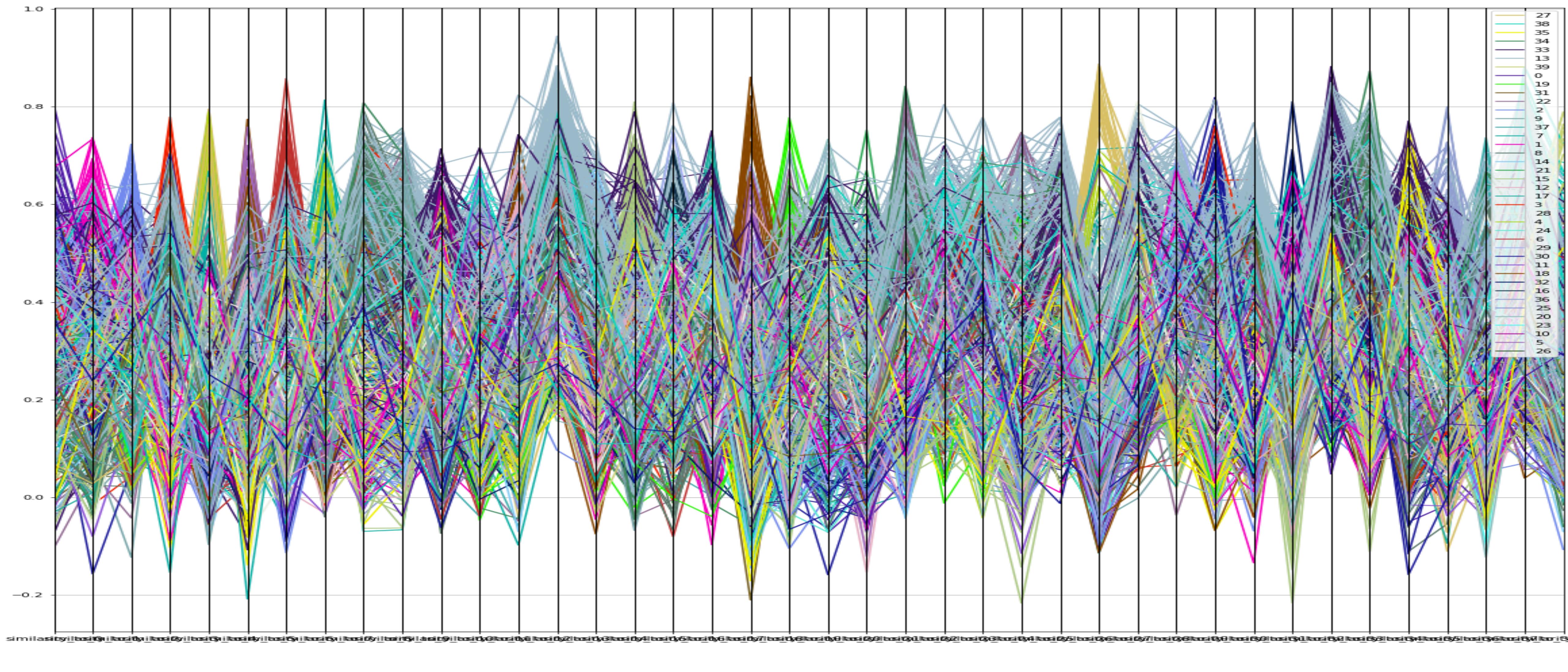
**Silhouette analysis for KMeans clustering on sample data with n\_clusters = 40**



The visualization of the clustered data.



# 10k notícias em 40 dimensões

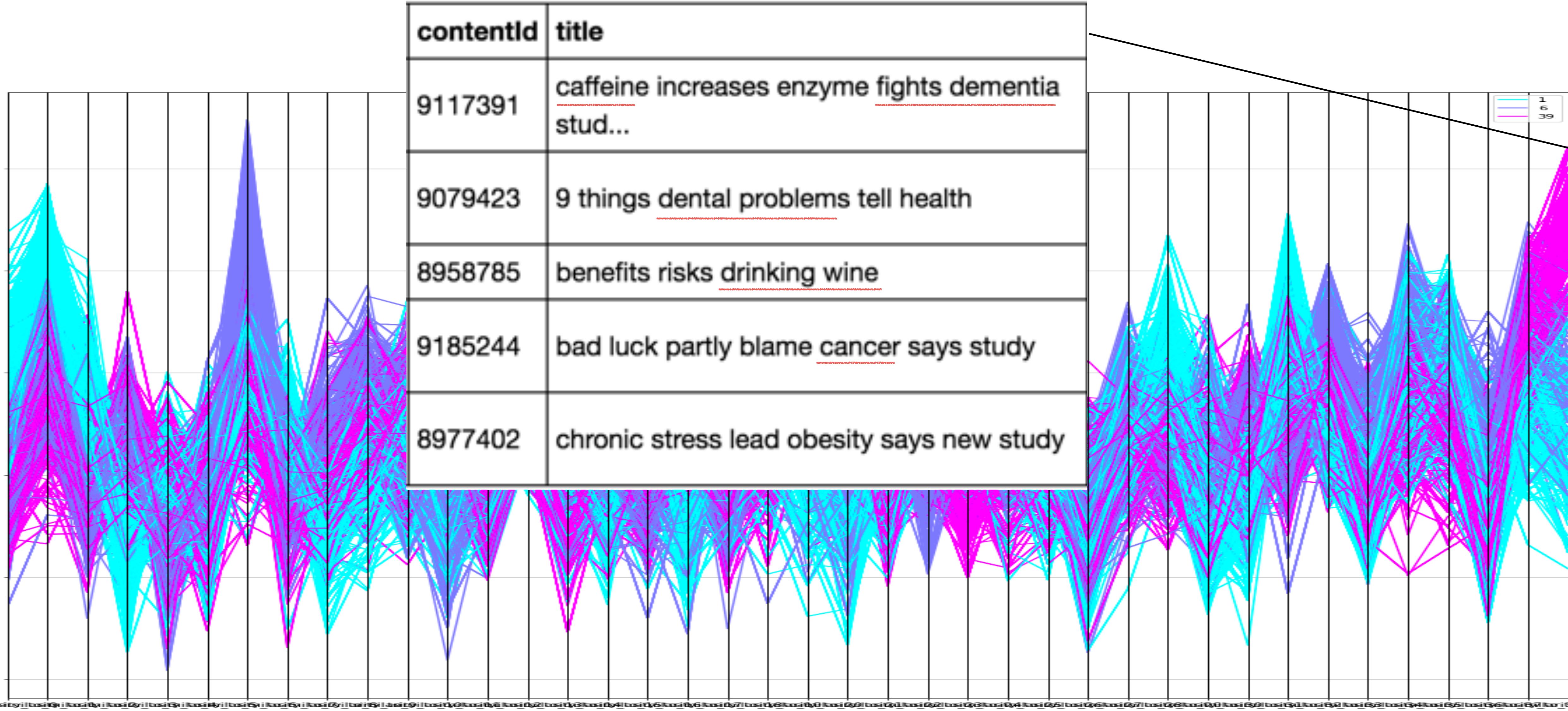


# Dimensão 0: “economia/desemprego”

contentId	title
9092218	new bill could help companies pay employee stu...
9116669	us adds 235k jobs unemployment rate falls 4.7 ...
9083135	us jobless claims drop lowest level since 1973
9083100	us applications unemployment benefits drop low...
9136106	fed raises rate sees hikes us economy improves

