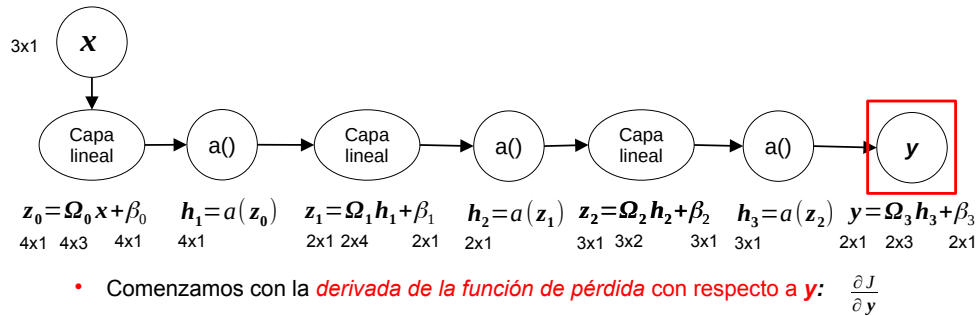


## 5.2 Optimización con redes profundas

- 5.2.1 Backpropagation
- 5.2.2 Inicialización de parámetros
- 5.2.3 Regularización

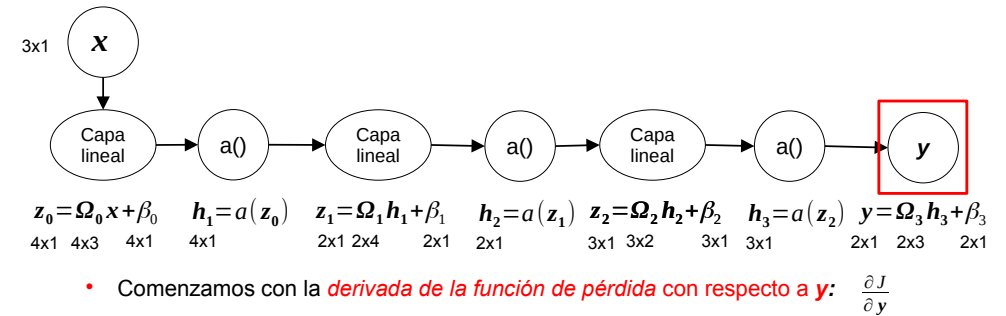
### Ejemplo



## 5.2 Optimización con redes profundas

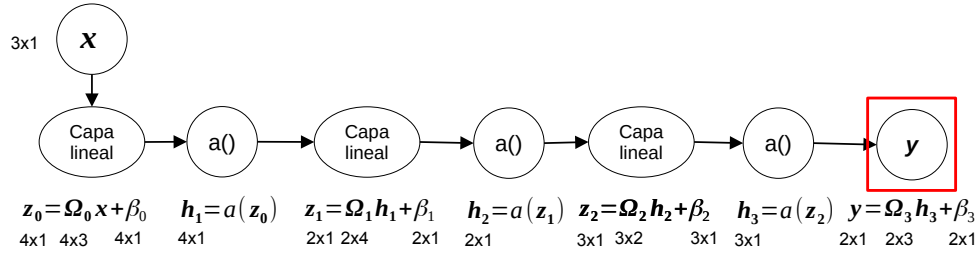
- 5.2.1 Backpropagation
- 5.2.2 Inicialización de parámetros
- 5.2.3 Regularización

### Ejemplo



$$\frac{\partial z_k}{\partial h_k} = \Omega_k^T \quad \frac{\partial a(z_k)}{\partial z_k} = \begin{bmatrix} I(z_{k,1} > 0) & 0 & 0 \\ 0 & I(z_{k,2} > 0) & 0 \\ 0 & 0 & I(z_{k,3} > 0) \end{bmatrix} \quad \frac{\partial z_k}{\partial \Omega_k} = h_k^T \quad \frac{\partial z_k}{\partial \beta_k} = I$$

## Ejemplo



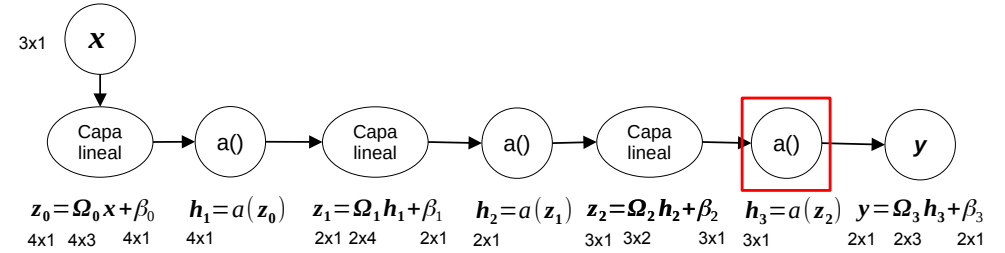
- Comenzamos con la **derivada de la función de pérdida con respecto a y**:  $\frac{\partial J}{\partial y}$
- Y vamos hacia atrás en por la red:

$$\frac{\partial L}{\partial h_3} = \frac{\partial y}{\partial h_3} \cdot \frac{\partial L}{\partial y} = \Omega_3^T \cdot \frac{\partial L}{\partial y} \quad \frac{\partial L}{\partial \beta_3} = \frac{\partial y}{\partial \beta_3} \cdot \frac{\partial L}{\partial y} = \frac{\partial L}{\partial y}$$

Derivadas de la función de pérdida con respecto a los parámetros

$$\frac{\partial z_k}{\partial h_k} = \Omega_k^T \quad \frac{\partial a(z_k)}{\partial z_k} = \begin{bmatrix} I(z_{k,1} > 0) & 0 & 0 \\ 0 & I(z_{k,2} > 0) & 0 \\ 0 & 0 & I(z_{k,3} > 0) \end{bmatrix} \quad \frac{\partial z_k}{\partial \Omega_k} = h_k^T \quad \frac{\partial z_k}{\partial \beta_k} = I$$

## Ejemplo

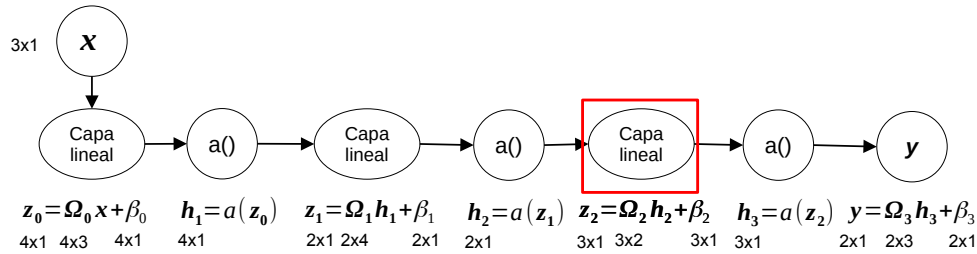


- Comenzamos con la **derivada de la función de pérdida con respecto a y**:  $\frac{\partial J}{\partial y}$
- Y vamos hacia atrás en por la red:

$$\frac{\partial L}{\partial z_2} = \frac{\partial h_3}{\partial z_2} \cdot \frac{\partial L}{\partial h_3} = \begin{bmatrix} I(z_{3,1} > 0) & 0 & 0 \\ 0 & I(z_{3,2} > 0) & 0 \\ 0 & 0 & I(z_{3,3} > 0) \end{bmatrix} \cdot \frac{\partial L}{\partial h_3}$$

$$\frac{\partial z_k}{\partial h_k} = \Omega_k^T \quad \frac{\partial a(z_k)}{\partial z_k} = \begin{bmatrix} I(z_{k,1} > 0) & 0 & 0 \\ 0 & I(z_{k,2} > 0) & 0 \\ 0 & 0 & I(z_{k,3} > 0) \end{bmatrix} \quad \frac{\partial z_k}{\partial \Omega_k} = h_k^T \quad \frac{\partial z_k}{\partial \beta_k} = I$$

