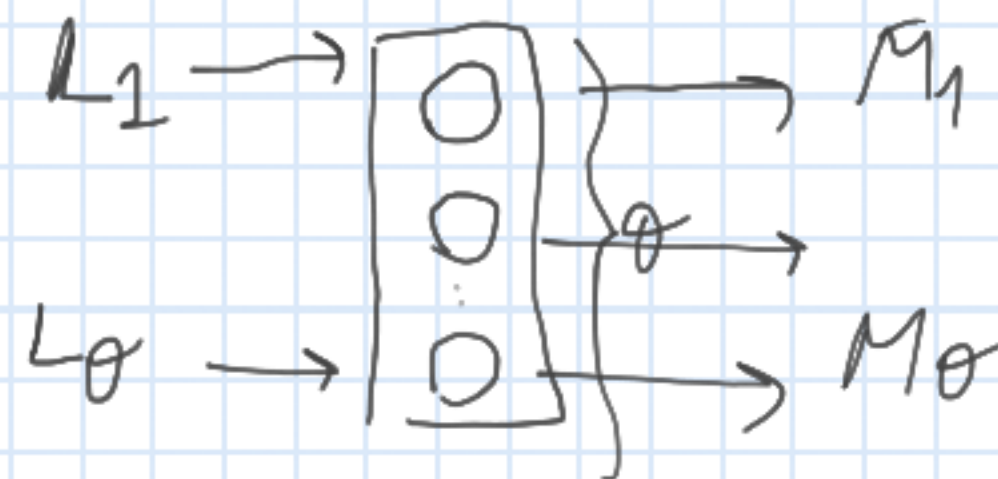
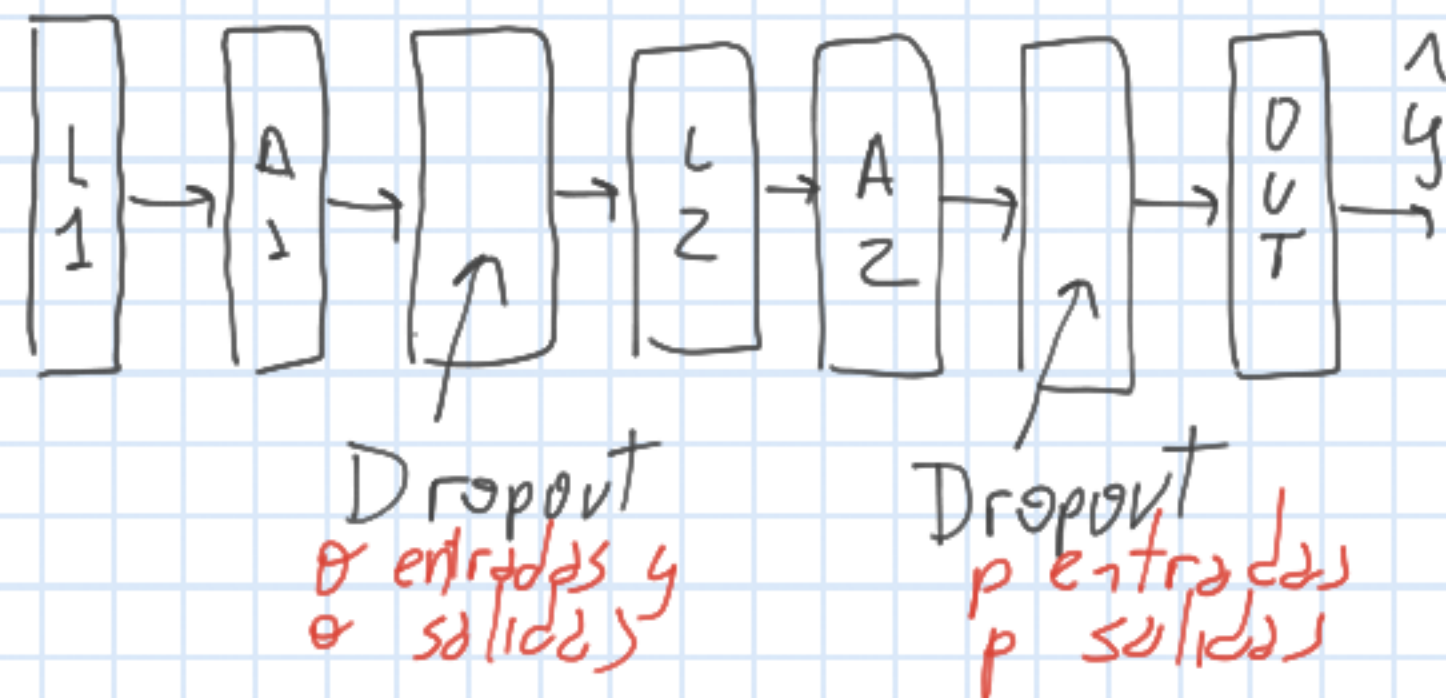
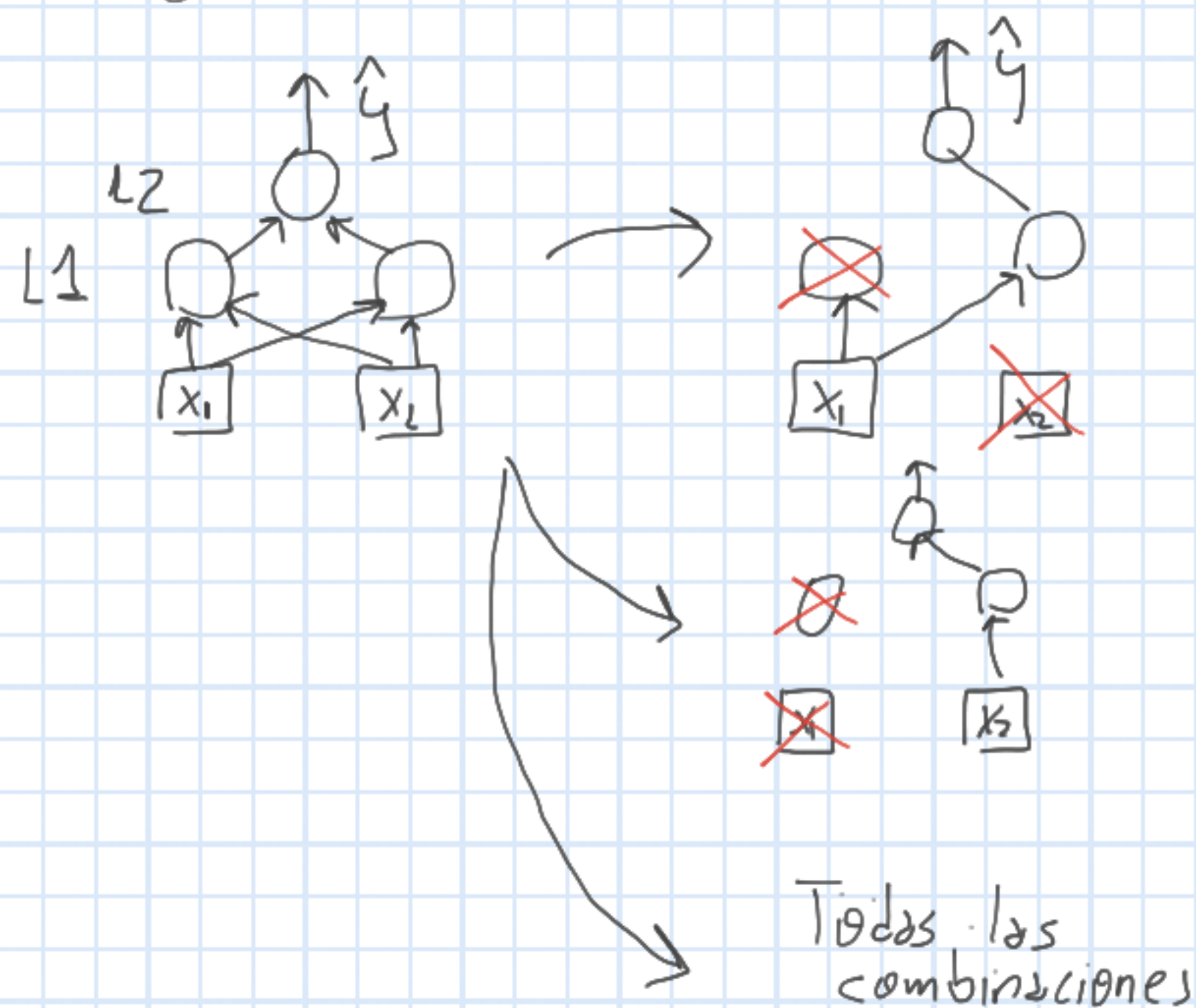


