

Clase 16: Redes Neuronales

Mauricio Castro C.
mcastro@mat.uc.cl

Departamento de Estadística, Pontificia Universidad Católica de Chile

TÓPICOS APLICADOS EN ESTADÍSTICA

Segundo Semestre 2022





Redes Neuronales Multi-Capa

Perceptrones multicapa

- **Minsky and Papert (1969):** Perceptrons. An introduction to computational geometry.



Perceptrones multicapa

- ▶ **Minsky and Papert (1969):** Perceptrons. An introduction to computational geometry.
- ▶ Las limitaciones de los perceptrones pueden superarse considerando **capas + transformaciones no lineales**.



Perceptrones multicapa

- ▶ **Minsky and Papert (1969):** Perceptrons. An introduction to computational geometry.
- ▶ Las limitaciones de los perceptrones pueden superarse considerando **capas + transformaciones no lineales**.
- ▶ Estas propuestas no se podían ejecutar debido a las limitaciones computacionales.



Perceptrones multicapa

- ▶ **Minsky and Papert (1969):** Perceptrons. An introduction to computational geometry.
- ▶ Las limitaciones de los perceptrones pueden superarse considerando **capas + transformaciones no lineales**.
- ▶ Estas propuestas no se podían ejecutar debido a las limitaciones computacionales.
- ▶ **Backpropagation algorithm:** programación dinámica con dos fases, **forward + backward**.

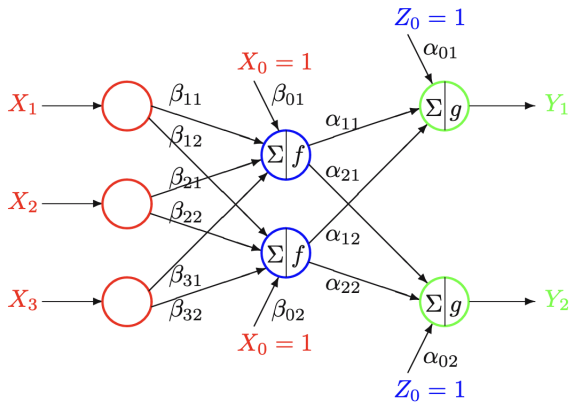


Perceptrones multicapa



- ▶ **Minsky and Papert (1969):** Perceptrons. An introduction to computational geometry.
- ▶ Las limitaciones de los perceptrones pueden superarse considerando **capas + transformaciones no lineales**.
- ▶ Estas propuestas no se podían ejecutar debido a las limitaciones computacionales.
- ▶ **Backpropagation algorithm:** programación dinámica con dos fases, forward + backward.
- ▶ **Programación dinámica:** método de optimización matemática + programación que resuelve un problema grande dividiéndolo en pequeños subproblemas de forma recursiva.

Perceptrones multicapa



input layer

hidden layer

output layer

Perceptrones multicapa

- Vector de entrada: $X = (X_1, \dots, X_r)^T$.



Perceptrones multicapa

- ▶ **Vector de entrada:** $X = (X_1, \dots, X_r)^\top$.
- ▶ **Vector de salida:** $Y = (Y_1, \dots, Y_s)^\top$.



Perceptrones multicapa

- ▶ **Vector de entrada:** $\mathbf{X} = (X_1, \dots, X_r)^\top$.
- ▶ **Vector de salida:** $\mathbf{Y} = (Y_1, \dots, Y_s)^\top$.
- ▶ **Mapeo:** no-lineal.



Perceptrones multicapa

- ▶ **Vector de entrada:** $X = (X_1, \dots, X_r)^\top$.
- ▶ **Vector de salida:** $Y = (Y_1, \dots, Y_s)^\top$.
- ▶ **Mapeo:** no-lineal.
- ▶ La capa oculta y los vectores de salida se llaman **nodos** o **neuronas**.