## Ejemplo



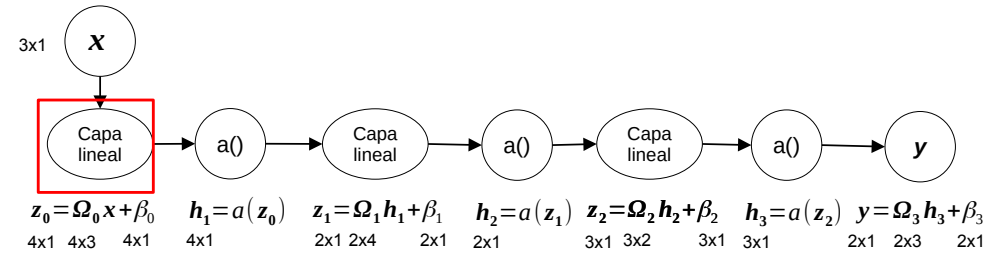
- Comenzamos con la **derivada de la función de pérdida con respecto a y**:  $\frac{\partial J}{\partial y}$
- Y vamos hacia atrás en por la red:

$$\frac{\partial L}{\partial h_2} = \frac{\partial z_2}{\partial h_2} \cdot \frac{\partial L}{\partial z_2} = \Omega_2^T \cdot \frac{\partial L}{\partial z_2} \quad \frac{\partial L}{\partial \beta_2} = \frac{\partial z_2}{\partial \beta_2} \cdot \frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial z_2}$$

Derivadas de la función de pérdida con respecto a los parámetros

$$\frac{\partial z_k}{\partial h_k} = \Omega_k^T \quad \frac{\partial a(z_k)}{\partial z_k} = \begin{bmatrix} I(z_{k,1} > 0) & 0 & 0 \\ 0 & I(z_{k,2} > 0) & 0 \\ 0 & 0 & I(z_{k,3} > 0) \end{bmatrix} \quad \frac{\partial z_k}{\partial \Omega_k} = h_k^T \quad \frac{\partial z_k}{\partial \beta_k} = I$$

## Ejemplo



- Comenzamos con la **derivada de la función de pérdida con respecto a y**:  $\frac{\partial J}{\partial y}$
- Y vamos hacia atrás en por la red:

$$\frac{\partial L}{\partial x} = \frac{\partial z_0}{\partial x} \cdot \frac{\partial L}{\partial z_0} = \Omega_0^T \cdot \frac{\partial L}{\partial z_0} \quad \frac{\partial L}{\partial \beta_0} = \frac{\partial z_0}{\partial \beta_0} \cdot \frac{\partial L}{\partial z_0} = \frac{\partial L}{\partial z_0}$$

Derivadas de la función de pérdida con respecto a los parámetros

$$\frac{\partial z_k}{\partial h_k} = \Omega_k^T \quad \frac{\partial a(z_k)}{\partial z_k} = \begin{bmatrix} I(z_{k,1} > 0) & 0 & 0 \\ 0 & I(z_{k,2} > 0) & 0 \\ 0 & 0 & I(z_{k,3} > 0) \end{bmatrix} \quad \frac{\partial z_k}{\partial \Omega_k} = h_k^T \quad \frac{\partial z_k}{\partial \beta_k} = I$$

## 5.2 Optimización con redes profundas

- 5.2.1 Backpropagation
- 5.2.2 Inicialización de parámetros
- 5.2.3 Normalización
- 5.2.4 Regularización

### 5.2.1 Inicialización de parámetros

- Necesidad de la inicialización
- Inicialización de Xavier y de He

### Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- ¿Cómo inicializamos los sesgos y los pesos?
- Equivalente a elegir un punto de arranque de la optimización en el regresor lineal/Gabor que hemos visto en el Tema 2.

### Inicialización: Pregunta 1

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- ¿Podemos inicializarlos todos a 0?

## Inicialización: Pregunta 2

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Pensemos en una red con 2 unidades ocultas. Si entrenamos la red hasta el final e intercambiamos los pesos de entrada y salida entre las dos unidades ¡Tendremos el mismo resultado! (simetría de la red).
  - ¿Qué pasaría si inicializamos los pesos de las dos unidades ocultas con los mismos valores?

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$

preactivación i-ésima en esa capa:

$$\mathbf{z}_k(i) = \sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j) + \cancel{\boldsymbol{\beta}_k(i)}$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

preactivación i-ésima en esa capa:

$$\mathbf{z}_k(i) = \sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j) + \boldsymbol{\beta}_k(i)$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$

preactivación i-ésima en esa capa:

$$\mathbf{z}_k(i) = \sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j)$$

## Repaso esperanza matemática: $E[\ ]$

$$E[g[x]] = \int_{-\infty}^{+\infty} g[x] Pr(x) dx$$

- **Interpretación:** es el valor promedio de  $g[x]$  cuando tenemos en cuenta la probabilidad de  $x$
- En el caso discreto tenemos  $g[x]$  para cada valor discreto de  $x$ :

$$E[g[x]] = \sum_{i=1}^N g[x_i] Pr(x_i)$$

## Repaso esperanza matemática: $E[\ ]$

- Propiedades de la Esperanza:

$$\begin{aligned} \mathbb{E}[k] &= k \\ \mathbb{E}[k \cdot g[x]] &= k \cdot \mathbb{E}[g[x]] \\ \mathbb{E}[f[x] + g[x]] &= \mathbb{E}[f[x]] + \mathbb{E}[g[x]] \\ \mathbb{E}[f[x]g[y]] &= \mathbb{E}[f[x]]\mathbb{E}[g[y]] \quad \text{if } x, y \text{ independent} \end{aligned}$$

## Repaso esperanza matemática: $E[\ ]$

Function $g[\bullet]$	Expectation
$x$	mean, $\mu$
$x^k$	$k$ th moment about zero
$(x - \mu)^k$	$k$ th moment about the mean
$(x - \mu)^2$	variance
$(x - \mu)^3$	skew
$(x - \mu)^4$	kurtosis

**Table B.1** Special cases of expectation. For some functions  $g[x]$ , the expectation  $\mathbb{E}[g[x]]$  is given a special name. Here we use the notation  $\mu_x$  to represent the mean with respect to random variable  $x$ .