* Regularización Dropout



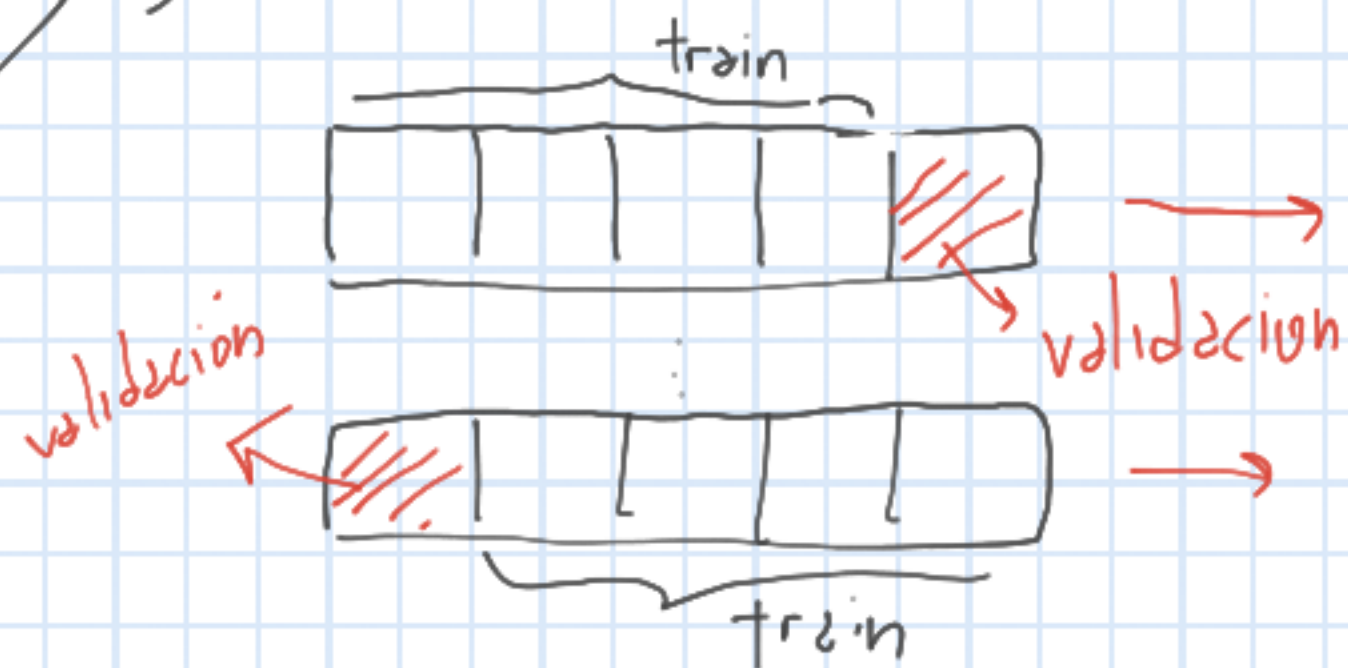
$P_i \sim \text{Bernoulli}(p)$
 P_i son iid
 $p \in [0, 1]$

$$M_i = \begin{cases} L_i & \text{if } P_i = 1 \\ 0 & \text{if } P_i = 0 \end{cases}$$

