

## The 46th International School for Young Astronomers ISYA 2025: Ecuador



In[.]:= nbook = EvaluationNotebook[];

## Demo ChatGPT

**El material se puede encontrar en:**

[https://github.com/josemramirez/  
ChatGPT curso 2025 Mathematica v1](https://github.com/josemramirez/ChatGPT_curso_2025_Mathematica_v1)

PDF:

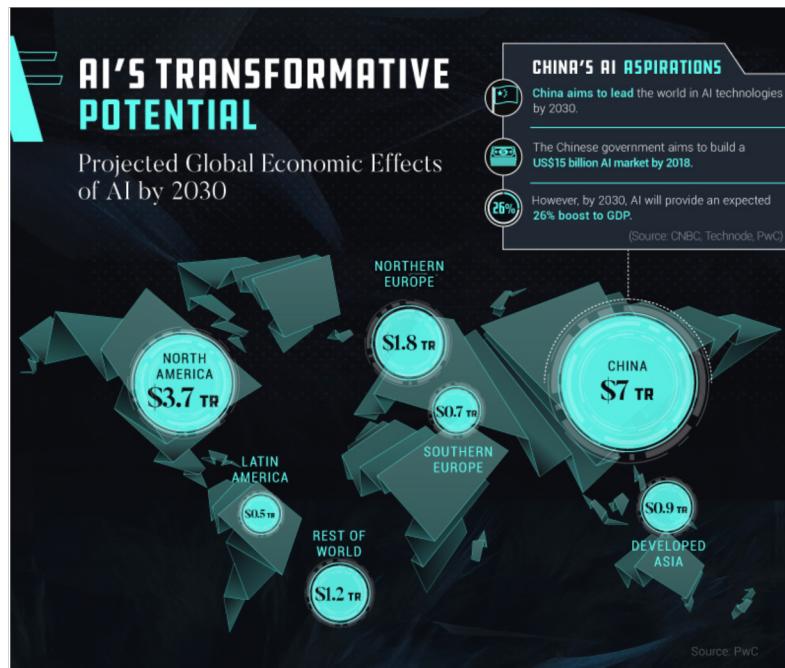
[https://drive.google.com/file/d/17oK2jn7bpv\\_fQYRBtahXYvLb7I2Xk4s4/view?  
usp=sharing](https://drive.google.com/file/d/17oK2jn7bpv_fQYRBtahXYvLb7I2Xk4s4/view?usp=sharing)



Out[·]=

◀ prev - + prox ▶

# ¿Dónde estamos? 101



Crédito: <https://www.visualcapitalist.com/economic-impact-artificial-intelligence-ai/>

Out[...]=

◀ prev - + prox ▶



Crédito: <https://www.visualcapitalist.com/economic-impact-artificial-intelligence-ai/>

Out[•]=

◀ prev – + prox ▶

# Inteligencia Artificial

Redes Neuronales,  
Natural Language Processing (NLP)  
Un curso introductorio 101



Demo de yachayGPT

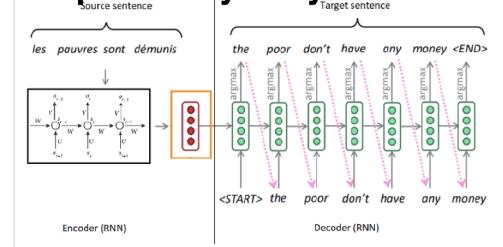


<https://t.me/yachayGPT4bot>

José Manuel Ramírez, ISYA Julio 2025



<https://t.me/yachayGPT4bot>



Out[-]=

◀ prev - + prox ▶

# Generative Adversarial Nets (GANs)

## Generative Adversarial Nets

Ian J. Goodfellow,<sup>1</sup> Jean Pouget-Abadie,<sup>2</sup> Mehdi Mirza,<sup>1</sup> Bing Xu,<sup>1</sup> David Warde-Farley,<sup>1</sup> Sherjil Ozair,<sup>1</sup> Aaron Courville,<sup>1</sup> Yoshua Bengio<sup>1</sup>  
<sup>1</sup>Département d'information et recherche opérationnelle  
<sup>2</sup>Université de Montréal  
 Montréal, QC H3C 3J7

### Abstract

We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model  $G$  that captures the data distribution, and a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than  $G$ . The training procedure for  $G$  is to maximize the probability of  $D$  making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary functions  $\mathcal{G}$  and  $\mathcal{D}$ , a unique solution exists with  $G$  outputting data drawn from the data distribution and  $D$  equal to  $\frac{1}{2}$  everywhere. In the case where  $G$  and  $D$  are defined by multilayer perceptrons, the entire system can be trained with backpropagation. There is no need for any Markov chains or unrolled approximate inference networks during either training or generation of samples. Experiments demonstrate the potential of the framework through qualitative and quantitative evaluation of the generated samples.

### 1 Introduction

The promise of deep learning is to discover rich, hierarchical models [2] that represent probability distributions over the kinds of data encountered in artificial intelligence applications, such as natural images, audio waveforms containing speech, and symbols in natural language corpora. So far, the most striking successes in deep learning have involved discriminative models, usually those that map a high-dimensional, rich sensory input to a class label [14, 22]. These striking successes have primarily been achieved by using neural networks with rectified linear units (ReLU) [19, 9, 10] which have a particularly well-behaved gradient. Deep generative models have had less of an impact, due to the difficulty of approximating many intractable probabilistic computations that arise in maximum likelihood estimation and related strategies, and due to difficulty of leveraging the benefits of piecewise linear units in the generative context. We propose a new generative model estimation procedure that sidesteps these difficulties.<sup>1</sup>

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

---

```

for number of training iterations do
    for  $k$  steps do
        • Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_{\theta}(z)$ .
        • Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
        • Update the discriminator by ascending its stochastic gradient:
            
$$\nabla_{\theta_D} \frac{1}{m} \sum_{i=1}^m [\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)})))]$$
.
    end for
    • Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_{\theta}(z)$ .
    • Update the generator by descending its stochastic gradient:
            
$$\nabla_{\theta_G} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)})))$$
.
end for
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

```

---

#### 4.1 Global Optimality of $p_{\theta} = p_{\text{data}}$

We first consider the optimal discriminator  $D$  for any given generator  $G$ .

**Proposition 1.** For  $G$  fixed, the optimal discriminator  $D$  is

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\theta}(\mathbf{x})} \quad (2)$$

*Proof.* The training criterion for the discriminator  $D$ , given any generator  $G$ , is to maximize the quantity  $V(G, D)$

$$V(G, D) = \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{x}} p_{\theta}(\mathbf{x}) \log(1 - D(G(\mathbf{x}))) d\mathbf{x} \\ = \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) + p_{\theta}(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x} \quad (3)$$

For any  $(a, b) \in \mathbb{R}^2 \setminus \{(0, 0)\}$ , the function  $y \rightarrow a \log(y) + b \log(1 - y)$  achieves its maximum in  $[0, 1]$  at  $\frac{a}{a+b}$ . The discriminator does not need to be defined outside of  $\text{Supp}(p_{\text{data}}) \cup \text{Supp}(p_{\theta})$ , concluding the proof.  $\square$

Note that the training objective for  $D$  can be interpreted as maximizing the log-likelihood for estimating the conditional probability  $P(Y=y|\mathbf{x})$ , where  $Y$  indicates whether  $\mathbf{x}$  comes from  $p_{\text{data}}$  (with  $y=1$ ) or from  $p_{\theta}$  (with  $y=0$ ). The minimax game in Eq. 1 can now be reformulated as:

**Theorem 1.** The global minimum of the virtual training criterion  $C(G)$  is achieved if and only if  $p_{\theta} = p_{\text{data}}$ . At that point,  $C(G)$  achieves the value  $-\log 4$ .

*Proof.* For  $p_{\theta} = p_{\text{data}}$ ,  $D_G^*(\mathbf{x}) = \frac{1}{2}$ , (consider Eq. 2). Hence, by inspecting Eq. 4 at  $D_G^*(\mathbf{x}) = \frac{1}{2}$ , we find  $C(G) = \log \frac{1}{4} + \log \frac{1}{2} = -\log 4$ . To see that this is the best possible value of  $C(G)$ , reached only for  $p_{\theta} = p_{\text{data}}$ , observe that

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}}[-\log 2] + \mathbb{E}_{\mathbf{x} \sim p_{\theta}}[-\log 2] = -\log 4$$

and that by subtracting this expression from  $C(G) = V(D_G^*, G)$ , we obtain:

$$C(G) = -\log(4) + KL\left(p_{\text{data}} \middle\| \frac{p_{\text{data}} + p_{\theta}}{2}\right) + KL\left(p_{\theta} \middle\| \frac{p_{\text{data}} + p_{\theta}}{2}\right) \quad (5)$$

where  $KL$  is the Kullback–Leibler divergence. We recognize in the previous expression the Jensen–Shannon divergence between the model's distribution and the data generating process:

$$C(G) = -\log(4) + 2 \cdot JSD(p_{\text{data}} || p_{\theta}) \quad (6)$$

Since the Jensen–Shannon divergence between two distributions is always non-negative and zero only when they are equal, we have shown that  $C^* = -\log(4)$  is the global minimum of  $C(G)$  and that the only solution is  $p_{\theta} = p_{\text{data}}$ , i.e., the generative model perfectly replicating the data generating process.  $\square$

#### 4.2 Convergence of Algorithm 1

**Proposition 2.** If  $G$  and  $D$  have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given  $G$ , and  $p_{\theta}$  is updated so as to improve the criterion

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}}[\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_{\theta}}[\log(1 - D_G^*(\mathbf{x}))]$$

then  $p_{\theta}$  converges to  $p_{\text{data}}$ .

*Proof.* Consider  $V(G, D) = U(p_{\theta}, D)$  as a function of  $p_{\theta}$  as done in the above criterion. Note that  $U(p_{\theta}, D)$  is convex in  $p_{\theta}$ . The subderivatives of a sum of convex functions include the derivative of the function at the point where the maximum is attained. In other words, if  $f(x) = \sup_{\alpha \in A} f_{\alpha}(x)$  and  $f_{\alpha}(x)$  is convex in  $x$  for every  $\alpha$ , then  $\partial f(x) \in \partial f$  if  $\beta = \arg \sup_{\alpha \in A} f_{\alpha}(x)$ . This is equivalent to computing a gradient descent update for  $p_{\theta}$  at the optimal  $D$  given the corresponding  $G$ .  $\sup_{p_{\theta}} U(p_{\theta}, D)$  is convex in  $p_{\theta}$  with a unique global optima as proven in Thm 1, therefore with sufficiently small updates of  $p_{\theta}$ ,  $p_{\theta}$  converges to  $p_{\text{data}}$ , concluding the proof.  $\square$

In practice, adversarial nets represent a limited family of  $p_{\theta}$  distributions via the function  $G(\mathbf{z}; \theta_G)$ , and we optimize  $\theta_G$  rather than  $p_{\theta}$  itself. Using a multilayer perceptron to define  $G$  introduces multiple critical points in parameter space. However, the excellent performance of multilayer perceptrons in practice suggests that they are a reasonable model to use despite their lack of theoretical guarantees.

## ChatGPT ¿Por qué funciona?

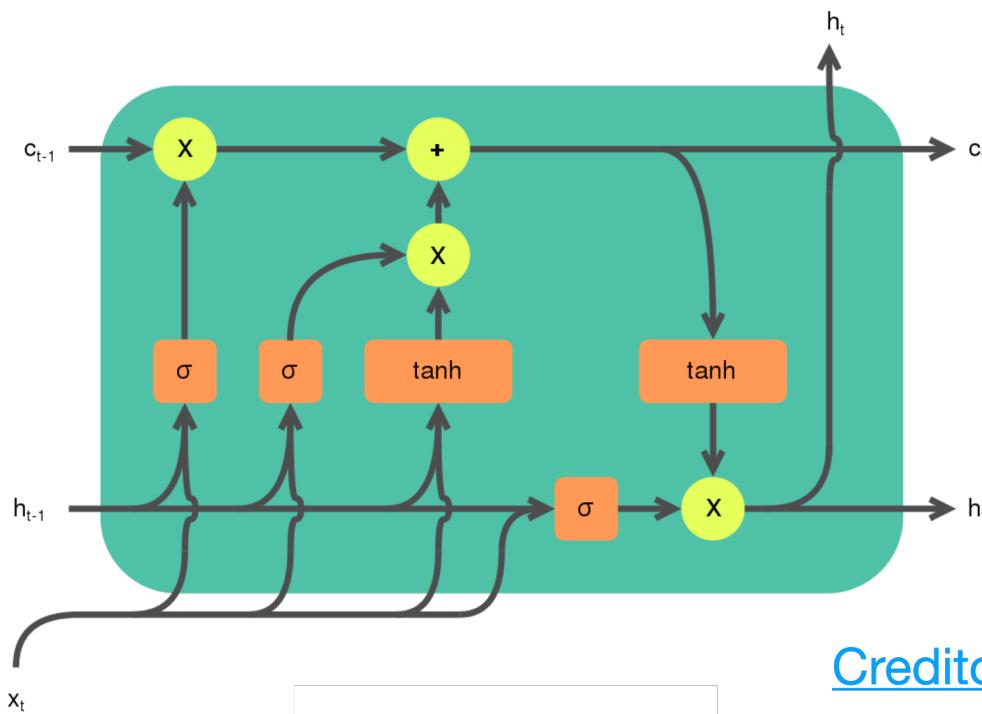
Zero2Hero (v1)

by JM Ramirez, Yachay 2025. (GPTZ2H)

basado en los escritos de:

<https://writings.stephenwolfram.com/2023/02/what-is-chatgpt-doing-and-why-does-it-work/>

# Es solo agregar una palabra a la vez!

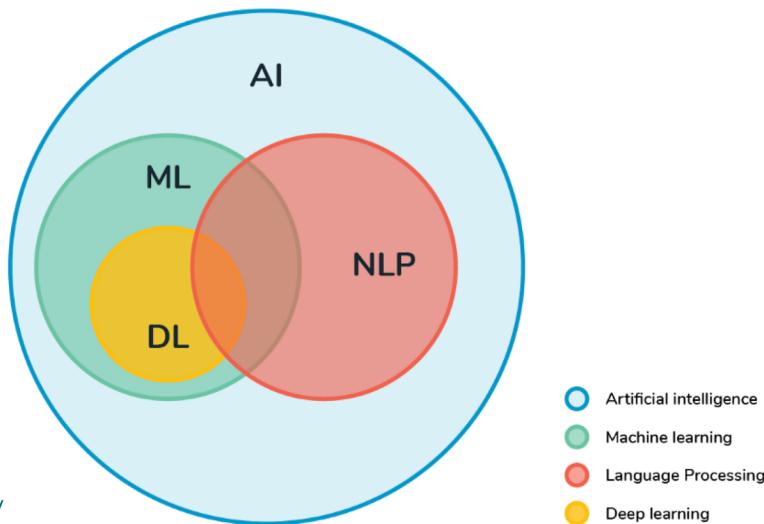


Credito: Wiki

Out[•]=

◀ prev – + prox ▶

# ¿Dónde estamos?



Credito: <https://becominghuman.ai/>

## ChatGPT: modelos LLMs

- Machine Learning ([Stanford](#))
  - DeepLearning ([DeepLearning.AI](#))
  - Bayesian Statistics ([HSE - Rusia](#))
  - Alfa Tester ([DeepLearning.AI](#))
  - Business Analytics ([Cambridge](#))