Perceptrones multicapa

- ▶ **Vector de entrada:** $\mathbf{X} = (X_1, \dots, X_r)^\top$.
- ▶ **Vector de salida:** $\mathbf{Y} = (Y_1, \dots, Y_s)^\top$.
- ▶ **Mapeo:** no-lineal.
- ▶ La capa oculta y los vectores de salida se llaman **nodos** o **neuronas**.
- ▶ La figura anterior posee 2 capas (la oculta + output), un vector de entrada de dimensión $r = 3$, un vector de salida de dimensión $s = 2$ y $t = 2$ nodos en la capa oculta.



Perceptrones multicapa

- Regresión múltiple: $s = 1$.



Perceptrones multicapa

- ▶ Regresión múltiple: $s = 1$.
- ▶ Regresión multivariada: $s > 1$.



Perceptrones multicapa

- ▶ Regresión múltiple: $s = 1$.
- ▶ Regresión multivariada: $s > 1$.
- ▶ Clasificación binaria: $s = 1$.



- ▶ Regresión múltiple: $s = 1$.
- ▶ Regresión multivariada: $s > 1$.
- ▶ Clasificación binaria: $s = 1$.
- ▶ Clasificación multiclase: $s = K - 1$, donde K es el número de clases.



Problema de una sola capa oculta

- Suponga una red de dos capas (una oculta + output).



Problema de una sola capa oculta

- ▶ Suponga una red de dos capas (una oculta + output).
- ▶ X_m , $m = 1, 2, \dots, r$ **nodos de entrada**.



Problema de una sola capa oculta

- ▶ Suponga una red de dos capas (una oculta + output).
- ▶ X_m , $m = 1, 2, \dots, r$ **nodos de entrada**.
- ▶ Z_j , $j = 1, \dots, t$ **nodos ocultos**.



Problema de una sola capa oculta

- ▶ Suponga una red de dos capas (una oculta + output).
- ▶ X_m , $m = 1, 2, \dots, r$ **nodos de entrada**.
- ▶ Z_j , $j = 1, \dots, t$ **nodos ocultos**.
- ▶ Y_k , $k = 1, 2, \dots, s$ **nodos de salida**.



Problema de una sola capa oculta

- ▶ Suponga una red de dos capas (una oculta + output).
- ▶ X_m , $m = 1, 2, \dots, r$ **nodos de entrada**.
- ▶ Z_j , $j = 1, \dots, t$ **nodos ocultos**.
- ▶ Y_k , $k = 1, 2, \dots, s$ **nodos de salida**.
- ▶ β_{mj} el peso de la relación entre X_m y Z_j , con sesgo β_{0j} .



Problema de una sola capa oculta

- ▶ Suponga una red de dos capas (una oculta + output).
- ▶ X_m , $m = 1, 2, \dots, r$ **nodos de entrada**.
- ▶ Z_j , $j = 1, \dots, t$ **nodos ocultos**.
- ▶ Y_k , $k = 1, 2, \dots, s$ **nodos de salida**.
- ▶ β_{mj} el peso de la relación entre X_m y Z_j , con sesgo β_{0j} .
- ▶ α_{jk} el peso de la relación entre Z_j y Y_k , con sesgo α_{0k} .

Problema de una sola capa oculta

► $X = (X_1, \dots, X_r)^T$.



Problema de una sola capa oculta

► $X = (X_1, \dots, X_r)^\top.$

► $Z = (Z_1, \dots, Z_t)^\top.$



Problema de una sola capa oculta

- ▶ $\mathbf{X} = (X_1, \dots, X_r)^\top$.
- ▶ $\mathbf{Z} = (Z_1, \dots, Z_t)^\top$.
- ▶ $U_j = \beta_{0j} + \mathbf{X}^\top \boldsymbol{\beta}_j$ y $V_k = \alpha_{0k} + \mathbf{Z}^\top \boldsymbol{\alpha}_k$.