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$

preactivación  $i$ -ésima en esa capa:

$$\mathbf{z}_k(i) = \sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j)$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$

Elevando al cuadrado en ambos lados:

$$\mathbf{z}_k(i)^2 = \left( \sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j) \right)^2$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$

$$E[\mathbf{z}_k(i)^2] = E\left[\left(\sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j)\right)^2\right]$$

La varianza es la esperanza de  $(x - \text{media})^2$  – pero aquí la media es 0

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$

Aplicando la esperanza en ambos lados:

$$E[\mathbf{z}_k(i)^2] = E\left[\left(\sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j)\right)^2\right]$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$

varianza de la preactivación i-ésima en esa capa:

$$\sigma_{\mathbf{z}_k(i)}^2 = E\left[\left(\sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j)\right)^2\right]$$

La varianza es la esperanza de  $(x - \text{media})^2$  – pero aquí la media es 0

## Varianza de la preactivación

$$\sigma_{z_k(i)}^2 = E \left[ \left( \sum_{j=1}^{D_h} w_{ij} \mathbf{h}_k(j) \right)^2 \right]$$

– Recordemos que:

$$(a+b)^2 = a^2 + 2ab + b^2$$

$$(a+b+c)^2 = a^2 + b^2 + c^2 + 2bc + 2ca + 2ab$$

$\vdots$

$$\left( \sum_i a_i \right)^2 = \sum_i a_i^2 + 2 \sum_{i < j} a_i a_j$$

## Varianza de la preactivación

$$\sigma_{z_k(i)}^2 = E \left[ \sum_{j=1}^{D_h} (w_{ij} \mathbf{h}_k(j))^2 + 2 \cdot \sum_{r < t} (w_{ir} \mathbf{h}_k(r) w_{it} \mathbf{h}_k(t)) \right]$$

– Recordemos que:

$$(a+b)^2 = a^2 + 2ab + b^2$$

$$(a+b+c)^2 = a^2 + b^2 + c^2 + 2bc + 2ca + 2ab$$

$\vdots$

$$\left( \sum_i a_i \right)^2 = \sum_i a_i^2 + 2 \sum_{i < j} a_i a_j$$

## Varianza de la preactivación

$$\sigma_{z_k(i)}^2 = E \left[ \sum_{j=1}^{D_h} (w_{ij} \mathbf{h}_k(j))^2 + 2 \cdot \sum_{r < t} (w_{ir} \mathbf{h}_k(r) w_{it} \mathbf{h}_k(t)) \right]$$

– Y con las propiedades de la esperanza:

$$\mathbb{E}[k] = k$$

$$\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$$

$$\Rightarrow \mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$$

$$\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]] \quad \text{if } x, y \text{ independent}$$

## Varianza de la preactivación

$$\sigma_{z_k(i)}^2 = \sum_{j=1}^{D_h} E[(w_{ij} \mathbf{h}_k(j))^2] + 2 \cdot \sum_{r < t} E[w_{ir} \mathbf{h}_k(r) w_{it} \mathbf{h}_k(t)]$$

– Y con las propiedades de la esperanza:

$$\mathbb{E}[k] = k$$

$$\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$$

$$\Rightarrow \mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$$

$$\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]] \quad \text{if } x, y \text{ independent}$$

## Varianza de la preactivación

$$\sigma_{z_k(i)}^2 = \sum_{j=1}^{D_h} E[(w_{ij} \mathbf{h}_k(j))^2] + 2 \cdot \sum_{r < t} E[w_{ir} \mathbf{h}_k(r) w_{it} \mathbf{h}_k(t)]$$

- Y con las propiedades de la esperanza:

$$\begin{aligned} \mathbb{E}[k] &= k \\ \mathbb{E}[k \cdot g[x]] &= k \cdot \mathbb{E}[g[x]] \\ \mathbb{E}[f[x] + g[x]] &= \mathbb{E}[f[x]] + \mathbb{E}[g[x]] \\ \Rightarrow \mathbb{E}[f[x]g[y]] &= \mathbb{E}[f[x]]\mathbb{E}[g[y]] \quad \text{if } x, y \text{ independent} \end{aligned}$$

## Varianza de la preactivación

$$\sigma_{z_k(i)}^2 = \sum_{j=1}^{D_h} E[w_{ij}^2] E[\mathbf{h}_k(j)^2] + 2 \cdot \sum_{r < t} E[w_{ir}] E[\mathbf{h}_k(r)] E[w_{it}] E[\mathbf{h}_k(t)]$$

- Y con las propiedades de la esperanza:

$$\begin{aligned} \mathbb{E}[k] &= k \\ \mathbb{E}[k \cdot g[x]] &= k \cdot \mathbb{E}[g[x]] \\ \mathbb{E}[f[x] + g[x]] &= \mathbb{E}[f[x]] + \mathbb{E}[g[x]] \\ \Rightarrow \mathbb{E}[f[x]g[y]] &= \mathbb{E}[f[x]]\mathbb{E}[g[y]] \quad \text{if } x, y \text{ independent} \end{aligned}$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$

$$\sigma_{z_k(i)}^2 = \sum_{j=1}^{D_h} E[w_{ij}^2] E[\mathbf{h}_k(j)^2] + 2 \cdot \sum_{r < t} E[w_{ir}] E[\mathbf{h}_k(r)] E[w_{it}] E[\mathbf{h}_k(t)]$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$

$$\sigma_{z_k(i)}^2 = \sum_{j=1}^{D_h} E[w_{ij}^2] \cdot E[\mathbf{h}_k(j)^2]$$

La varianza es la esperanza de  $(x - \text{media})^2$  – pero aquí la media es 0



## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$

$$\sigma_{z_k(i)}^2 = \sum_{j=1}^{D_h} \sigma_{\Omega}^2 \cdot E[\mathbf{h}_k(j)^2]$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$
- Varianza de las activaciones en la capa k es  $\sigma_k^2$

$$\sigma_{z_k(i)}^2 = \sigma_{\Omega}^2 \sum_{j=1}^{D_h} E[\mathbf{h}_k(j)^2]$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$

$$\sigma_{z_k(i)}^2 = \sigma_{\Omega}^2 \sum_{j=1}^{D_h} E[\mathbf{h}_k(j)^2]$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} \mathbf{z}_k &= \mathbf{\Omega}_k \mathbf{h}_k + \boldsymbol{\beta}_k \\ &= \mathbf{\Omega}_k \mathbf{a}[\mathbf{z}_{k-1}] + \boldsymbol{\beta}_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\boldsymbol{\beta}_k = \mathbf{0}$
- Pesos,  $\mathbf{\Omega}_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_{\Omega}^2$
- Varianza de las activaciones en la capa k es  $\sigma_k^2$

$$\sigma_{z_k(i)}^2 = \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \sigma_k^2$$