José M. Ramírez

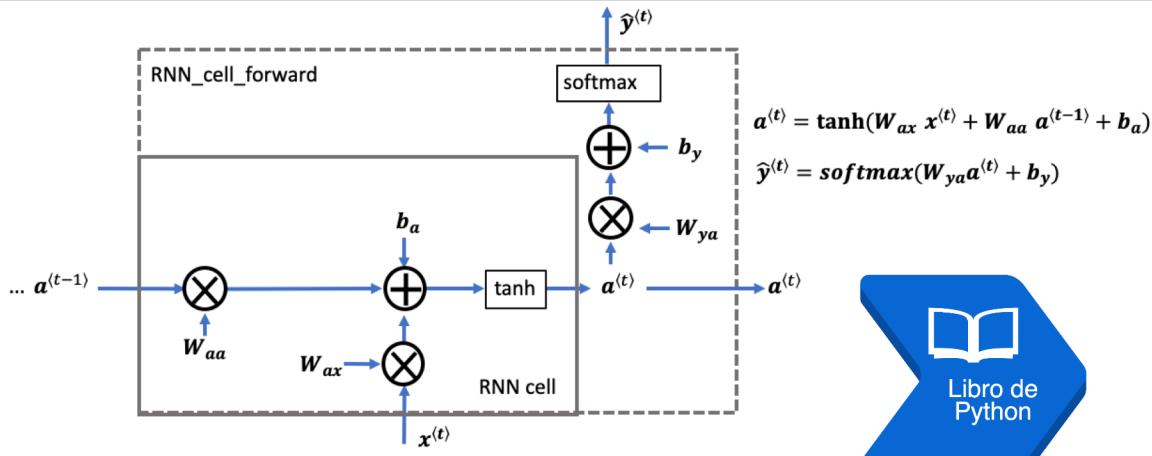


DeepLearning.ai



# Una celda RNN (101)

x



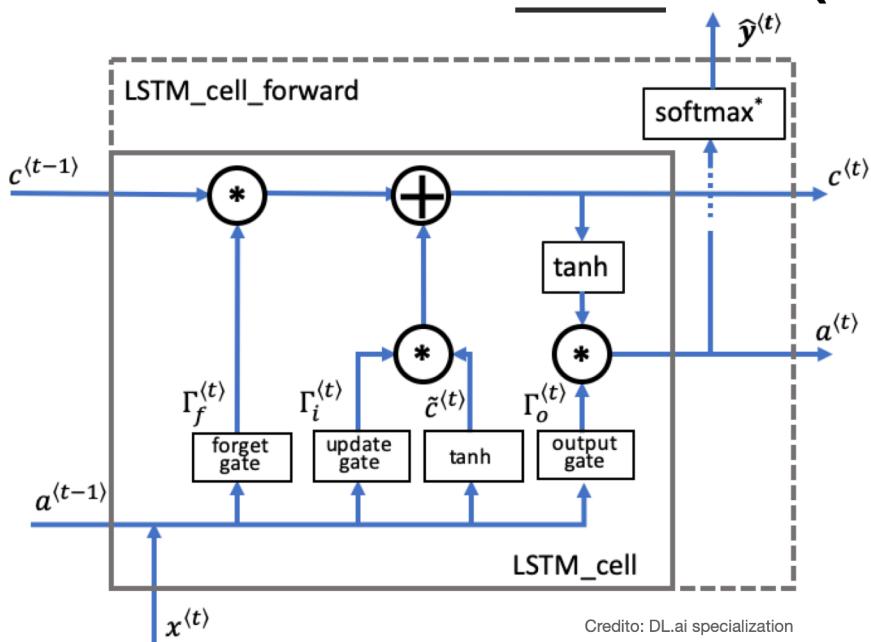
Crédito: DL.ai specialization

Out[1]:

◀ prev -
+ prox ▶

The best thing about AI is its ability to \_\_\_\_\_

## Una celda LSTM (101)



◀ prev – + prox ▶

# NN & A.I (v1.0)

Notas:

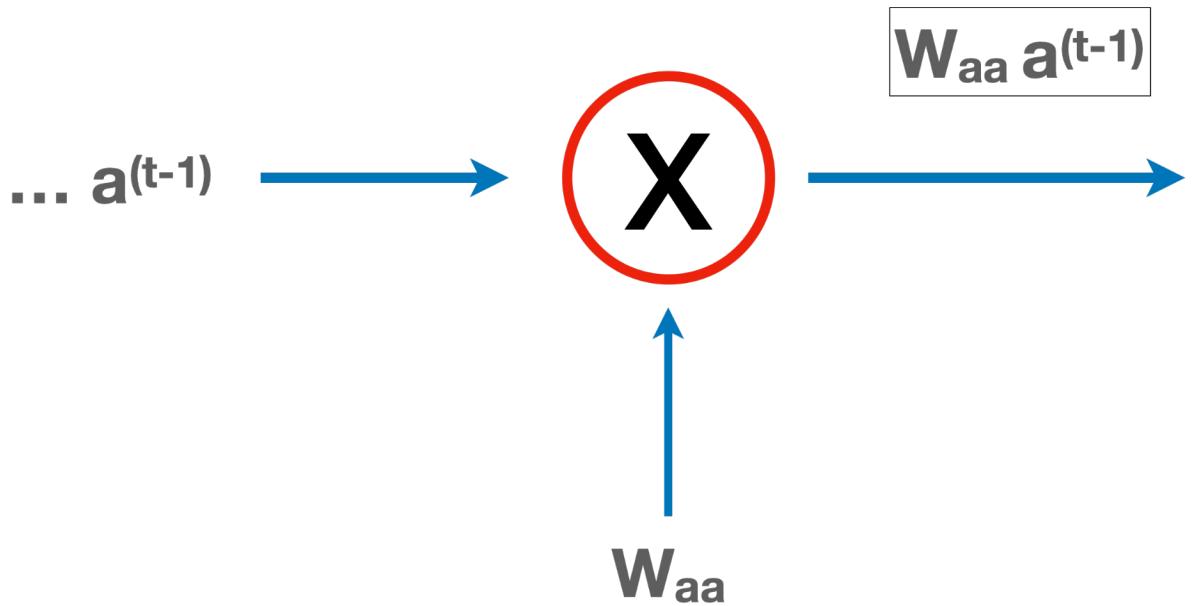
- Un poco de álgebra de vectores.
- Un poco de jerga computacional: Arreglos, Loops, Funciones, Argumentos, Variables, Inputs, Outputs.
- Láminas en Español, pero Libro de Python mezclado (estilo libre).
- Visualización. 2D Plots, 3D Plots.
- Un poco de ajustes de polinomios a datos experimentales.

José Manuel Ramírez, Junio 2025

...

Out[...]=

◀ prev - + prox ▶



Out[...]=

◀ prev - + prox ▶

# NN & A.I (v1.0)

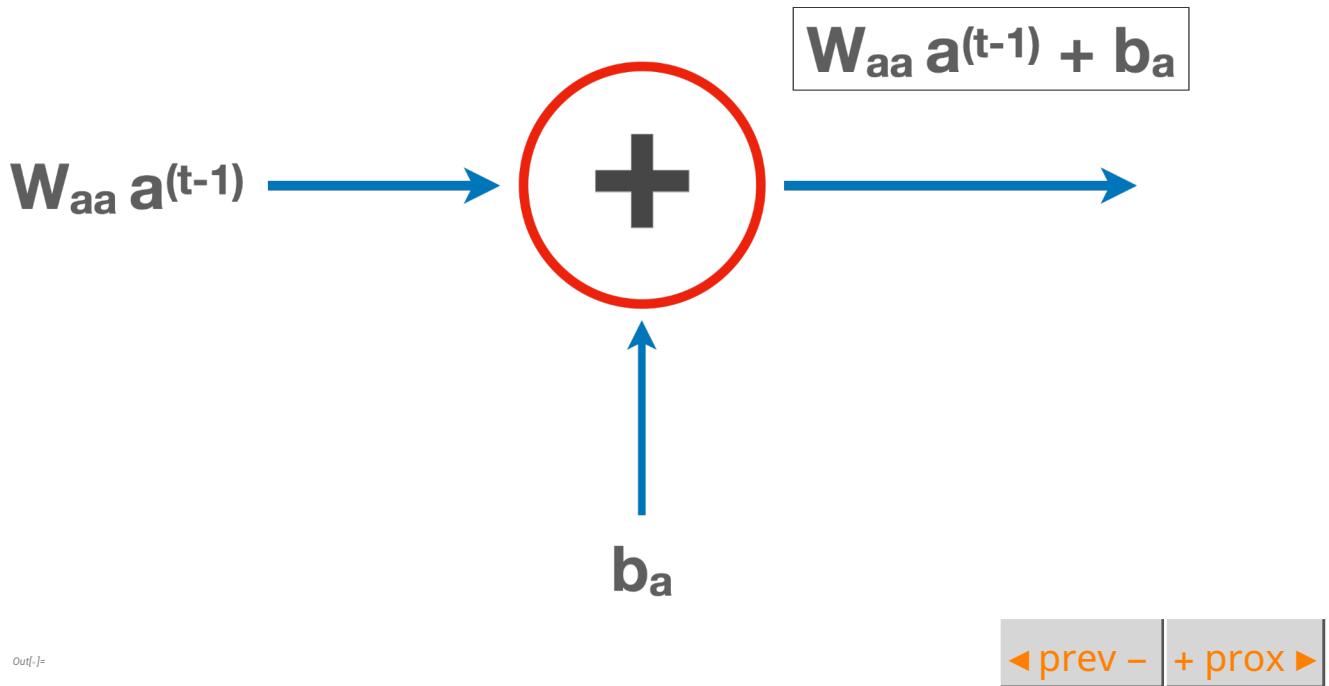
Notas:

- **Un poco de inteligencia artificial: Machine Learning, Deep Learning, NLP=Natural Language Processing.**
- **Input, output, activation functions, Matrices, operaciones con vectores.**
- **Celdas RNN, LSTM, GRU.**
- **“Entrenamiento” de una red neuronal.**

José Manuel Ramírez, Junio 2025

Out[*-*]=

◀ prev - + prox ▶



# NN & A.I (v1.0)

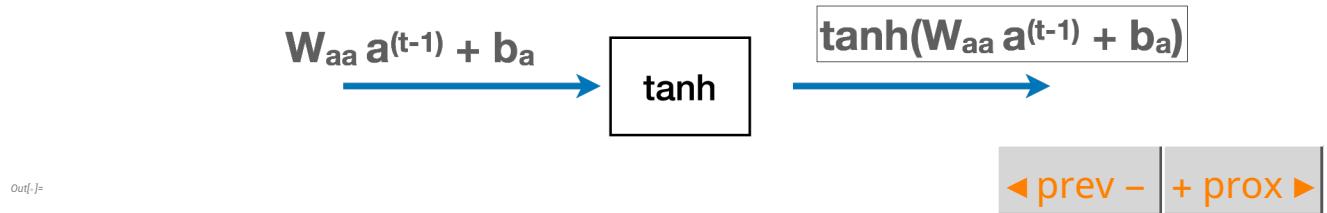
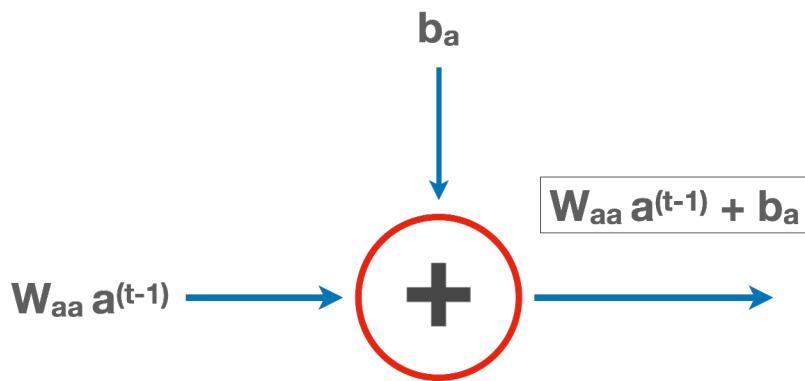
Material:

- Una función  $\tanh(x)$ .
- Recurrent Neural Network (RNN).
- LSTM cell.
- Modelar la salida de un solo LSTM.
- Crear Una “red” de LSTM para predecir texto.
- Escritura de un código de Python para realizar la tarea!
- Comentarios finales.

José Manuel Ramírez, Junio 2025

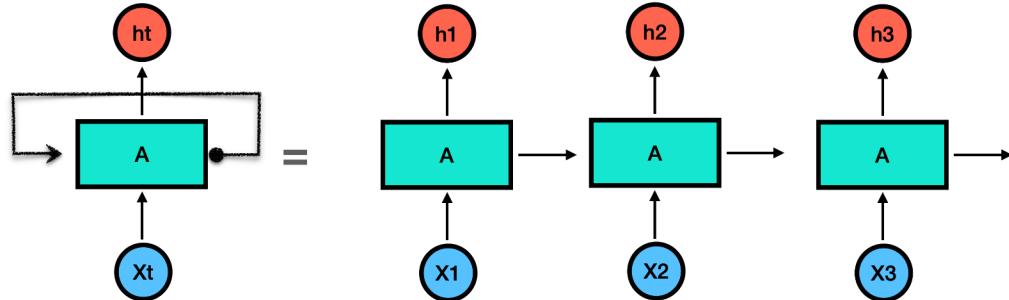
Out[*-*]=

◀ prev - + prox ▶



x

## Una red Neuronal ‘Recurrente’



An ‘unrolled’ Recurrent Neural Network  
(RNN)

Out[-]:=

◀ prev - + prox ▶

## Terminology Alert!!! Layers

(Oh no! It's another  
**Terminology Alert!!!**  
**Double Ugh!!**)

1

Neural networks are organized in **Layers**. Usually, a neural network has multiple **Input Nodes** that form an **Input Layer**...

*Out[.,]=*

2

...and usually there are multiple **Output Nodes** that form an **Output Layer**...

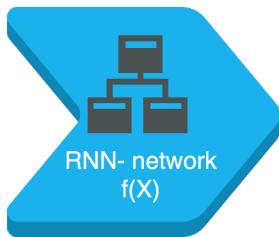
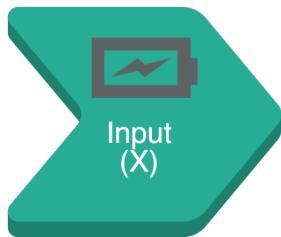
◀ prev – + prox ▶



The best thing about AI is its ability to \_\_\_\_\_

## Modelo Básico 101

x



### ● Datos a Modelar

Necesitamos elegir los datos de donde queremos sacar nuestro modelo.

### ● Celdas RNN

Caja negra que va a procesar la entrada: una función!

### ● Predicciones (y)

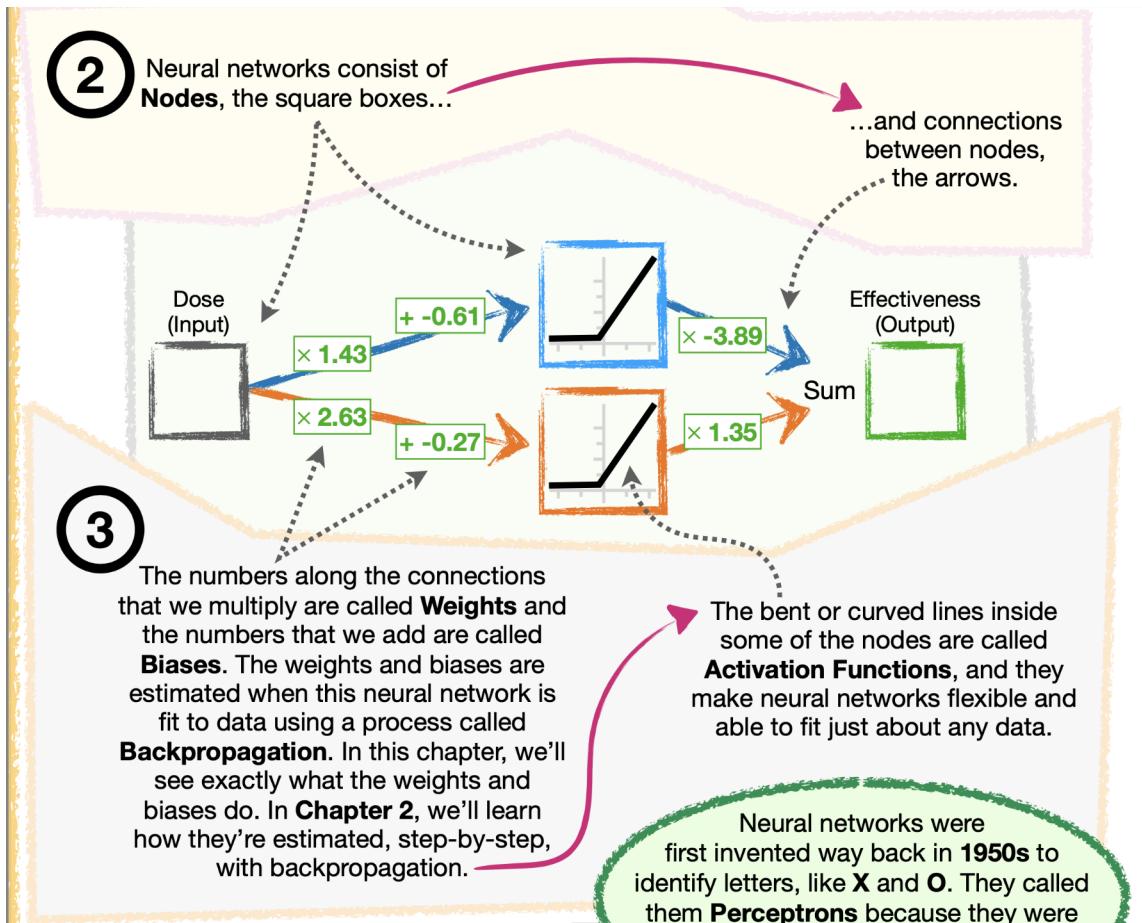
A una caja negra le inyectas una entrada, esta la procesa y sale algo!

