

# 2. Entrenamiento Regularización Optimización

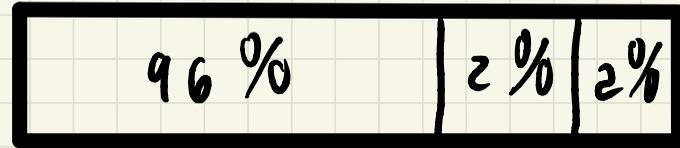
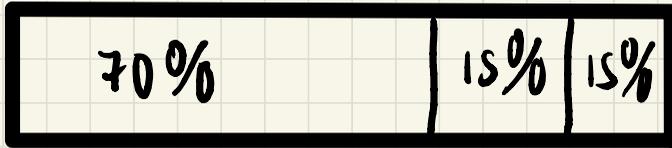
Marzo 2024

# ① ENTRENAMIENTO : DATOS

- LAS MUESTRAS DE DESARROLLO (VALIDACION) y PRUEBA SON cada vez menores:

ADICIONALMENTE

DL



- Como vemos más adelante lo más importante es que la muestra de desarrollo y prueba venga de la misma distrib.

↓  
se necesitan muchos datos para entrenar y error fuera de muestra se estima con 2% de muchos datos

② ENTRENAMIENTO: ANÁLISIS

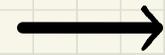
SESGO (APROXIMACION) vs

VARIANZA (ESTIMACION)

ENTRENAMIENTO

SESGO ALTO

PROBLEMAS EN ENTRENAMIENTO



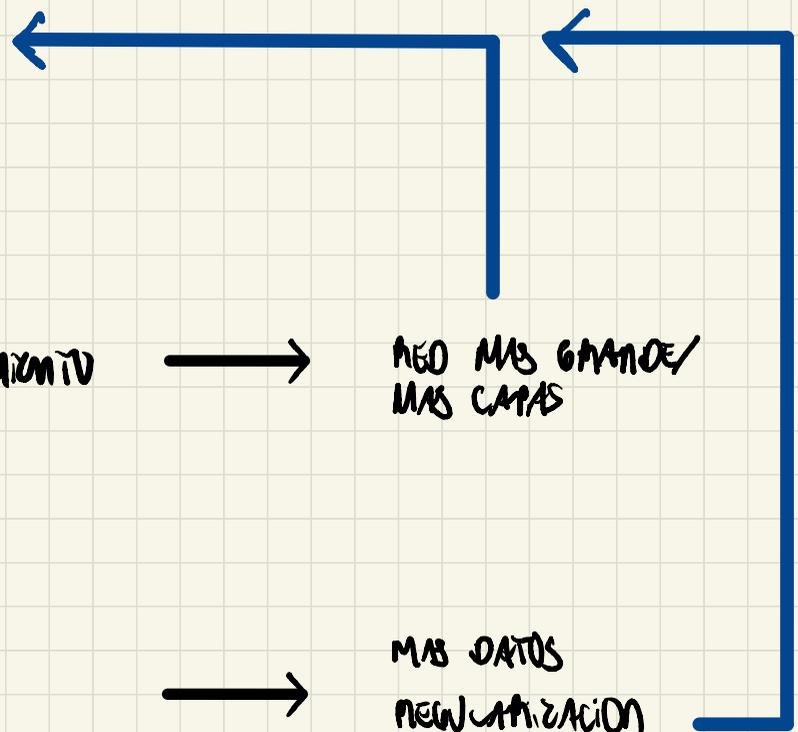
RED MAS GRANDE/  
MAS CAPAS

VARIANZA ALTA

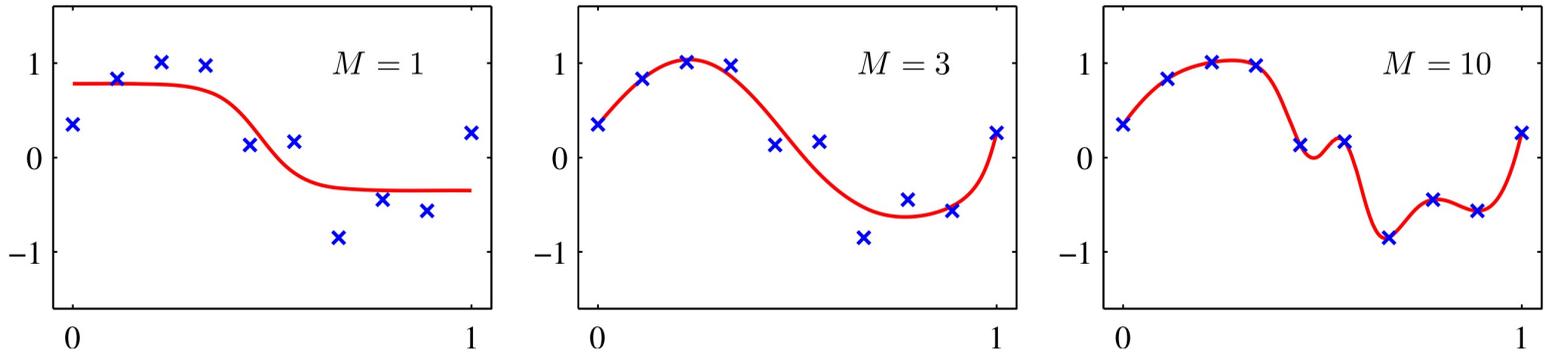
PROBLEMAS EN MUESTRAS  
DEV / TEST



MAS DATOS  
REGULARIZACION  
RED MAS GRANDE

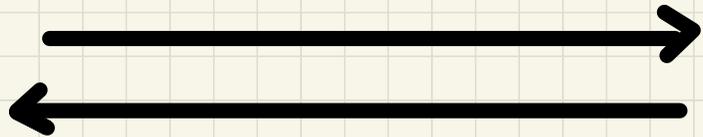


# EJEMPLO : SESGO vs. VARIANZA



**Figure 5.9** Examples of two-layer networks trained on 10 data points drawn from the sinusoidal data set. The graphs show the result of fitting networks having  $M = 1, 3$  and  $10$  hidden units, respectively, by minimizing a sum-of-squares error function using a scaled conjugate-gradient algorithm.