## Inicialización

- Las preactivaciones de una capa lineal:

$$\begin{aligned} z_k &= \Omega_k h_k + \beta_k \\ &= \Omega_k a[z_{k-1}] + \beta_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\beta_k = 0$
- Pesos,  $\Omega_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_\Omega^2$
- Varianza de las activaciones en la capa k es  $\sigma_k^2$

$$\sigma_{z_k(i)}^2 = D_h \cdot \sigma_\Omega^2 \cdot \sigma_k^2$$

### 5.2.1 Inicialización de parámetros

- Necesidad de la inicialización
- Inicialización de Xavier y de He**

Understanding the difficulty of training deep feedforward neural networks. Xavier Glorot, Yoshua Bengio. AISTATS 2010: 249-256

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification  
Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun. ICCV 2015

## Inicialización

- Las preactivaciones de una capa lineal:

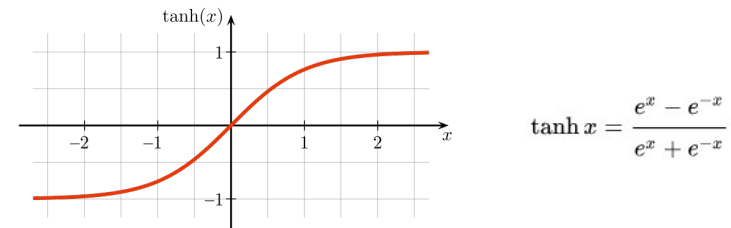
$$\begin{aligned} z_k &= \Omega_k h_k + \beta_k \\ &= \Omega_k a[z_{k-1}] + \beta_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\beta_k = 0$
- Pesos,  $\Omega_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_\Omega^2$
- Varianza de las activaciones en la capa k es  $\sigma_k^2$
- Todas las preactivaciones de la capa k comparten varianza:

$$\sigma_{z_k}^2 = D_h \cdot \sigma_\Omega^2 \cdot \sigma_k^2$$

## Activación tanh

- La función de activación tanh tiene la forma:



- Si suponemos que al comienzo del entrenamiento las preactivaciones siempre se encuentran en la parte lineal:

$$\begin{aligned} \tanh[z_k] = z_k &\Rightarrow \sigma_{z_k}^2 = \sigma_{k+1}^2 \Rightarrow \sigma_{k+1}^2 = D_h \cdot \sigma_\Omega^2 \cdot \sigma_k^2 \\ &\Rightarrow \sigma_{z_{k-1}}^2 = \sigma_k^2 \Rightarrow \sigma_{z_k}^2 = D_h \cdot \sigma_\Omega^2 \cdot \sigma_{z_{k-1}}^2 \end{aligned}$$

## Inicialización: Pregunta 3

- Las preactivaciones de una capa lineal:

$$\begin{aligned} z_k &= \Omega_k h_k + \beta_k \\ &= \Omega_k a[z_{k-1}] + \beta_k \end{aligned}$$

- Establecer todos los sesgos a 0:  $\beta_k = 0$
- Pesos,  $\Omega_k$ , con distribución Normal: **media 0** y **varianza**  $\sigma_\Omega^2$ 
  - ¿Qué pasará al movernos por la red si la varianza de los pesos de cada capa es tal que  $D_h \cdot \sigma_\Omega^2 < 1$ ?
  - ¿Qué pasará al movernos por la red si la varianza es tal que  $D_h \cdot \sigma_\Omega^2 > 1$ ?

$$\sigma_{k+1}^2 = D_h \cdot \sigma_\Omega^2 \cdot \sigma_k^2$$

## Relación entre varianzas (asume tanh)

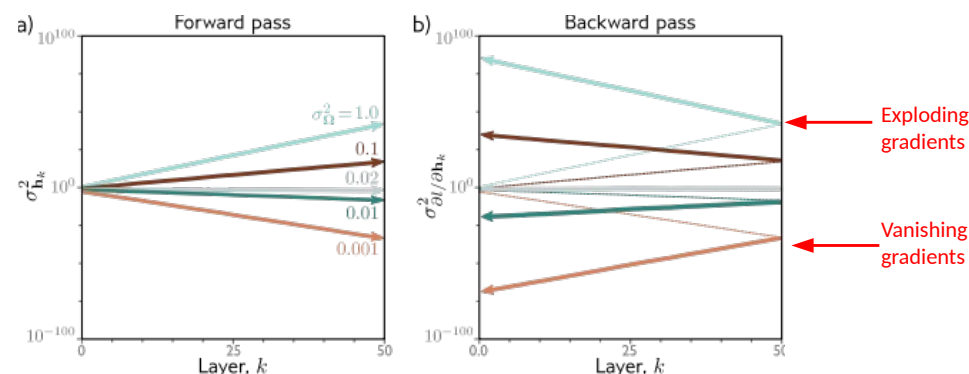
- La relación entre la varianza de las activaciones en la capa  $k$ , y la varianza de las activaciones en la capa  $k+1$  es:

$$\sigma_{k+1}^2 = D_h \cdot \sigma_\Omega^2 \cdot \sigma_k^2$$

- La relación entre la varianza de los gradientes en la capa oculta  $k+1$  y la varianza de los gradientes en capa  $k$  es:

$$\sigma_k^2 = D_{h'} \cdot \sigma_\Omega^2 \cdot \sigma_{k+1}^2$$

## Problemas numéricos con el gradiente



**Figure 7.4** Weight initialization. Consider a deep network with 50 hidden layers and  $D_h = 100$  hidden units per layer. The network has a 100 dimensional input  $\mathbf{x}$  initialized with values from a standard normal distribution, a single output fixed at  $y = 0$ , and a least squares loss function. The bias vectors  $\beta_k$  are initialized to zero and the weight matrices  $\Omega_k$  are initialized with a normal distribution with mean zero and five different variances  $\sigma_\Omega^2 \in \{0.001, 0.01, 0.02, 0.1, 1.0\}$ . a)

## Inicialización de Xavier

- Forward pass:** se busca que la varianza de las activaciones en la capa  $k+1$  **sea igual** que la varianza de las activaciones en la capa  $k$ :

$$\sigma_{k+1}^2 = D_h \cdot \sigma_\Omega^2 \cdot \sigma_k^2 \Rightarrow \boxed{\sigma_\Omega^2 = \frac{1}{D_h}}$$

N° de unidades en la capa  $k$

- Backward pass:** se busca que la varianza de los gradientes en la capa  $k$  **sea igual** que la varianza de los gradientes en la capa  $k+1$ :

$$\sigma_k^2 = D_{h'} \cdot \sigma_\Omega^2 \cdot \sigma_{k+1}^2 \Rightarrow \boxed{\sigma_\Omega^2 = \frac{1}{D_{h'}}}$$

N° de unidades en la capa  $k+1$

## Inicialización de Xavier

- Si  $D_h \neq D_{h'}$  entonces podemos usar  $(D_h + D_{h'})/2$  como número de unidades:

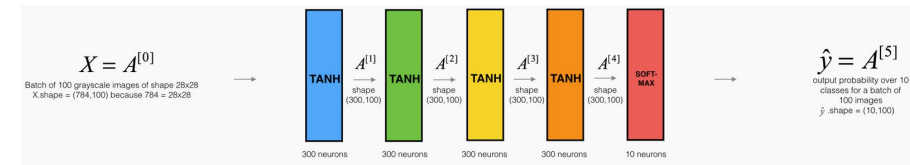
$$\sigma_{\Omega}^2 = \frac{1}{\frac{D_h + D_{h'}}{2}} = \frac{2}{D_h + D_{h'}}$$

Nº de entradas + Nº de salidas de la capa K

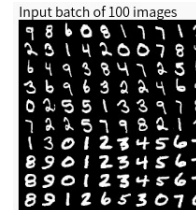
- Si esta varianza se utiliza en una distribución uniforme,  $U(a, -a)$  que tiene varianza  $a^2/3 \rightarrow$  tendremos que usar una uniforme:

$$U\left(-\sqrt{\frac{6}{D_h + D_{h'}}}, \sqrt{\frac{6}{D_h + D_{h'}}}\right)$$

## Otras Inicializaciones vs Xavier

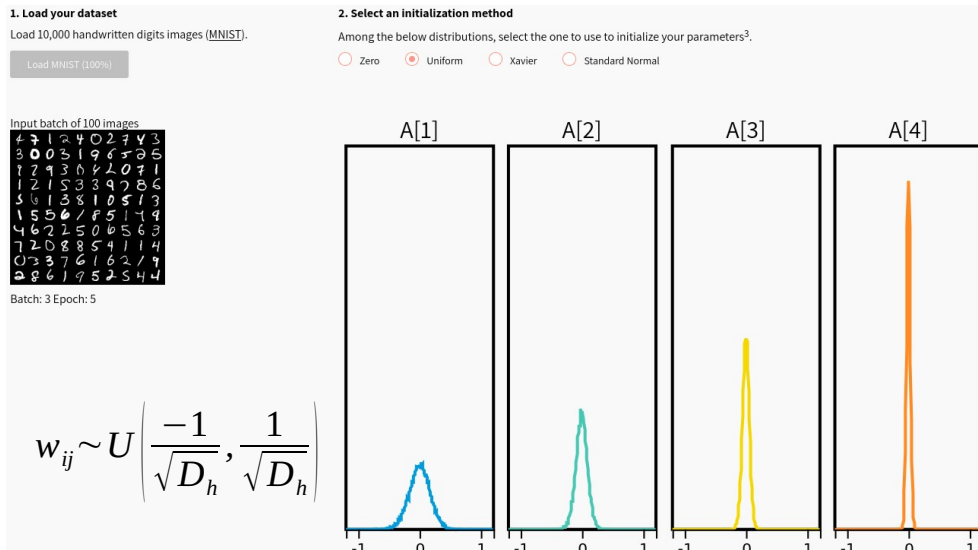


10,000 handwritten digits images (MNIST).