Problema de una sola capa oculta

- ▶ $\mathbf{X} = (X_1, \dots, X_r)^\top$.
- ▶ $\mathbf{Z} = (Z_1, \dots, Z_t)^\top$.
- ▶ $U_j = \beta_{0j} + \mathbf{X}^\top \beta_j$ y $V_k = \alpha_{0k} + \mathbf{Z}^\top \alpha_k$.
- ▶ Sean $f_j(\cdot)$ y $g_k(\cdot)$ funciones de activación de la capa oculta y la capa de salida, con $j = 1, \dots, t$ y $k = 1, \dots, s$. Entonces,

$$\begin{aligned} Z_j &= f_j(U_j), \quad j = 1, 2, \dots, t \\ \mu_k(\mathbf{X}) &= g_k(V_k), \quad k = 1, 2, \dots, s, \end{aligned}$$

con $\beta_j = (\beta_{1j}, \dots, \beta_{rj})^\top$ y $\alpha_k = (\alpha_{1k}, \dots, \alpha_{tk})^\top$.

Problema de una sola capa oculta

- Usando notación matricial, tenemos que para el k -ésimo nodo de salida, con $k = 1, \dots, s$:



Problema de una sola capa oculta

- Usando notación matricial, tenemos que para el k -ésimo nodo de salida, con $k = 1, \dots, s$:

$$Y_k = \mu_k(\mathbf{X}) + \epsilon_k,$$

$$\text{con } \mu_k(\mathbf{X}) = g_k \left(\alpha_{0k} + \sum_{j=1}^t \alpha_{jk} f_j \left(\beta_{0j} + \sum_{m=1}^r \beta_{mj} X_m \right) \right).$$



Problema de una sola capa oculta

- ▶ Usando notación matricial, tenemos que para el k -ésimo nodo de salida, con $k = 1, \dots, s$:

$$Y_k = \mu_k(\mathbf{X}) + \epsilon_k,$$

$$\text{con } \mu_k(\mathbf{X}) = g_k \left(\alpha_{0k} + \sum_{j=1}^t \alpha_{jk} f_j \left(\beta_{0j} + \sum_{m=1}^r \beta_{mj} X_m \right) \right).$$

- ▶ $f_j(\cdot)$'s son usualmente **logística** o **tanh** (sigmoideas).



Problema de una sola capa oculta

- ▶ Usando notación matricial, tenemos que para el k -ésimo nodo de salida, con $k = 1, \dots, s$:

$$Y_k = \mu_k(\mathbf{X}) + \epsilon_k,$$

$$\text{con } \mu_k(\mathbf{X}) = g_k \left(\alpha_{0k} + \sum_{j=1}^t \alpha_{jk} f_j \left(\beta_{0j} + \sum_{m=1}^r \beta_{mj} X_m \right) \right).$$

- ▶ $f_j(\cdot)$'s son usualmente **logística** o **tanh** (sigmoideas).
- ▶ $g_k(\cdot)$'s son usualmente **lineales** (regresión) o **sigmoideas** (clasificación).



Problema de una sola capa oculta

- ▶ Usando notación matricial, tenemos que para el k -ésimo nodo de salida, con $k = 1, \dots, s$:

$$Y_k = \mu_k(\mathbf{X}) + \epsilon_k,$$

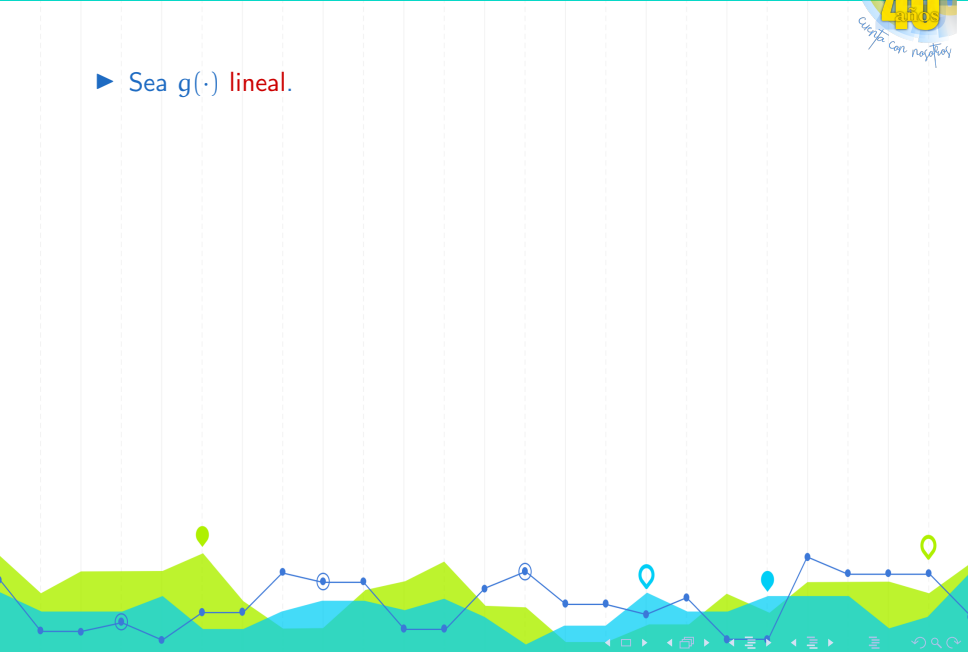
$$\text{con } \mu_k(\mathbf{X}) = g_k \left(\alpha_{0k} + \sum_{j=1}^t \alpha_{jk} f_j \left(\beta_{0j} + \sum_{m=1}^r \beta_{mj} X_m \right) \right).$$

- ▶ $f_j(\cdot)$'s son usualmente **logística** o **tanh** (sigmoideas).
- ▶ $g_k(\cdot)$'s son usualmente **lineales** (regresión) o **sigmoideas** (clasificación).
- ▶ ϵ_k es un ruido Gaussiano con media 0 y varianza σ_k .



Problema de una sola capa oculta + $s = 1$

- Sea $g(\cdot)$ **lineal**.



Problema de una sola capa oculta + $s = 1$

- ▶ Sea $g(\cdot)$ **lineal**.
- ▶ Sea $f_j(\cdot) = f(\cdot)$, $\forall j$ *i.e.*, **todas iguales y sigmoideas**.



Problema de una sola capa oculta + $s = 1$

- ▶ Sea $g(\cdot)$ **lineal**.
- ▶ Sea $f_j(\cdot) = f(\cdot)$, $\forall j$ **i.e., todas iguales y sigmoideas**.
- ▶ Entonces, $Y = \mu(\mathbf{X}) + \epsilon$ con



Problema de una sola capa oculta + $s = 1$

- ▶ Sea $g(\cdot)$ **lineal**.
- ▶ Sea $f_j(\cdot) = \mathbf{f}(\cdot)$, $\forall j$ *i.e.*, **todas iguales y sigmoideas**.
- ▶ Entonces, $Y = \mu(\mathbf{X}) + \epsilon$ con

$$\mu(\mathbf{X}) = \alpha_0 + \sum_{j=1}^t \alpha_j \mathbf{f} \left(\beta_{0j} + \sum_{m=1}^r \beta_{mj} X_m \right).$$



Problema de una sola capa oculta + $s = 1$

- ▶ Sea $g(\cdot)$ **lineal**.
- ▶ Sea $f_j(\cdot) = \mathbf{f}(\cdot)$, $\forall j$ **i.e., todas iguales y sigmoideas**.
- ▶ Entonces, $Y = \mu(\mathbf{X}) + \epsilon$ con

$$\mu(\mathbf{X}) = \alpha_0 + \sum_{j=1}^t \alpha_j \mathbf{f} \left(\beta_{0j} + \sum_{m=1}^r \beta_{mj} X_m \right).$$

- ▶ La red anterior equivale al **perceptron de una sola capa**.