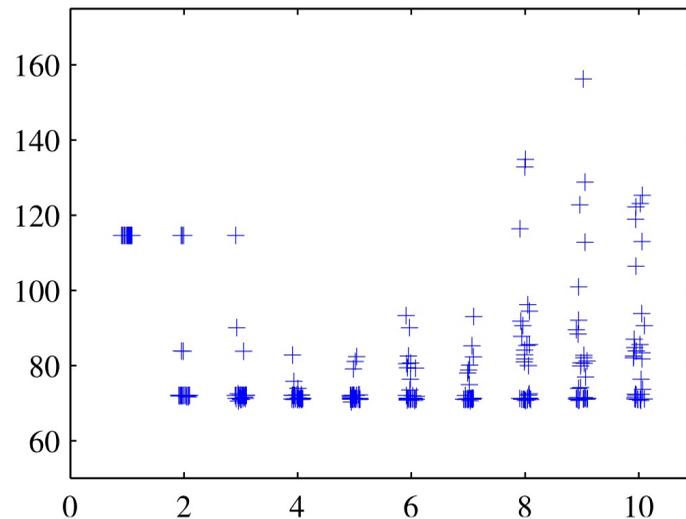
SESGO



complejidad  
varianza

- LA relación entre el número de unidades ocultas  $M$  y el error fuera de muestra es difícil de descubrir por los óptimos locales

**Figure 5.10** Plot of the sum-of-squares test-set error for the polynomial data set versus the number of hidden units in the network, with 30 random starts for each network size, showing the effect of local minima. For each new start, the weight vector was initialized by sampling from an isotropic Gaussian distribution having a mean of zero and a variance of 10.



### ③ ENTWICKLUNG: DATA AUGMENTATION

- rotations / translations / reflections
- zooming
- adding noise

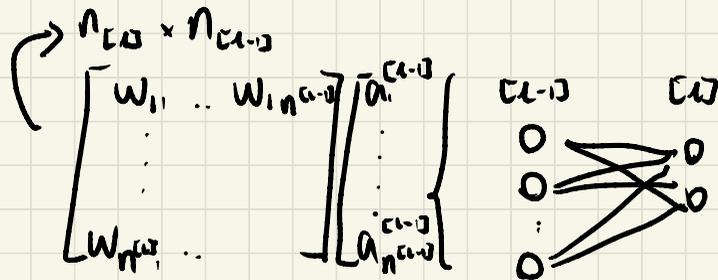
# ① REGULARIZACIÓN: WEIGHT DELAY

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)})$$
$$= \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

donde:

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l-1)}} (w_{ij}^{(l)})^2$$

↓  
lo mismo que  
la norma  $L_2$



- cuando hacemos backpropagation:

$$dW^{[L]} \leftarrow \text{backprop} + \frac{\lambda}{m} W^{[L]}$$

minor

- luego el algoritmo de gradiente descendente:

$$W^{[L]} \leftarrow W^{[L]} - d \, dW_{\text{orig.}} - d \frac{\lambda}{m} W^{[L]}$$

↓  
gradiente  
original sin regularización

$$W^{[L]} \leftarrow \underbrace{\left(1 - \frac{\lambda d}{m}\right)}_{\text{weight decay}} W^{[L]} - d \, dW_{\text{orig.}}$$

weight decay

- Regularización funciona por la razón que ya vimos en el curso (i.e. selección de variable suave)

- Adicionalmente si  $\lambda \uparrow$ ,  $\|w\|_2 \downarrow 0$

$\Rightarrow \|z\| \downarrow 0$  y  $g(z)$  es casi una función lineal (e.g.  $y = \tanh$  o  $\text{ReLU}$ ) que es una función muy simple.

REGULARIZACIÓN: DESDE UN PUNTO DE VISTA BAYESIANO

- consideramos un modelo lineal de funciones base:

$$y(x, w) = w^T \phi(x)$$

donde:

$$\phi(x) = (1, \phi_1(x), \dots, \phi_{n-1}(x))^T$$

- las funciones  $\phi_i$  se llaman funciones base.
  - supongamos ahora que  $\hat{y} = y(x, w) + \epsilon$  donde  $\epsilon \sim N(0, \sigma^{-1})$
  - luego la distribución de probabilidad:
- $$p(\hat{y} | x, w, \beta) = N(\hat{y} | y(x, w), \beta^{-1})$$

- Suponiendo independencia condicional de los observadores la log verosimilitud es:

$$\frac{1}{2} \ln(\rho) - \frac{1}{2} \ln(2\pi) - \rho \sum_{i=1}^n (y_i - \underbrace{w^T \phi(x_i)}_{\hat{y}_i})^2$$

y replicamos el resultado que maximizar la verosimilitud es lo mismo que minimizar el error cuadrático

- introducción de prior:  $\rho \rightarrow \rho$  conocido

- verosimilitud:  $p(\hat{y} | x, w, \rho) = N(\hat{y} | y(x, w), \rho^{-1})$

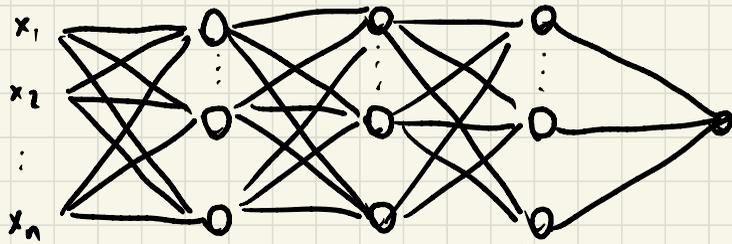
- prior conjugada:  $p(w | d) = N(w | 0, d^{-1} I)$

- LA posterior también es normal:  $\hookrightarrow$  por simplicidad centrado en cero

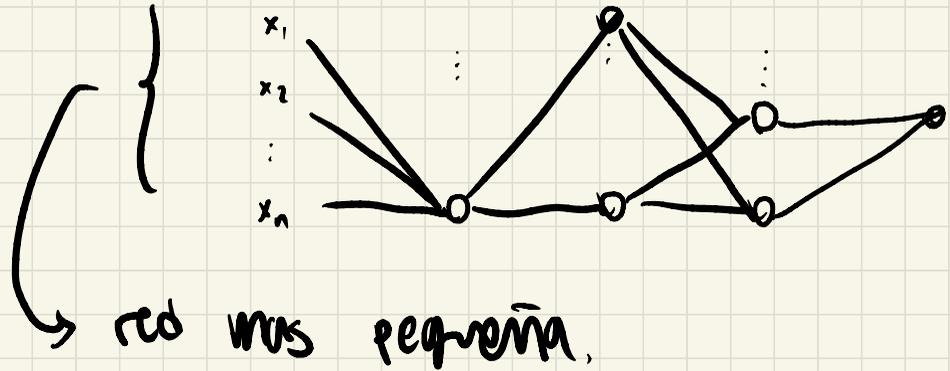
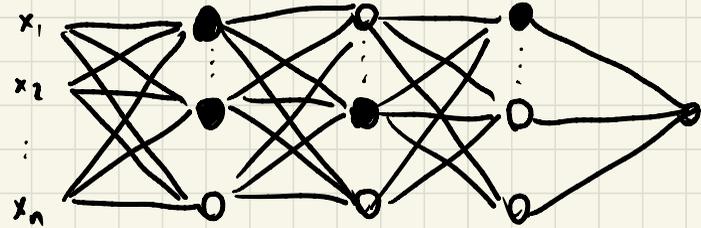
(i.e. posterior de todos los datos observados)

$$\log p(w | y) = -\frac{\rho}{2} \sum_{i=1}^m (y_i - w^T \phi(x_i))^2 - \frac{d}{2} w^T w + \text{const.}$$