Transformers:

<https://www.jeremyjordan.me/transformer-architecture/>

Out[...]=

◀ prev - + prox ▶



Out[...]=

◀ prev -
+ prox ▶

```
In[1]:= model =  
NetModel[{"GPT2 Transformer Trained on WebText Data", "Task" → "LanguageModeling"}]  
Out[1]= NetChain[  
In[2]:= model["The best thing about AI is its ability to", {"TopProbabilities", 5}]  
Out[2]= { do → 0.0288513, understand → 0.0307801, make → 0.0319072, predict → 0.0349754,  
learn → 0.0445319}  
In[3]:= Dataset[ReverseSort[Association[%]], ItemDisplayFunction → (PercentForm[#, 2] &)]  
Out[3]= 
```

## Si aplicamos el modelo

```
In[4]:= NestList[StringJoin[#, model[#], "Decision"]] &,  
"The best thing about AI is its ability to", 7]  
Out[4]= {The best thing about AI is its ability to,  
The best thing about AI is its ability to learn,  
The best thing about AI is its ability to learn from,  
The best thing about AI is its ability to learn from experience,  
The best thing about AI is its ability to learn from experience.,  
The best thing about AI is its ability to learn from experience. It,  
The best thing about AI is its ability to learn from experience. It's,  
The best thing about AI is its ability to learn from experience. It's not}
```

■ Si seguimos con Temperatura Zero:

...

```
In[1]:= StringReplace[Nest[StringJoin[#, model[#], "Decision"]], &,
  "The best thing about AI is its ability to", 132], "\n" .. → " "]
```

Out[1]= The best thing about AI is its ability to learn from experience. It's not just a matter of learning from experience, it's learning from the world around you. The AI is a very good example of this. It's a very good example of how to use AI to improve your life. It's a very good example of how to use AI to improve your life. The AI is a very good example of how to use AI to improve your life. It's a very good example of how to use AI to improve your life. The AI is a very good example of how to use AI to improve your life. It's a very good example of how to use AI to improve your life.

...

NOT GREAT...

## Pero con temperature de 0.8

```
In[1]:= NestList[StringJoin[#, model[#,{ "RandomSample", "Temperature" → .8}]] &,
  "The best thing about AI is its ability to", 7]
Out[1]= {The best thing about AI is its ability to,
  The best thing about AI is its ability to provide,
  The best thing about AI is its ability to provide a,
  The best thing about AI is its ability to provide a certain,
  The best thing about AI is its ability to provide a certain level,
  The best thing about AI is its ability to provide a certain level of,
  The best thing about AI is its ability to provide a certain level of privacy,
  The best thing about AI is its ability to provide a certain level of privacy when}
```

... Mira como queda

```
In[2]:= Table[Nest[StringJoin[#, model[#,{ "RandomSample", "Temperature" → .8}]] &,
  "The best thing about AI is its ability to", 7], 5]
Out[2]= {The best thing about AI is its ability to be able to estate    ly more,
  The best thing about AI is its ability to understand, to tell you which functions,
  The best thing about AI is its ability to apply usually significant cognitive
  procedures to complex,
  The best thing about AI is its ability to learn. That is the only thing,
  The best thing about AI is its ability to solve problems in the natural world."}
```

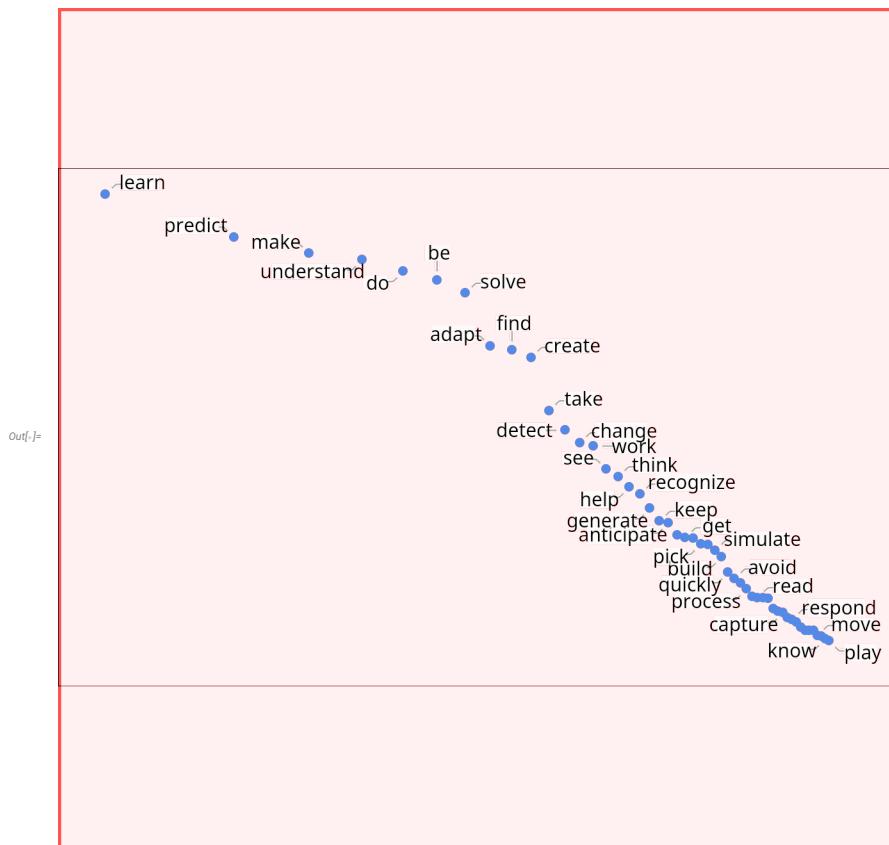
■  $n^{-1}$  “Power law”

...

• • •

## ListLogLogPlot

```
Take[ReverseSort[model[{"The best thing  
about AI is its ability to",  
"Probabilities"}]], 50],  
ImageSize \[Rule] {700, 700},  
FrameStyle \[Rule]  
Directive[FontFamily \[Rule] "Times"],  
LabelStyle \[Rule] Directive[Black, 16],  
Frame \[Rule] True, PlotStyle \[Rule] Thick]
```



Un poco mejor que el caso de temperatura ZERO (GPT-2):

```
In[1]:= StringReplace[Nest[StringJoin[#, model[#, {"RandomSample", "Temperature" → .8}]]] &,
  "The best thing about AI is its ability to", 100], "\n" .. → " "]
```

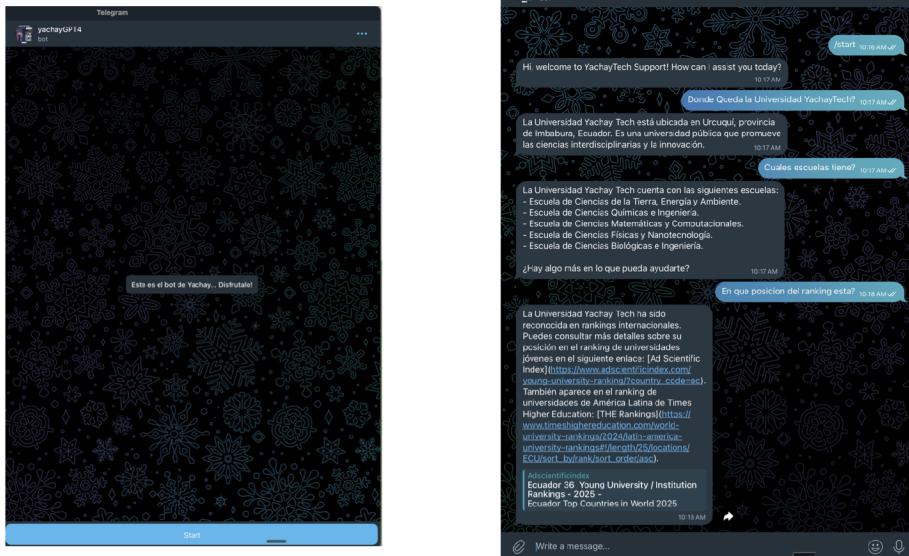
Out[1]= The best thing about AI is its ability to find its way around and work so much harder to understand the realm. The ability to do so is also what allows it to know the ten most common problems in your business (or, better yet, beyond it), and it has many other skills of its own staff. So you'll need to be extremely worrying about your AI to actually learn to interact with it. And in martyrdom, it does tend to be a bit cruel: you'll need an AI that doesn't want to

■ Con GPT-4 (@Telegram. Pedirle a la audiencia que pueda, usar el bot)



Demo: <https://t.me/yachayGPT4bot>

## Demo ChatGPT





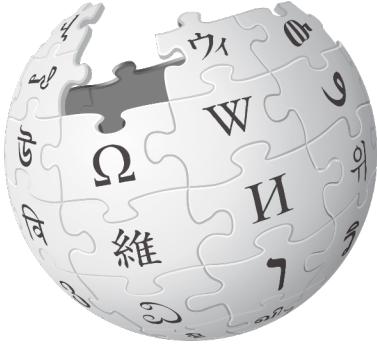
## ¿De donde provienen esas probabilidades?

---

Un problema mucho mas sencillo!  
solo contemos las letras de un artículo en Wikipedia...

Articulos con "cats"

# Contemos en Wikipedia



```
In[1]:= LetterCounts[WikipediaData["cats"]]  
Out[1]= <| e → 4345, a → 3485, t → 3411, i → 2818, s → 2632, n → 2530, o → 2476, r → 2180,  
h → 1607, l → 1598, c → 1433, d → 1346, m → 1023, u → 932, g → 768, f → 743,  
p → 664, y → 618, b → 515, w → 475, v → 403, k → 196, T → 121, x → 94, C → 82,  
A → 76, I → 68, S → 59, F → 43, z → 42, H → 35, E → 33, D → 30, M → 29, B → 29,  
N → 27, q → 23, P → 21, G → 20, j → 19, L → 17, o → 17, U → 16, R → 15, W → 13,  
K → 11, J → 7, Z → 5, V → 4, Q → 2, í → 2, ä → 2, á → 2, Y → 1, è → 1, ž → 1,  
d' → 1, ö → 1, ð → 1, ī → 1, ī → 1 |>
```

Con dogs?

```
In[1]:= LetterCounts[WikipediaData["dogs"]]

Out[1]= <| e → 3908, a → 2736, o → 2606, i → 2556, t → 2521, s → 2399, n → 2336, r → 1864,
d → 1581, h → 1460, l → 1355, c → 1082, g → 928, m → 858, u → 782, f → 662,
p → 636, y → 499, b → 462, w → 409, v → 405, k → 151, T → 90, C → 85, I → 80,
A → 73, x → 71, S → 65, D → 58, B → 47, P → 40, N → 38, z → 36, W → 27, H → 27,
q → 25, M → 24, L → 24, E → 23, F → 21, O → 21, U → 20, R → 18, G → 17, j → 17,
K → 10, J → 7, V → 7, Y → 4, é → 2, X → 1, á → 1, ó → 1, ê → 1 |>
```

En el Ingles!

```
In[2]:= English LANGUAGE [ character frequencies ]

Out[2]= {e → 12.7%, t → 9.06%, a → 8.17%, o → 7.51%, i → 6.97%, n → 6.75%,
s → 6.33%, h → 6.09%, r → 5.99%, d → 4.25%, l → 4.03%, c → 2.78%,
u → 2.76%, m → 2.41%, w → 2.36%, f → 2.23%, g → 2.02%, y → 1.97%,
p → 1.93%, b → 1.49%, v → 0.978%, k → 0.772%, j → 0.153%, x → 0.150%,
q → 0.0950%, z → 0.0740%}
```

Un ejemplo, si generamos una secuencia con estas probabilidades:

```
In[1]:= with[{freq =
Entity["Language", "English::385w8"][
EntityProperty[
"Language", "LetterFrequency"]
]}, SeedRandom[53424]; StringJoin[
RandomChoice[UnitConvert[Values[freq]] →
Keys[freq], 150]]]
```

Out[1]= rronoitadatcaeaesaotdoysaroioyiinnbantoioestlhddeocneooewceseciselnodrtrdgriscsatsepes :  
dcniouhoetsedehedslernevstotheindtbmnaohngotannbthrdthtonsiieldn

## ■ Podemos hacerlo con palabras

veamos

```
In[1]:= With[{freq = Entity["Language", "English::385w8"][
EntityProperty["Language", "LetterFrequency"]
]}, SeedRandom[34232]; StringJoin[
RandomChoice[Append[UnitConvert[Values[freq]
], .18] → Append[Keys[freq], " "], 150]]]
```

Out[1]= sd n oeiain satnwhoo eer rtr ofianordrenapwokom del oaas ill e h f  
rellptohltvoettseodtrncilntehtotrkthrslo hdaol n sriaefr hthehtn ld gpod a h y oi

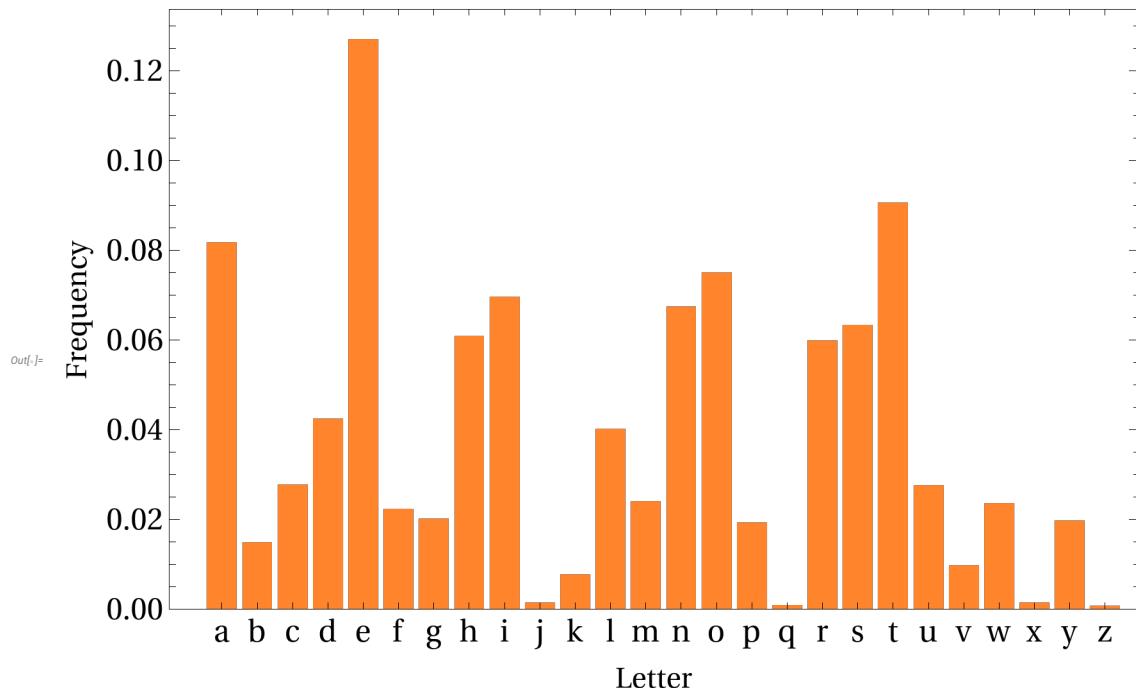
Forzarlas a parecerse al Ingles

```
In[1]:= SeedRandom[23424]; StringJoin[Riffle[Table[
RandomLettersWord[0], {25}], " "]]
```

Out[1]= ni hilwhuei kjtn isjd erogofnr n rwhwfao rcuw lis fahte uss cpnc nlu oe nusaetat  
llfo oeme rrhrtn xdses ohm oa tne ebedcon oarvthv ist

Las Probabilidades

```
In[1]:= BarChart[Normal[Values[KeySort@Entity["Language", "English::385w8"]][
  EntityProperty["Language", "LetterFrequency"]]],
  ChartLabels → Alphabet[], Frame → True,
  FrameLabel → {{"Frequency", FontSize → 24}, None},
  {{"Letter", FontSize → 24}, None},
  FrameTicksStyle → Directive[FontSize → 26],
  ChartStyle → Lookup[ChatTechColors["Data"], "Orange"],
  ChartBaseStyle → EdgeForm[], ImageSize → Large]
```