* Optimizadores

(1) Dataset, redes neuronales

→ Cross-Validation K-folds.



metrica-val 1

validation

metrica-val k

$$\text{metrica-val} = \frac{1}{k} \sum_{i=1}^k \text{metrica-val}_i$$

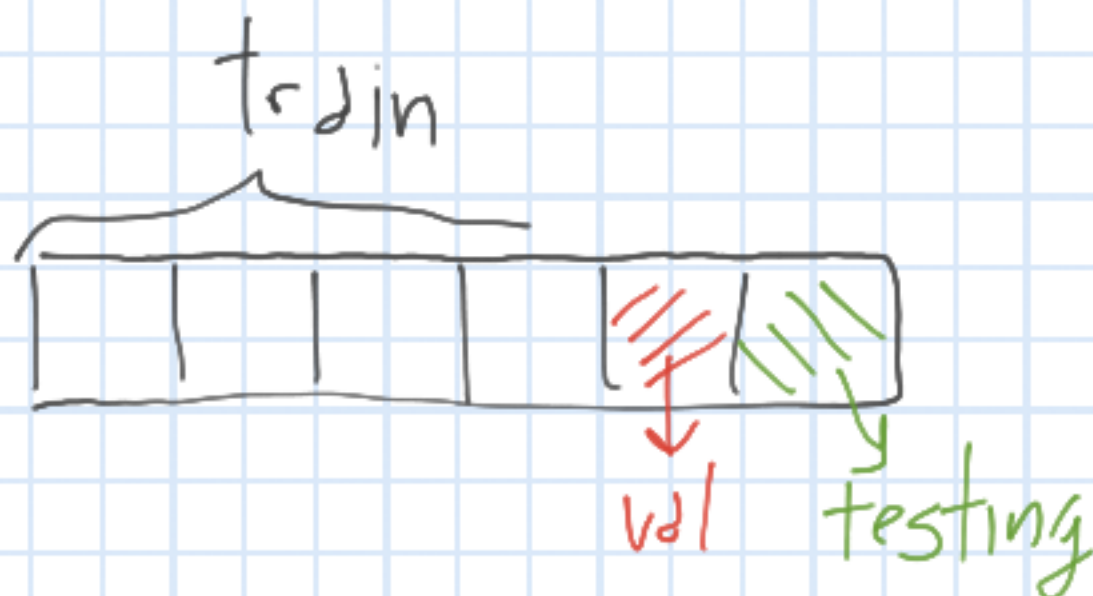
seleccionar HP modelo

me quedo con el mejor HP

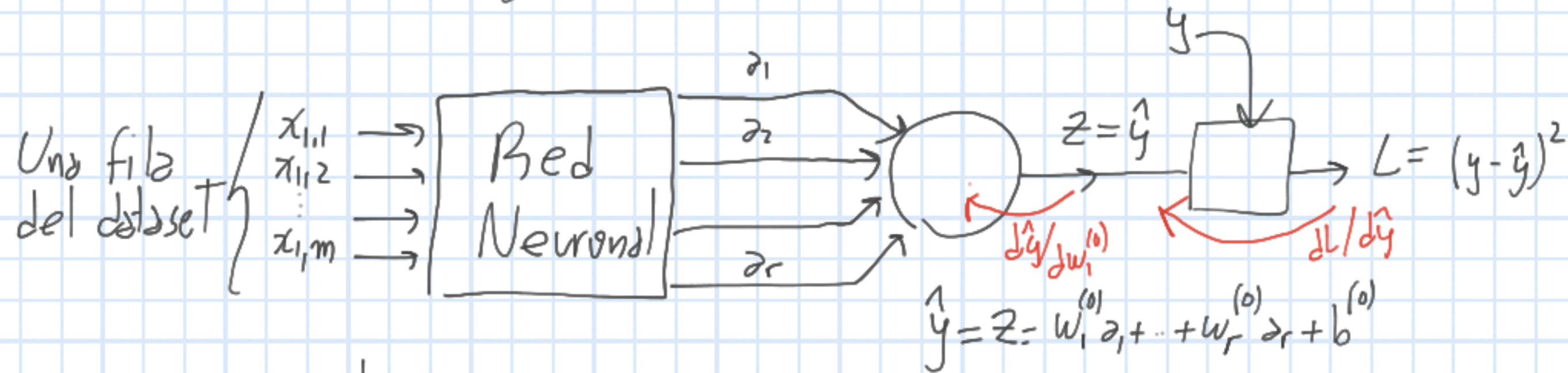


metrica final

En DL:



* (2) Data normalization



* No normalizar
entrada afecta
velocidad de
entrenamiento

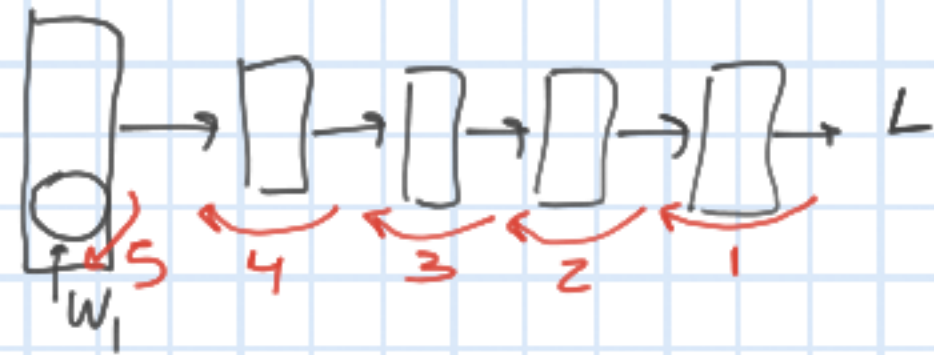
* No normalizar
la salida puede
generar gradientes
que explotan

$$\frac{dL}{dw_1^{(0)}} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dw_1^{(0)}} = -z(y - \hat{y}) \cdot z_1$$

$$w_1^{(0)} \leftarrow w_1^{(0)} - \alpha \frac{dL}{dw_1^{(0)}}$$

$$w_1^{(0)} \leftarrow w_1^{(0)} + \alpha z(y - \hat{y}) \boxed{z_1}$$

* (3) Vanishing gradients



$$\frac{dL}{dw_1} = (1) (2) (3) (4) (5)$$

usando como activación $\sigma(\cdot)$

tiende a cero

$$w_1 \leftarrow w_1 - \alpha \frac{dL}{dw_1} \rightarrow w_1 \leftarrow w_1 \quad (\text{no aprende})$$

aprox. 0

¿Cómo se resuelve? \rightarrow ReLU

* (4) Exploding gradients

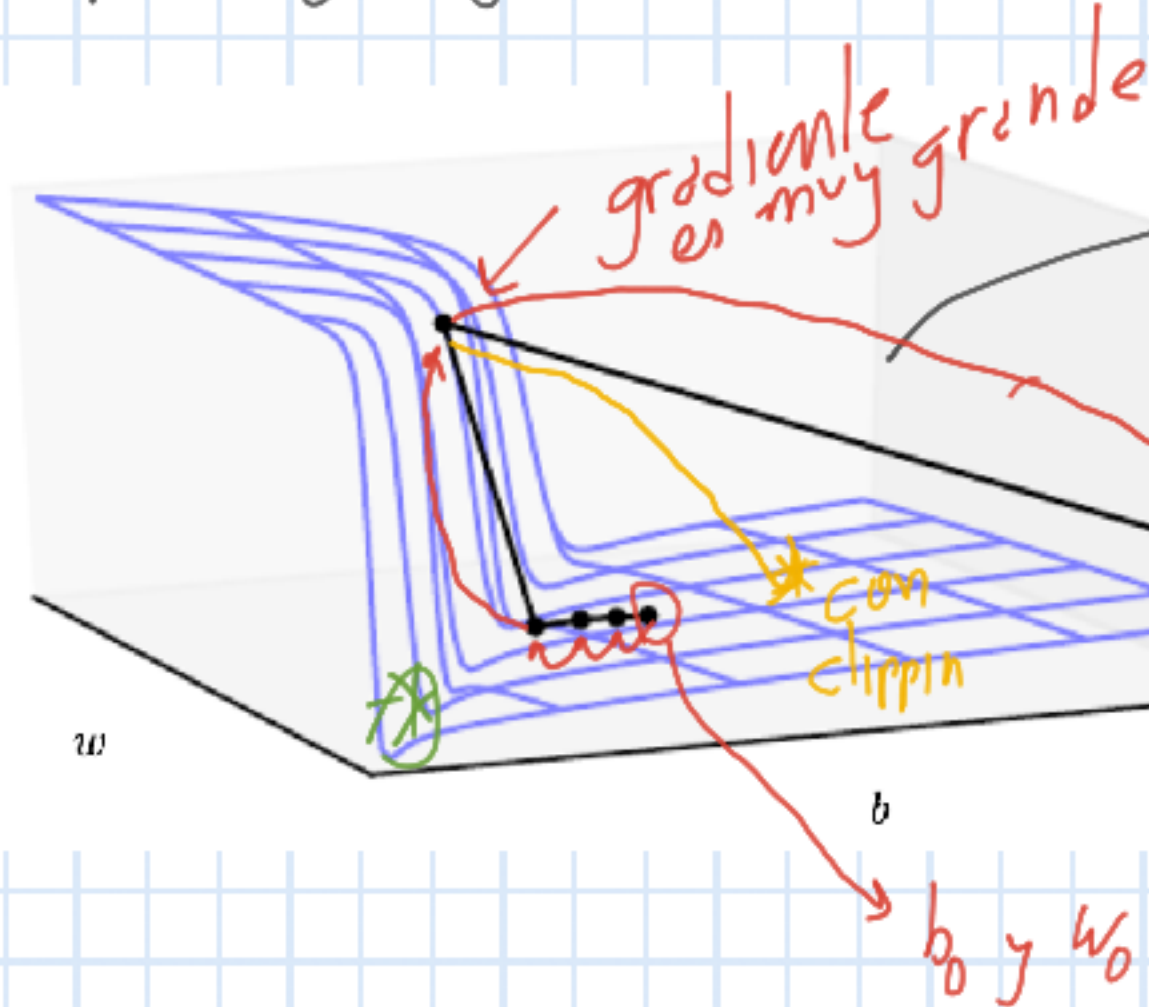
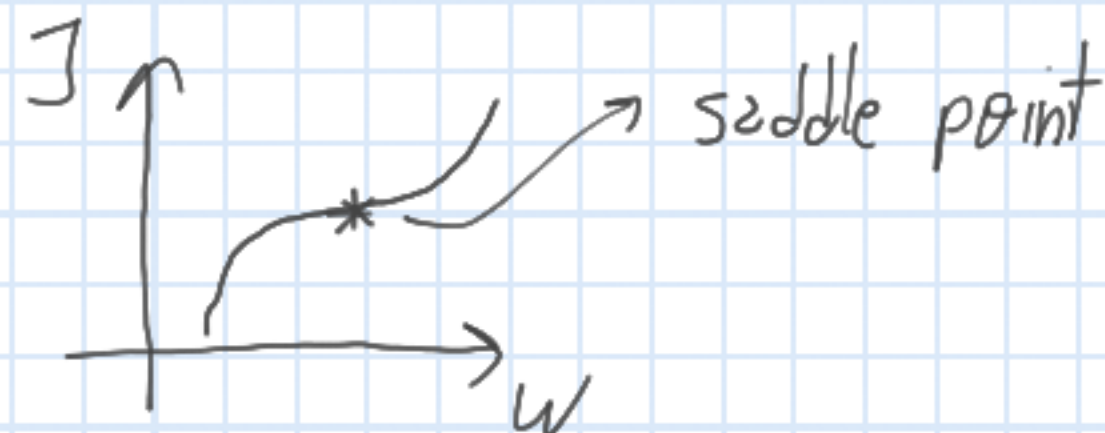
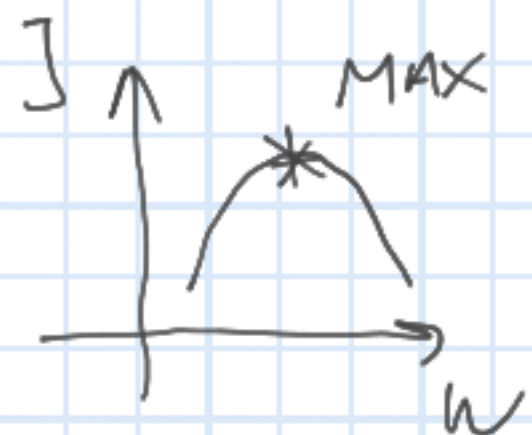
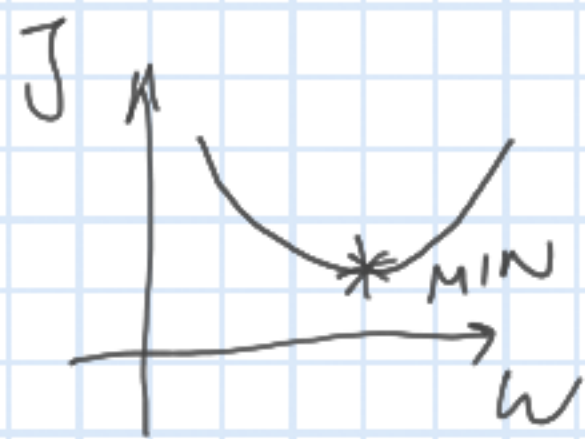