Problema de una sola capa oculta + $s = 1$

- ▶ Sea $g(\cdot)$ **lineal**.
- ▶ Sea $f_j(\cdot) = \mathbf{f}(\cdot)$, $\forall j$ **i.e., todas iguales y sigmoideas**.
- ▶ Entonces, $Y = \mu(\mathbf{X}) + \epsilon$ con

$$\mu(\mathbf{X}) = \alpha_0 + \sum_{j=1}^t \alpha_j \mathbf{f} \left(\beta_{0j} + \sum_{m=1}^r \beta_{mj} X_m \right).$$

- ▶ La red anterior equivale al **perceptron de una sola capa**.
- ▶ Si $\mathbf{f}(\cdot)$ y $g(\cdot)$ son lineales, la red es una combinación lineal de inputs.

Теорема. При любом целом $n \geq 2$ существуют такие определенные на единичном отрезке $E^1 = [0; 1]$ непрерывные действительные функции $\psi^{pq}(x)$, что каждая определенная на n -мерном единичном кубе E^n непрерывная действительная функция $f(x_1, \dots, x_n)$ представима в виде

$$f(x_1, \dots, x_n) = \sum_{q=1}^{q=2n+1} \chi_q \left[\sum_{p=1}^n \psi^{pq}(x_p) \right], \quad (1)$$

где функции $\chi_q(y)$ действительны и непрерывны.

A.N. Kolmogorov: On the representation of continuous functions of many variables by superposition of continuous functions of one variable and addition, Doklady Akademii Nauk SSSR, 1957, Volume 114, Number 5, 953–956



- Kolmogorov establece que podrían aproximarse todas las funciones continua a través de una red de dos capas, incluyendo una oculta con un gran número de nodos ocultos, outputs y pesos.



- ▶ Kolmogorov establece que podrían aproximarse todas las funciones continua a través de una red de dos capas, incluyendo una oculta con un gran número de nodos ocultos, outputs y pesos.
- ▶ Mejor la aproximación, más nodos ocultos son requeridos.



- ▶ Kolmogorov establece que podrían aproximarse todas las funciones continua a través de una red de dos capas, incluyendo una oculta con un gran número de nodos ocultos, outputs y pesos.
- ▶ Mejor la aproximación, más nodos ocultos son requeridos.
- ▶ El resultado se conoce como el **Teorema Universal de Aproximación**.



- ▶ Kolmogorov establece que podrían aproximarse todas las funciones continua a través de una red de dos capas, incluyendo una oculta con un gran número de nodos ocultos, outputs y pesos.
- ▶ Mejor la aproximación, más nodos ocultos son requeridos.
- ▶ El resultado se conoce como el **Teorema Universal de Aproximación**.
- ▶ El teorema anterior dice que podemos aproximar cualquier función arbitraria con una red neuronal de una capa oculta.



- ▶ Kolmogorov establece que podrían aproximarse todas las funciones continua a través de una red de dos capas, incluyendo una oculta con un gran número de nodos ocultos, outputs y pesos.
- ▶ Mejor la aproximación, más nodos ocultos son requeridos.
- ▶ El resultado se conoce como el **Teorema Universal de Aproximación**.
- ▶ El teorema anterior dice que podemos aproximar cualquier función arbitraria con una red neuronal de una capa oculta.
- ▶ Sin embargo, no nos dice si esta aproximación es la mejor.



Problema de más de una capa oculta

► Sea $\mu(\mathbf{X}) = g(\alpha_0 + \mathbf{A}f(\beta_0 + \mathbf{B}\mathbf{X}))$.



Problema de más de una capa oculta

- ▶ Sea $\mu(\mathbf{X}) = g(\alpha_0 + \mathbf{A}f(\beta_0 + \mathbf{B}\mathbf{X}))$.
- ▶ $\mathbf{B} = (\beta_{ij})$ matriz de pesos entre los nodos de entrada y la capa oculta de dimensión $t \times r$.