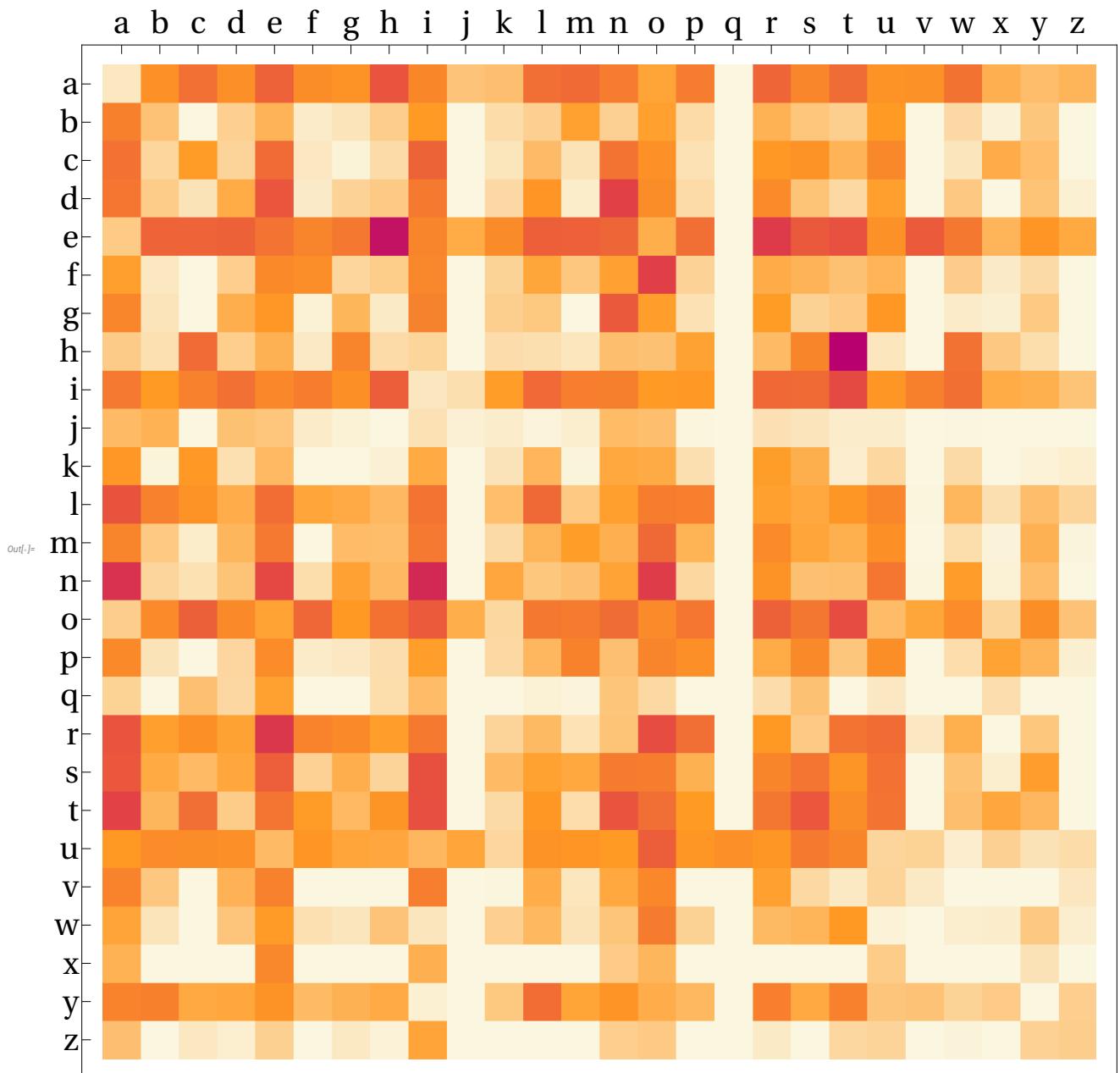


## ■ La probabilidad de Pares

...

2 —grams (nombre de la teoría)

```
In[1]:= MatrixPlot[Transpose@With[{data = 
  Outer[Lookup[Lookup[data, Key@{#1}], Key@{#2}, 0] &, Alphabet[], Alphabet[]]}, 
  FrameTicks → {{#, None}, {None, #}} &[Table[{i, FromLetterNumber[i]}, {i, 26}]], 
  FrameTicksStyle → Directive[FontSize → 30], 
  ColorFunction → ChatTechColors["Orange"], ImageSize → Full]]
```



Veamos que la q es ZERO excepto para la "u" !!! .. interesante.

Usando 2-grams

```
In[1]:= SeedRandom[23424]; StringJoin[Riffle[Table[
  RandomLettersWord[1],
  {25}], " "]]
```

```
Out[1]= olur weatmai froginda wha rcothe vo os olev cit ato peme trito refla arn edon
          parutast cuchrhedar ts we tas pof ten sath stik tea
```

Y con n-grams:

```
In[1]:= SeedRandom[234098234];
Grid[Table[
{style[Text@n, Lighter@Gray],
StringJoin[Riffle[
Table[RandomLettersWord[n],
{12}], " "]]}, {n, 0, 5}],
FrameStyle → LightGray,
Frame → All, Alignment → Left,
Spacings → {1, 0.5}]]
```

0	on gxeeetowmt tsifhy ah aufnsoc ior oia itlt bnc tu ih uls
1	ri io os ot timumumoi gymyestit ate bshe abol viowr wotybeat mecho
2	wore hi usinallistin hia ale warou pothe of premetra bect upo pr
3	qual musin was witherins wil por vie surgedygua was suchinguary outheydays       theresist
4	stud made yello adenced through theirs from cent intous wherefo proteined screa
5	special average vocab consumer market prepara injury trade consa usually speci       utility

Con millones de palabras

```
In[1]:= SeedRandom[2134]; StringJoin[Riffle[Table[
  WeightedRandomWord[], {30}], " "]]
```

```
Out[1]= from not emergency tree of is its are and their the rose and default cold statuesque
had worth director orientation have possibly one the the used be said that little
```

No tiene sentido!!!

## Con n-grams, empezando por la palabra cats:

...

```
In[1]:= Column[{"cat through shipping variety is made the aid emergency can the",
  "cat for the book flip was generally decided to design of",
  "cat at safety to contain the vicinity coupled between electric public",
  "cat throughout in a confirmation procedure and two were difficult music",
  "cat on the theory an already from a representation before a"},  

  Spacings → 1, Dividers → {None, Thread[Range[2, 5] → LightGray]}]
```

cat through shipping variety is made the aid emergency can the

cat for the book flip was generally decided to design of

Out[1]= cat at safety to contain the vicinity coupled between electric public

cat throughout in a confirmation procedure and two were difficult music

cat on the theory an already from a representation before a

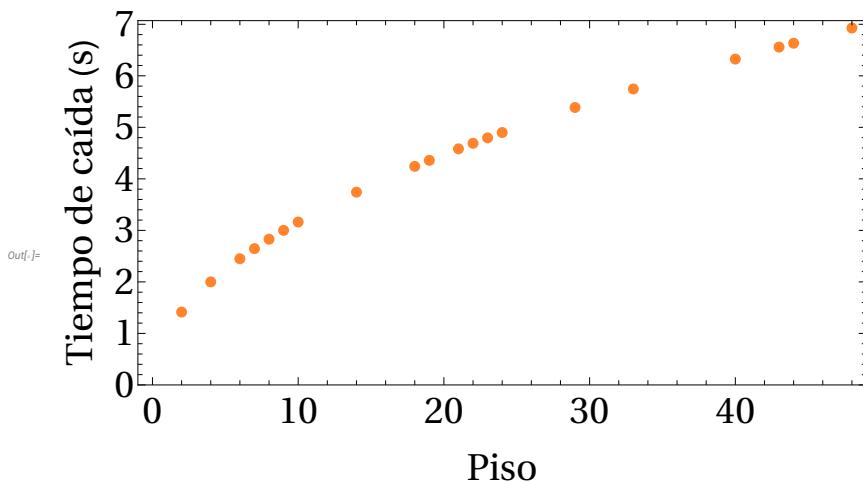


# ¿Qué es un modelo?

■ Como Galileo, cuánto dura la caída de una bola desde la torre de Pizza?

Tenemos estos datos

```
In[1]:= SeedRandom[34535]; ListPlot[{#, Sqrt[#]} & /@ Sort[
  RandomSample[Range[50], 20]], Frame -> True, PlotRange -> {0, Sqrt[50]},
  FrameLabel -> {{Style["Tiempo de caída (s)", FontSize -> 30], None},
    {Style["Piso", FontSize -> 30], None}},
  FrameTicksStyle -> Directive[FontSize -> 30],
  PlotStyle -> {ChatTechColors["Data"]["Orange"],
    PointSize[.015]}, AspectRatio -> 1/2]
```

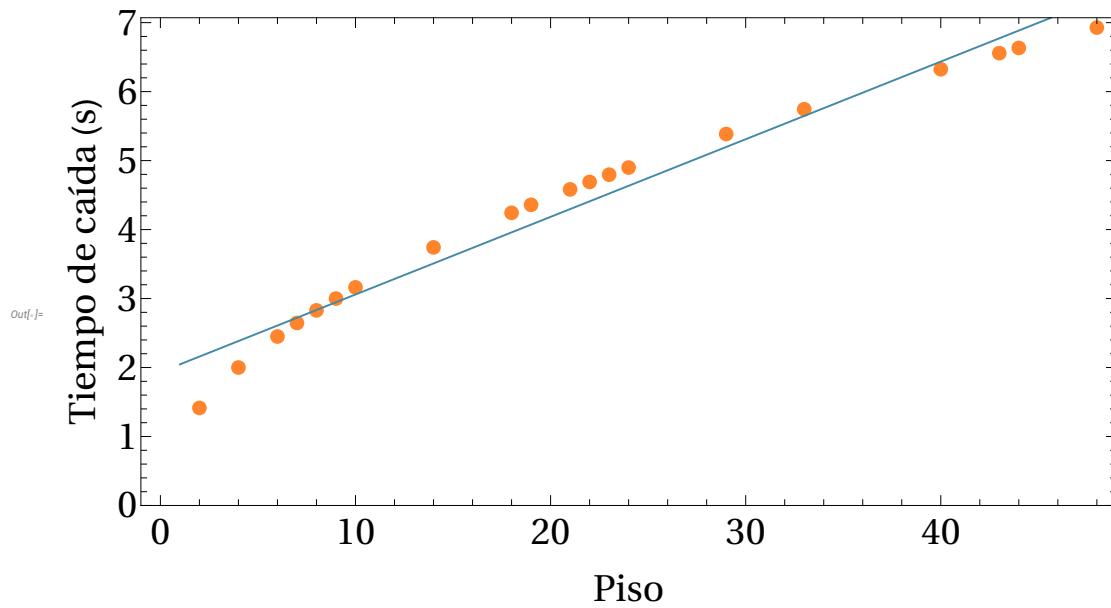


Simplifiquemos:

```

With[{data = (SeedRandom[34 535];
{#, Sqrt[#]} & /@ Sort[RandomSample[Range[50], 20]])},
Show[ListPlot[data, Frame → True,
FrameLabel → {{Style["Tiempo de caída (s)", FontSize → 30], None},
{Style["Piso", FontSize → 30], None}},
FrameTicksStyle → Directive[FontSize → 30],
PlotRange → {0, Sqrt[50]}, PlotStyle → {ChatTechColors["Data"]["Orange"],
PointSize[.015]}, AspectRatio → 1/2],
Plot[Evaluate[Fit[data, {1, x}, x]], {x, 1, 50},
PlotRange → {0, Sqrt[50]},
PlotStyle → ChatTechColors["Data"]["Teal"]]]]

```

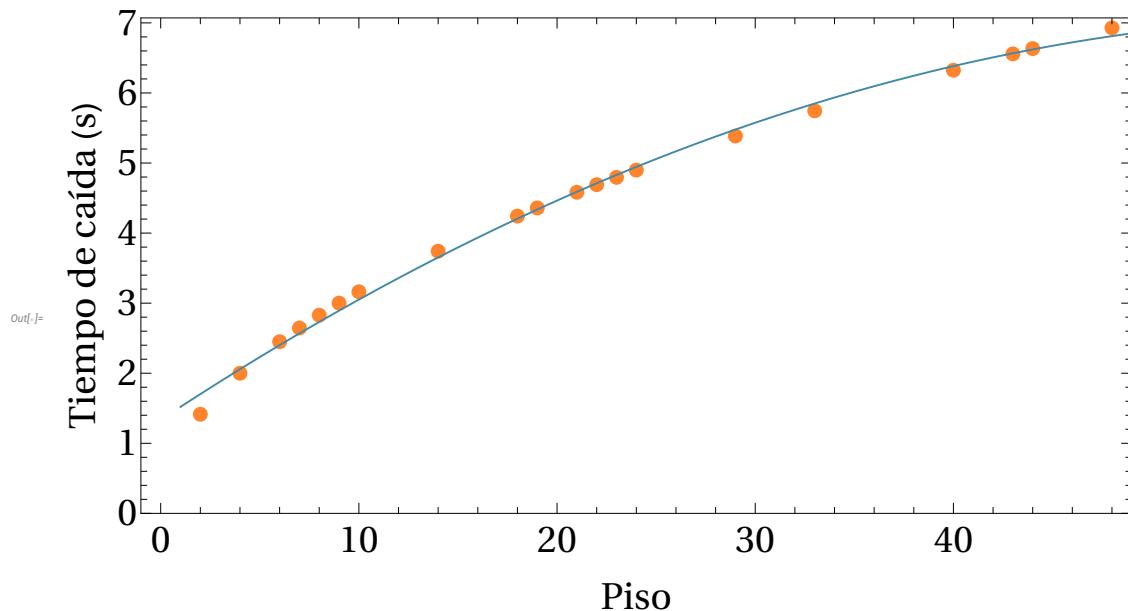


Ajustamos!!

```

With[{data = (SeedRandom[34 535];
  {#, Sqrt[#]} & /@ Sort[RandomSample[Range[50], 20]]),
 Show[ListPlot[data, Frame → True,
  FrameLabel → {{Style["Tiempo de caída (s)", FontSize → 30], None},
   {Style["Piso", FontSize → 30], None}},
  FrameTicksStyle → Directive[FontSize → 30],
  PlotRange → {0, Sqrt[50]}, PlotStyle → {ChatTechColors["Data"]["Orange"],
   PointSize[.015]}, AspectRatio → 1/2],
 Plot[Evaluate[Fit[data, {1, x, x^2}, x]], {x, 1, 50}],
 PlotRange → {0, Sqrt[50]},
 PlotStyle → ChatTechColors["Data"]["Teal"]]]
]

```

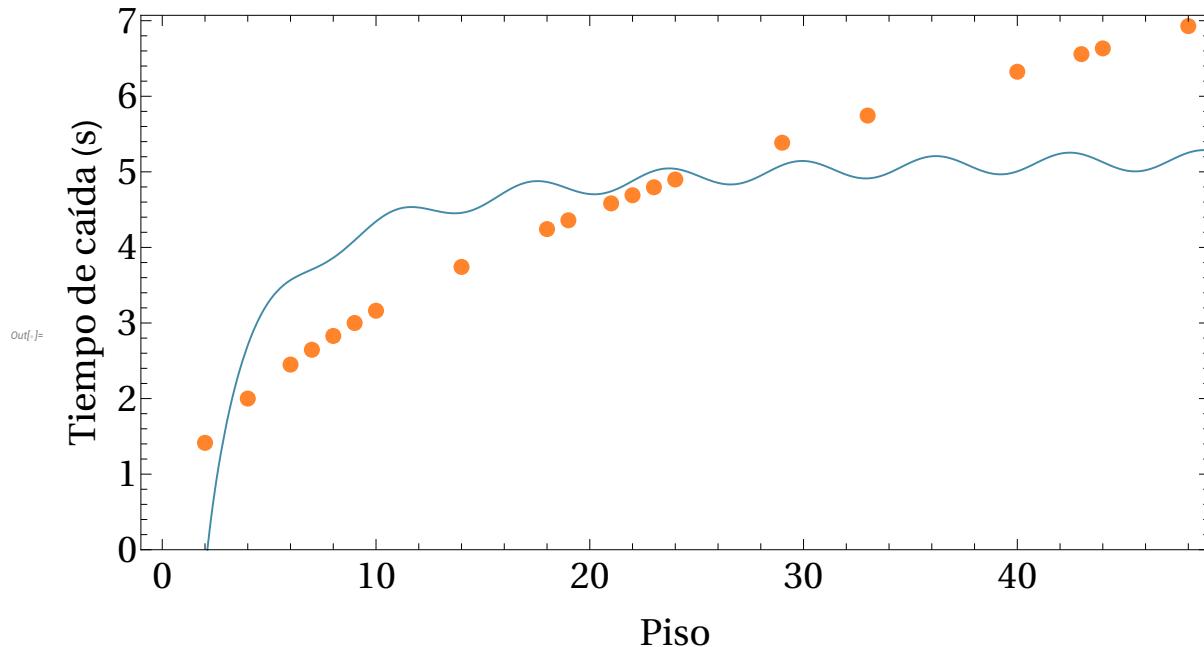


Un mal modelo

```

With[{data = (SeedRandom[34535];
  {#, Sqrt[#]} & /@ Sort[RandomSample[Range[50], 20]]}),
 Show[ListPlot[data, Frame → True,
  FrameLabel → {{Style["Tiempo de caída (s)", FontSize → 30], None},
   {Style["Piso", FontSize → 30], None}},
  FrameTicksStyle → Directive[FontSize → 30],
  PlotRange → {0, Sqrt[50]}, PlotStyle → {ChatTechColors["Data"]["Orange"],
   PointSize[.015]}, AspectRatio → 1/2],
 Plot[Evaluate[(a + b/x + c Sin[x]) /. FindFit[data, a + b/x + c Sin[x], {a, b, c}, x]],
 {x, 1, 50}, PlotRange → {0, Sqrt[50]},
 PlotStyle → ChatTechColors["Data"]["Teal"]]
]]

```



ChatGPT lo hace con 175 —  
millardos de Parámetros



## Modelos para tareas muy HUMANAS !

Modelo de como funciona el Cerebro?