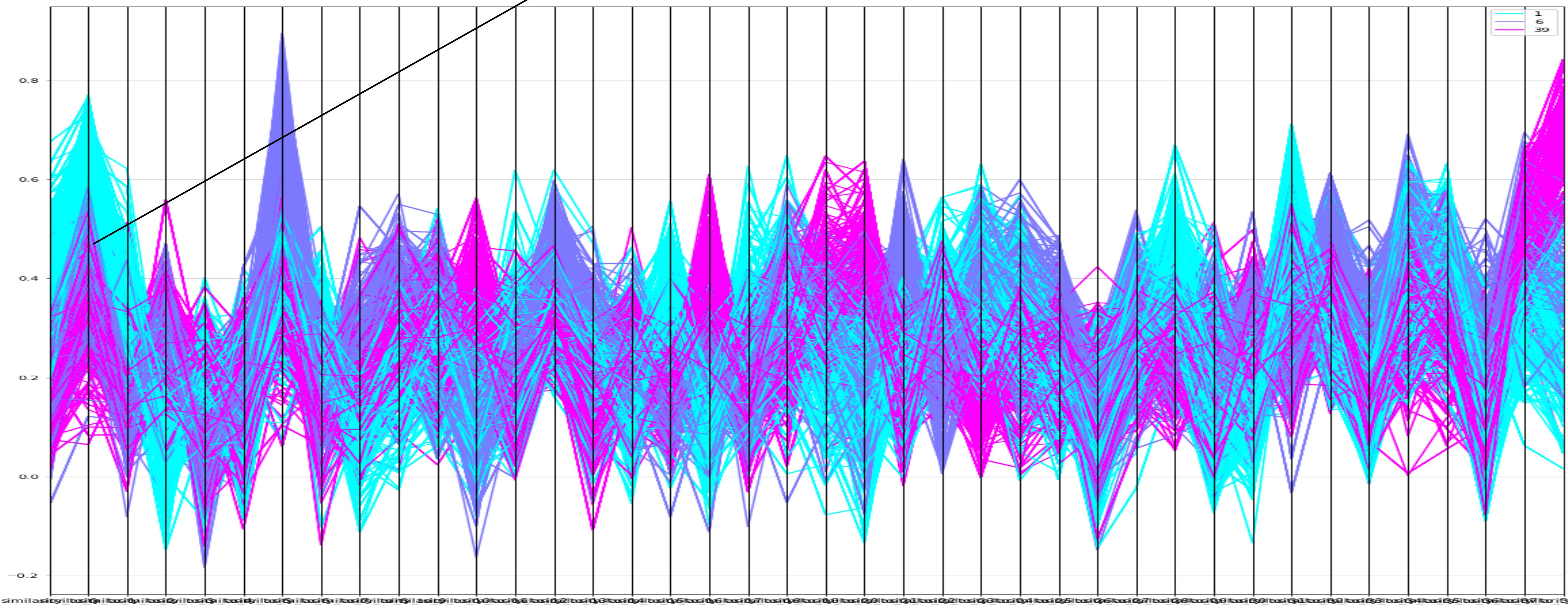


# Dimensão 39: “dicas de saúde”

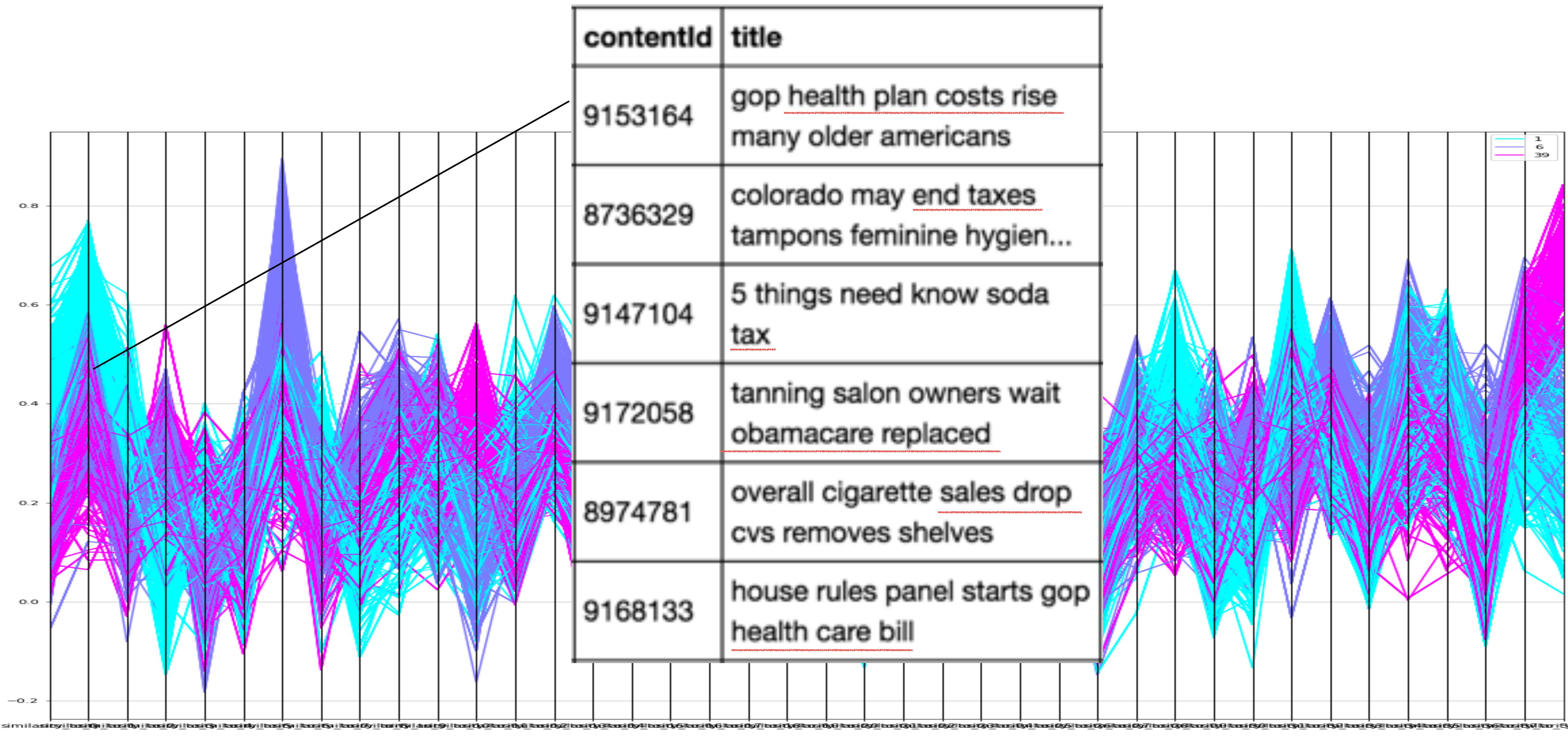


# Artigos “magenta” próximo à “economia”?

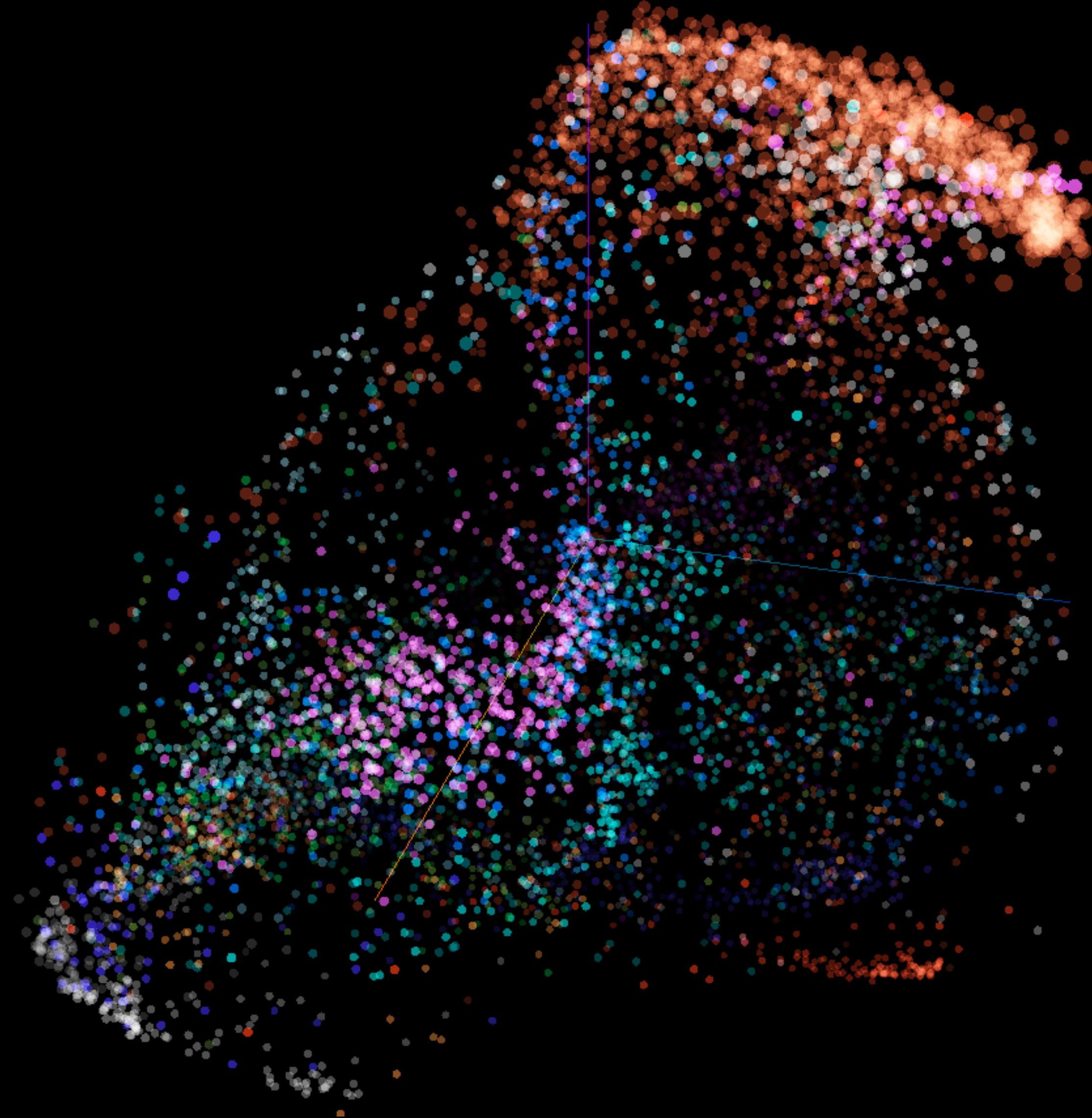
?



# Artigos “magenta” próximo a “economia”?



# Visualização com T-SNE



malaysia says autopsy shows nerve agent effects ^

contentId 8983442  
title malaysia says autopsy shows nerve agent effects  
section \$news  
subsection \$\$national\_news  
keywords #autopsy #dead #half\_brother #kim\_jong\_un #malaysia #nerve\_agent #results  
cluster 5

## malaysia says autopsy shows nerve agent effects

- woman escapes illegal alien kidnapping attempt near mexico border
- widow orlando pulse nightclub shooter released trial
- priest arrested plot poison head georgian church
- lawyer malaysia may compromised kim jong nam case
- airline pilot wife found dead toxicology report pending
- police maryland girl planned school attack shotgun
- woman tries shoot drone filming house
- police pakistan shrine custodian kills 20 stabbing attack
- british parliament attack know
- police release 12 people arrested connection london attack
- heavily armed teen arrested 2 wounded shooting french high school
- teen asylum seeker id suspect norway explosive case
- dallas officials seek details 156 emergency alarms malfunctioned
- russia challenges trump say would syria
- us military confirms syria bombing denies mosque hit
- navy missiles light skies northern california
- albuquerque un peacekeeper kidnapped africa
- shark kills 17 year old surfer southern australia

## NOWCAST WCVB On Demand

≡  
MENU

# O livro que os bancos não querem que você leia



Advertisement

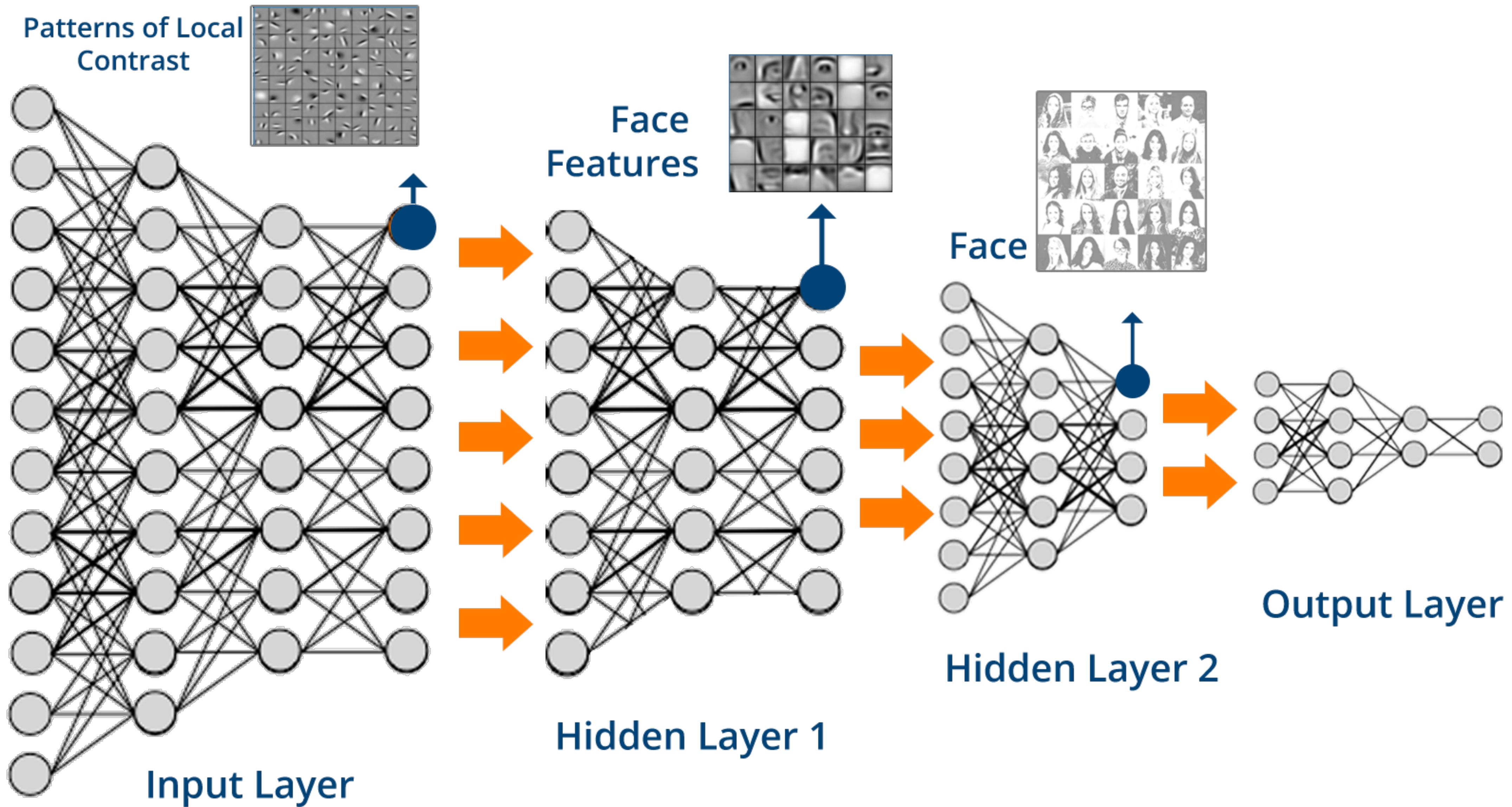


TITLE	DATE	CONTENT	DISPLAY
<b>Profile in Courage Award: Obama reflects on courage, health care</b>	May 7, 2017 11:21 PM	Video	Standard Article
<b>Obama talks about the courage required to pass the ACA</b>	May 7, 2017 11:37 PM	Video	Standard Article
<b>Trooper chases drunken driving suspect who rear-ended cruiser</b>	May 7, 2017 11:10 PM	Default	Standard Article
<b>Is that Facebook post claiming 'Friends' is returning real?</b>	May 7, 2017 10:11 PM	Default	Standard Article
<b>Twin sisters give birth to baby boys on same day</b>	May 8, 2017 1:24 AM	Video	Standard Article
<b>1 dead, 1 in serious condition after double stabbing in Milford</b>	May 7, 2017 11:42 PM	Default	Standard Article
<b>Police Pull Over Couple To Reveal Pregnancy Announcement</b>	May 7, 2017 2:37 PM	Video	Standard Article
<b>Obama reflects on health care, meaning of courage in award acceptance speech</b>	May 7, 2017 11:32 PM	Default	Standard Article
<b>Teacher sells plasma just to make enough for his family of four</b>	May 7, 2017 4:52 PM	Video	Standard Article

View Interactive Radar

Kelly Ann Cicalese

# **Parte #2 - Deep Learning**



# Motivações para usar DNN

- V1 captura de cada usuário a intensidade de seu interesse por cada tópico latente
- Queremos capturar a dinâmica comportamental do usuário em diferentes contextos:
  - Temporal: Dia da semana, horário, dia do mês, dia do ano, etc...
  - Assuntos: Sequência dos últimos artigos visitados e seus assuntos.
  - Localização: Em casa? No trabalho? Passeando?