## ② REGULARIZATION : DROPOUT



↓  
Keep neuron  
with prob  $p$



- en el momento de implementar es importante rescaladar las funciones de activacion porq̃e:

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$



si esta se reduce en  $p$  entonces  $z^{[L]}$  tambien se reduce por un factor de  $p$ . LA IDEA DE invariant dropout es rescaladar asi:

$$z^{[L]} = \frac{w^{[L]} a^{[L-1]}}{p} + b^{[L]}$$

ESTO PARA DE MANTENER EL VALOR ESPERADO DE  $z^{[L]}$  constante

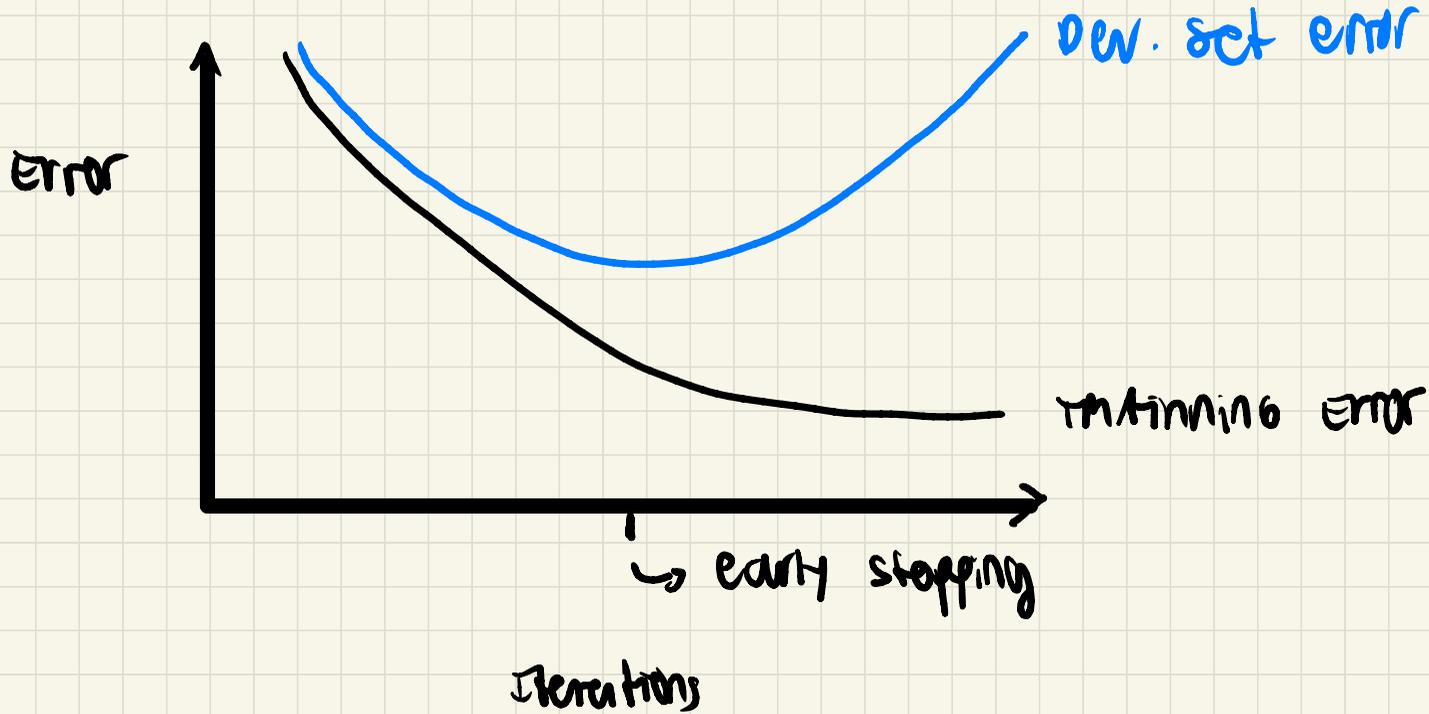
## Observaciones

- Dropout distribuye la importancia de las neuronas impidiendo que un peso tome demasiada importancia
- en principio se puede usar una probabilidad dropout en cada capa (e.g. si hay muchos pesos en una capa)
- se puede hacer incluso en la capa de entrada (no muy común)
- cuando se hacen predicciones no se aleatoriza.

- un problema es que la función de costo no es decreciente durante el entrenamiento. La función no está bien definida. Luego se puede verificar el código primero sin dropout y después se regulariza.

③

# EARLY STOPPING



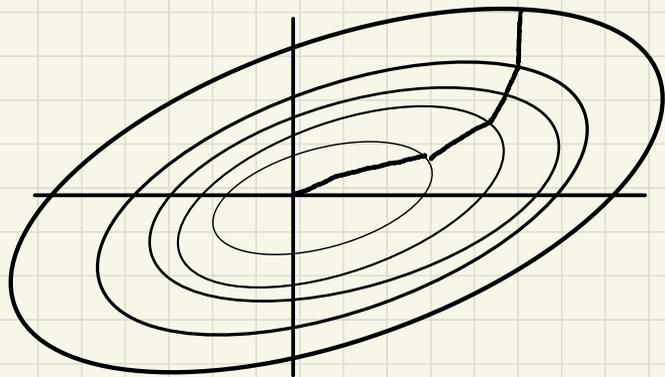
- un problema con esta estrategia es una mezcla de los problemas:  $\left. \begin{array}{l} \text{optimización de } J(w, b) \\ \text{Avoid overfitting.} \end{array} \right\}$

cuando se hace early stopping se detra de optimizar buscando reducir overfitting. separar los dos problemas de manera como orthogonalización.