Comencemos por una tarea sencilla :  
"Reconocimiento de imágenes"

In[1]:=

```
First /@ RandomSample[ResourceData["MNIST"], 10]
```

Out[1]= {6, 2, 4, 2, 5, 1, 7, 1,  
3, 8}

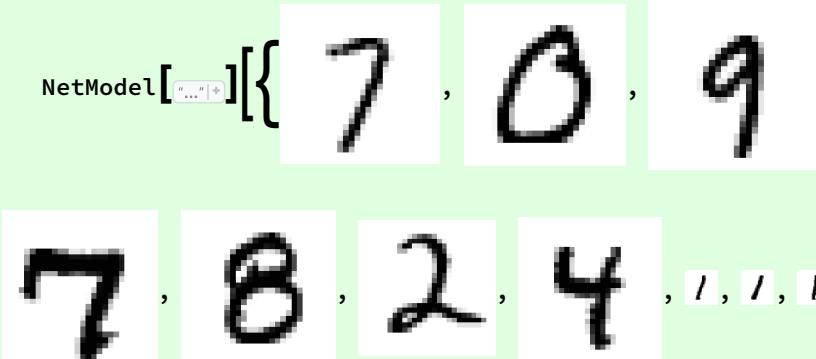
## Para cada imagen:

Como humanos podemos reconocer esto!

**¿Podemos construir un modelo que  
sepa a que numero representa  
una imagen?**

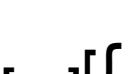
Si !

Tratemos a la función como una caja negra:

```
In[1]:= 
Style[
  NetModel["..."]@{7, 0, 9, 7, 8, 2, 4, 1, 1, 1},
  , 70]
]
```

$\{7, 0, 9, 7, 8, 2, 4, 1, 1, 1\}$

**Que ocurre?:** Vamos a progresivamente BLUR las imagenes

```
In[1]:= Style[  
  NetModel[""], {     , 70]  
{2, 2, 2, 1, 1, 1, 1, 1, 1}
```



## -- Redes Neuronales --

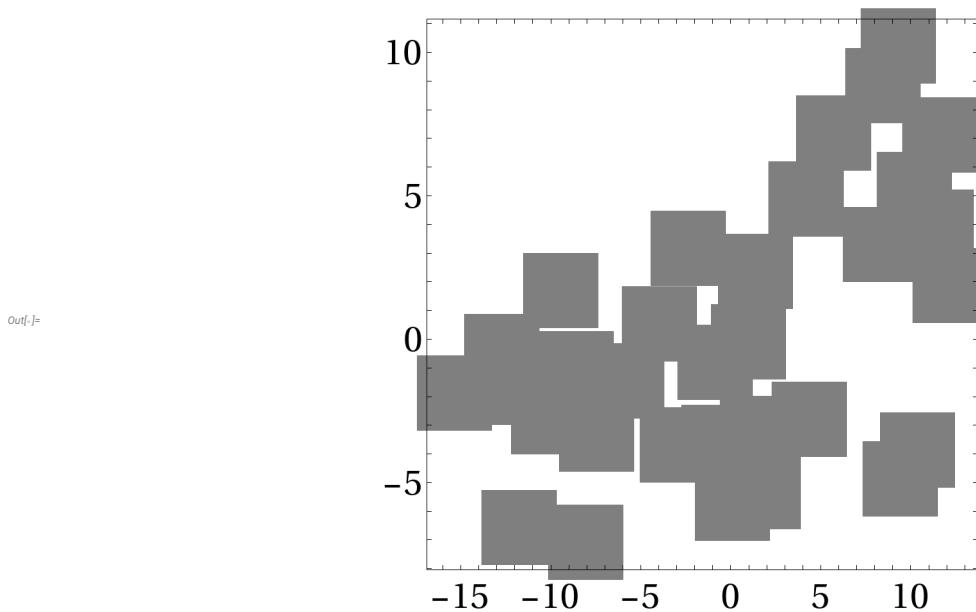
1940 - Como funciona el cerebro!

En el cerebro, despues de que los fotones chocan el ojo:



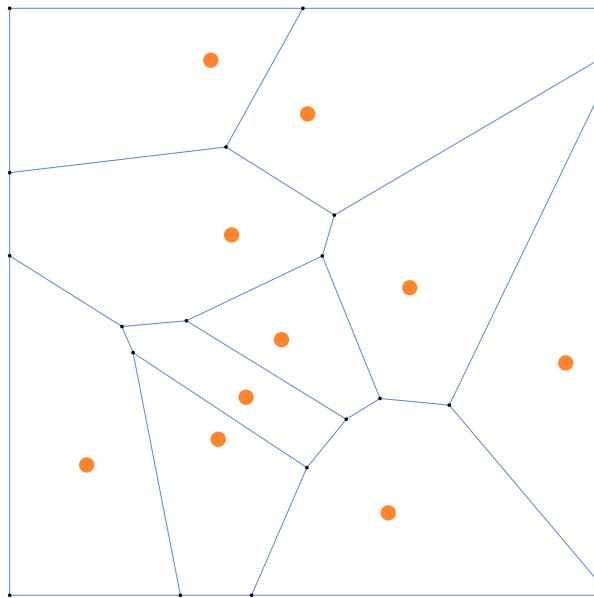
No hay teoría. Solo funcionó en 1998!! -- Noción del atractor: queremos que los 1 estén juntos..

```
In[1]:= FeatureSpacePlot[Catenate[Cases[ResourceData["MNIST"],  
      (x_ → ##) :> (x → SetAlphaChannel[x, ColorNegate[x]]), {1}, 16] & /@ {1, 2}],  
      RandomSeeding → 5, Frame → True, FrameTicksStyle → Directive[FontSize → 30]]
```



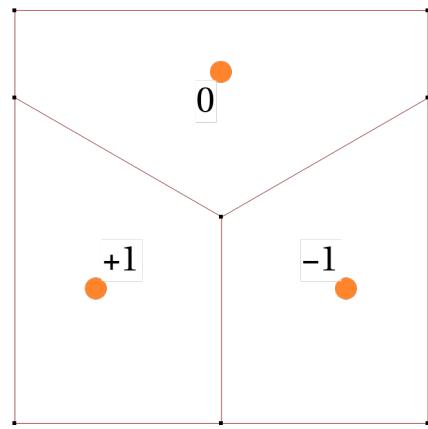
Imagine que ...

```
In[1]:= SeedRandom[3 423 424];
With[{pts = RandomReal[{{-.9, .9}, {10, 2}}]},
  Show[VoronoiMesh[pts, {{-1, 1}, {-1, 1}}},
    PlotTheme → "Lines", BaseStyle → GrayLevel[.7]],
  Graphics[Style[Point[pts], PointSize[.025],
    ChatTechColors["Data"][[Orange]]]
  ]]]]
```



Como hacemos de una red neuronal realice la tarea  
de “Reconocer imágenes” ?

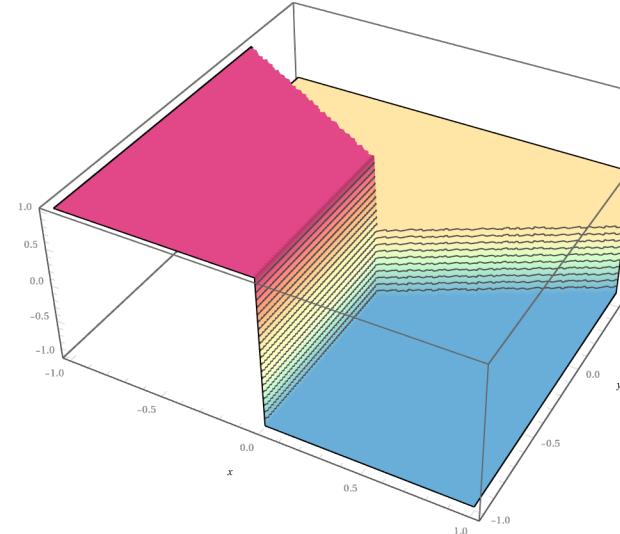
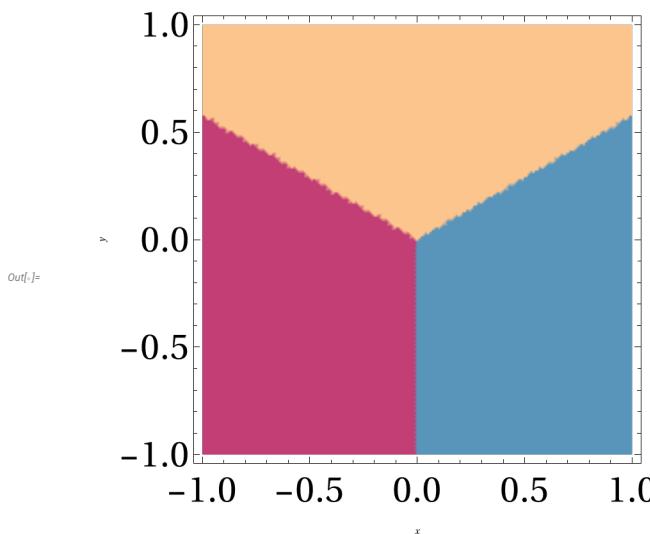
```
In[1]:= SeedRandom[3 423 424];
With[{pts = .7 CirclePoints[3]},
Show[VoronoiMesh[pts, {{-1, 1}, {-1, 1}}], PlotTheme -> "Lines",
BaseStyle -> GrayLevel[.7]], ListPlot[MapIndexed[Labeled[\!\(\#2\), Style[
Switch[First[\!\(\#2\)], 1, Style["-1", 30], 2, Style["0", 30], 3, Style["+1", 30]
], Black]] &, pts], PlotStyle -> Directive[PointSize[.05],
ChatTechColors["Data"]["Orange"]]]]]]
```



■ Dado {x,y} en donde cae?

Una función así:

```
In[1]:= With[{nf = Nearest[Thread[7 CirclePoints[3] → {-1, 0, 1}]]},
GraphicsRow[{DensityPlot[nf[{x, y}], {x, -1, 1}, {y, -1, 1},
ColorFunction → Function[ChatTechColors["Normal"][[2 # - 1]]],
PlotPoints → 50, FrameTicksStyle → Directive[FontSize → 30], FrameLabel → (Style[#, Italic] & /@ {"x", "y"}), ImageSize → 310],
Plot3D[nf[{x, y}], {x, -1, 1}, {y, -1, 1}, PlotPoints → 50, MeshFunctions → (#3 &), Mesh → 20,
TicksStyle → Gray,
ColorFunction → Function[ChatTechColors["Normal"][[2 # 3 - 1]]], SphericalRegion → False, AxesLabel → (Style[#, Italic] & /@ {"x", "y"}),
BarLegend[{Function[ChatTechColors["Normal"][[#]], {-1, 1}], LegendMarkerSize → 150}],
}, Spacings → {25, 5}]]
```



Como?

```
In[1]:= Clear[FormulaNeuralNetworkGraph]
FormulaNeuralNetworkGraph[layerCounts : {_Integer, _Integer, _Integer}] :=
Block[{gr1, gr2, gr21, gr3, gr4, gr, bc},
gr1 = IndexGraph[CompleteGraph[Take[layerCounts, 2]]];
gr2 =
Graph[Map[(layerCounts[[1]] + #) → (layerCounts[[1]] + layerCounts[[2]] + #) &,
Range[layerCounts[[2]]]]];
gr3 = IndexGraph[CompleteGraph[Take[layerCounts, -2]],
layerCounts[[1]] + layerCounts[[2]] + 1];
bc = layerCounts[[1]] + 2 * layerCounts[[2]];
gr4 = Graph[Map[(bc + #) → (bc + layerCounts[[3]] + #) &,
Range[layerCounts[[3]]]], VertexLabels → "Name"]];
```

```

gr = GraphUnion[gr1, gr2, gr3, gr4];

Graph[gr,
  GraphLayout → {"MultipartiteEmbedding", "VertexPartition" → {layerCounts[[1]],
    layerCounts[[2]], layerCounts[[2]], layerCounts[[3]], layerCounts[[3]]}}];

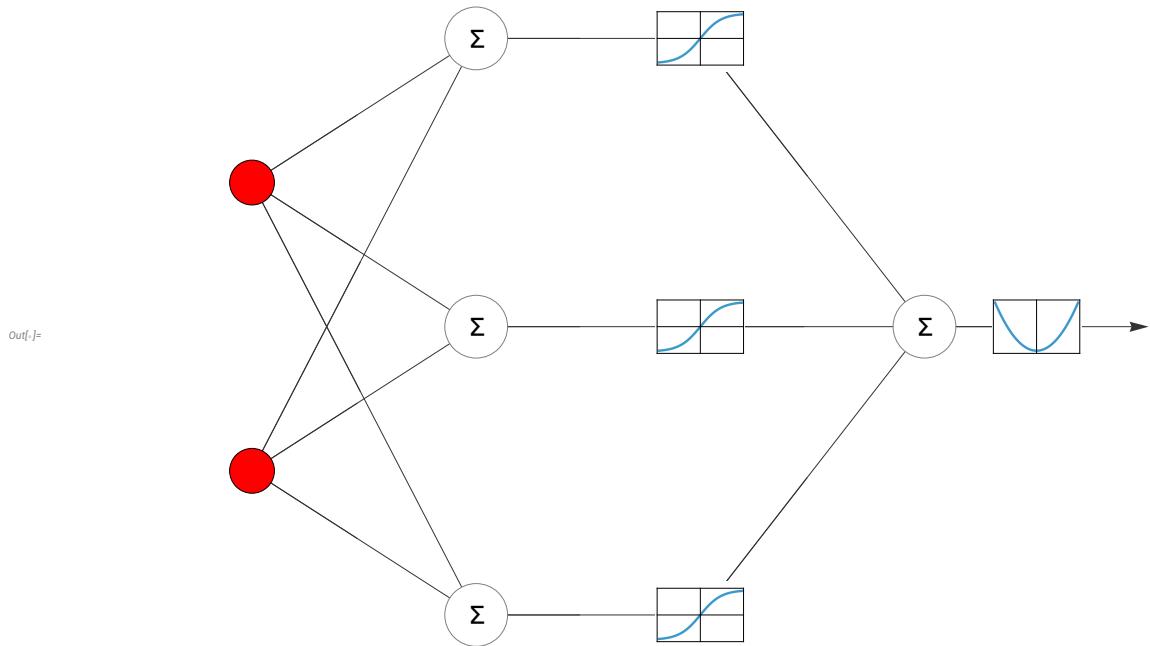
Clear[FormulaNeuralNetworkGraphPlot]
Options[FormulaNeuralNetworkGraphPlot] = Options[Graphics];

FormulaNeuralNetworkGraphPlot[
  layerCounts : {__Integer, __Integer, __Integer}, func1_, opts : OptionsPattern[]] :=
  FormulaNeuralNetworkGraphPlot[layerCounts, func1, # &, opts];

FormulaNeuralNetworkGraphPlot[layerCounts : {__Integer, __Integer, __Integer},
  func1_, func2_, opts : OptionsPattern[]] :=
  Block[{pl0pts, grFunc1, grFunc2, gr, vNames, vCoords, vNameToCoordsRules, edgeLines},
    pl0pts = {PlotTheme → "Default", Axes → True, Ticks → False,
      Frame → True, FrameTicks → False, ImageSize → Small};
    grFunc1 = Plot[func1[x], {x, -2, 2}, Evaluate[pl0pts]];
    grFunc2 = Plot[func2[x], {x, -2, 2}, Evaluate[pl0pts]];
    gr = FormulaNeuralNetworkGraph[layerCounts];
    vNames = VertexList[gr];
    vCoords = VertexCoordinates /. AbsoluteOptions[gr, VertexCoordinates];
    vNameToCoordsRules = Thread[vNames → vCoords];
    edgeLines = Arrow @ ReplaceAll[List @@@ EdgeList[gr], vNameToCoordsRules];
    Graphics[{Arrowheads[0.02], GrayLevel[0.2], edgeLines, EdgeForm[Black],
      FaceForm[Red], Map[Disk[#, 0.1] &, vCoords[[1 ;; -layerCounts[[1]] - 1]]];
      Black, Map[{EdgeForm[Gray], FaceForm[White],
        Disk[#, 0.14], Text[Style["\u2211", 16, Bold], #]} &,
        Join[vCoords[[layerCounts[[1]] + 1 ;; layerCounts[[1]] + layerCounts[[2]]]],
        vCoords[[-2 layerCounts[[1]] ;; -layerCounts[[1]] - 1]]], Map[
        {EdgeForm[None], FaceForm[White], Rectangle[#+{0.2, 0.15}, #+{0.2, 0.15}],
        Inset[grFunc1, #1, Center, 0.4]} &, vCoords[[Total[layerCounts[[1 ;; 2]]] + 1
        ;; Total[layerCounts[[1 ;; 2]] + layerCounts[[2]]]]], Map[{EdgeForm[None],
        FaceForm[White], Rectangle[#+{0.2, 0.15}, #+{0.2, 0.15}]]}]}];
  
