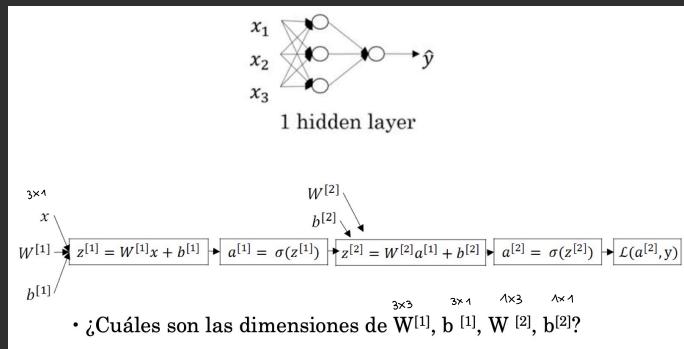
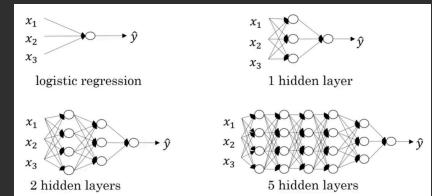
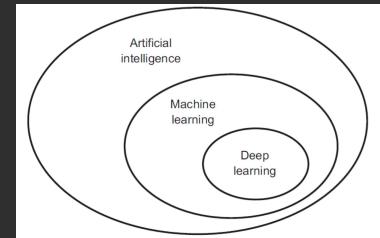
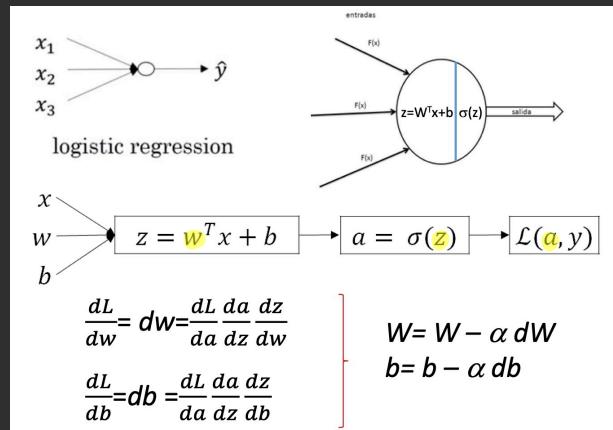


- "Deep" not because it understands deeper, but because the model has more hidden layers
- automatize the "feature engineering"
- it needs powerful hardware, lots of data

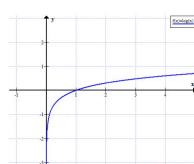


$(m_l \times m_{l-1})$ Without transposing

Loss function logistic classification

$$L(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

- Si $y=1$ $L(y, \hat{y}) = -\log \hat{y}$: minimizar $L(y, \hat{y}) \rightarrow$ maximizar \hat{y}
- Si $y=0$ $L(y, \hat{y}) = -\log (1-\hat{y})$: minimizar $L(y, \hat{y}) \rightarrow$ minimizar \hat{y}



$$y=0 \Rightarrow L(y, \hat{y}) = -\log(1-\hat{y})$$

$$\hat{y}=0 \Rightarrow L(y, \hat{y}) = \infty$$

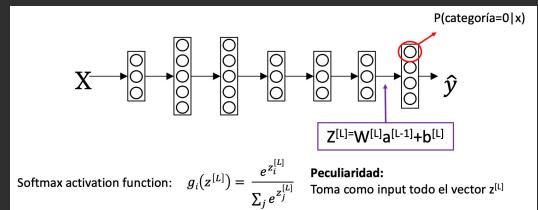
$$\hat{y}=1 \Rightarrow L(y, \hat{y}) \approx \infty$$

$$y=1 \Rightarrow L(y, \hat{y}) = -\log \hat{y}$$

$$\hat{y}=0 \Rightarrow L(y, \hat{y}) \approx \infty$$

$$\hat{y}=1 \Rightarrow L(y, \hat{y}) = 0$$

Logistic is for binary classification while softmax for multiple classification.



Generalized Loss Function - Cross-Entropy

$$L(y, \hat{y}) = -\sum_j y_j \log \hat{y}_j \quad \text{Cross-Entropy}$$

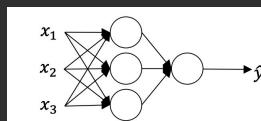
to minimize L we need to maximize the probability of belonging to the correct class.