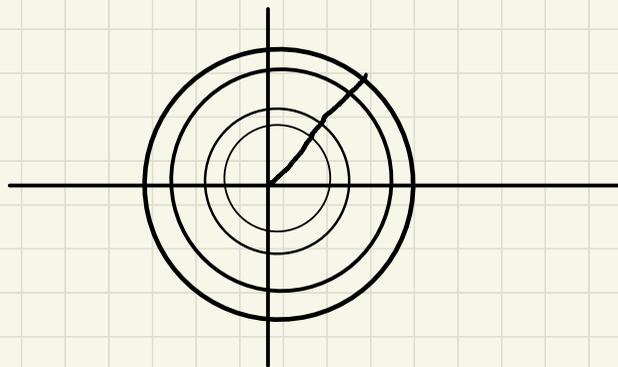
- una mejor alternativa pero más costosa computacionalmente es hacer una única mezcla o usar el dev set para elegir hiperparámetros. Por ejemplo de regularización  $L^2$ .

# ① OPTIMIZACIÓN : NORMALIZACIÓN, IMP/EXPOSION, GRADO CHECKING

- si se normalizan los datos gradiente descendiente funciona mejor
- $f: \mathbb{R}^3 \rightarrow \mathbb{R}$



SIN NORMALIZAM



NORMALIZADO

## - IMPULSION \ EXPUSION DE GRADIENTES

CONSIDERE ESTA RED:



$g(z) = z$ ,  $b = 0$  en todos los capas

$$\Rightarrow \hat{y} = w^{(1)} \dots w^{(L)} x$$

dependiendo de  $w_{ij}^{(L)} \gtrless 1$  por muy poco, puede crecer o decrecer exponencialmente.

- una forma de corregir esto es minimizando mejor los pesos.

$$\text{como } z = w_1 x_1 + w_2 x_2 \dots + w_n x_n + b$$

entre unos elementos en la suma potencialmente más grande  $z$ .

Entonces se puede normalizar  $w$  de tal forma que la suma se "controle":  $\text{var}(w) = \frac{1}{n}$ .

EN LA PRÁCTICA LO MEJOR ES  $\text{var}(w) = \sqrt{\frac{2}{n}}$   
LUEGO SE IMPLEMENTA COMO:

$$w^{[k]} = \text{random} * \sqrt{\frac{2}{n^{[k-1]}}}$$

Esta es la forma más común de hacerlo cuando las funciones de activación son relu. Para otros funciones se usan otras normalizaciones.

## Gradient checking

tomar  $w^{L1}, \dots, w^{L2}, b^{L1}, \dots, b^{L2}$  y poner como un solo parámetro  $\theta$ .

Escribir el costo como función de  $\theta : J(\theta_1, \dots, \theta_n)$

verificar:

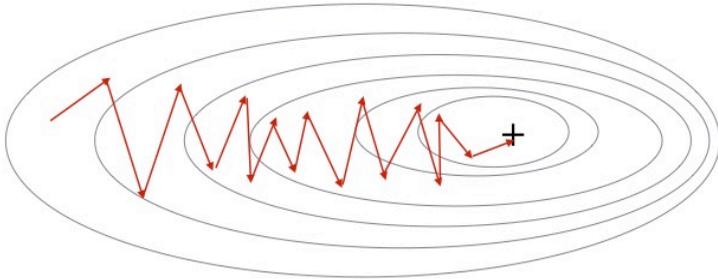
$$d\theta[i] \cong d\theta_{\text{approx}} = \text{derivada numérica usando la secante.}$$

por ejemplo:

$$\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} < \epsilon$$

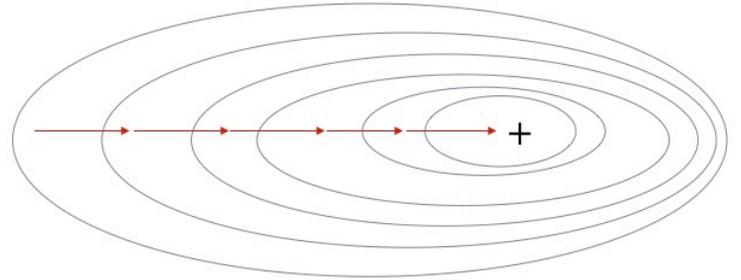
## ② GRADIENT DESCENT : STOCHASTIC vs. BATCH

Stochastic Gradient Descent



ACTUALIZACION RÁPIDA  
DE PARÁMETROS

Gradient Descent



ACTUALIZACION LENTA  
DE PARÁMETROS  
NO APROVECHA VELOCIZACIÓN

# ③ BATCH vs Mini BATCH GRADIENT DESCENT

$X =$	64 training examples	64 training examples	64 training examples	...	...	...	64 training examples	<64 training examples
$Y =$	64 training examples	64 training examples	64 training examples	...	...	...	64 training examples	<64 training examples
	$\underbrace{\hspace{10em}}$ mini_batch 1    mini_batch 2    mini_batch 3			...	...	...	$\underbrace{\hspace{10em}}$ mini_batch $\lfloor m/64 \rfloor$	$\underbrace{\hspace{10em}}$ mini_batch $\lfloor m/64 \rfloor + 1$

$$X = \begin{bmatrix} X^{(1)} & & & X^{(m)} \end{bmatrix}$$

$\underbrace{\hspace{15em}}_{X^{(1)}} \quad \underbrace{\hspace{15em}}_{X^{(2)}} \quad \dots \quad \underbrace{\hspace{15em}}_{X^{(m)}}$

$$Y = \begin{bmatrix} y^{(1)} & & & y^{(m)} \end{bmatrix}$$

$\underbrace{\hspace{15em}}_{y^{(1)}} \quad \dots \quad \underbrace{\hspace{15em}}_{y^{(m)}}$

MÉTODO DE PARTICIONES

$$X = \begin{pmatrix} x_0^{(1)} & x_0^{(2)} & \dots & x_0^{(m-1)} & x_0^{(m)} \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m-1)} & x_1^{(m)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{12286}^{(1)} & x_{12286}^{(2)} & \dots & x_{12286}^{(m-1)} & x_{12286}^{(m)} \\ x_{12287}^{(1)} & x_{12287}^{(2)} & \dots & x_{12287}^{(m-1)} & x_{12287}^{(m)} \end{pmatrix}$$

$$X = \begin{pmatrix} x_0^{(5)} & x_0^{(16)} & \dots & x_0^{(2)} & x_0^{(m-1)} \\ x_1^{(5)} & x_1^{(16)} & \dots & x_1^{(2)} & x_1^{(m-1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{12286}^{(5)} & x_{12286}^{(16)} & \dots & x_{12286}^{(2)} & x_{12286}^{(m-1)} \\ x_{12287}^{(5)} & x_{12287}^{(16)} & \dots & x_{12287}^{(2)} & x_{12287}^{(m-1)} \end{pmatrix}$$

$$Y = \begin{pmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m-1)} & y^{(m)} \end{pmatrix}$$

$$Y = \begin{pmatrix} y^{(5)} & y^{(16)} & \dots & y^{(2)} & y^{(m-1)} \end{pmatrix}$$

METODO DE  
ALFA TOMIZACION

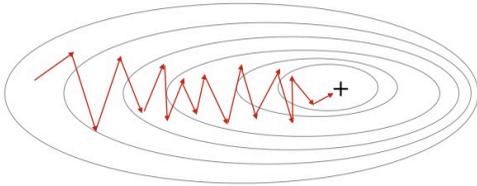
- en cada miniBATCH se estima todo: Forward, backward, costo, etc y se calcula nuevos w, b. repitiendo este proceso con todos los minibatches se actualiza repetidamente los w, b.