```

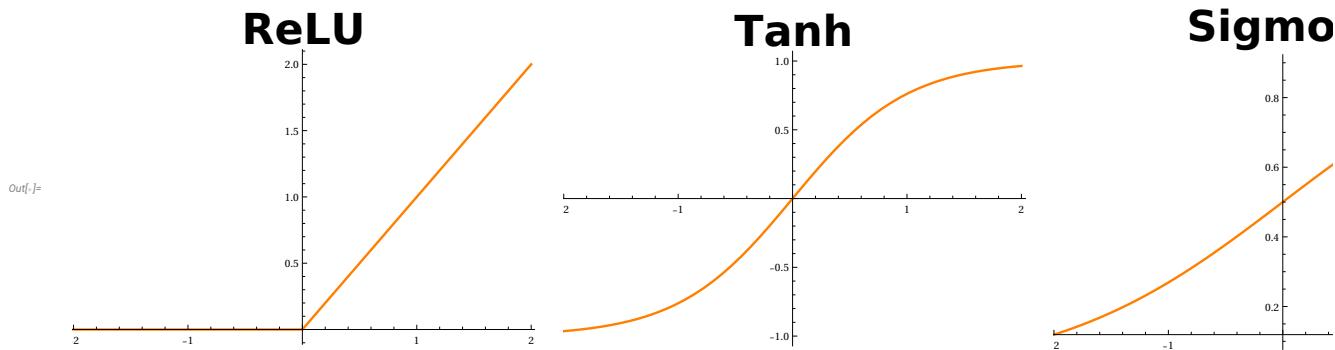
```
# + {0.2, 0.15}], Inset[grFunc2, #1, Center, 0.4]} &, MapThread[  
Mean@*List, {vCoords[[-2 layerCounts[-1] ;; -layerCounts[-1] - 1],  
vCoords[[-layerCounts[-1] ;; -1]]]}], opts]];
```

```
In[4]:= FormulaNeuralNetworkGraphPlot[  
{2, 3, 1},  
Tanh, #^2 &, ImageSize -> 800]
```



# Funciones de Activación

```
GraphicsRow[
 Labeled[With[{f = ElementwiseLayer[\#, "Input" → "Real"]}, Plot[f[x], {x, -2, 2},
 PlotStyle → Directive[Thick, Orange], (*Thicker lines*)
 FrameTicksStyle → Directive[FontSize → 30], ImageSize → 400]],
 Style[\#, FontSize → 36, Bold], Top] & /@ {"ReLU", Tanh, "Sigmoid"}, 
 ImageSize → Large]]
```



Una función matematica cualquiera:

```
In[1]:= NetToFunction[NetRankData2D["Result431"], {x, y}, {w, b}][{-1, 1}] /. Ramp → f // TraditionalForm
```

Out[1]//TraditionalForm=

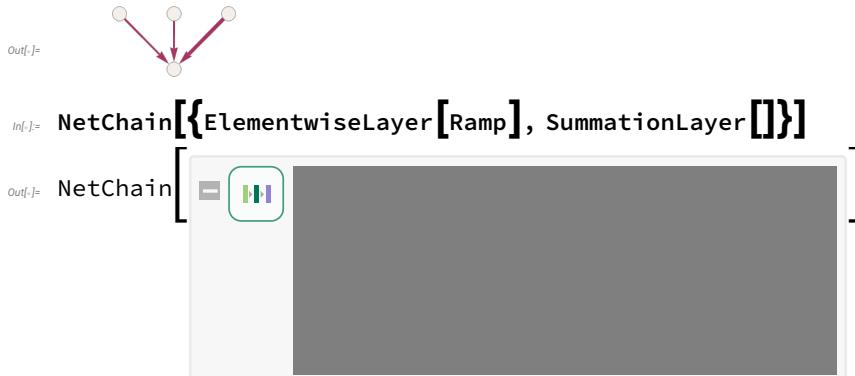
$$w_{511} f(w_{311} f(b_{11} + x w_{111} + y w_{112}) + w_{312} f(b_{12} + x w_{121} + y w_{122}) + w_{313} f(b_{13} + x w_{131} + y w_{132}) + w_{314} f(b_{14} + x w_{141} + y w_{142}) + b_{31}) + w_{512} f(w_{321} f(b_{11} + x w_{111} + y w_{112}) + w_{322} f(b_{12} + x w_{121} + y w_{122}) + w_{323} f(b_{13} + x w_{131} + y w_{132}) + w_{324} f(b_{14} + x w_{141} + y w_{142}) + b_{32}) + w_{513} f(w_{331} f(b_{11} + x w_{111} + y w_{112}) + w_{332} f(b_{12} + x w_{121} + y w_{122}) + w_{333} f(b_{13} + x w_{131} + y w_{132}) + w_{334} f(b_{14} + x w_{141} + y w_{142}) + b_{33}) + b_{51}$$

In[2]:= "LayersList"

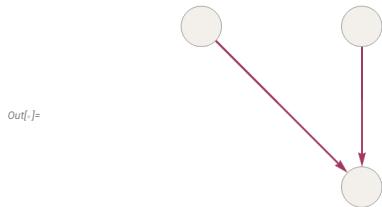
Out[2]= LayersList

The neural net of ChatGPT also just corresponds to a mathematical function like this—but effectively with billions of terms.

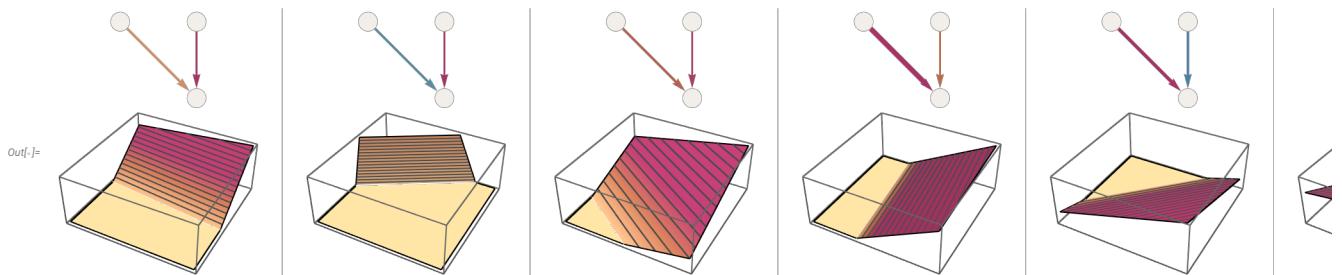
```
In[1]:= NetGraphPlot[LinearLayer["Weights" → {{1, 1, 3}}, "Biases" → {1}],
  PlotRange → {{-1.14, 6.8}, Automatic}]
```



```
In[3]:= NetGraphPlot[
  LinearLayer["Weights" → {{1, 1}}, "Biases" → {1}],
  PlotRange → {{-.9, 1.8}, Automatic}]
```

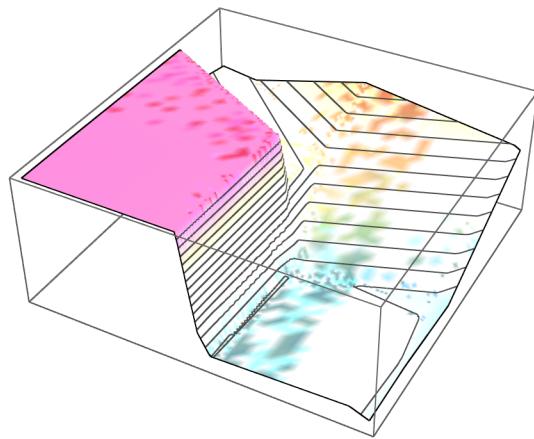


```
In[1]:= Module[{netelem3D},
  netelem3D[{a_, b_, c_}] :=
  {NetGraphPlot[LinearLayer["Weights" → {{a, b}}, "Biases" → {c}],
    PlotRange → {{-0.8, 1.8}, Automatic}],
   Plot3D[Ramp[{x, y}.{a, b} + c], {x, -1, 1}, {y, -1, 1}, PlotPoints → 40,
    ColorFunction → Function[ChatTechColors["Normal"][[#3]]], Ticks → None,
    MeshFunctions → {#3 &}, ColorFunctionScaling → False, ClippingStyle → None]};
  GraphicsGrid[
  netelem3D[#] & /@ {2{.1, .8, 0}, 3{-.3, .5, -.3}, 3{.2, .3, .3}, 5{.8, .1, .3},
    4{.5, -.5, .3}, 2{.2, -.7, .3}} // Transpose, Spacings → {80, 0},
  ImageSize → 500, Dividers → {{False, {True}, False}}, FrameStyle → LightGray]
  ]
```



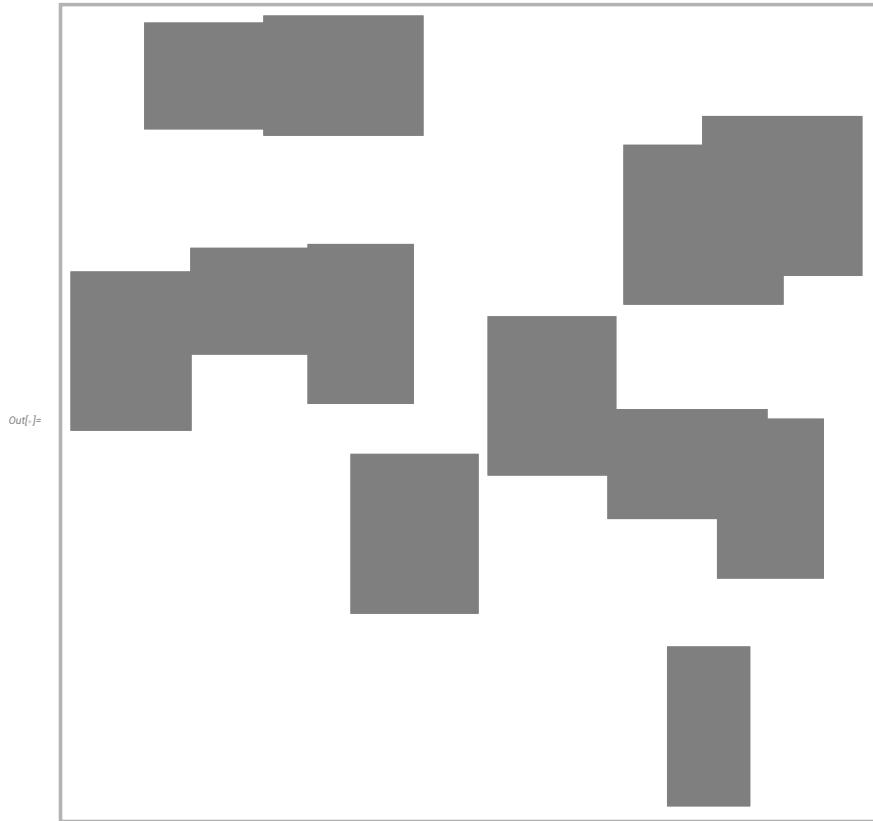
Que hay de redes mas grande? Bueno sale algo asi:

```
In[1]:= GraphicsRow[ { DensityAndContourPlot3D[  
    PointRankData2D["Result431"],  
    PlotPoints -> 40], ImageSize -> 500 } ]
```



Una tarea mas elaborada.

In[...]:= **Framed[FeatureSpacePlot[ $\{\dots\}$ , LabelingSize  $\rightarrow$  60], FrameStyle  $\rightarrow$  {GrayLevel[0, .3]}]**



Diferentes capas:

*In[1]:=*

```
GraphicsGrid[
 Partition[Image /@ Take[NetModel["Wolfram ImageIdentify Net V1"], 1][[1]], 10]]
```

*Out[1]=*

*In[2]:=*

```
GraphicsGrid[Partition[
 Take[Image /@ Take[NetModel["Wolfram ImageIdentify Net V1"], 10][[1]], 50], 10]]
```



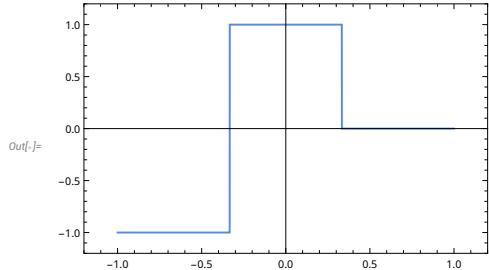
# Machine Learning y entrenamiento de redes neuronales:

## 1. Como se entrena una red neuronal?

Lo que queremos es encontrar WEIGHTS! (w)

Vamos a tratar de encontrar una red que reproduzca esta función:

```
In[1]:= Plot[Piecewise[{{{-1, x < -1/3}, {1, -1/3 < x < 1/3}, {0, x > 1/3}}, {x, -1, 1}], Exclusions -> None, Frame -> True, PlotRangePadding -> .2]
```



Necesitamos una red que luzca mas o menos asi:

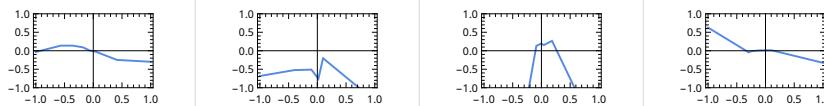
```
In[1]:= NetGraphPlot[PiecewiseNet[{4, 3, 1}],  
EdgeStyle → GrayLevel[.1, 1.],  
"AddArrows" → True,  
ImageSize → 380]
```

Out[1]=

## Cuales Weights usamos?:

```
GraphicsGrid[  
With[{net = NetInitialize[PiecewiseNet[{4, 3, 1}], RandomSeeding → ##,  
Method → {"Random", "Biases" → NormalDistribution[0, .2]}]},  
{{NetGraphPlot[net, ImageSize → {140}], Plot[net[x], {x, -1, 1}, Frame → True,  
PlotRange → {-1, 1}, ImageSize → {140}]} & /@ {4, 6, 8, 10}} // Transpose,  
ImageSize → 700, Alignment → Center, Dividers → {{False, {True}, False}},  
FrameStyle → LightGray, Spacings → {50, -20}]
```

Out[1]=



Vemos que ninguno lo reproduce bien!

Conseguimos mejores weights y:

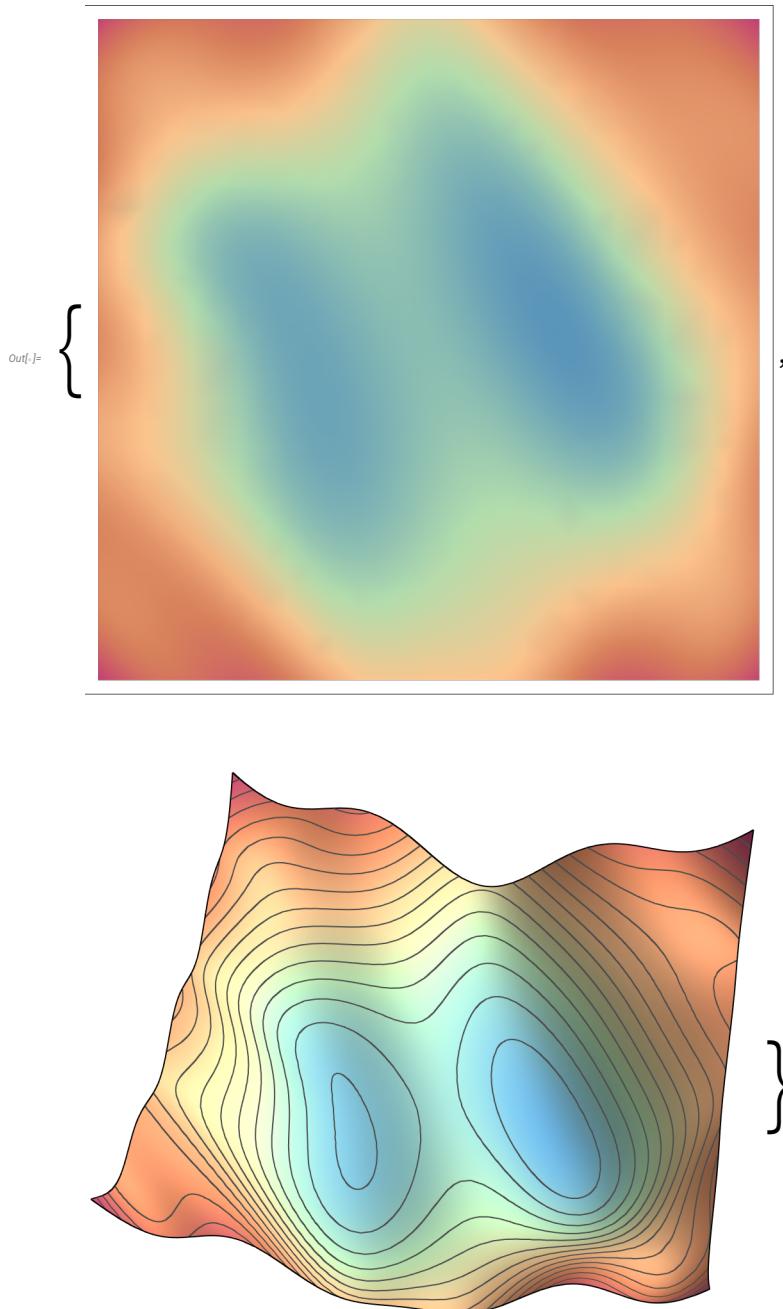
```
In[1]:= (** GraphicsGrid[ [ PiecewiseData1D ] ,
  ImageSize → 600, Dividers → { { False, { True }, False } } ,
  FrameStyle → LightGray, Spacings → { 45, Automatic } ] **)
```