Cost Function

$$\text{Cost} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{y}_i)$$

Activation Function

Without it we could only learn hyperplane decision boundaries.

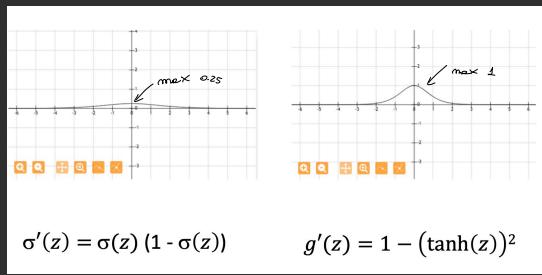
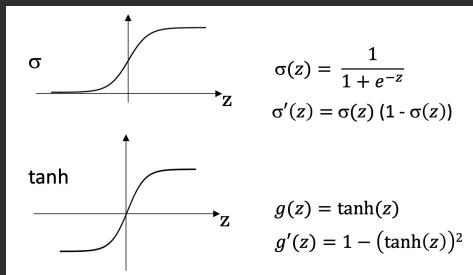


Dado un cierto x :

$$\begin{aligned} z^{[1]} &= W^{[1]}x + b^{[1]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) = z^{[1]} \\ z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= g^{[2]}(z^{[2]}) = z^{[2]} \end{aligned}$$

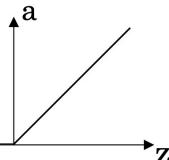
$$\begin{aligned} a^{[1]} &= z^{[1]} = w^{[1]}x + b^{[1]} \\ a^{[1]} &= z^{[1]} = (w^{[1]}_1 x_1 + w^{[1]}_2 x_2 + b^{[1]}) \\ a^{[2]} &= w^{[2]}(w^{[1]}_1 x_1 + b^{[1]}) + b^{[2]} \\ a^{[2]} &= (w^{[2]}_1 w^{[1]}_1 x_1 + b^{[2]}) + (w^{[2]}_2 w^{[1]}_2 x_2 + b^{[2]}) \\ &= \underline{w^T x + b^T} \end{aligned}$$

!Sería equivalente a tener un único nodo!



Vanish gradient

- having lots of hidden layers and using sigmoid could bring to not learn anything because $\sigma(x)$ has 0.25 as maximum, so multiplying lots of small numbers let the gradient tends to 0.
- In the hidden layers, sigmoid is not useful.



ReLU

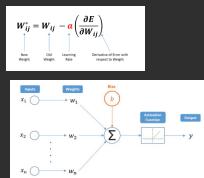
$$f(x) = \max(0, x)$$

- fast
- no gradient vanishing
- not differentiable in 0 . For convention is chosen 0 or 1 in that point.

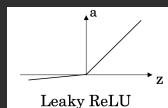
Dying ReLU

- high LR \Rightarrow weights could become negatives

- high negative bias \Rightarrow input relu negative



Solve the problem



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases}$$

How to use activation function in deep neural networks

Input layer: tanh, σ

Hidden layers: ReLU, leaky ReLU

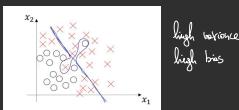
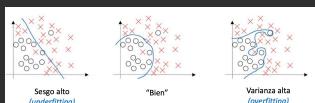
Output layer: binary classification \Rightarrow tanh, σ