# Referência: Survey de Jul/2017

Deep Learning based Recommender System: A Survey and New Perspectives • 35:5

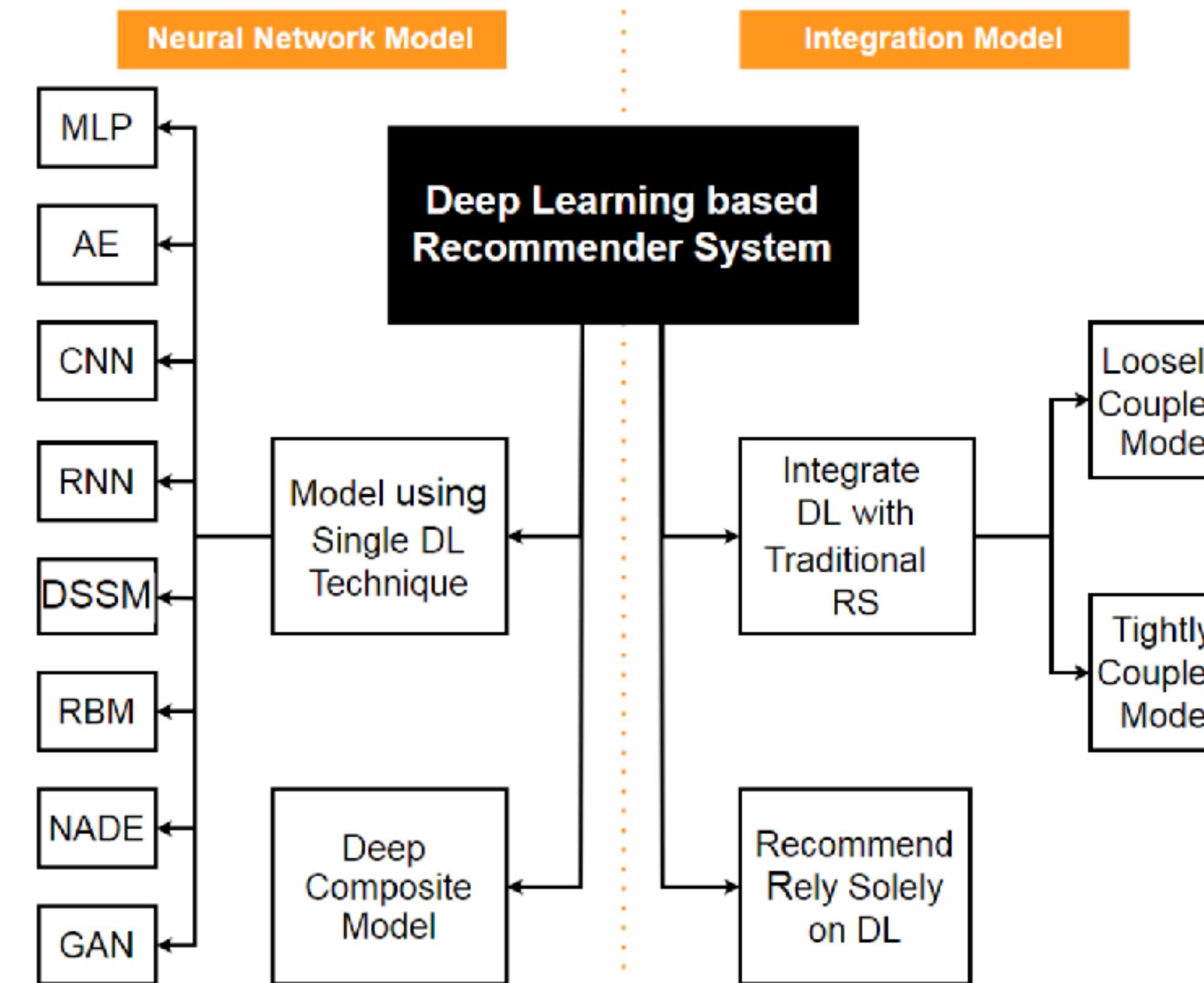
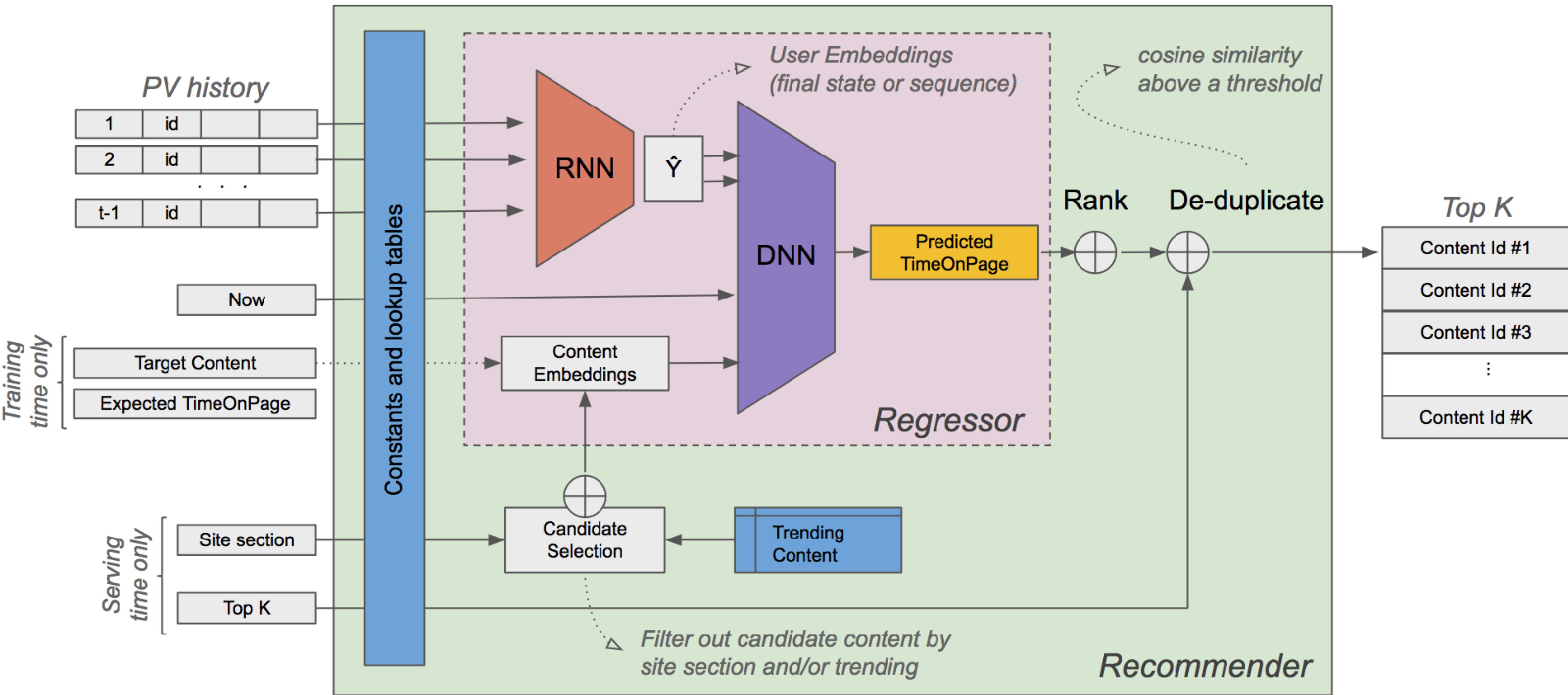


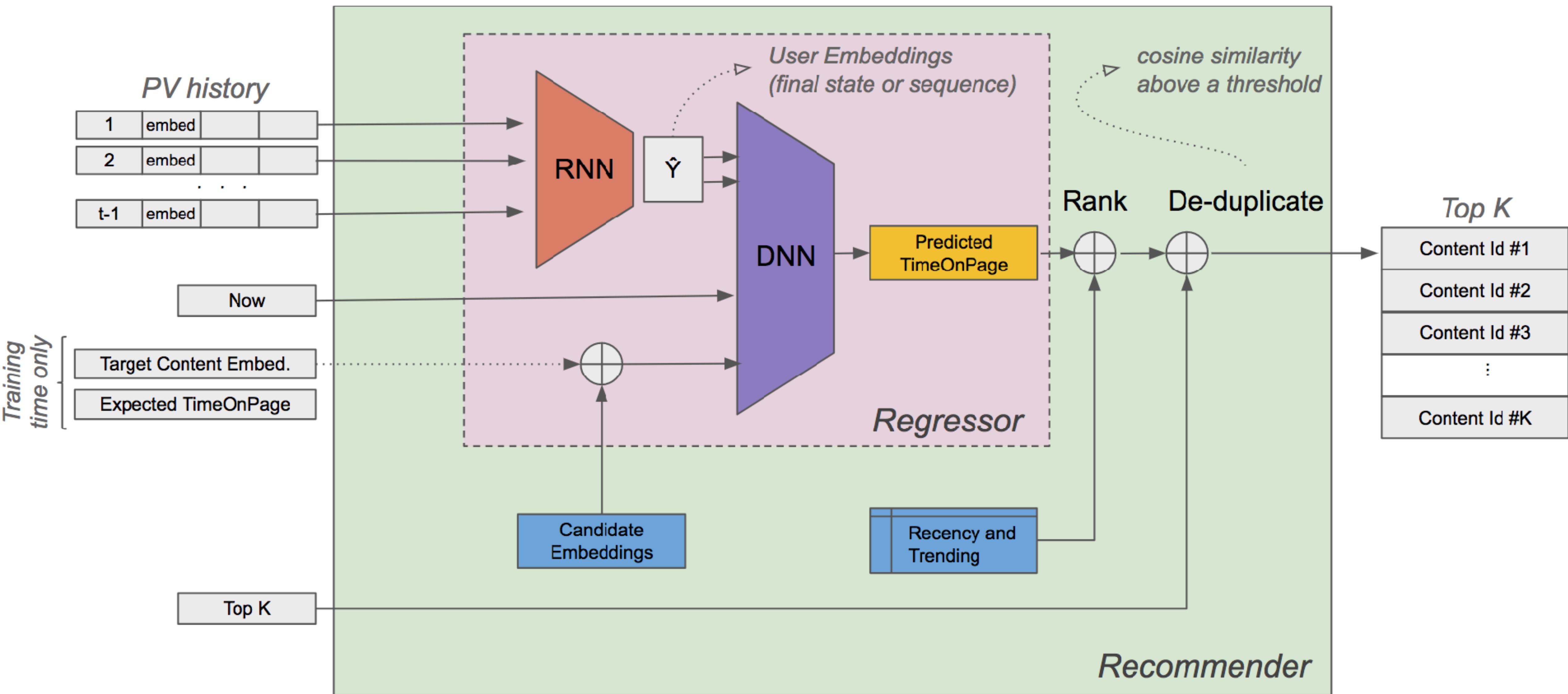
Fig. 1. Two-dimension scheme for classification of deep learning based recommender system. The left part illustrates the neural network models, and the right part illustrates the integration models.