Loss function (L2).

```
In[2]:= GraphicsGrid[ [ PiecewiseData1D , "ShowLoss" → True ], ImageSize → 700,
 Alignment → Center, Dividers → { { False, { True }, False } } , FrameStyle → LightGray, Spacings → { 15, -10 } ]
```

Como un ejemplo agarremos una función que solo sea función de dos variables w1 y w2. Entonces tenemos una función loss que a su vez es función de w1 y w2 y lucir así:

```
In[3]:= Module[ { bound = Pi , fun },
 fun = x ^ 2 + y ^ 2 + Sin[ (x - 1 / 2) (y - 1 / 2) ] + 2 Sin[ 2 (x + 1 / 2) + (y + 1 / 2) ];
 { DensityPlot[ fun, { x, -bound, bound }, { y, -bound, bound },
 ColorFunction → Function[
 [ ChatTechColors [ "Normal" ] [ 2 # - 1 ] ],
 FrameTicks → None, ImageSize → Large ],
 Plot3D[ fun, { x, -bound, bound }, { y, -bound, bound },
 ColorFunction → Function[
 [ ChatTechColors [ "Normal" ] [ 2 # 3 - 1 ] ],
 Ticks → None,
 MeshFunctions → { # 3 & },
 PlotPoints → 50,
 PerformanceGoal → "Quality",
 SphericalRegion → False,
 Boxed → False, Axes → False,
 ViewPoint → { 0.33, -2.01, 2.70 }, ImageSize → Large ] ]
 }
```



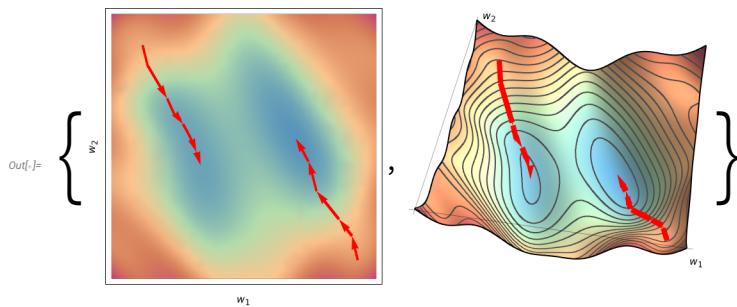
Técnicas para minimizar:

```
In[7]=
Module[{bound = Pi, fun, paths, IterateDescent, OverlapSplits},
  OverlapSplits[splits_] := Append[Map[
    Append[First[#], First[Last[#]]] &,
    Partition[splits, 2, 1]], Last[splits]];
  
```

```

IterateDescent[step_][pt_] := Plus[pt, Times[
  step, ReplaceAll[
    -D[fun, #] & /@ {x, y}, Thread[{x, y} → pt]]];
fun = x^2 + y^2 + Sin[(x - 1/2)(y - 1/2)] + 2 Sin[2(x + 1/2) + (y + 1/2)];
paths = TimeConstrained[FixedPointList[
  IterateDescent[0.12], N@#], 2, Nothing
] & /@ {.8{-3, 3}, .9{3, -3}};
paths = OverlapSplits[SplitBy[#, Floor[EuclideanDistance[#[[1]], #]] &]] & /@ paths;
Show[#, ImageSize → {Automatic, 250}] & /@
{Show[DensityPlot[fun, {x, -bound, bound}, {y, -bound, bound},
  FrameLabel → (Text /@ {"w1", "w2"}),
  ColorFunction → Function[
  ChatTechColors["Normal"][[2 # - 1]],
  FrameTicks → None],
  Graphics[{Thick, Red,
  Arrowheads[0.05],
  Map[Arrow[#] &, paths, {2}]}]
], Show[Plot3D[fun, {x, -bound, bound}, {y, -bound, bound},
  ColorFunction → Function[
  ChatTechColors["Normal"][[2 #3 - 1]],
  Ticks → None,
  MeshFunctions → {#3 &},
  PlotPoints → 50,
  Axes → {True, True, False},
  AxesOrigin → {-Pi, -Pi, 18},
  AxesLabel → (Text /@ {"w1", "w2"}),
  PerformanceGoal → "Quality",
  SphericalRegion → False,
  Boxed → False,
  ViewPoint → {0.33, -2.01, 2.70}],
  Graphics3D[{Thick, Red, Arrowheads[.02],
  Map[Arrow[Append[#, (fun + 2)/. Thread[{x, y} → (#)]]
  ] & /@ #] &, paths, {2}]}]
]}
]

```

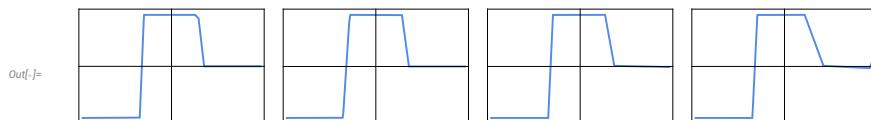


Podemos usar cálculo:

Derivadas para funciones con 175 millardos de párametros!

Diferentes redes dando los mismos resultados:

```
In[.] =
GraphicsRow[Plot[\#x], {x, -1, 1},
  Frame → True,
  FrameTicks → None,
  ImageSize → 100] & /@ PiecewiseData1D["AlternativeTrainings"],
ImageSize → 700]
```



Pero si pedimos valores extrapolados, podemos tener diferentes resultados:

```
In[.] =
GraphicsRow[
Plot[\#x], {x, -3, 3}, PlotRange → {-2, 2}, Frame → True, FrameTicks → None,
Prolog → {GrayLevel[.8, .2], Rectangle[{-3, -3}, {-1, 3}],
Rectangle[{1, -3}, {3, 3}]}, ImageSize → 200] & /@
PiecewiseData1D["AlternativeTrainings"], ImageSize → 750]
```





## El arte de entrenar una red neuronal -

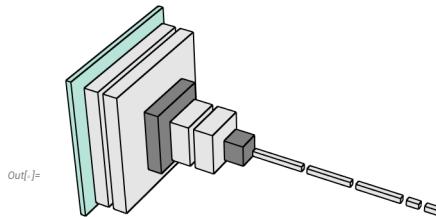
- 1) Arquitectura
  - 2) data
  - 3) Se pueden usar redes ya entrenadas y agregar cosas.
- 

Parece que es mejor end-to-end, dejando descubrir los pasos intermedios etc, Y componentes sencillos

---

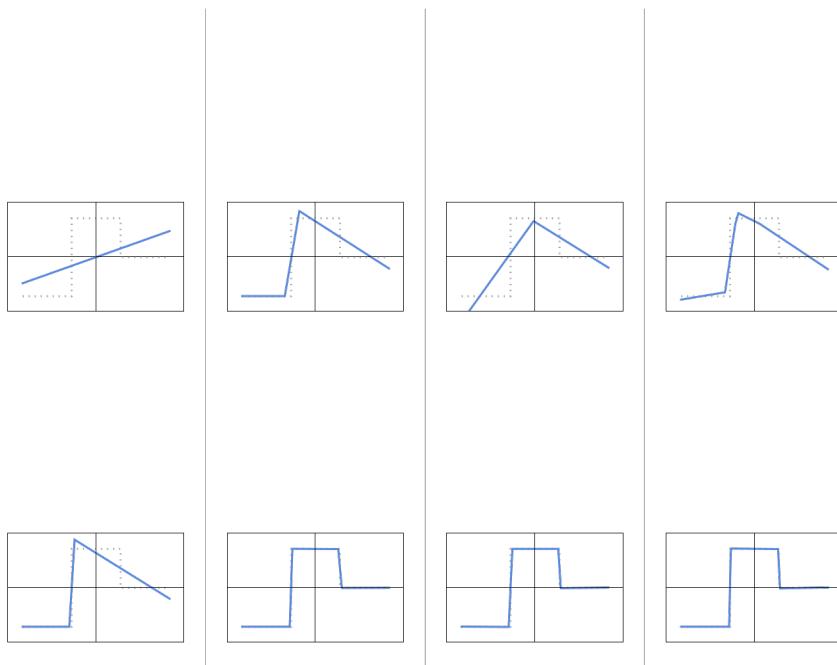
Estamos lidiando con arreglo de numeros. Y en el andar podemos rearreglar, cambiar de forma, etc, Por ejemplo en el reconocimiento de digitos empezamos por una imagen 2D y terminamos con output de 1D:

In[1]:= `[*] NetBlockPlot3D [NetModel["LeNet Trained on MNIST Data"]]`



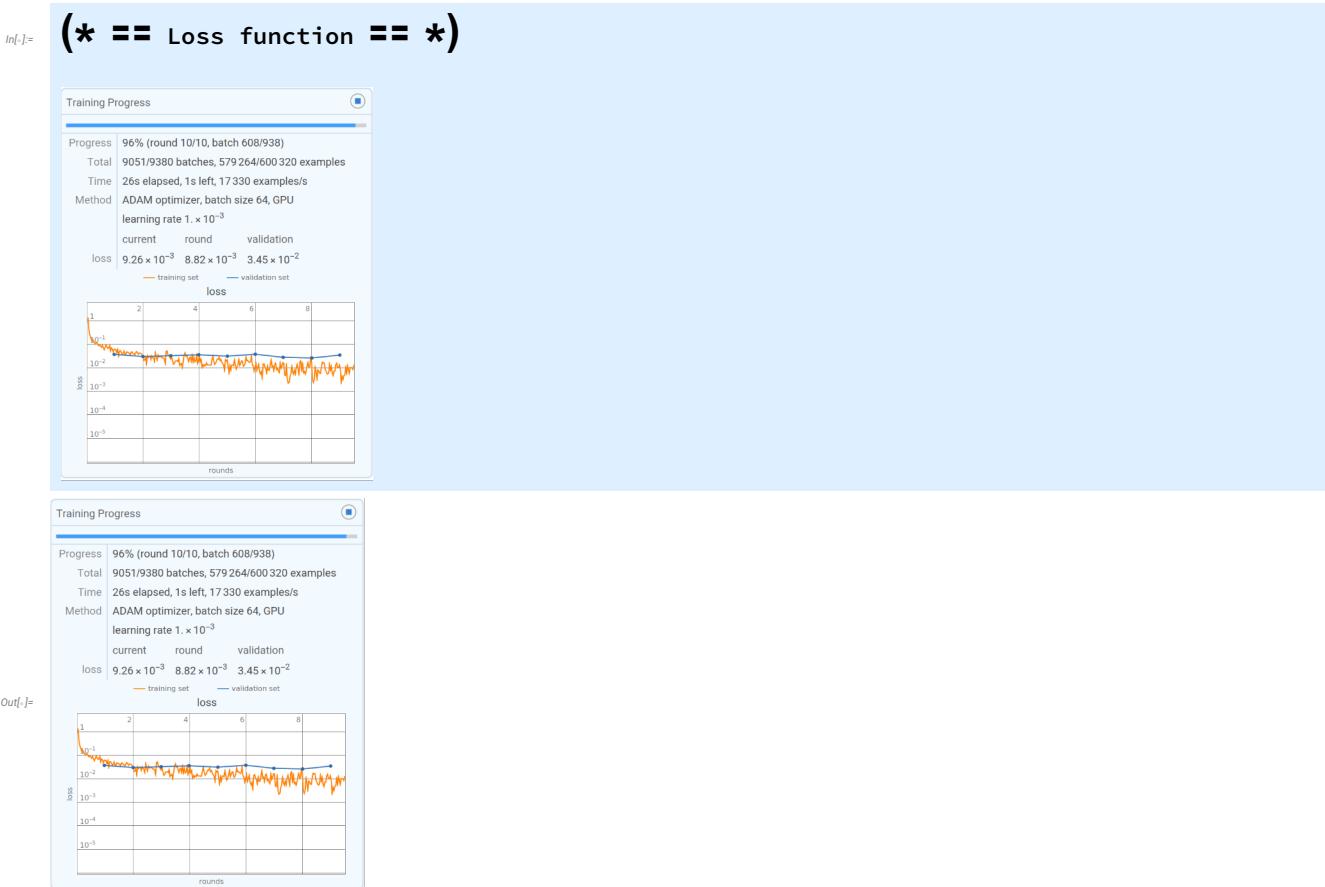
Tratar de reproducir la función escalon de arriba

```
In[1]:= GraphicsGrid[Flatten[Transpose /@ Partition[MapThread[
{[n] NetGraphPlot [#, #[2],
PlotRange -> {{-1.5, 1.5}, All}],
Show[Plot[[n] PiecewiseData1D ["TargetFunction"][[x], {x, -1, 1}, Frame -> True,
PlotRange -> {{-1.2, 1.2}, {-1.4, 1.4}}, FrameTicks -> None,
PlotStyle -> Directive[GrayLevel[.3, .5], Dotted], Exclusions -> None],
Plot[#[x], {x, -1, 1}, Frame -> True, FrameTicks -> None]]]&,
{[n] PiecewiseData1D ["SmallNetsV2"], {}, {}, {}, {AspectRatio -> 1,
PlotRange -> {{-2.5, 1.5}, All}}, {}, {}, {}, {}}], 4], 1],
Dividers -> {{False, {True}, False}}, FrameStyle ->
LightGray,
Spacings -> {80, -20}, ImageSize ->
700,
Alignment -> {Center, {Bottom, Center}},
AspectRatio ->
1]
]
```



“Transfer Learning” es importante!

Lo importante es ver que la función Loss decrezca!





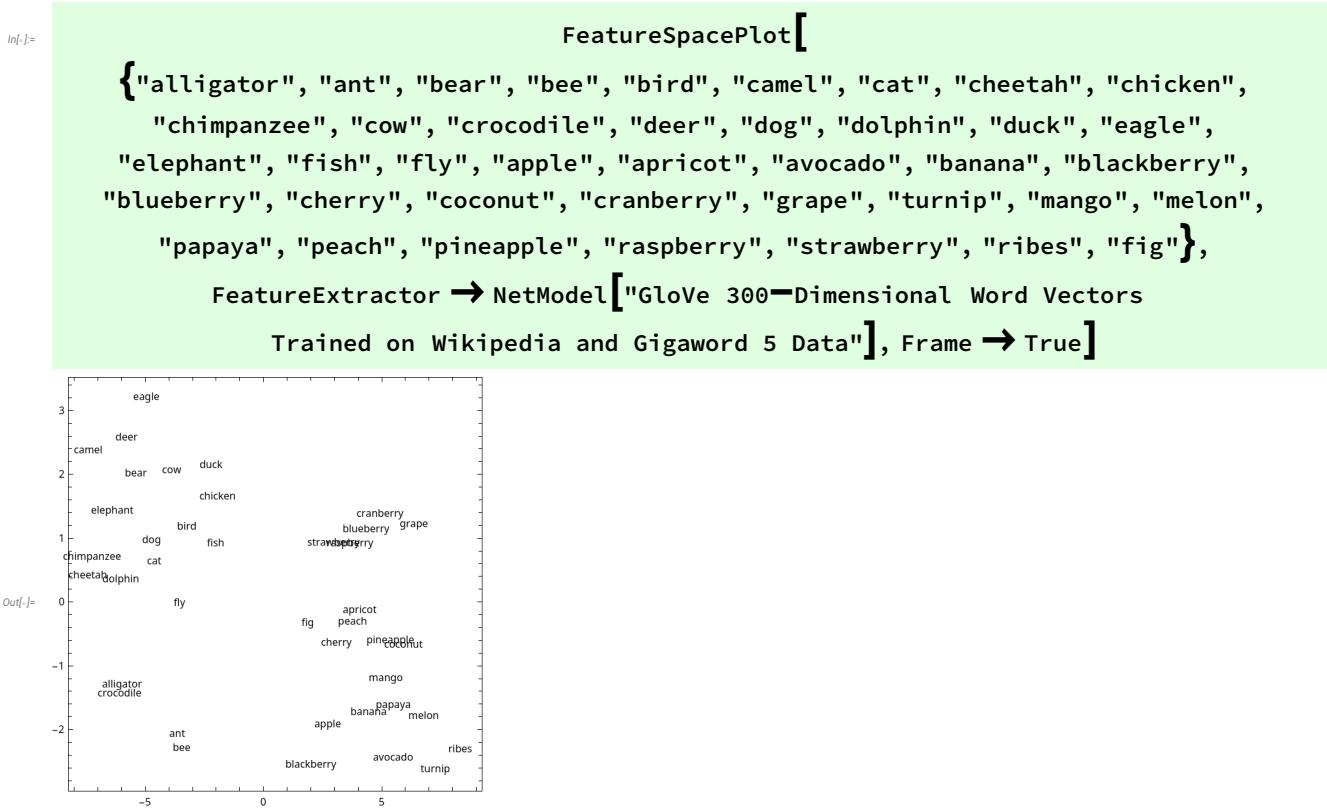
Seguro que una red lo suficientemente grande puede hacer lo que sea!