- luego en contraste a batch gradient descent una sola pasada por todos los datos permite actualizar una sola vez los parámetros
- con mini batch gradient descent cuando se hace una pasada completa de los datos (llamado época) se actualizan varias veces los parámetros.
- al hacer esto con varias épocas se llega mucho más rápido a los parámetros en minimización el costo.

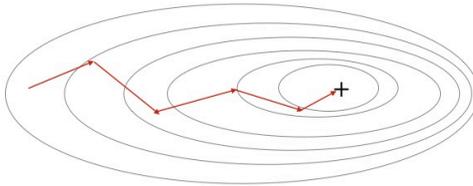
MINIBATCHES  
TAMAÑO 1



Stochastic Gradient Descent



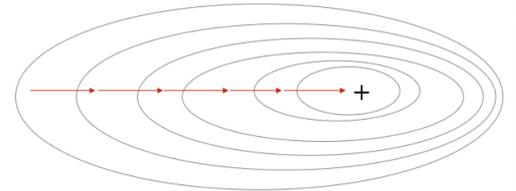
Mini-Batch Gradient Descent



BATCH



Gradient Descent





#### ④ exponential weighted gradient descent

- observation  $\{\theta_k\}$
- exponentially weighted average  $v_t = \beta v_{t-1} + (1 - \beta)\theta_k$
- correction de sesgo: al comienzo  $v_t$  puede estar sistemáticamente por debajo de  $\theta_k$  ( $\theta_k \geq 0$ )  
y para corregir esto en primeras iteraciones:

$$v_t \rightarrow \frac{v_t}{1 - \beta^t}$$

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t$$

$$v_1 = \beta v_0 + (1-\beta) \theta_1$$

$$v_2 = \beta^2 v_0 + \beta(1-\beta) \theta_1 + (1-\beta) \theta_2$$

$$v_3 = \cancel{\beta^3 v_0} + \beta^2 (1-\beta) \theta_1 + \beta(1-\beta) \theta_2 + (1-\beta) \theta_3, \quad v_0 = 0$$

$$v_t = \left[ \sum_{i=1}^t \beta^{t-i} \theta_i \right] (1-\beta)$$

$$v_t = \theta \left[ \sum_{i=1}^t \beta^{t-i} \right] (1-\beta), \quad \text{si } \theta_i = \theta$$

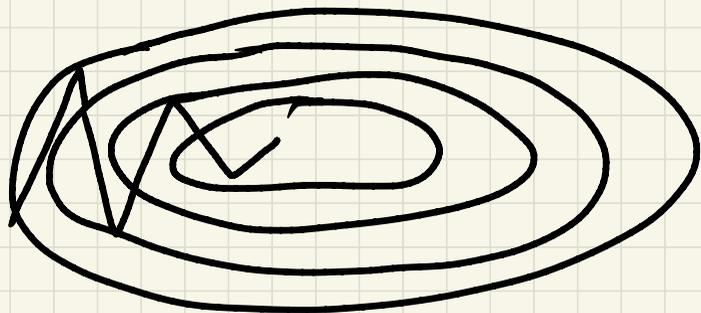
$$= \theta \left[ \frac{1 - \beta^t}{1 - \beta} \right] (1-\beta)$$

$$= \theta (1 - \beta^t)$$

NOTA : TODO LOS SIGUIENTES ALGORITMOS SE  
IMPLEMENTAN EN CADA MINI BATCH.

⑤

## Gradient descent with momentum.



- los movimientos en  $y$  no son eficientes

- en el eje  $x$  si son más eficientes.

- el método de momentum consiste en suavizar los gradientes en cada iteración utilizando el método anterior de "pesos exponenciales"

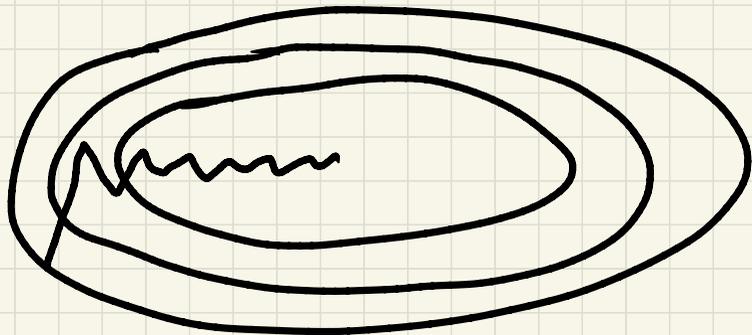
la idea es disminuir  $y$  y aumentar  $x$ .

$$V_{dw} = \beta V_{dw} + (1-\beta)dW$$

$$V_{db} = \beta V_{db} + (1-\beta)db$$

se puede incluir  
el ajuste por sesgo.

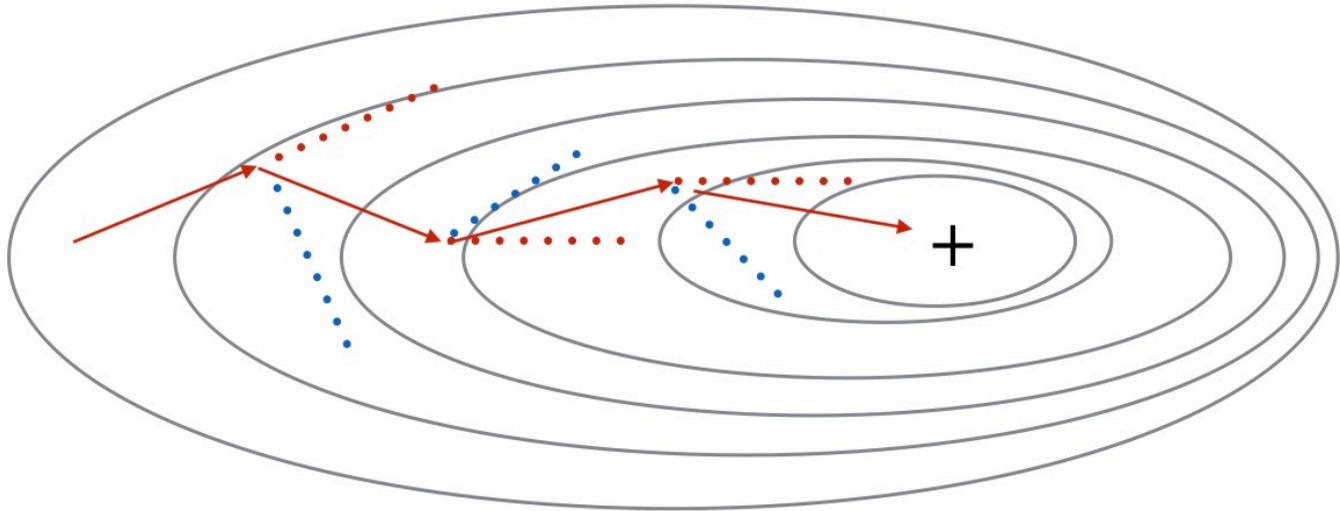
¿por que funciona? Disminuye el recorrido en el etc y.