- Neural Autoregressive Distribution Estimation (NADE) [57, 108] is an unsupervised neural network built atop autoregressive model and feedforward neural network. It is a tractable and efficient estimator for modelling data distributions and densities.

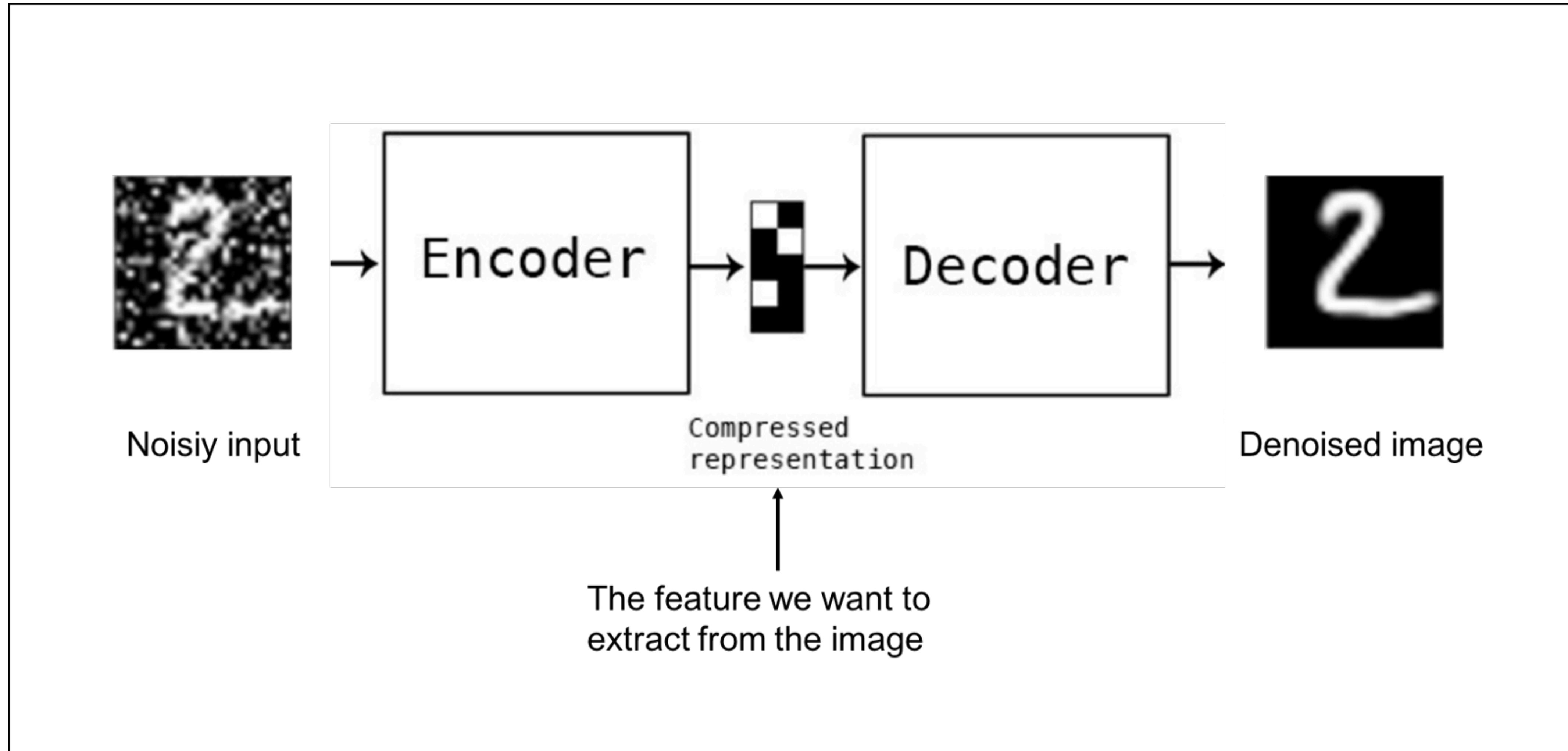
# Proposta Inicial



# Proposta após feedback do Google

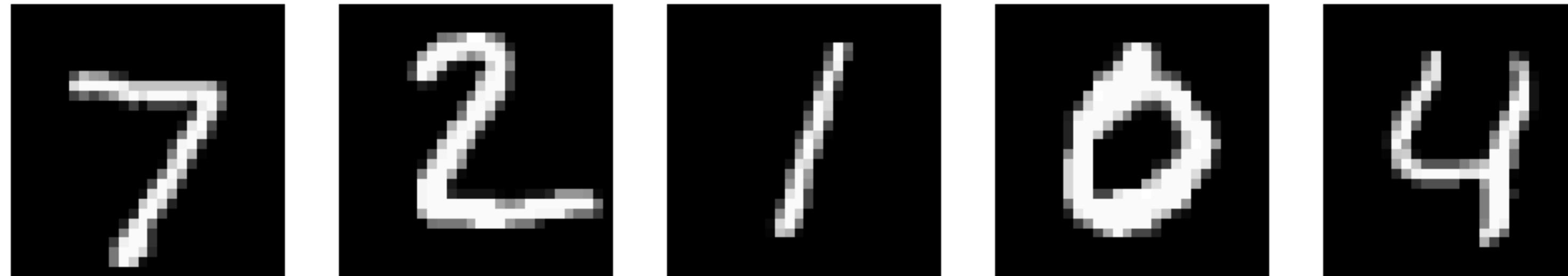