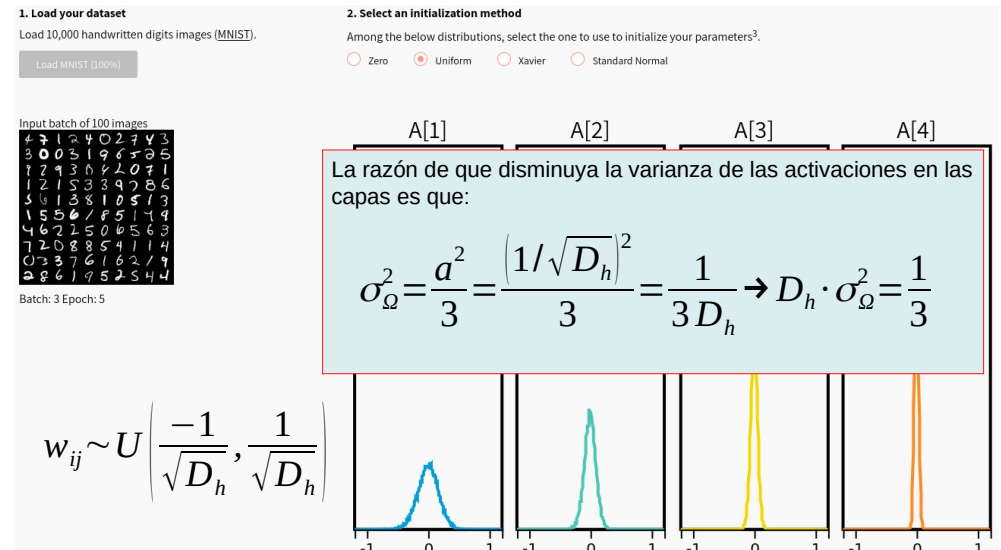


<https://www.deeplearning.ai/ai-notes/initialization/index.html>

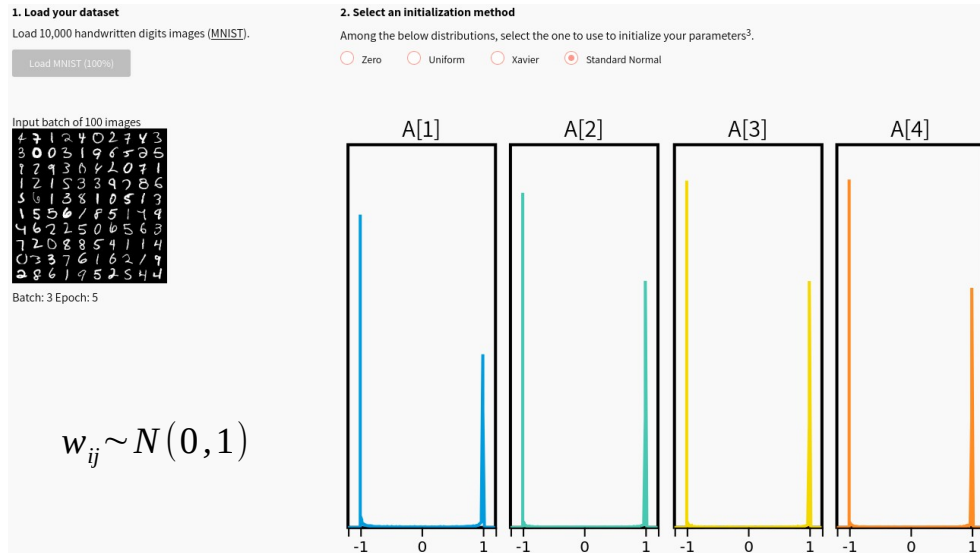
## Otras Inicializaciones vs Xavier



## Otras Inicializaciones vs Xavier



## Otras Inicializaciones vs Xavier

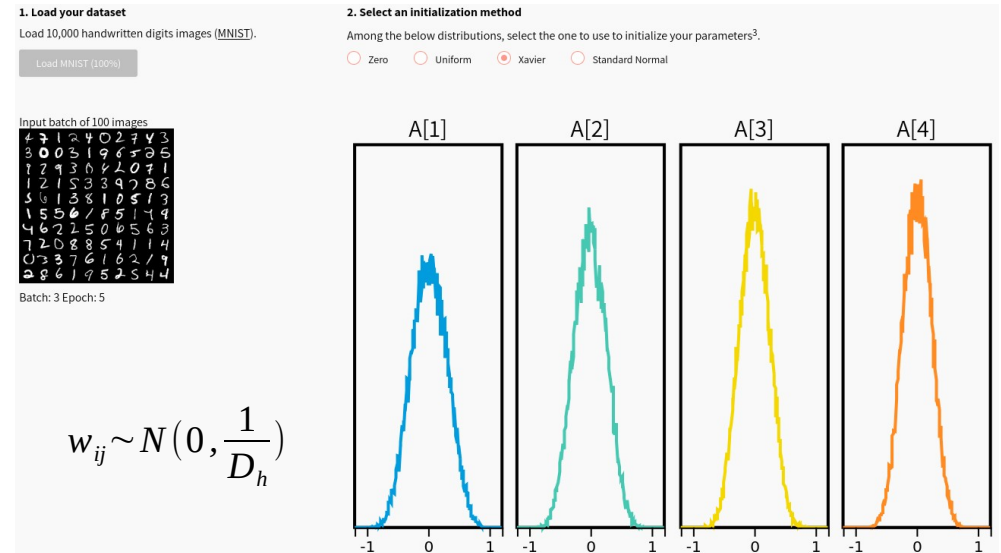


## Relación entre varianzas (asume ReLU)

- ¿Y qué pasa si tenemos ReLU como función de activación?

$$\begin{aligned}
 \sigma_{z_k}^2 &= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} E[\mathbf{h}_k(j)^2] \\
 &= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} E[\text{ReLU}(\mathbf{z}_{k-1}(j))^2] \\
 &= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \int_{-\infty}^{\infty} \text{ReLU}(\mathbf{z}_{k-1}(j))^2 \text{Pr}(\mathbf{z}_{k-1}(j)) d\mathbf{z}_k(j) \\
 &= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \int_0^{\infty} \mathbf{z}_{k-1}(j)^2 \text{Pr}(\mathbf{z}_{k-1}(j)) d\mathbf{z}_k(j) \\
 &= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \frac{\sigma_{z_{k-1}}^2}{2} = \frac{D_h \sigma_{\Omega}^2 \sigma_{z_{k-1}}^2}{2}
 \end{aligned}$$

## Otras Inicializaciones vs Xavier



## Relación entre varianzas (asume ReLU)

- La relación entre la varianza de las preactivaciones en la capa k, y la varianza de las preactivaciones en la capa k-1 es:

$$\sigma_{z_k}^2 = \frac{1}{2} \cdot D_h \cdot \sigma_{\Omega}^2 \cdot \sigma_{z_{k-1}}^2$$

- La relación entre la varianza de los gradientes de las preactivaciones en la capa k-1 y la varianza de los gradientes

$$\sigma_{z_{k-1}}^2 = \frac{1}{2} \cdot D_h \cdot \sigma_{\Omega}^2 \cdot \sigma_{z_k}^2$$

## Inicialización de He (asume ReLU)

- Forward pass:** se busca que la varianza de las preactivaciones en la capa k **sea igual** que la varianza de las preactivaciones en la capa k-1:

$$\sigma_{z_k}^2 = \frac{1}{2} \cdot D_h \cdot \sigma_{\Omega}^2 \cdot \sigma_{z_{k-1}}^2 \rightarrow \boxed{\sigma_{\Omega}^2 = \frac{2}{D_h}} \leftarrow \begin{array}{l} \text{N}^\circ \text{ de unidades} \\ \text{en la capa k} \end{array}$$

- Backward pass:** se busca que la varianza de los gradientes en la capa k-1 **sea igual** que la varianza de los gradientes en la capa k:

$$\sigma_{z_{k-1}}^2 = \frac{1}{2} \cdot D_{h'} \cdot \sigma_{\Omega}^2 \cdot \sigma_{z_k}^2 \rightarrow \boxed{\sigma_{\Omega}^2 = \frac{2}{D_{h'}}} \leftarrow \begin{array}{l} \text{N}^\circ \text{ de unidades} \\ \text{en la capa k} \end{array}$$

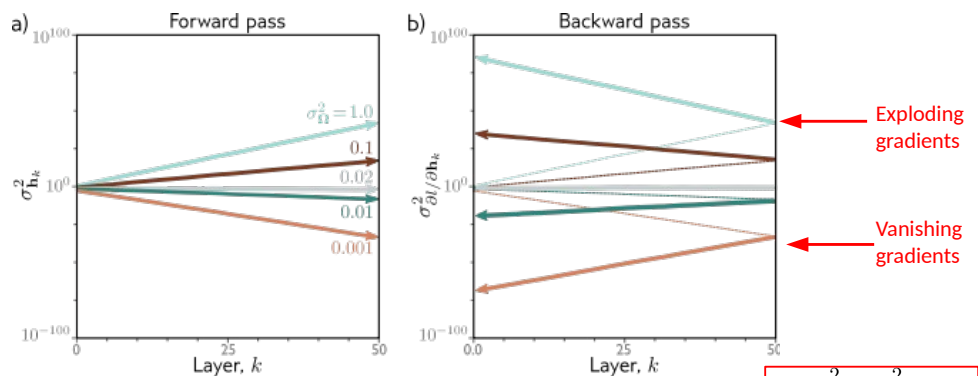
## Inicialización de He (asume ReLU)

- Si  $D_h \neq D_{h'}$ , entonces podemos usar  $(D_h + D_{h'})/2$  como número de unidades:

$$\sigma_{\Omega}^2 = \frac{2}{\frac{D_h + D_{h'}}{2}} = \frac{4}{\boxed{D_h + D_{h'}}}$$

Nº de entradas + Nº de salidas de la capa K

## Problemas numéricos con el gradiente



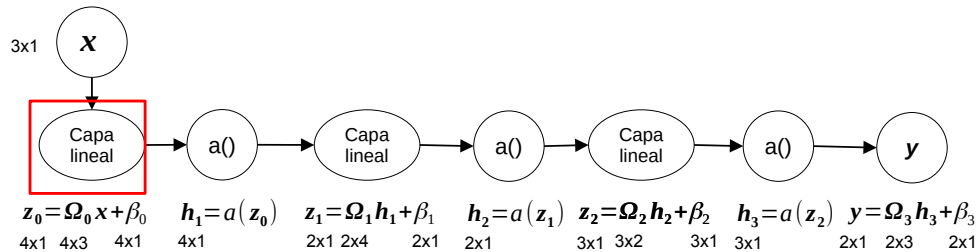
**Figure 7.4** Weight initialization. Consider a deep network with 50 hidden layers and  $D_h = 100$  hidden units per layer. The network has a 100 dimensional input  $\mathbf{x}$  initialized with values from a standard normal distribution, a single output fixed at  $y = 0$ , and a least squares loss function. The bias vectors  $\beta_k$  are initialized to zero and the weight matrices  $\Omega_k$  are initialized with a normal distribution with mean zero and five different variances  $\sigma_{\Omega}^2 \in \{0.001, 0.01, 0.02, 0.1, 1.0\}$ . a)

$$\sigma_{\Omega}^2 = \frac{2}{D_h} = \frac{2}{100} = 0.02$$

## 5.2 Optimización con redes profundas

- 5.2.1 Backpropagation
- 5.2.2 Inicialización de parámetros
- 5.2.3 Normalización
- 5.2.4 Regularización

## El problema de los valores grandes ...

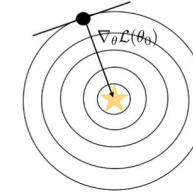


- Recordemos:  $\frac{\partial L}{\partial \Omega_0} = \frac{\partial L}{\partial z_0} \cdot \frac{\partial z_0}{\partial \Omega_0} = \frac{\partial L}{\partial z_0} \cdot x^T$
- ¿Qué pasa si en  $x = (x_1, x_2, x_3)^T$  uno de los  $x_i$  tiene una magnitud mucho mayor que los otros?