Matematica no trivial es un ejemplo de tarea dificil!

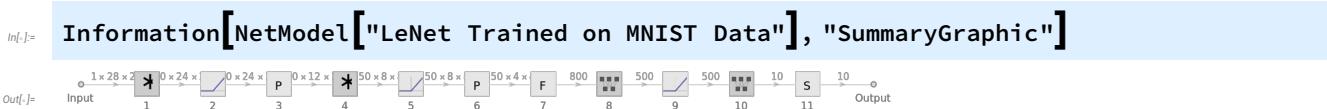


## El concepto de Embeddings -

Las redes neuronales como estan concebidas justo ya trabajan con numeros! Y por lo tanto tenemos que transformar el texto a numero! ... Como un “**meaning space**”



En la tarea de reconocimiento de imagen, tenemos 11 capas:



Por ejemplo le damos un 4:

```
In[3]:= SoftmaxLayer[]@Drop[NetModel["LeNet Trained on MNIST Data"], -1][4]
```

$$\left\{ 1.42072 \times 10^{-22}, 7.69859 \times 10^{-14}, 1.96532 \times 10^{-16}, 5.55234 \times 10^{-21}, 1., 8.33845 \times 10^{-14}, 6.89742 \times 10^{-17}, 6.52287 \times 10^{-19}, 6.51468 \times 10^{-12}, 1.9751 \times 10^{-14} \right\}$$

Que pasa si tomamos la salida justo antes del SoftMax?:

## Una capa anterior a softmax:

```
In[7]:= Drop[NetModel["LeNet Trained on MNIST Data"], -1][4]  
Out[7]= { -26.134, -6.02347, -11.994, -22.4684, 24.1717, -5.94363, -13.0411, -17.7021,  
-1.58528, -7.38389}
```

O sea, un 4 lo podemos representar como una imagen de una sola fila:

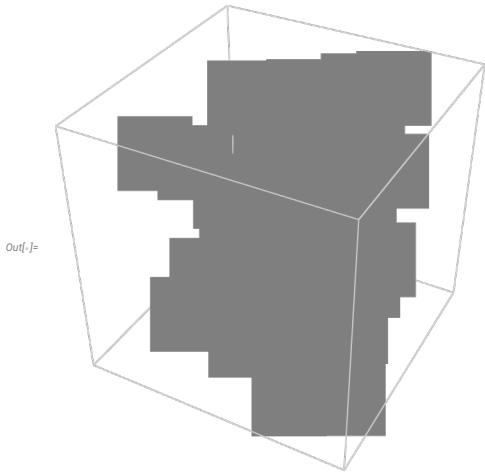
```
#> ArrayPlot[{{Drop[NetModel["LeNet Trained on MNIST Data"], -1][#]} / 30,
  ImageSize -> 160, Frame -> None, Mesh -> True,
  ColorFunction -> ChatTechColors["OrangeShort"]}] & /@ {4, 4, 4, 4, 8, 8}

{4 -> , 4 -> , 4 -> , 4 -> ,
  8 -> , 8 -> 

```

Tratemos de construir un “espacio de imagen” como el meaning space del embedding para palabras, en 3D

```
In[7]:= FeatureSpacePlot3D[SeedRandom[123];
  # \[Function] SetAlphaChannel[#, ColorNegate[#]] & @@@@ RandomSample[ResourceData["MNIST"], 50],
  FeatureExtractor \[Function] NetDrop[NetModel["LeNet Trained on MNIST Data"], -3]]
```

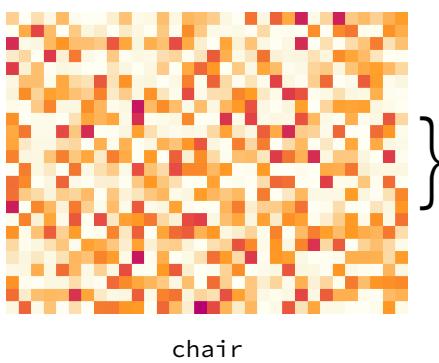
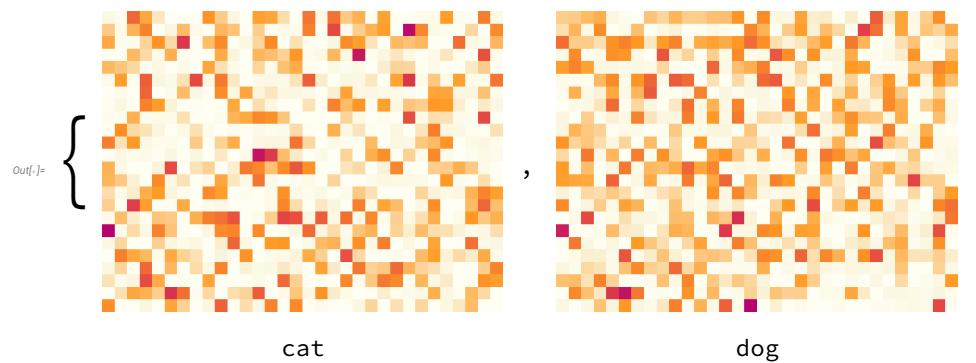




La tarea que nos trajo aquí ! ... “word prediction”.  
Imagina que tenemos “the \_\_\_ cat”.

Como podemos trabajar esa caracterización de las palabras. Por ejemplo:

```
In[7]:= Labeled[ArrayPlot[Partition[
  First[NetExtract[NetModel["GPT2 Transformer Trained on WebText Data"],
    {"embedding", "embeddingtokens"}]][NetExtract[
      NetModel["GPT2 Transformer Trained on WebText Data"], "Input"]][#]]], 32],
  Frame \[Rule] None, ImageSize \[Rule] 350, ColorFunction \[Rule] ChatTechColors["Orange"]],
  #]& /@ {"cat", "dog", "chair"}]
```



Lo bueno es que tenemos una forma de poner a las palabras en una forma AMIGABLE para una red neuronal.

ChatGPT lida con “tokens”. Unidades de lenguaje.

# Adentro de ChatGPT -- Una red GPT-3 con 175 millardos de parametros!

Su elemento de arquitectura mas destacable es una cosa que se llama “Transformer”

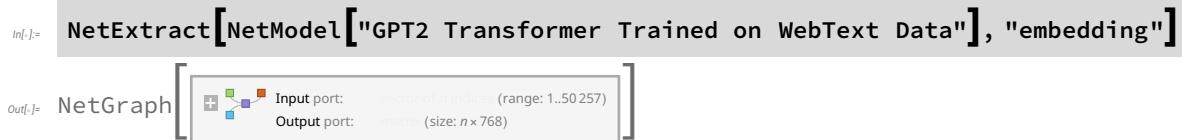
Introducen la noción de “Attention” y la idea de prestar mas atención a un sector de la secuencia de token que de las demás partes. Recuerden, agregar el siguiente “token”!

## Tres pasos:

- 1) Agarra el texto y hace el embedding de los tokens.
- 2) Los trata como lo hace cualquier red neuronal.
- 3) Genera un arreglo de 50,000 numeros y le asigna probabilidades.

Hay un monton de detalles en la manera que la red esta creada.

(1) El modulo de Embeddings. Aqui GPT-2

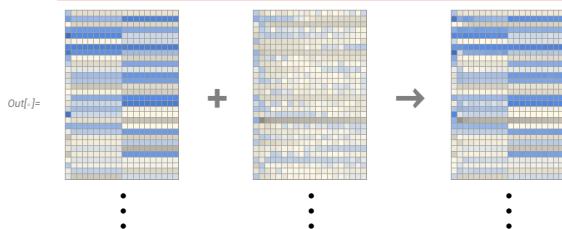


`n` es el numero de tokens dado como “input”. `[x]` 768 es la dimensión del arreglo con la que estamos representando cada token.

El **SI** es importante. Es una secuencia de enteros de **POSICIONES** para saber donde estaba cada palabra o token. Y de estos enteros sacar otro vector embedding!

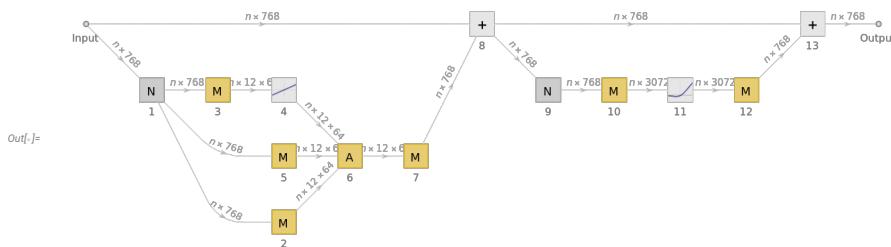
Aqui lo que hace el modulo de Embedding para la secuencia:  
“hello my dear how are you hope everything is all right”

```
In[.]:= Module[
{matrixPlot},
matrixPlot[matrix_?MatrixQ]:= Labeled[
MatrixPlot[Take[Transpose@matrix, 30],
Mesh→True, ImageSize→100, Frame→None], "⋮"];
Row[Riffle[
matrixPlot/@Values@NetModel["GPT2 Transformer Trained on WebText Data"][[1,
"hello hello hello hello hello hello hello hello",
"hello hello bye bye bye bye bye bye bye bye",
{NetPort[{"embedding", "embeddingtokens"}],
NetPort[{"embedding", "embeddingpos"}], NetPort[{"embedding"}]}],
Style[#, 18, Gray, Bold]&/@{{"+", "→"}}, Spacer[5]]]
]
```



## (2) El bloque de atención (attention)

```
In[.]:= (* == Un bloque de attention: 12 para GPT-2 y 96 para GPT-3 == *)
Information[NetFlatten@NetExtract[
NetModel["GPT2 Transformer Trained on WebText Data"],
>{"decoder", 1}], "FullSummaryGraphic"]
```



Los 12 attention heads lucen asi:

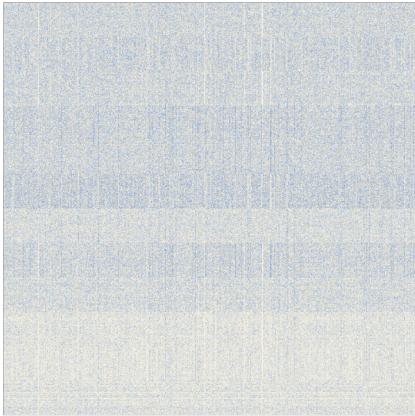
Despues de procesar: y como (3) lo pasamos por una red complemento conectada!

como una red cualquiera!

In[1]:=

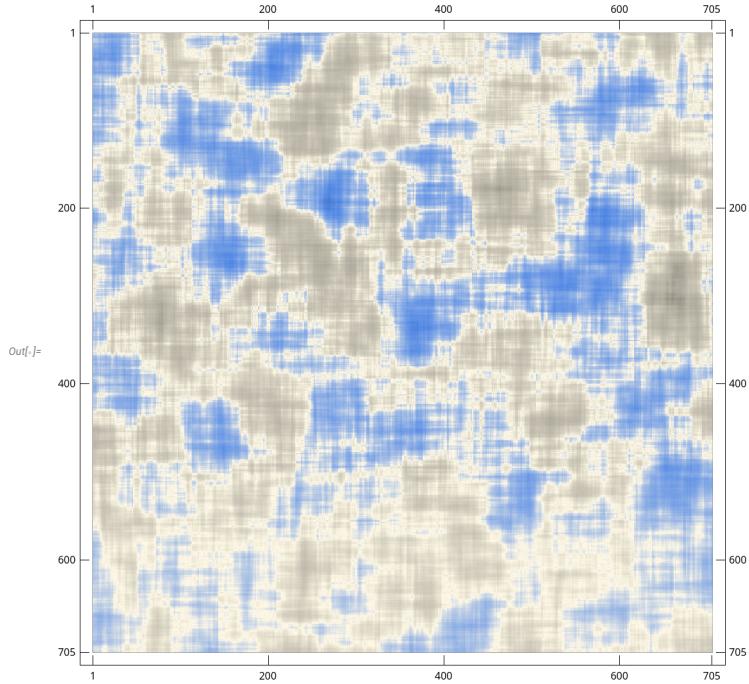
```
MatrixPlot[Normal@NetExtract[
  NetModel["GPT2 Transformer Trained on WebText Data"],
  {"decoder", 1, "1", "attention", 1(*key*), "Net", "Weights"}],
  Frame → False, FrameTicks → None]
```

Out[1]:=



Tomando un average que se mueve 64x64:

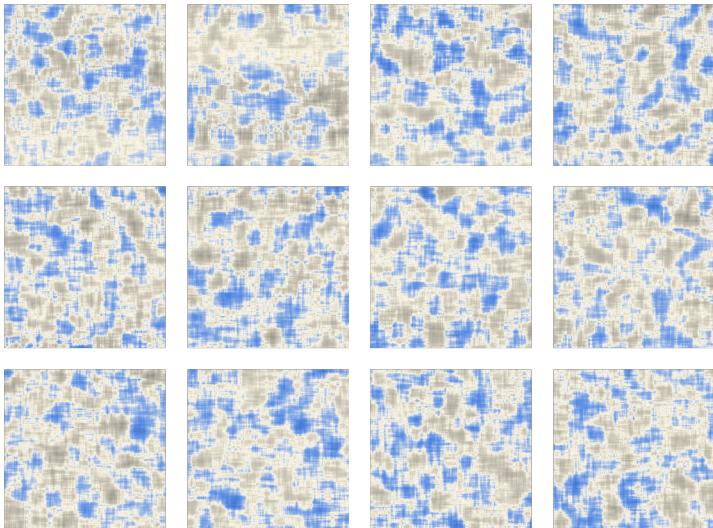
```
In[1]:= MatrixPlot[
  MovingAverage[MovingAverage[#, 64] & /@(
    Normal @ NetExtract[
      NetModel["GPT2 Transformer Trained
        on WebText Data"], {
        "decoder", 1, "1", "attention", 1
        (*key*), "Net", "Weights"}], 64],
  Frame → True, ImageSize → Large]
```



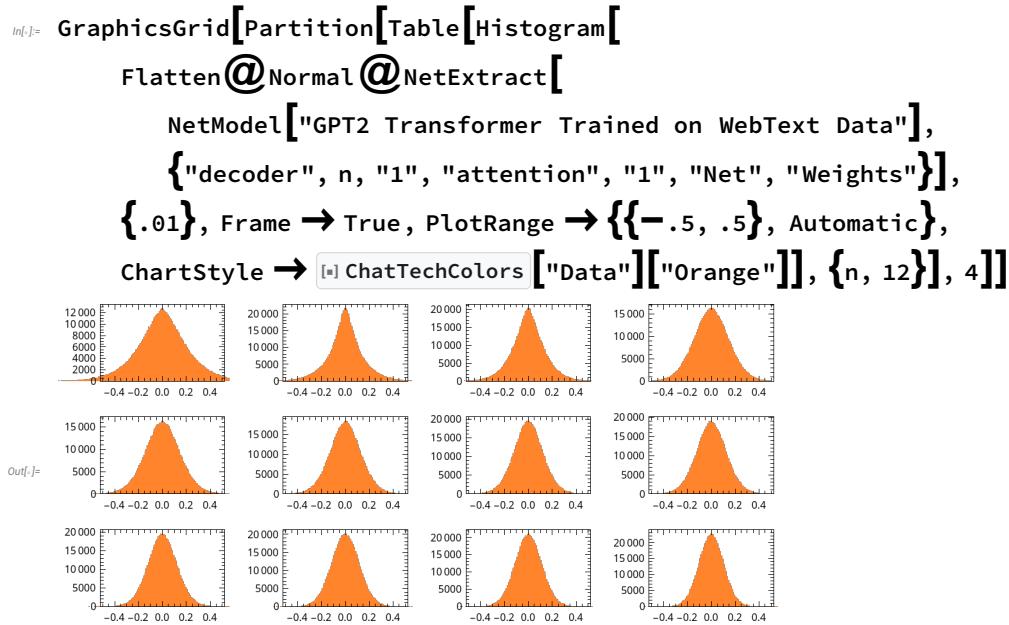
# Que determina esa estructura? El lenguaje humano en ultima instancia!

Los averages que se mueven para todos los 12 heads:

```
In[1]:= GraphicsGrid[Partition[Table[MatrixPlot[
  MovingAverage[MovingAverage[#, 64] & /@ Normal @ NetExtract[
    NetModel["GPT2 Transformer Trained on WebText Data"],
    {"decoder", n, "1", "attention", 1(*key*), "Net", "Weights"}], 64],
  Frame → False], {n, 12}], 4]]
```



Aunque luzcan parecidas, sus distribuciones:



Que hizo el transformer? Bueno transformó el embedding original en otro embedding!

(último paso) Toma estos embeddings y asigna probabilidades!

No hay looping back, ni repetidos reprocesamientos.. nada .. solo forward.  
Bueno también lee hacia atrás para saber qué palabra colocar.