# Deep Autoencoder

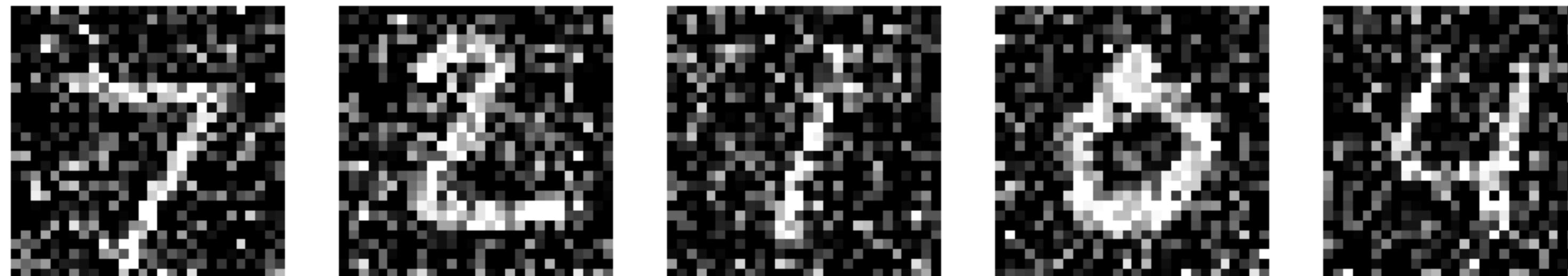


# Deep Autoencoder

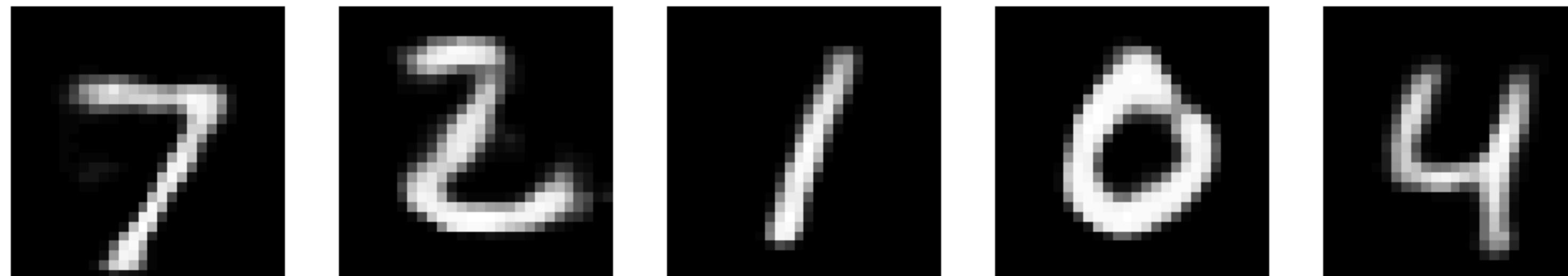
Original Images



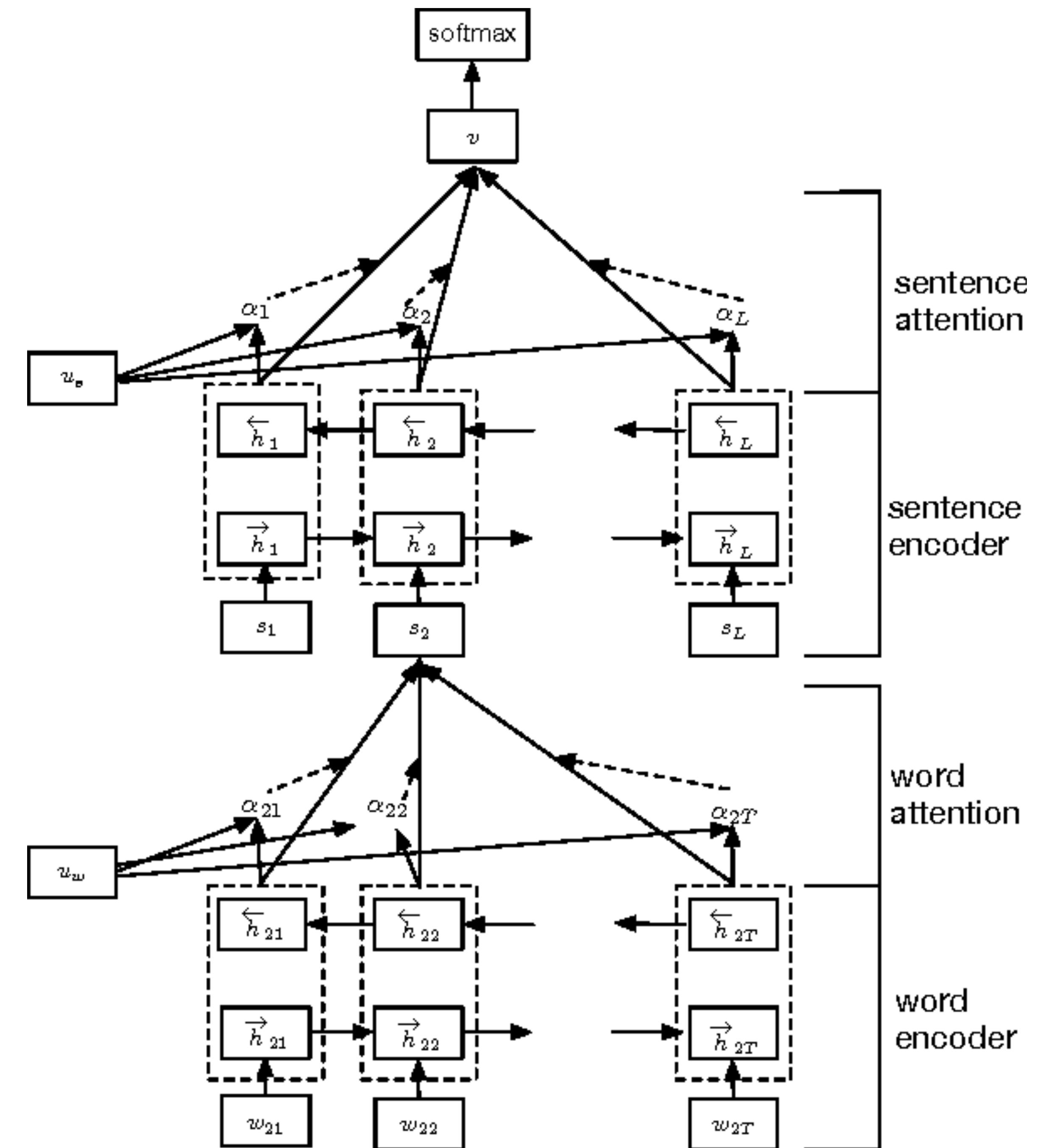
Noisy Input



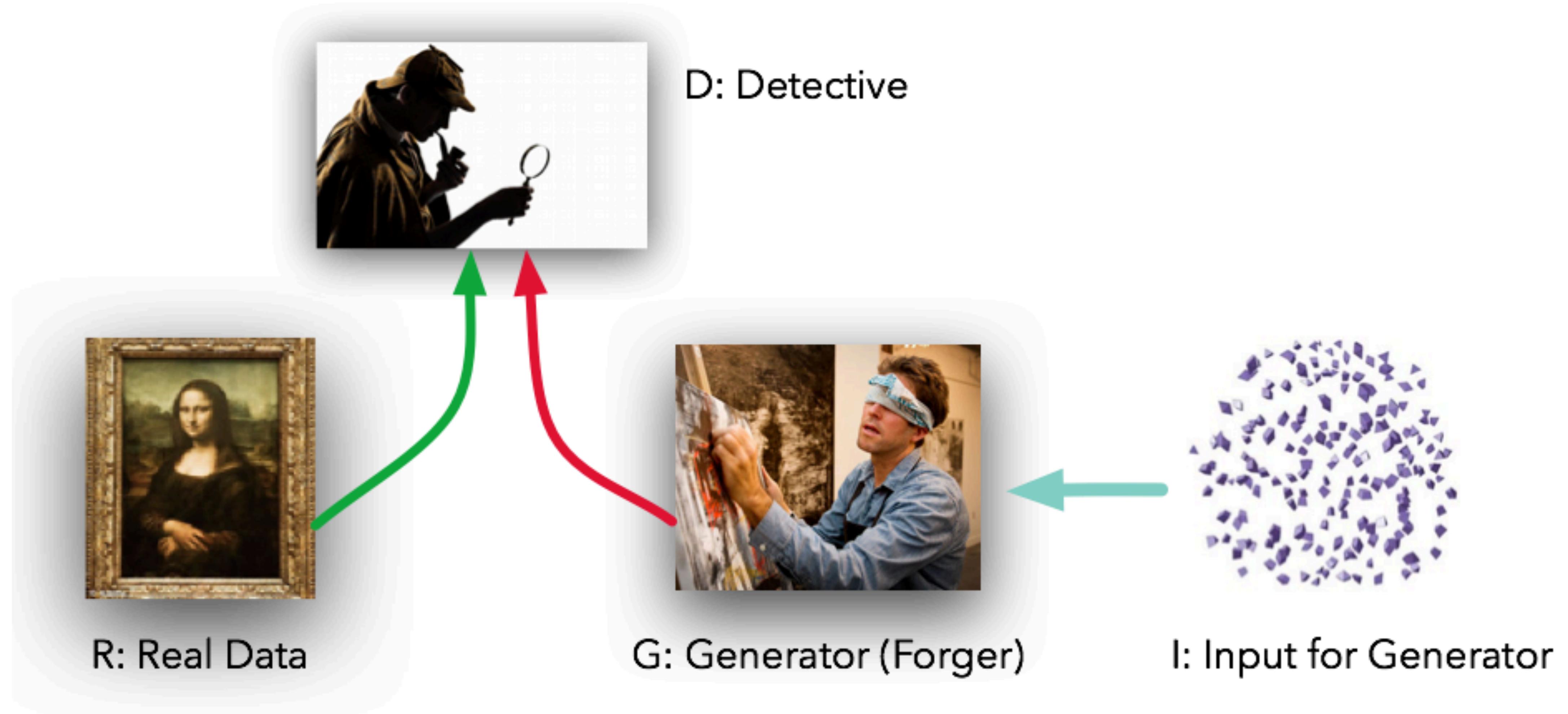
Autoencoder Output



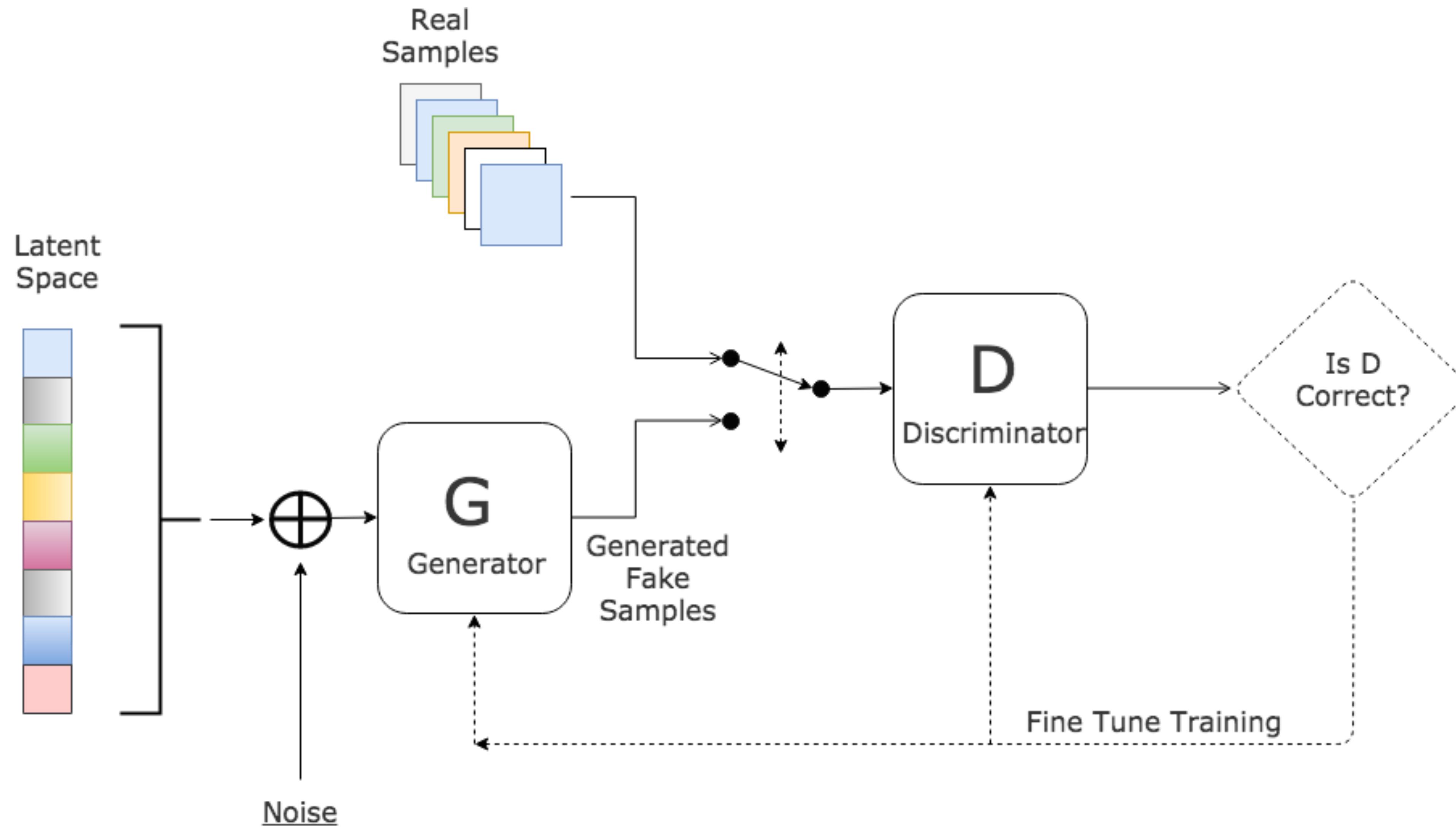