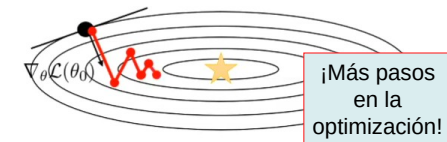
Las componentes en el gradiente con respecto a  $\Omega_0$  tendrán magnitudes muy diferentes (son el resultado de un producto con  $x^T$ )

## El problema de los valores grandes ...

- Con magnitudes iguales en los gradientes tendremos una función de coste a optimizar:



- Con magnitudes muy diferentes en los gradientes tendremos:



## El problema de los valores grandes ...

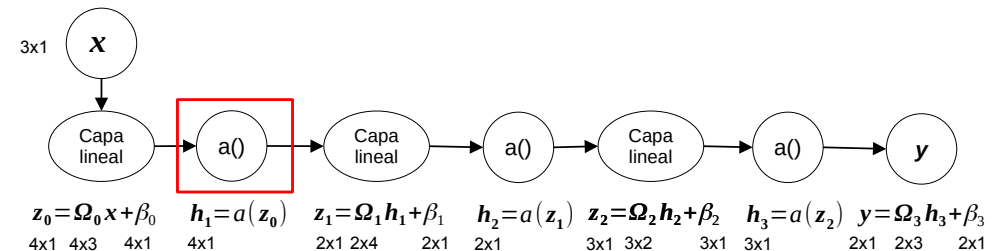
- Estandarizar entradas y etiquetas (regresión) con  $\mu=0$  y  $\sigma=1$ :



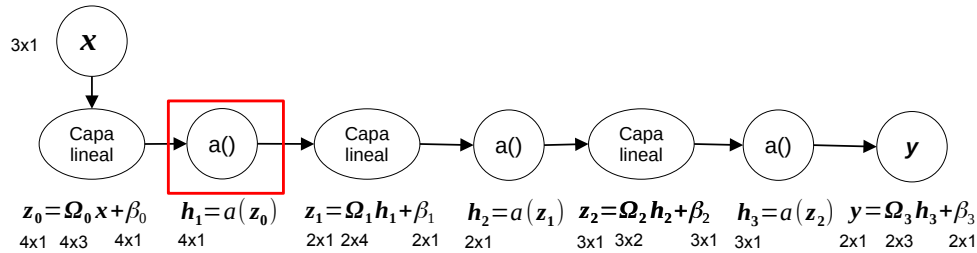
- Para hacer  $\mu=0$ :  $\bar{x}_i = x_i - E[x]$   $E[x] \approx \frac{1}{N} \cdot \sum_{i=1}^N x_i$
- Para  $\sigma=1$ :  $\bar{x}_i = \frac{x_i - E[x]}{\sqrt{E[(x - E[x])^2]}}$

## El problema de los valores grandes ...

- Estandarizar entradas y etiquetas (regresión) con  $\mu=0$  y  $\sigma=1$
- ¿Qué pasa si empezamos a tener activaciones con magnitudes muy diferentes en cualquier capa?



## ¿Estandarizar las activaciones?



- Media y desviación típica del componente  $i$  del vector de activaciones  $\mathbf{h}_i$  (N datos de entrenamiento):

$$\mu_i(i) = \frac{1}{N} \cdot \sum_{j=1}^N h_i(i, j) \quad \sigma_i(i) = \sqrt{\frac{1}{N} \cdot \sum_{j=1}^N (h_i(i, j) - \mu_i(i))^2}$$

- La activación normalizada:  $\tilde{h}_i(i) = \frac{h_i(i) - \mu_i(i)}{\sigma_i(i)}$
- Y la preactivación siguiente:  $\mathbf{z}_1 = \Omega_1 \tilde{\mathbf{h}}_1 + \beta_1$

Dependen de  $\Omega_0$  y  $\beta_0$  y ¡hay que recalcular media y varianza en cada paso de cálculo de gradiente!

## Batch normalization

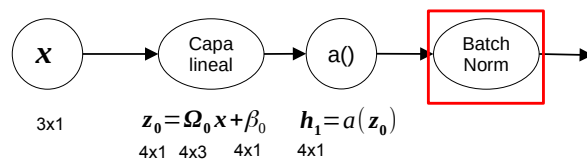
- Media y varianza del componente  $i$  del vector de activaciones  $\mathbf{h}_i$  calculado con los B datos del mini-batch:

$$\mu_i(i) \approx \frac{1}{B} \cdot \sum_{j=1}^B h_i(i, j) \quad \sigma_i(i) \approx \sqrt{\frac{1}{B} \cdot \sum_{j=1}^B (h_i(i, j) - \mu_i(i))^2}$$

- La activación normalizada se modifica con dos parámetros que se aprenden:

$$\tilde{h}_i(i) = \frac{h_i(i) - \mu_i(i)}{\sigma_i(i)} \cdot \gamma + \beta$$

## Batch normalization



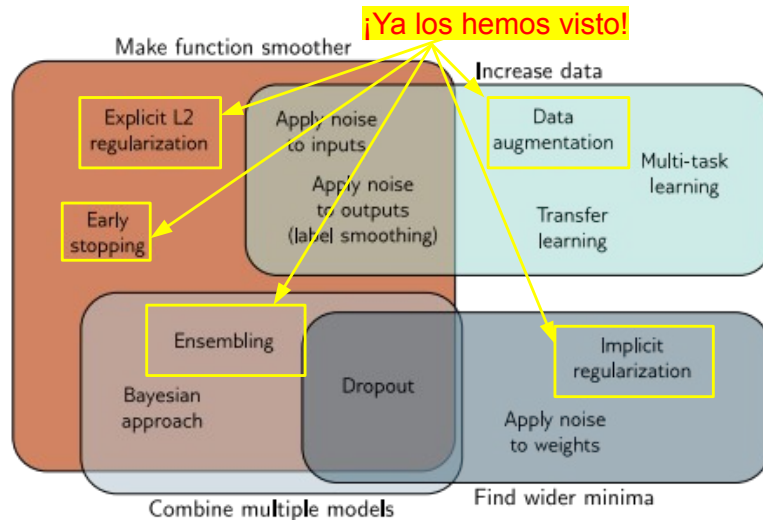
- Hay una capa adicional de normalización del mini-batch (batch norm)
- Sus parámetros se entrenan con backpropagation.
- ¿Qué se hace en inferencia con la media y la varianza de cada capa?