gráfico $J(w, b)$

perdi toda la optimización

\rightarrow gradient clipping
(limita el máx valor que puede tomar)

Optimización de funciones



$$(1D) \quad \frac{dJ}{dw} = 0$$

(n dim)

$$\bar{\nabla}_w(J) = \begin{bmatrix} \partial J / \partial w_1 \\ \vdots \\ \partial J / \partial w_n \end{bmatrix} = \bar{0}$$

modelos
lineales

$\bar{\nabla}(J) = \bar{0}$
tiene forma
cerrada

modelos
más complejos

$\bar{\nabla}(J) = \bar{0}$
no tiene
forma
cerrada

$$L: \mathbb{R}^1 \rightarrow \mathbb{R}^n$$

n params desconocidos

$$\bar{\nabla}(L) = \begin{bmatrix} \partial L / \partial w_1 \\ \vdots \\ \partial L / \partial w_n \end{bmatrix} \in \mathbb{R}^{n \times 1}$$

$$H(L) = \begin{bmatrix} \partial^2 L / \partial w_1^2 & \partial^2 L / \partial w_1 \partial w_2 & \dots & \partial^2 L / \partial w_1 \partial w_n \\ \vdots & \vdots & \ddots & \vdots \\ \partial^2 L / \partial w_n \partial w_1 & \dots & \dots & \partial^2 L / \partial w_n^2 \end{bmatrix}$$

necesito
algoritmos de optimización
(SGD)