# HANs: Hierarchical Attention Network



# GAN: Generative Adversarial Network (2013)



# GAN: Generative Adversarial Network (2013)

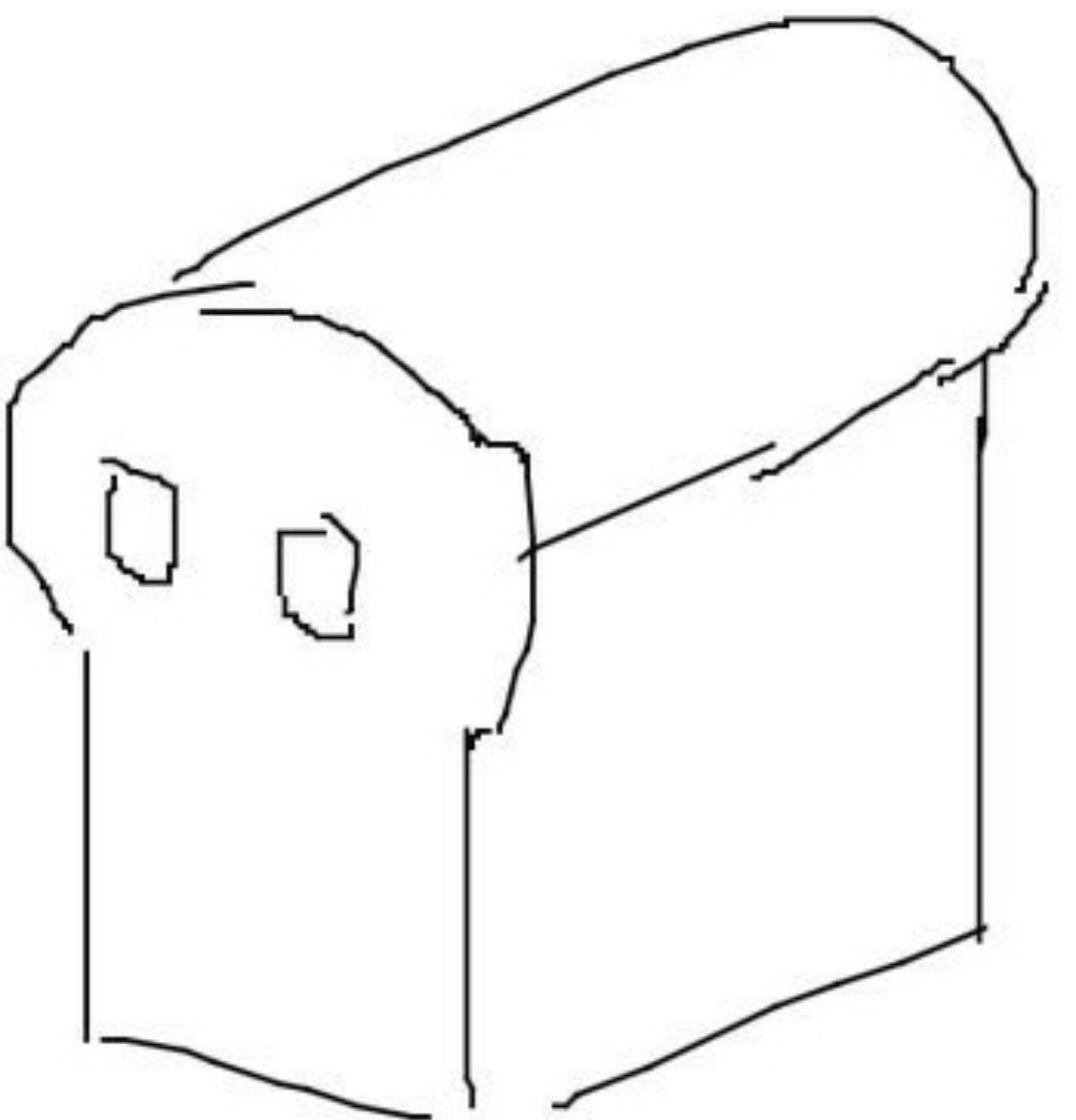


# **Espaço Latente “Gato”**

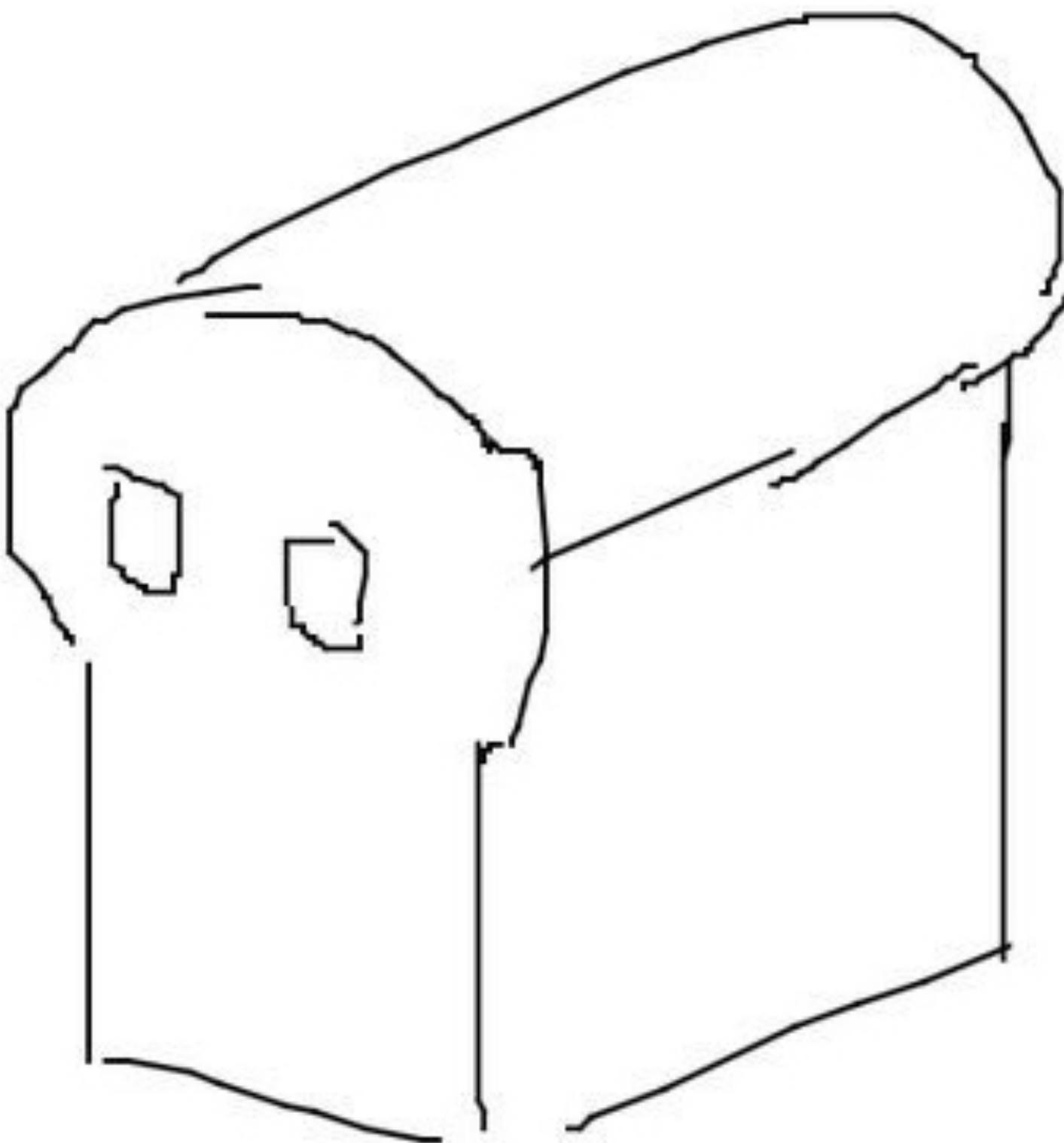
# Gato fake



# Pão de forma => ?



Pão de forma => “Pão Gato” :-)



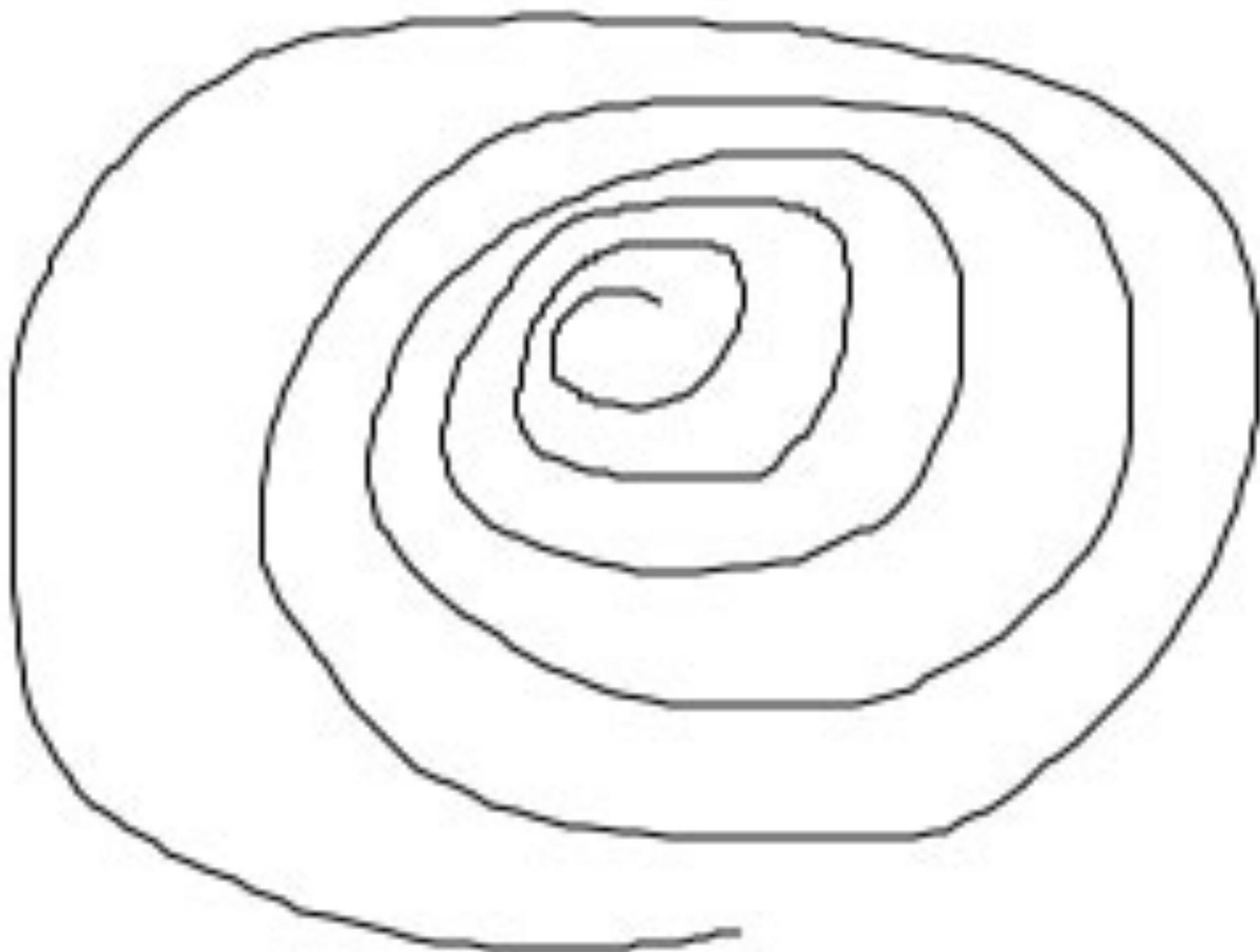
# Beholder (Dungeons & Dragons)