$$W = W - dV_{dw}$$

$$b = b - dV_{db}$$

# Momentum



## ⑤ PMS POP

- PARA DISMINUIR LAS OSCILACIONES en  $Y$  y  
AUMENTARLAS en  $X$ :

$$- S_{dw} = \beta S_{dw} + (1-\beta) dW * dW$$

$$S_{db} = \beta S_{db} + (1-\beta) db * db$$

$$- w = w - d \frac{dW}{\sqrt{S_{dw}}}$$

$$- b = b - d \frac{db}{\sqrt{S_{db}}}$$

} PARA EVITAR DIVIDIR POR  
CERO USUALMENTE SE

$$\text{MODIFICA } S_{dw} \leftarrow S_{dw} + \epsilon$$

$$S_{db} \leftarrow S_{db} + \epsilon$$

## ⑥ ADAM OPTIMIZATION ALGORITHM (ADAPTIVE MOMENT ESTIMATION)

ESTE ALGORITMO COMBINA MOMENTUM Y RMS PROP

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

$\left\{ \begin{array}{l} V_{dw}, V_{db} \text{ actualizar con momentum y parámetro } \beta_1 \\ S_{dw}, S_{db} \text{ " " RMS prop y " } \beta_2 \end{array} \right.$

→ hacer corrección de sesgos de ambos

Actualización final:

$$w \leftarrow w - \alpha \frac{V_{dw}}{\sqrt{S_{dw}}}, \quad b \leftarrow b - \alpha \frac{V_{db}}{\sqrt{S_{db}}}$$