multiclass classification \Rightarrow softmax

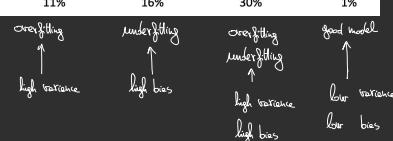
Dataset

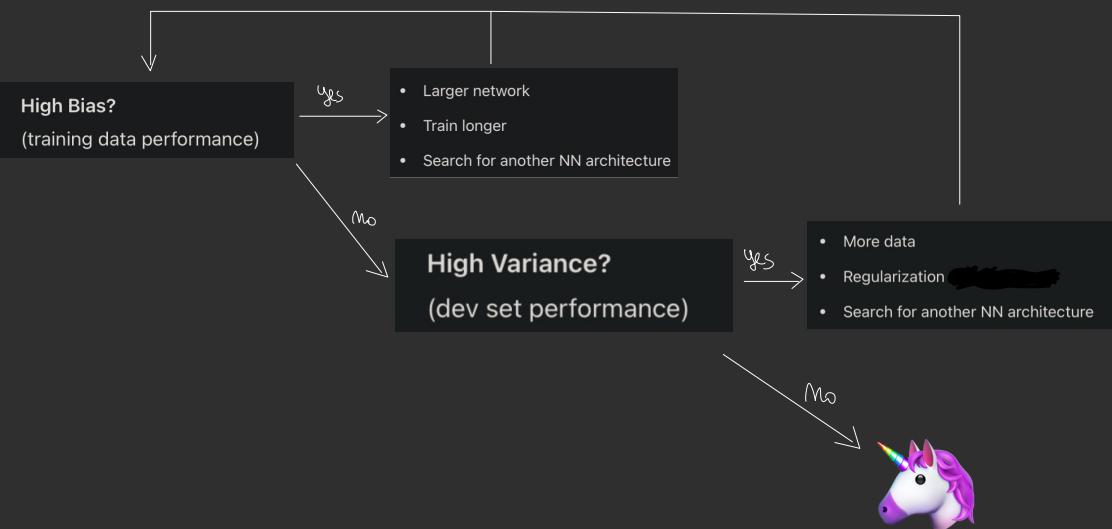
Training set	Dev set	Test set
Antes de Big Data: 70% (train) / 30% (test) 60% (train) / 20% (dev) / 20% (test)		~ 10.000
Big Data: 98% (train) / 1% (dev) / 1% (test)		$\sim 1.000.000$ 10000!

Bias \ Variance



Error humano ~ 0%	Train set error	1%	15%	15%	0.5%
	Dev set error	11%	16%	30%	1%





Regularization

- If weight are smaller, the network is less specialized

1 L1 Regularization in Neural Networks

- Adds the sum of the absolute values of the weights to the loss function:
$$L = \text{Loss} + \lambda \sum |w_i|$$
- Encourages sparsity by forcing some weights to become exactly zero.
- Can help prune unnecessary connections in a neural network.

Effects in NN:

- Reduces complexity by eliminating some weights.
- Leads to a sparser network (can be useful for feature selection).
- May cause instability when training deep networks.

2 L2 Regularization in Neural Networks

- Adds the sum of the squared values of the weights to the loss function:
$$L = \text{Loss} + \lambda \sum w_i^2$$
- Forces weights to be small but not necessarily zero.
- Helps distribute learning across all neurons.

Effects in NN:

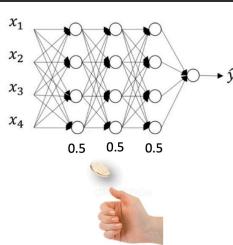
- Prevents overfitting by discouraging large weights.
- Works well in deep networks where all features contribute.
- Doesn't enforce sparsity but smoothens learning.

3 Key Differences in Neural Networks

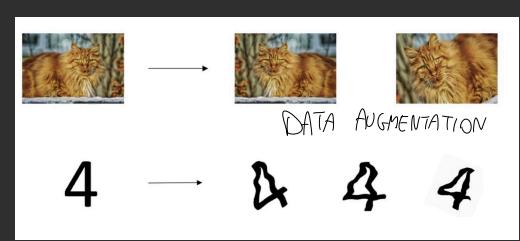
Feature	L1 Regularization (Lasso)	L2 Regularization (Ridge)
Effect on weights	Some weights become zero (sparse network)	Weights become small but not zero
Sparsity	Yes (reduces unnecessary connections)	No (keeps all connections)
Best use case	When you want to prune neurons/connections	When you want to stabilize learning
Training effect	Can make training unstable	More stable training

Conclusion:

- L1 removes unimportant weights (good for reducing model size).
- L2 shrinks weights evenly (good for stabilizing training).
- For deep learning, L2 (weight decay) is more commonly used.



DROPOUT



4

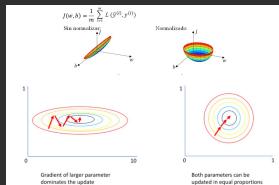
4

4

Normalize input

- help the training

$$z = \frac{x - \mu}{\sigma}$$



Weight initialization

- ReLU / Leaky ReLU → He Initialization ✓
- Sigmoid / Tanh → Xavier or Lecun ✓
- Multi layer → Orthogonal ✓
- Bias → Zero ✓

Gradient descent

1. Batch Gradient Descent

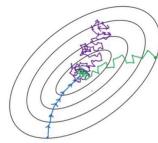
- Updates the model's weights after processing the entire dataset (one full epoch).
- Advantage:** More stable updates since they use the full dataset.
- Disadvantage:** Computationally expensive and slow for large datasets.
- Best for:** Small datasets where computational efficiency is not a concern.