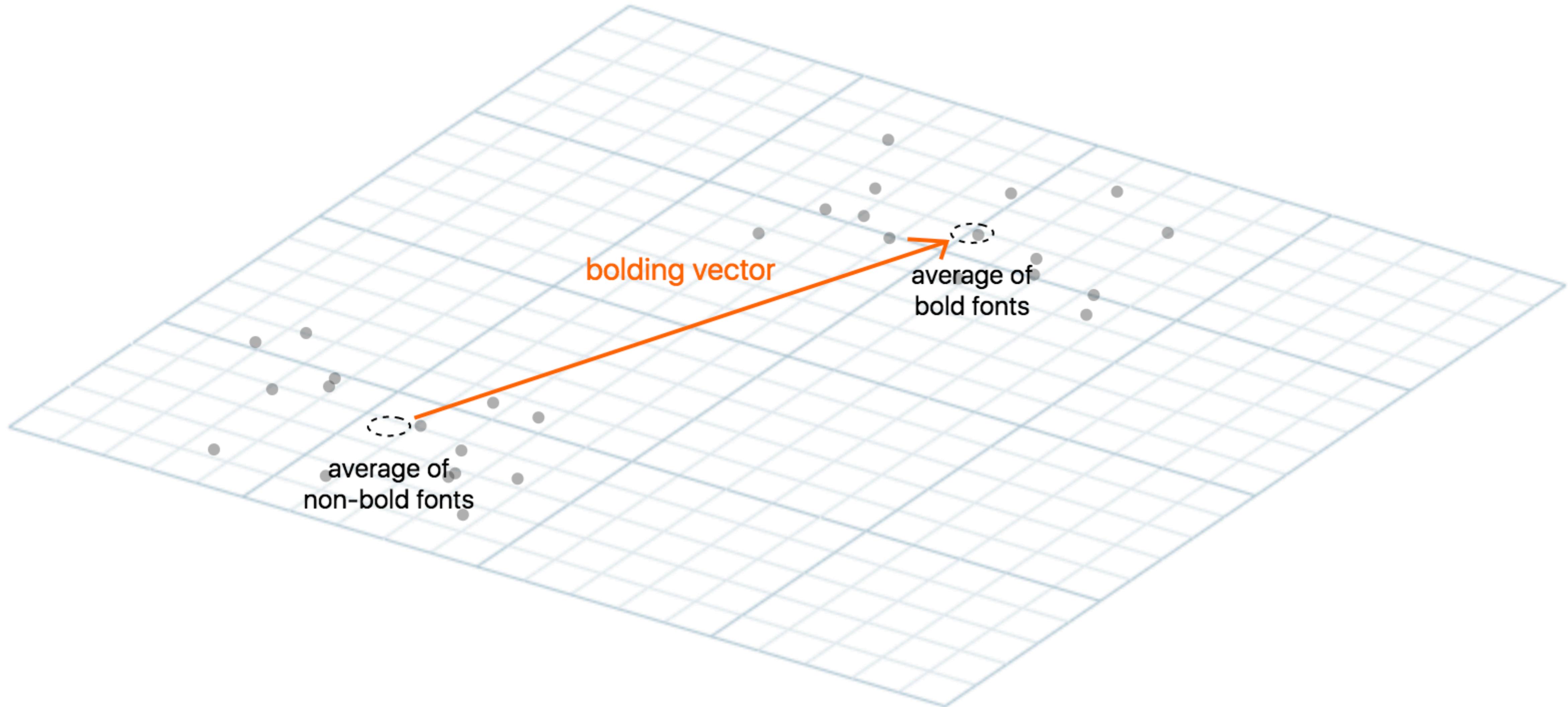
# Beholder Gato



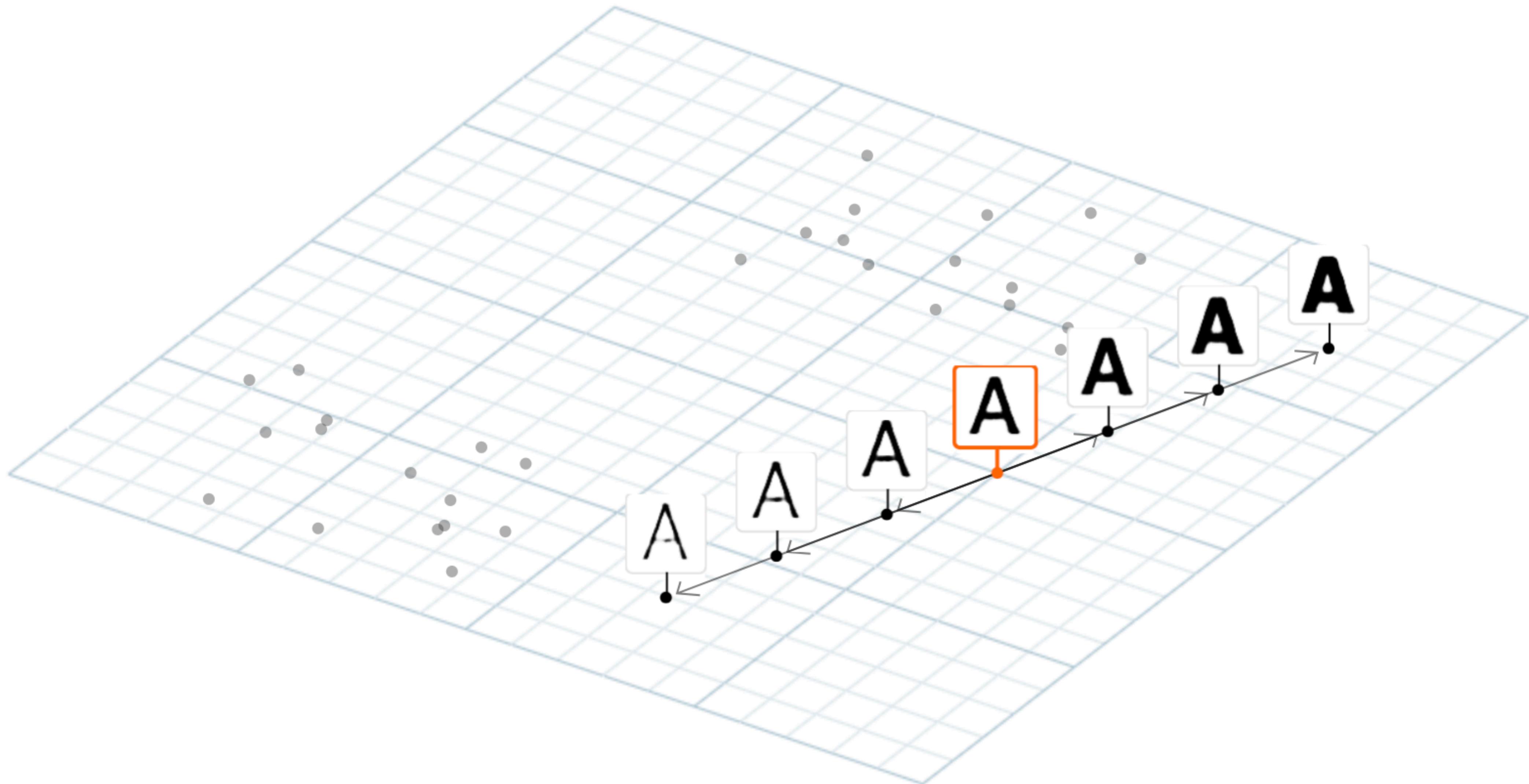
# Gato espiral



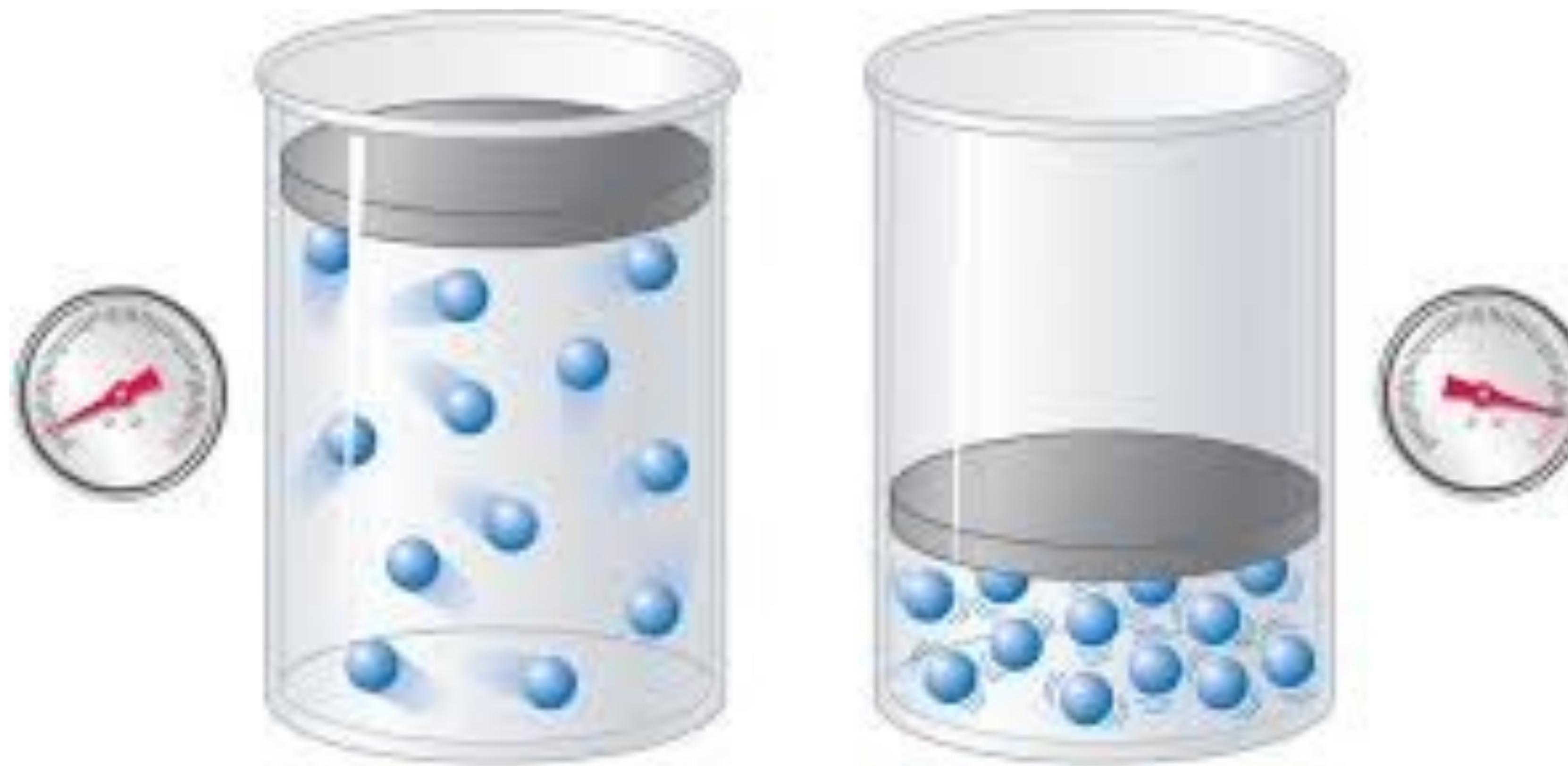
# Estilo de Fontes



# Estilo de Fontes



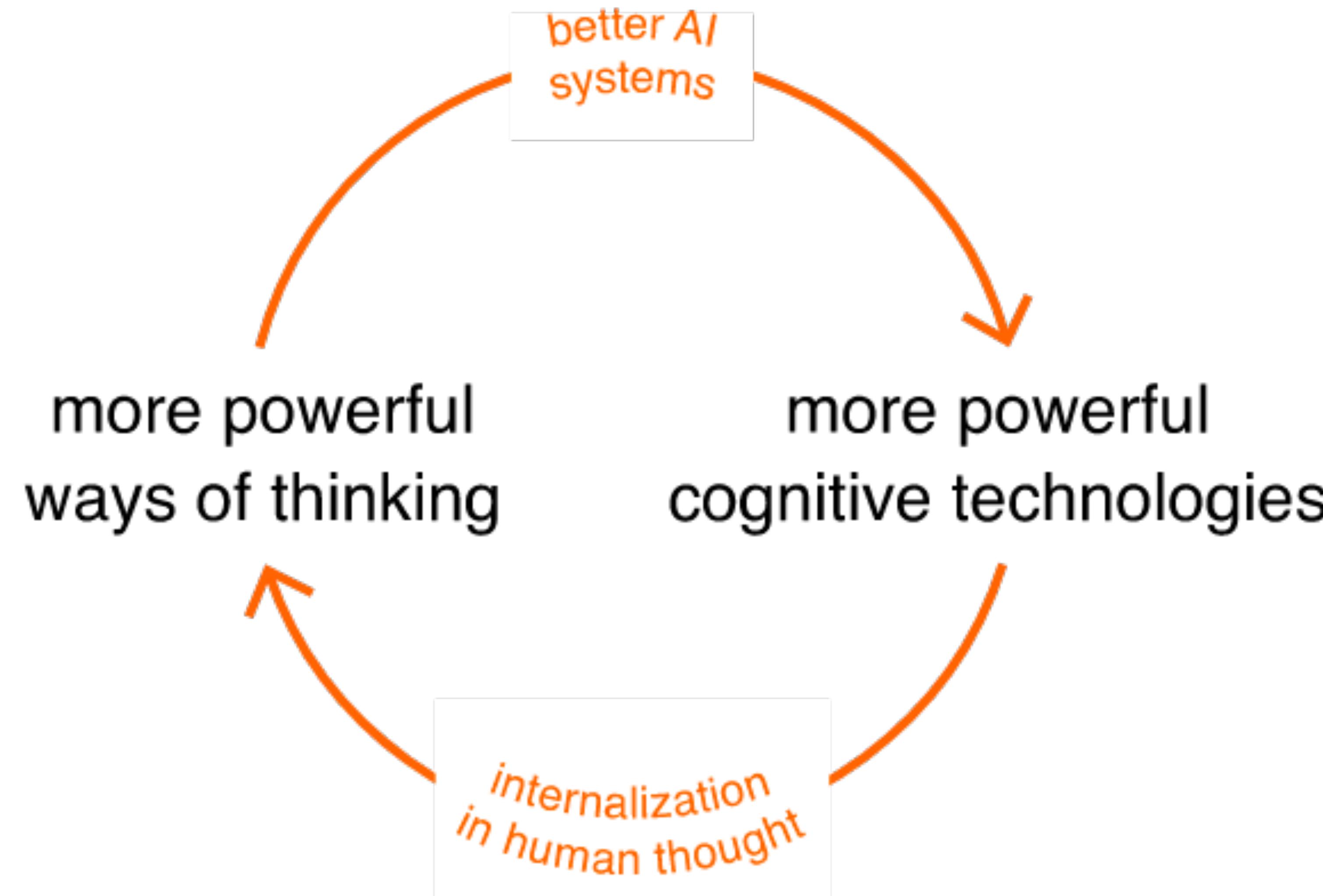
# Novas Primitivas



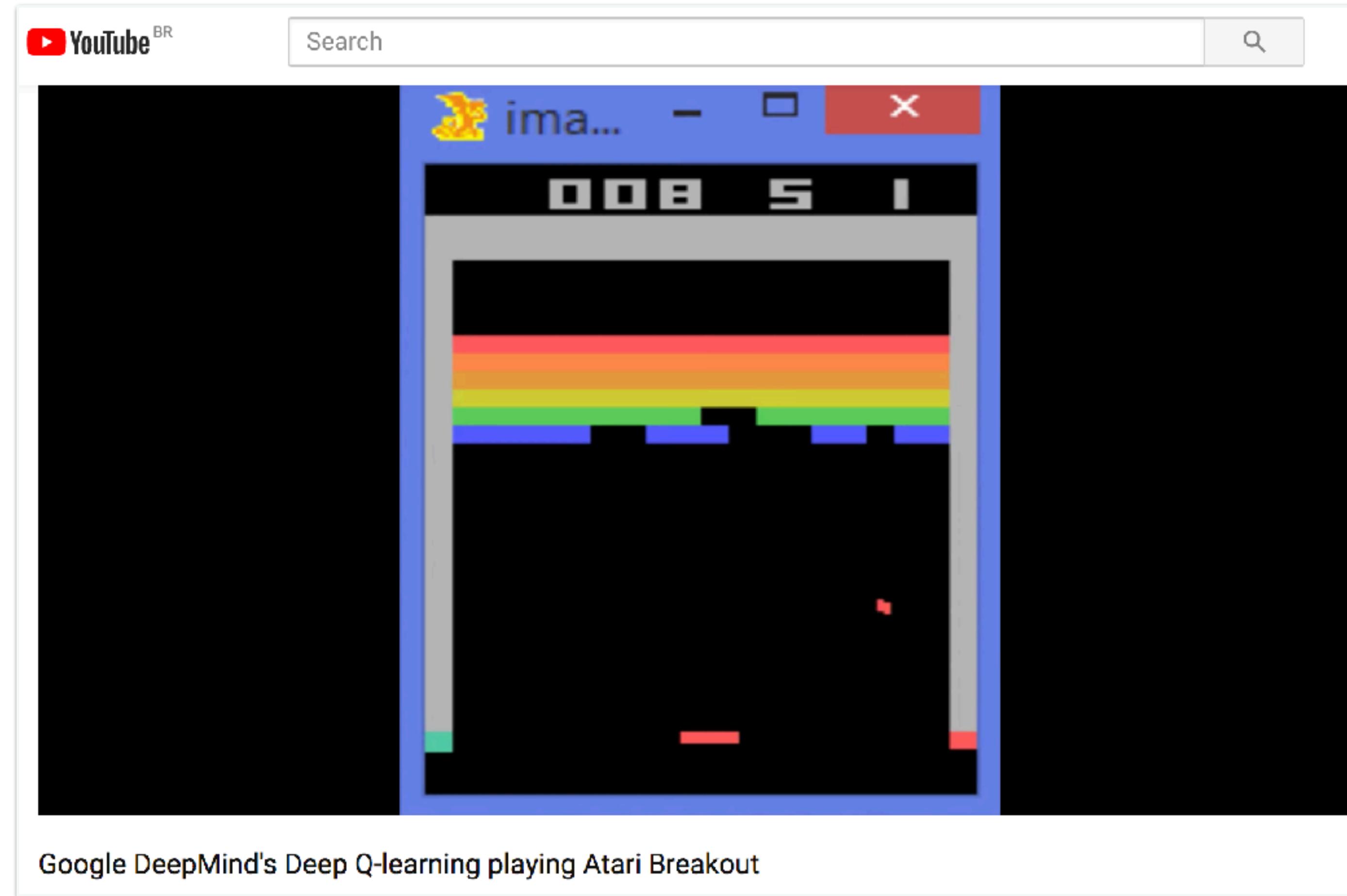
**(a) Low pressure**

**(b) High pressure**

# Singularidade

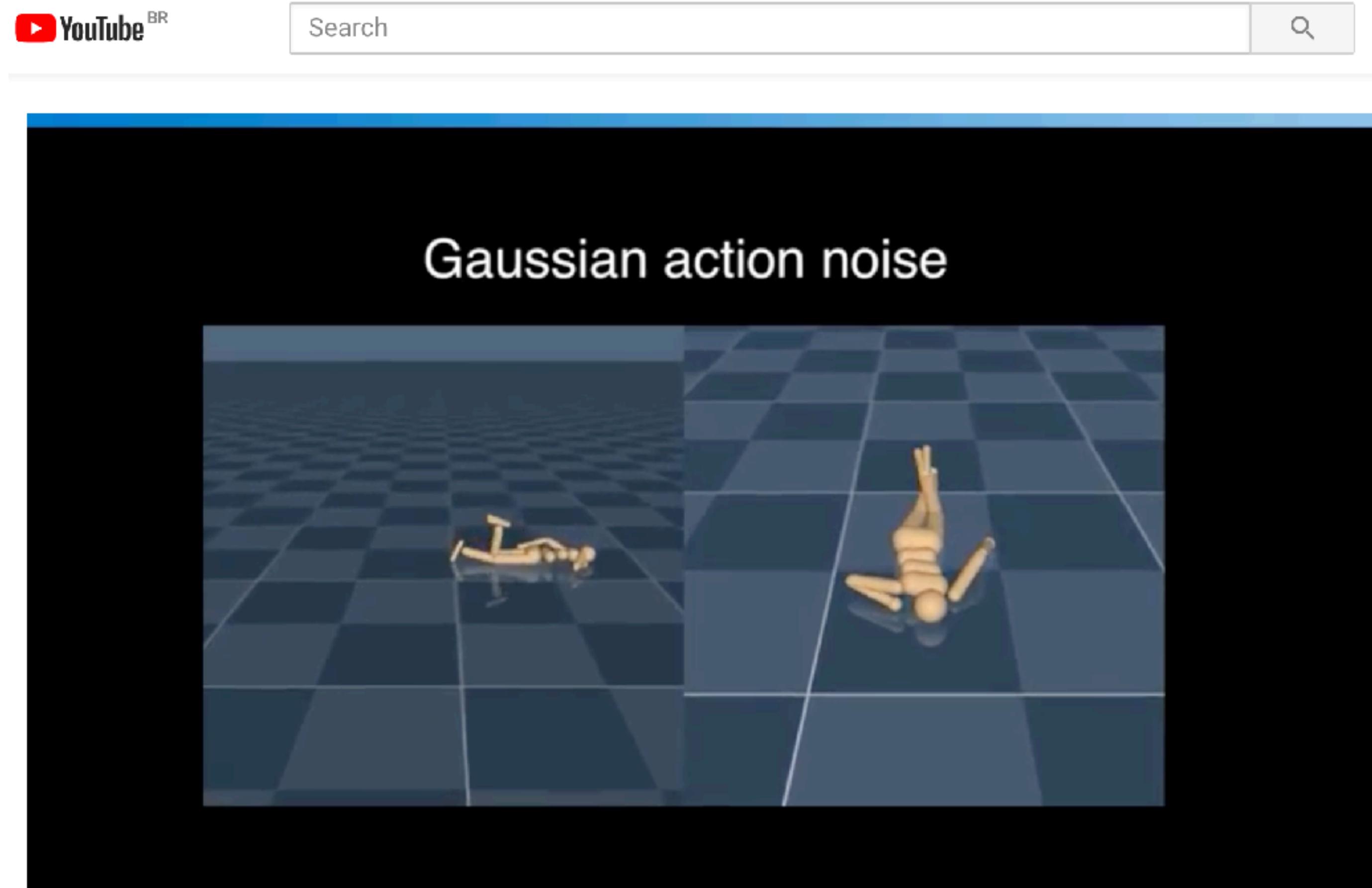


# Reinforcement Learning



<https://www.youtube.com/watch?v=V1eYniJ0Rnk>

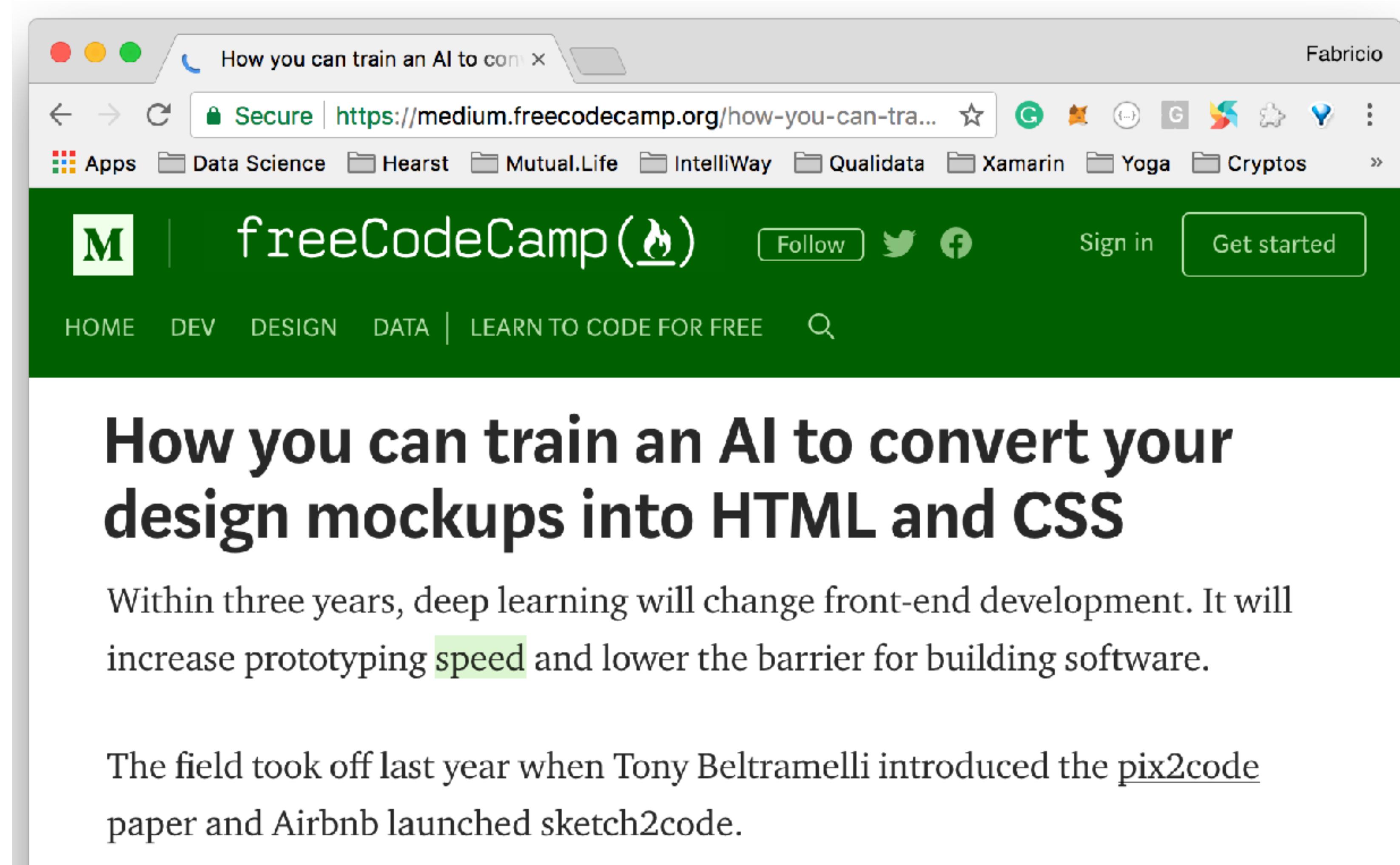
# Reinforcement Learning



RAAIS 2017 - Raia Hadsell, Senior Research Scientist at DeepMind

[https://www.youtube.com/watch?v=0e\\_uGa7ic74&feature=youtu.be&t=2065](https://www.youtube.com/watch?v=0e_uGa7ic74&feature=youtu.be&t=2065)

# Profissões impactadas: todas!



The screenshot shows a web browser window with a green header bar. The header bar includes a red circular icon, a yellow circular icon, a green circular icon, a blue refresh icon, the text "How you can train an AI to convert your design mockups into HTML and CSS", and the name "Fabricio". Below the header is a toolbar with icons for back, forward, search, and other functions. The main content area has a dark green background. At the top left is a white square icon with a black "M". To its right is the text "freeCodeCamp" followed by a small logo consisting of a flame-like shape above a downward-pointing arrow. To the right of the logo are buttons for "Follow", a Twitter icon, a Facebook icon, "Sign in", and a "Get started" button. Below this is a navigation bar with links for "HOME", "DEV", "DESIGN", "DATA", "LEARN TO CODE FOR FREE", and a magnifying glass icon for search. The main title of the article is "How you can train an AI to convert your design mockups into HTML and CSS". The text below the title reads: "Within three years, deep learning will change front-end development. It will increase prototyping speed and lower the barrier for building software." The text at the bottom states: "The field took off last year when Tony Beltramelli introduced the pix2code paper and Airbnb launched sketch2code."

How you can train an AI to convert your design mockups into HTML and CSS

Within three years, deep learning will change front-end development. It will increase prototyping speed and lower the barrier for building software.

The field took off last year when Tony Beltramelli introduced the pix2code paper and Airbnb launched sketch2code.