## 5.2 Optimización con redes profundas

- 5.2.1 Backpropagation
- 5.2.2 Inicialización de parámetros
- 5.2.3 Normalización
- 5.2.4 Regularización



## Métodos de regularización



## Recordatorio Bagging (Classifier Ensemble)

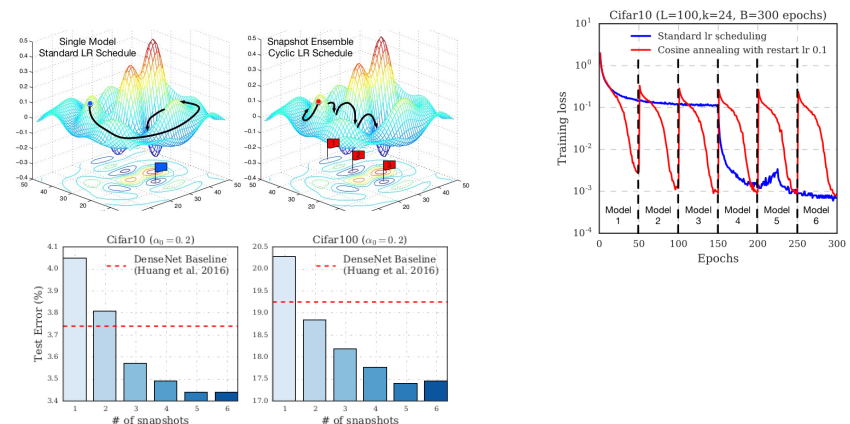
- Algoritmo **Bootstrapping Aggregating** (Leo Breiman 1996)
  - Para  $i = 1 \dots M$ 
    - ▶ Obtener  $n^* < n$  muestra de  $\mathcal{D}$  con reemplazamiento.
    - ▶ Aprender un clasificador  $C_i$  sobre la nueva muestra.
  - El clasificador final es una votación de los  $C_1, \dots, C_M$ .
- El clasificador resultante del Bagging:
  - Incrementa la estabilidad del clasificador
  - Reduce la varianza:
    - Mejora las técnicas de varianza alta (Redes de neuronas, Árboles de decisión)
    - No mejora las técnicas de varianza baja (Un clasificador lineal)

## Bagging de redes de neuronas

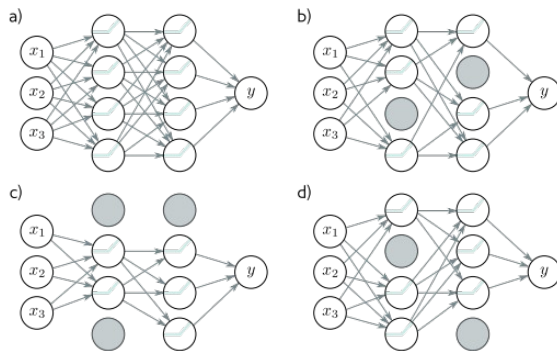
- En el entrenamiento de las redes **existen muchas fuentes de aleatoriedad**:
  - Inicialización aleatoria
  - Permutación aleatoria del mini-batch
  - Descenso de Gradiente Estocástico
- Se podría hacer Bagging de la misma red entrenada sobre los mismos datos cambiando la semilla del generador de números aleatorios:
  - Se necesitan M entrenamientos ¡muy costoso!

## Bagging de redes de neuronas

- Tomar parámetros de diferentes instantes de un único entrenamiento (snapshots) como los modelos a usar en Bagging



## Dropout: apagar neuronas aleatoriamente



Aleatoriamente algunas activaciones se igualan a 0 en la paso hacia adelante (forward pass)

### En el entrenamiento:

Para toda  $h_{ij}$ , cambiarla por  $h_{ij} \cdot m_{ij}$

$m_{ij} \sim \text{Bernoulli}(0.5)$  (1 con probabilidad 0.5, 0 en otro caso)

## Dropout: apagar neuronas aleatoriamente

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """

    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)
```

Andrej Karpathy

## Dropout en inferencia ...

### En el entrenamiento:

Para toda  $h_{ij}$ , cambiarla por  $h_{ij} \cdot m_{ij}$

$m_{ij} \sim \text{Bernoulli}(0.5)$  (1 con probabilidad 0.5, 0 en otro caso)

### En inferencia:

- Podríamos realizar la inferencia muchas veces con una máscaras de dropout generada cada vez y hacer Bagging explícito con los resultados.

¡Eso es costoso!

## Dropout en inferencia ...

### En el entrenamiento:

Para toda  $h_{ij}$ , cambiarla por  $h_{ij} \cdot m_{ij}$

$m_{ij} \sim \text{Bernoulli}(0.5)$  (1 con probabilidad 0.5, 0 en otro caso)

### En inferencia (solución 1):

- ¿Y si usamos todas las neuronas ( $m_{ij}=1$ ) en inferencia?
- Entrenamiento: en promedio  $\frac{1}{2}$  de las activaciones son 0 en cada capa lineal.
- Inferencia: Ninguna de ellas es 0, con lo que  $z_k(i) = \sum_j w_{kj} \cdot h_k(j)$  se multiplicarán por 2.
- Solución:**  $\Omega_k(j) = 1/2 \cdot \Omega_k(j)$  (dividir los pesos por

## Dropout en inferencia ... otra manera

En el entrenamiento (solución 2):

Para toda  $h_{ij}$ , cambiarla por  $h_{ij} \cdot m_{ij}$

$m_{ij} \sim \text{Bernoulli}(p)/p$  (1/p con probabilidad p, 0 en otro caso)

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

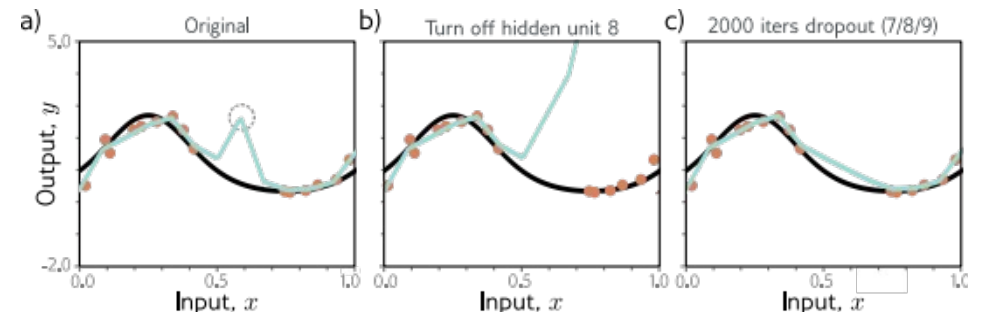
    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensemble forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    out = np.dot(W3, H2) + b3
```

test time is unchanged!

Andrej Karpathy

## Dropout



Puede eliminar puntos estimados que están lejos de los datos y no contribuyen a la pérdida de entrenamiento.