$$\bar{w} \leftarrow \bar{w} + \alpha \bar{\nabla}(J)$$

Punto: \bar{w}_0 , $J(\bar{w}_0)$ dato

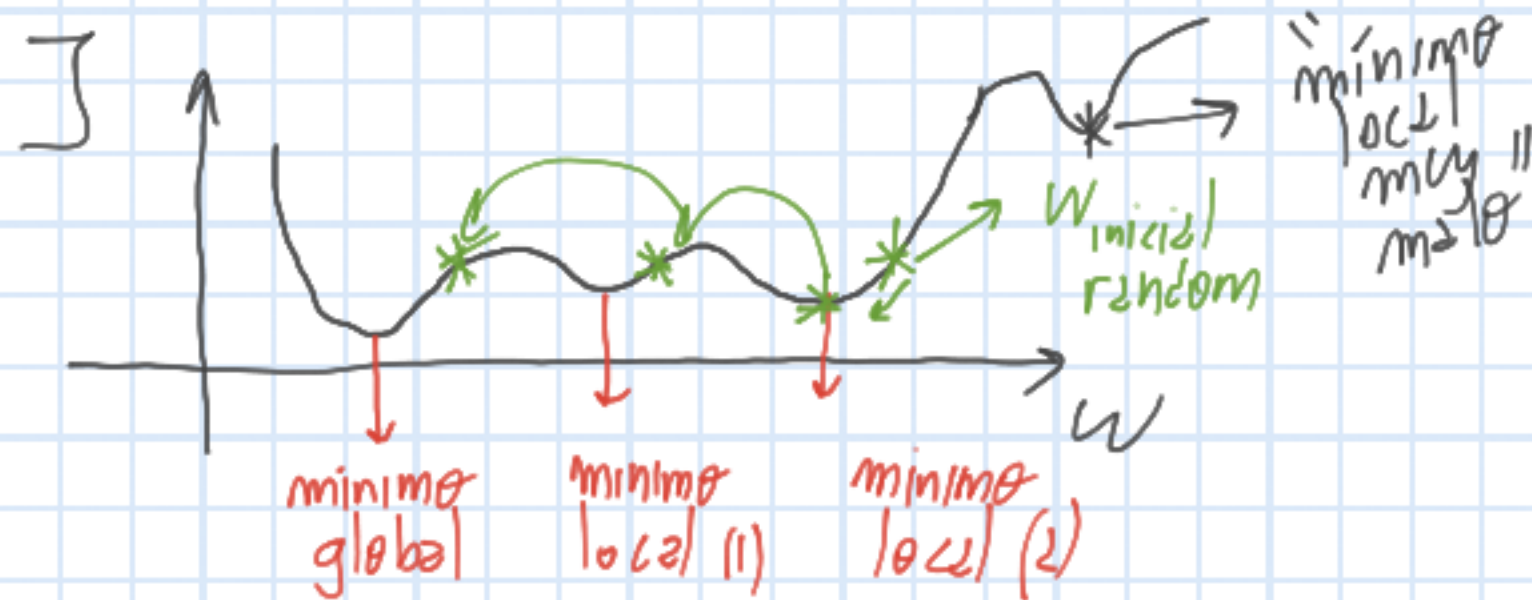
$$J(\bar{w}) \underset{\substack{\uparrow \\ \text{Taylor}}}{\approx} J(\bar{w}_0) + (\bar{w} - \bar{w}_0)^T \nabla J + \frac{1}{2} (\bar{w} - \bar{w}_0)^T H (\bar{w} - \bar{w}_0) + \dots$$

$$\underbrace{J(\bar{w} - \alpha \nabla J)}_{\substack{\text{eq. act.} \\ \text{paso}}} \approx J(\bar{w}_0) - \alpha \nabla J^T \nabla J + \frac{1}{2} \alpha^2 \nabla J^T H \nabla J + \dots$$

$$\rightarrow \text{Ejemplo} = J = \frac{1}{2} w^2 \quad \rightarrow \quad \begin{aligned} \partial J / \partial w &= w \\ \partial^2 J / \partial w^2 &= 1 \end{aligned}$$