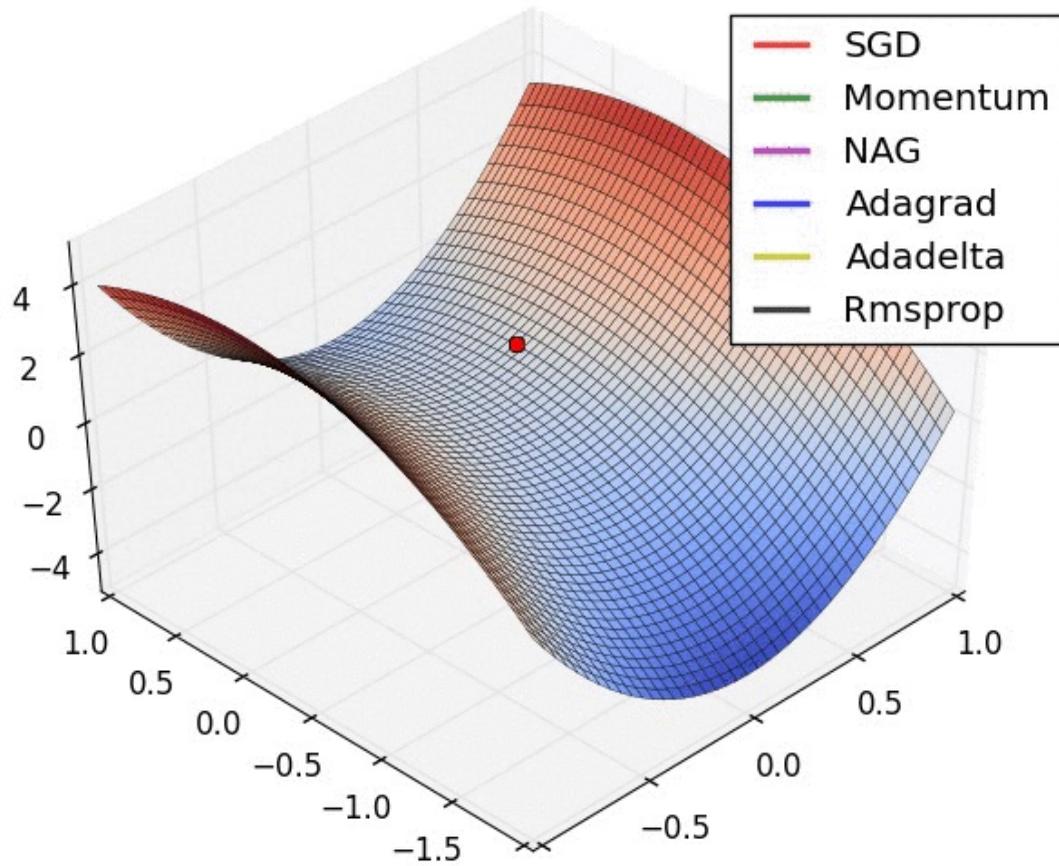
Iti per poi retro

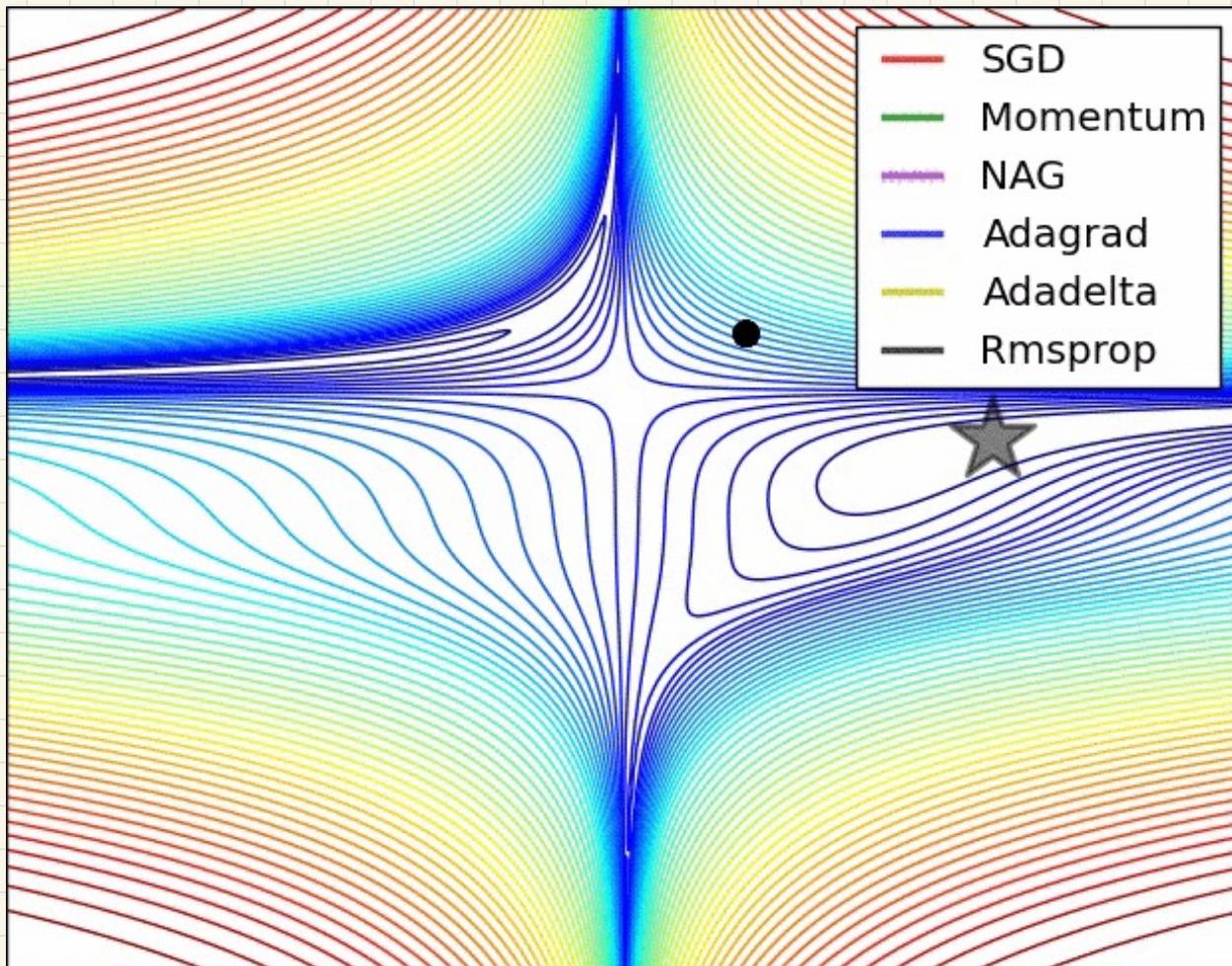
d long que equilibrio

$$\beta_1 \sim 0.9$$

$$\beta_2 \sim 0.999$$

$$\epsilon \sim 10^{-8}$$





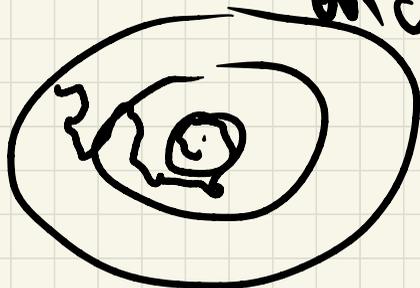
## ⑦ LEARNING RATE DECAY

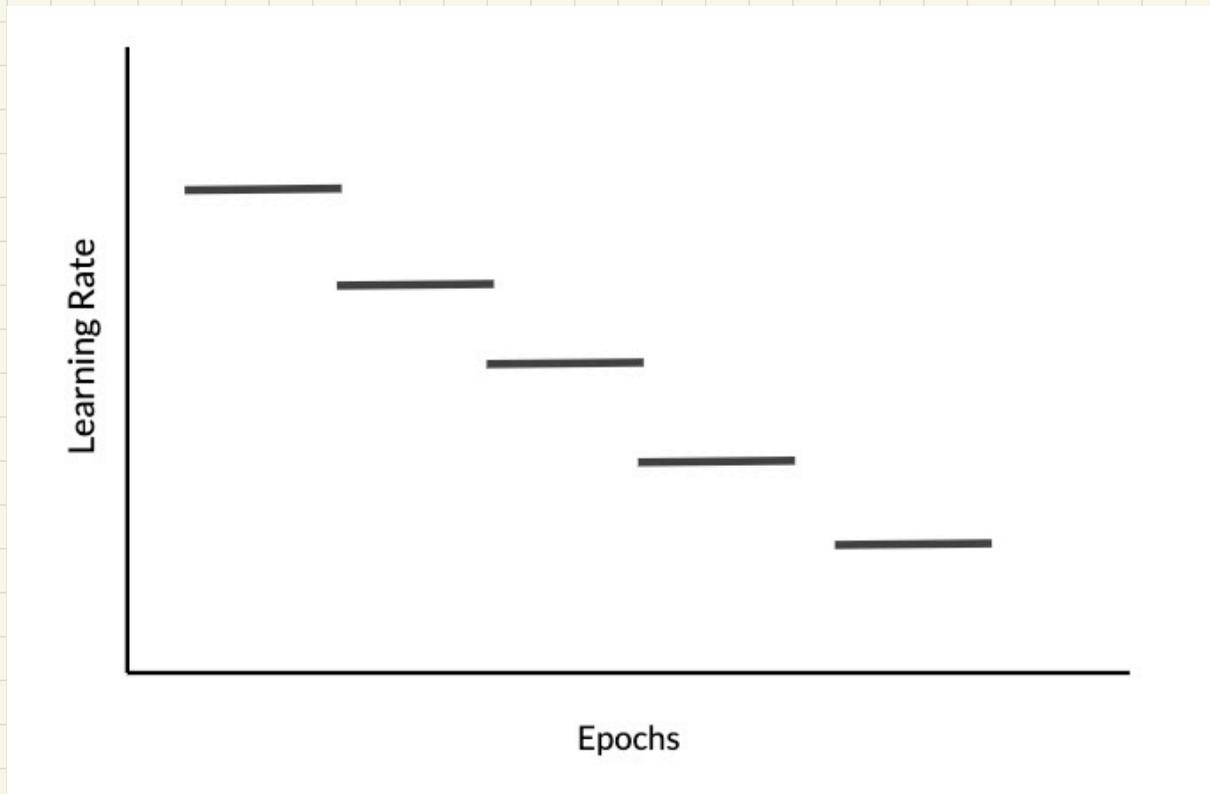
- Si  $d$  es grande se pueden dar saltos ineficientes cerca del óptimo
- una forma de irlo ajustando es:

$$d_t = \frac{1}{1 + \text{decay\_rate} \cdot \text{epoch\_number}} \times d_0$$

↓  
hiperparámetro

puede dar  
muchos saltos  
en cada mini  
batch





# ① CALIBRACION DE HIPERPARAMETROS

$\alpha$ ,  $\beta$ ,  $\beta_1$ ,  $\beta_2$ ,  $\epsilon$ ,  $L$ ,  $\eta$ , learning rate decay, minibatch size