2. Stochastic Gradient Descent (SGD)

- Updates the model's weights after each individual sample.
- Advantage:** Very fast updates, does not require processing the entire dataset.
- Disadvantage:** High variance in updates, may lead to noisy convergence.
- Possible Solution:** Decrease the learning rate over time using a decay function.
- Best for:** Large datasets where immediate updates are preferred.

3. Mini-Batch Gradient Descent

- Updates the model's weights after processing a small subset (mini-batch) of the dataset.
- Advantage:** Balances speed (like SGD) and stability (like Batch GD).
- Disadvantage:** Choosing the right batch size is crucial.
- Best for:** Most deep learning tasks since it leverages parallelization.



Stochastic gradient descent

Mini-Batch gradient descent:

Batch gradient descent:

Method	Update Frequency	Speed	Stability	Suitability
Batch GD	After full dataset	Slow	Stable	Small datasets
SGD	After each sample	Very fast	Noisy	Large datasets
Mini-Batch GD	After mini-batch	Fast	Balanced	Deep learning

Batch normalization

What is Batch Normalization?

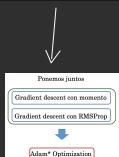
Batch Normalization (BN) is a technique that makes training faster and more stable by normalizing the values (activations) in a mini-batch.

Momentum / RMSProp / Adam

Optimizer	Key Idea	Strengths	Weaknesses
Momentum	Uses past gradients to accelerate learning	✓ Faster convergence ✓ Reduces oscillations	✗ Struggles with adaptive learning rates
RMSprop	Adapts learning rate for each parameter using a moving average of squared gradients	✓ Good for non-stationary problems ✓ Works well in deep networks	✗ Can become too aggressive in updates
Adam	Combines Momentum and RMSprop	✓ Best of both worlds ✓ Faster & more stable	✗ Uses more memory

$$\begin{array}{c} \text{last } | v_0 | \rightarrow | v_1 | \rightarrow | v_2 | \rightarrow \text{gradents} \\ \beta = 0.9 | v_0 | = 0.5 \\ \rightarrow \end{array} \quad \begin{array}{l} v_{dW} = \beta v_{dW} + (1 - \beta) dW \\ v_{db} = \beta v_{db} + (1 - \beta) db \end{array} \quad \begin{array}{l} W = W - \alpha v_{dW} \\ b = b - \alpha v_{db} \end{array}$$

$$\rightarrow \quad \begin{array}{l} s_{dW} = \beta s_{dW} + (1 - \beta) dW^2 \\ s_{db} = \beta s_{db} + (1 - \beta) db^2 \end{array} \quad \begin{array}{l} W = W - \alpha \frac{dW}{\sqrt{s_{dW}} + \epsilon} \leftarrow \text{Prevent division by 0} \\ b = b - \alpha \frac{db}{\sqrt{s_{db}} + \epsilon} \end{array}$$



$$\begin{aligned} v_{dW} &= \beta_1 v_{dW} + (1 - \beta_1) dW \\ v_{db} &= \beta_1 v_{db} + (1 - \beta_1) db \\ s_{dW} &= \beta_2 s_{dW} + (1 - \beta_2) dW^2 \\ s_{db} &= \beta_2 s_{db} + (1 - \beta_2) db^2 \\ W &= W - \alpha \frac{v_{dW}}{\sqrt{s_{dW}} + \epsilon} \\ b &= b - \alpha \frac{v_{db}}{\sqrt{s_{db}} + \epsilon} \end{aligned}$$

Momentum

RMSProp

Instead of using the same learning rate for all parameters, RMSProp scales it based on recent gradients.

• If ϵ_{dW} is big, the step is small.

Why Do We Need RMSProp?
Imagine you're walking down a mountain (minimizing loss), but the path is very uneven.

• Some steps are very steep (big gradients), some are flat.

• Some areas are flat (near the minimum).

• If you take big steps, you might fall (overshooting).

• If you take small steps, it will take forever to reach the bottom.

• RMSProp helps by updating your step size dynamically so it isn't always the same.

• In steep areas, you take smaller steps to avoid overshooting.

• In flat areas, you take bigger steps (or speed up learning).

The Problem: Walking Down a Mountain

- The terrain is rough – Some areas are steep, some are flat, and some have sharp drops.
- If you take big steps, you might trip – Standard Gradient Descent can overshoot the minimum.
- If you take small steps, it will take forever – RMSProp can slow down too much in flat areas.
- You also don't want to zig-zag too much – Momentum helps smooth your path.