Minimo global vs minimo local

DL \rightarrow la función J es no convexa

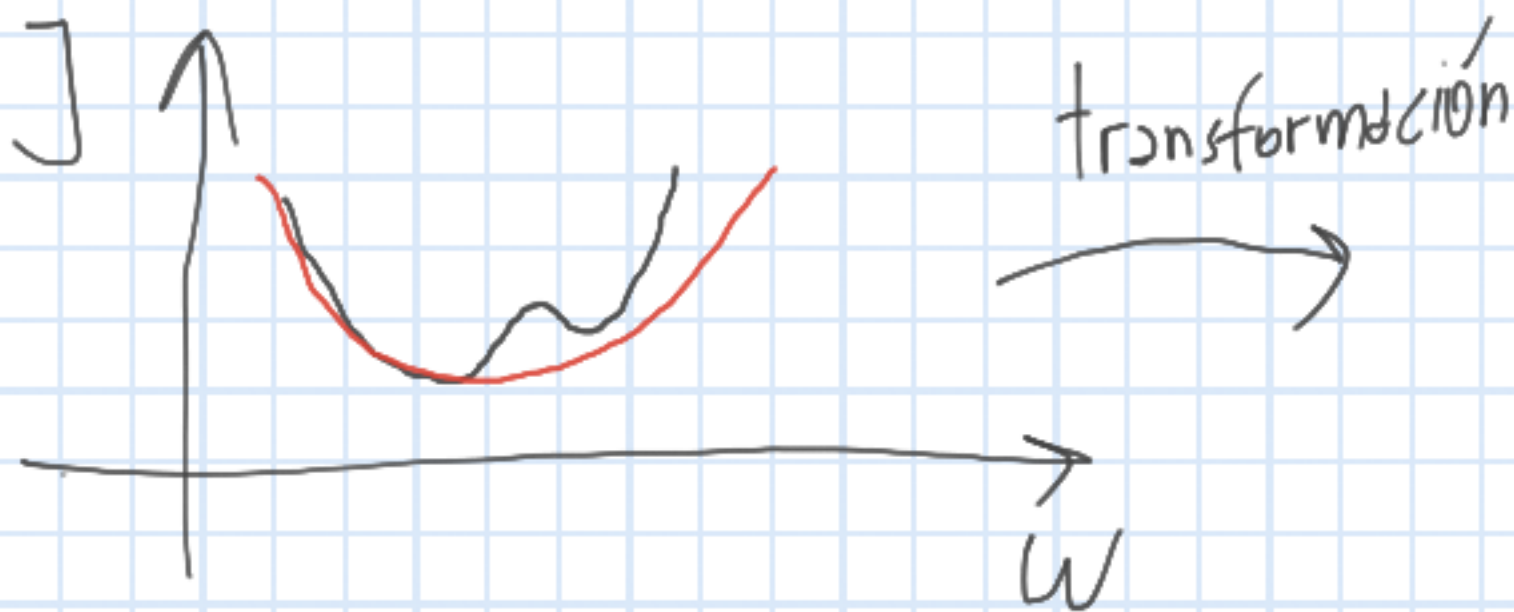


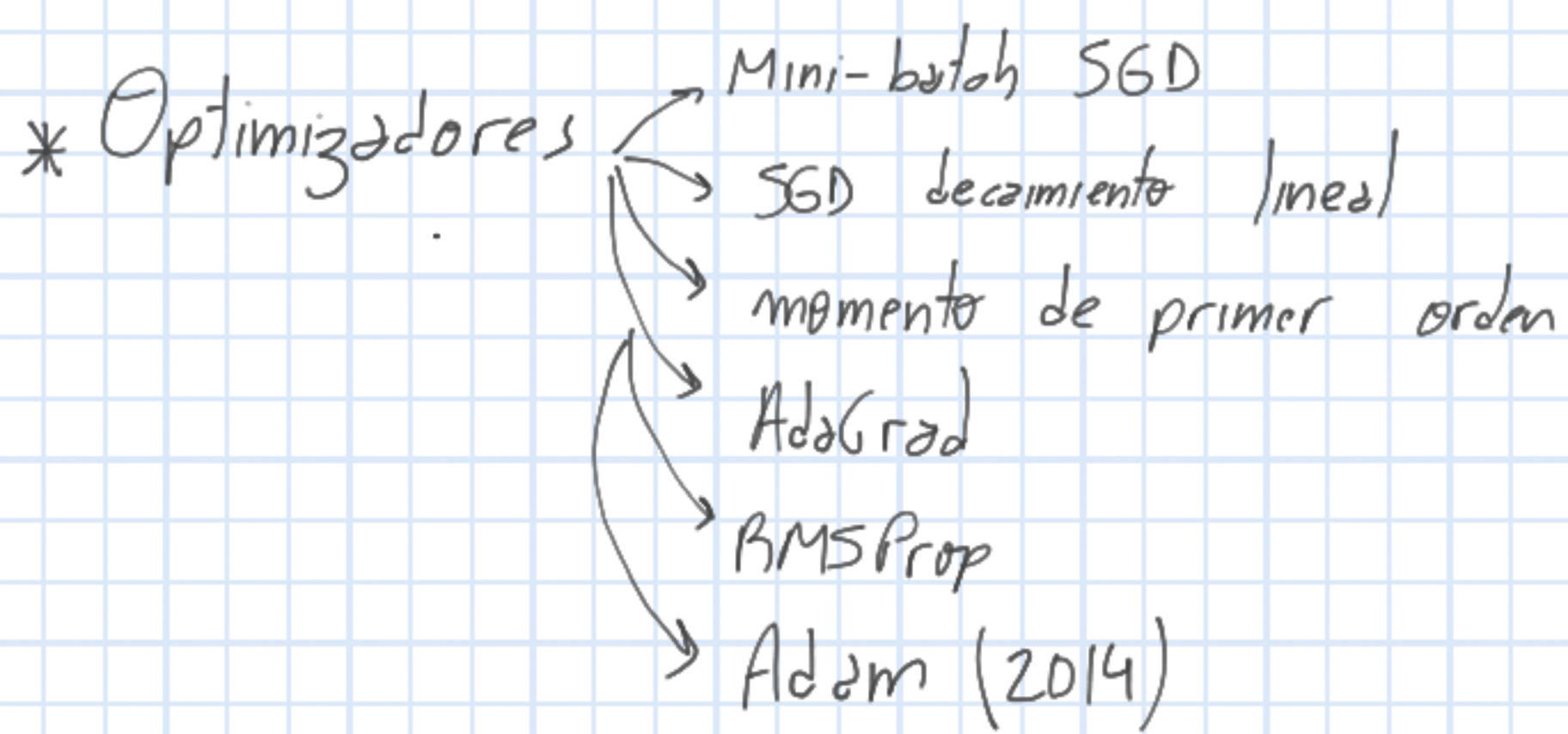
Pesos de la red neuronal
inician o con dist \mathcal{U} o \mathcal{N}

$$W_i \leftarrow \mathcal{U}\left(-\frac{1}{\sqrt{m}}, +\frac{1}{\sqrt{m}}\right)$$