1      2                      2      3                      3                      2

- RANKING DE MAS IMPORTANTES: 1, 2, 3

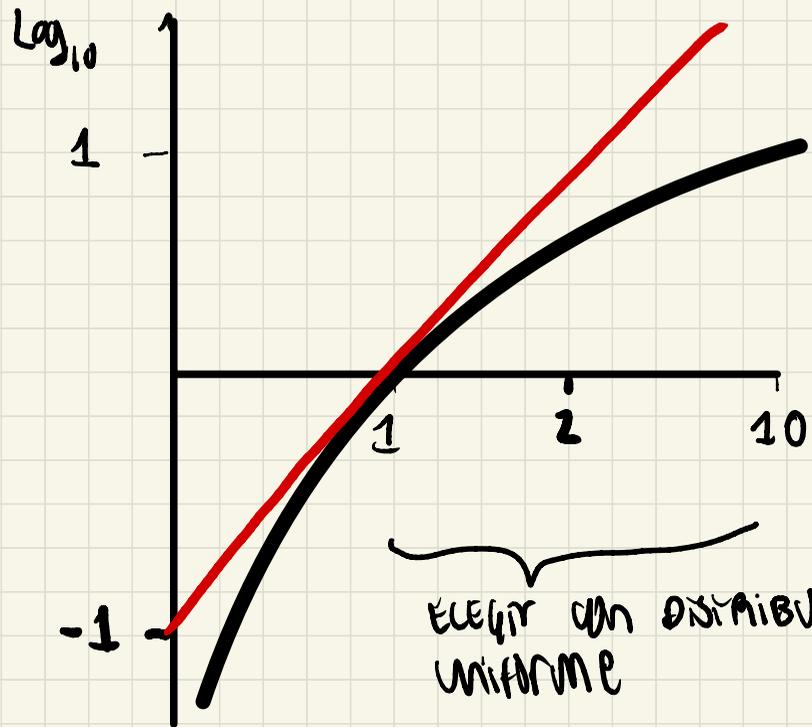
- ESTRATEGIAS CALIBRACION:

1 ELEGIR ALEATORIAMENTE

2 DE GUESO A FINO

3 ESCALA APROPIADA:

Si  $d \in [0, 1]$  por prob. mas cercano a cero  
MUESTREAR USANDO ESCALA LOGARITMICA.



ELEGIR una DISTRIBUCION UNIFORME

## ② BATCH NORMALIZATION

- ¿cómo normalizar en pasos intermedios para facilitar el aprendizaje de los parámetros del paso siguiente?

- Supongamos que estamos en capa  $l$

Definir  $\mu^{(l)} = \frac{1}{m} \sum_{i=1}^m z^{(l,i)}$

$$S^{(l)} = \frac{1}{m} \sum_{i=1}^m (z^{(l,i)} - \mu^{(l)})^2$$

$$z_{\text{norm}}^{(l,i)} = \frac{z^{(l,i)} - \mu^{(l)}}{\sqrt{S^{(l)} + \epsilon}}$$

- BATCH NORMALIZATION  $z^{(l,i)} = \gamma z_{\text{norm}}^{(l,i)} + \beta$

↓  
PARÁMETROS PARA APRENDER

- Observese que si:  $\left. \begin{array}{l} \gamma = \sqrt{\sigma^2 + \epsilon} \\ \rho = \mu \end{array} \right\} \Rightarrow \tilde{z}^{(i,j)} = z_{norm}^{(i,j)}$

Weyo poder calibrar  $\gamma, \rho$  permite controlar la media y varianza de las  $z$  normalizadas y que no sean cero y uno.

- PARA INCORPORAR en el entrenamiento de la red con cada mini-batch se puede hacer esto:

$$X^{(i,j)} \xrightarrow{w^{(i,j)}} z^{(i,j)} \xrightarrow{\gamma^{(i,j)} \rho^{(i,j)}} \tilde{z}^{(i,j)} \rightarrow g^{(i,j)}(\tilde{z}^{(i,j)}) \xrightarrow{w^{(i,j)}} z^{(i,j)} \rightarrow \dots$$

observese que no es necesario incluir el sesgo porque cuando se normaliza se elimina

## - Implementación de Gradient Descent

Loop  $t=1 \dots$  num iteraciones

datos  $X^{(t)}$  forward propagation

en cada capa reemplazar  $z^{(k)}$  por  $\hat{z}^{(k)}$

hacer Backprop:  $dW^{(k)}$ ,  ~~$dZ^{(k)}$~~ ,  $d\beta^{(k)}$ ,  $d\gamma^{(k)}$

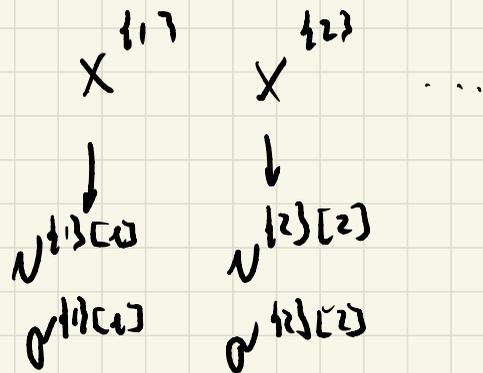
Actualizar parámetros de la forma estándar

end loop

- se puede combinar con momentum, rmsprop, Adam.

- ¿cuando se va a hacer la predicción con un ejemplo que  $N, \hat{N}$  se usa?

- si el modelo se entreno con min batches entonces todos tienen el mismo tamaño:



→  $N$  se calcula usando exponentially weighted average

→  $\hat{N}$

- este  $N, \hat{N}$  se usa en el momento de predecir un ejemplo:  $Z_{norm}$  se calcula usando este  $N, \hat{N}$  y con esto  $\hat{z}$  y se calcula toda la red.

## ① MULTICLASES

- sea  $C = \#$  clases

-  $n^{cls} = 4$

- función de pérdida:  $\ell(\hat{y}, y) = - \sum_{i=1}^C y_i \ln(\hat{y}_i)$

- función de costo:  $J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})$

- cuando no comencemos el backprop:  $\frac{\partial J}{\partial z^{cls}} = \delta z^{cls} = \hat{y} - y$  que es

- en vez de hacer a mano el resto la idea es usar un

framework