donde m es la cant. entradas

uniforme

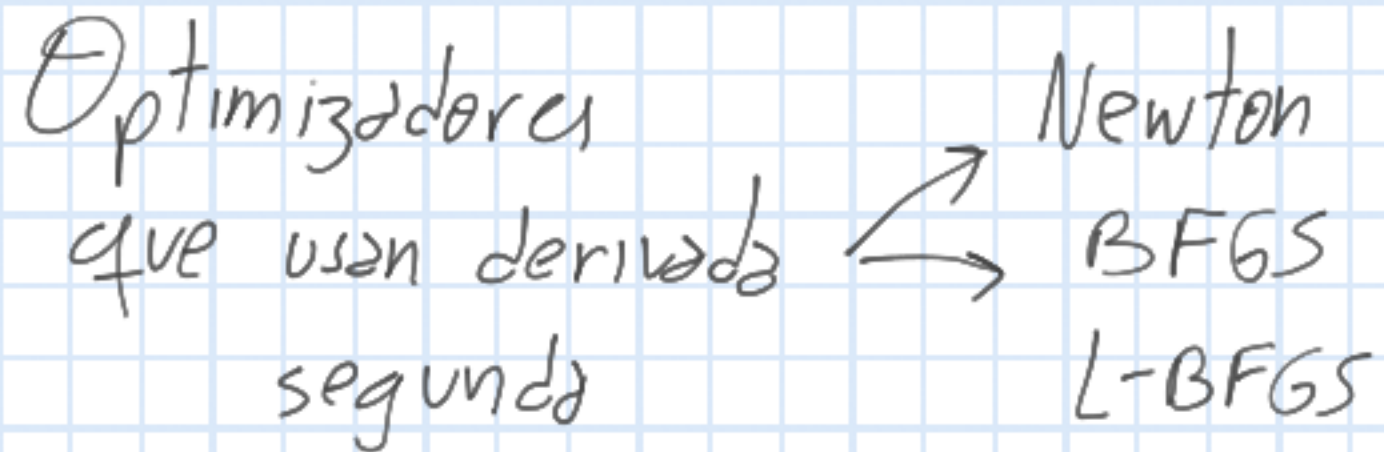




$$W \leftarrow W - \alpha \nabla J$$

siempre usando la derivada primera

distintas estrategias para los momentos



dependen de H

(1) Mini-Batch

For e in epochs:

For b in batch:

* forward...

$$G = \nabla_{\bar{w}} \left(\frac{1}{b} \sum_{i=1}^b \ell(f(\bar{x}_i, \bar{w}), y_i) \right)$$

$$\bar{w} \leftarrow \bar{w} - \alpha \bar{G}$$

es el
gradiente

(2) Mini-batch con decaimiento lineal

for $k=0$
for e in epoch:

for b in batch:

* forward...

$$\alpha_k = \left(1 - \frac{k}{T}\right) \alpha_0 + \frac{k}{T} \alpha_T$$

$$\bar{G} = \nabla_{\bar{w}} \left(\frac{1}{b} \sum \ell(\bar{g}_i, y_i) \right)$$

$$\bar{w} \leftarrow \bar{w} - \alpha_k \bar{G}$$

$$k = k + 1$$

$$k=0 \rightarrow \alpha_k = \alpha_0$$

$$k=T \rightarrow \alpha_k = \alpha_T$$

* First order momentum

for e in epoch:

for b in batch:

forward ...

$$\bar{G} = \nabla \left(\frac{1}{b} \sum_{i=1}^b \ell(\hat{y}_i, y_i) \right)$$

$$\hat{y}_i = f(\bar{x}_i, \bar{w})$$

$$\bar{V} \leftarrow \gamma \bar{V} + \alpha \bar{G}$$

$$\bar{w} \leftarrow \bar{w} - \bar{V}$$

$\gamma = 0 \rightarrow$ vanilla mini-batch

$\gamma \neq 0$

* Momento de segundo orden: AdaGrad

for e in epoch:

for b in batch:

forward

$$\bar{G} = \nabla \left(\frac{1}{b} \sum_{i=1}^b \ell(\hat{y}_i, y_i) \right)$$

$$\bar{F} \leftarrow \bar{F} + \bar{G} \odot \bar{G}$$

HP
(element
wise
multiplication)

$$\bar{G} \odot \bar{G} = \begin{bmatrix} \partial J / \partial w_1 & \partial J / \partial w_1 \\ \vdots & \vdots \\ \partial J / \partial w_n & \partial J / \partial w_n \end{bmatrix}$$

$$\bar{\Delta} \leftarrow - \frac{\alpha}{\sqrt{\bar{F}}} \odot \bar{G}$$

element
wise

$$\bar{W} \leftarrow \bar{W} + \bar{\Delta} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} - \begin{bmatrix} \alpha_1 & \partial L / \partial w_1 \\ \vdots & \vdots \\ \alpha_n & \partial L / \partial w_n \end{bmatrix}$$

BMS Prop

for e in epoch:
for b in batch:

forward

$$\bar{G} = \nabla \left(\frac{1}{b} \sum_{i=1}^b \ell(\hat{y}_i, y_i) \right)$$

$$\bar{\Gamma} = p \bar{\Gamma} + (1-p) \bar{G} \odot \bar{G}$$

$$\bar{\Delta} = -\alpha / \sqrt{\bar{\Gamma}} \odot \bar{G}$$

$$\bar{W} \leftarrow \bar{W} + \bar{\Delta}$$

$p=0 \rightarrow$ no tiene memoria

$p=1 \rightarrow$ usa mucha memoria

HPs



epoch	✓
batch	✓
p	✓
α	✓

Adam (2014)

for e in epoch:

for b in batch:

forward

$$\bar{G} = \nabla \left(\frac{1}{b} \sum_{i=1}^b \ell(\hat{y}_i, y_i) \right)$$

momento
1st order

$$\bar{V} \leftarrow \rho_1 \bar{V} + (1 - \rho_1) \bar{G}$$

HPs

↳

epoch
batch
 ρ_1
 ρ_2
 α

momento
2nd order

$$\bar{F} \leftarrow \rho_2 \bar{F} + (1 - \rho_2) \bar{G} \odot \bar{G}$$

$$\bar{\Delta} \leftarrow -\alpha / \sqrt{\bar{F}} \bar{V}$$

$$\bar{W} \leftarrow \bar{W} + \bar{\Delta}$$