Problema de más de una capa oculta

- ▶ Sea $\mu(\mathbf{X}) = g(\alpha_0 + \mathbf{A}f(\beta_0 + \mathbf{B}\mathbf{X}))$.
- ▶ $\mathbf{B} = (\beta_{ij})$ matriz de pesos entre los nodos de entrada y la capa oculta de dimensión $t \times r$.
- ▶ $\mathbf{A} = (\alpha_{jk})$ matriz de pesos entre la capa oculta y el output de dimensión $s \times t$.



Problema de más de una capa oculta

- ▶ Sea $\mu(\mathbf{X}) = g(\alpha_0 + \mathbf{A}\mathbf{f}(\beta_0 + \mathbf{B}\mathbf{X}))$.
- ▶ $\mathbf{B} = (\beta_{ij})$ matriz de pesos entre los nodos de entrada y la capa oculta de dimensión $t \times r$.
- ▶ $\mathbf{A} = (\alpha_{jk})$ matriz de pesos entre la capa oculta y el output de dimensión $s \times t$.
- ▶ $\beta_0 = (\beta_{01}, \dots, \beta_{0t})^\top$, $\alpha_0 = (\alpha_{01}, \dots, \alpha_{0s})^\top$, $\mathbf{f} = (f_1, \dots, f_t)^\top$ y $\mathbf{g} = (g_1, \dots, g_s)^\top$ (funciones de activación no lineales).



Problema de más de una capa oculta

- Un caso especial de lo anterior resulta cuando $f_j(\cdot)$ y $g_k(\cdot)$ son funciones identidad.



Problema de más de una capa oculta

- ▶ Un caso especial de lo anterior resulta cuando $f_j(\cdot)$ y $g_k(\cdot)$ son funciones identidad.
- ▶ Aquí lo anterior se reduce a $\mu(\mathbf{X}) = \mu + \mathbf{A}\mathbf{B}\mathbf{X}$ con $\mu = \alpha_0 + \mathbf{A}\beta_0$.



Problema de más de una capa oculta

- ▶ Un caso especial de lo anterior resulta cuando $f_j(\cdot)$ y $g_k(\cdot)$ son funciones identidad.
- ▶ Aquí lo anterior se reduce a $\mu(\mathbf{X}) = \mu + \mathbf{A}\mathbf{B}\mathbf{X}$ con $\mu = \alpha_0 + \mathbf{A}\beta_0$.
- ▶ Este modelo es conocido como el **modelo de regresión de rango reducido** (Izenman, 1975, Journal of Multivariate Analysis, 5:248-264).



Problema de más de una capa oculta

- ▶ Un caso especial de lo anterior resulta cuando $f_j(\cdot)$ y $g_k(\cdot)$ son funciones identidad.
- ▶ Aquí lo anterior se reduce a $\mu(\mathbf{X}) = \mu + \mathbf{ABX}$ con $\mu = \alpha_0 + \mathbf{A}\beta_0$.
- ▶ Este modelo es conocido como el **modelo de regresión de rango reducido** (Izenman, 1975, Journal of Multivariate Analysis, 5:248-264).
- ▶ La diferencia con los modelos de regresión tradicional es que el modelo de rango reducido impone que el rango de \mathbf{AB} sea menor que $\min\{r, s\}$.



Algoritmo Backpropagation - Fase Forward

- Los inputs ingresan a la red neuronal.



Algoritmo Backpropagation - Fase Forward

- ▶ Los inputs ingresan a la red neuronal.
- ▶ Dados los pesos actuales, se calculan las todas las componentes de las capas de la red.



Algoritmo Backpropagation - Fase Forward

- ▶ Los inputs ingresan a la red neuronal.
- ▶ Dados los pesos actuales, se calculan las todas las componentes de las capas de la red.
- ▶ El output se compara con el valor real.



Algoritmo Backpropagation - Fase Forward

- ▶ Los inputs ingresan a la red neuronal.
- ▶ Dados los pesos actuales, se calculan las todas las componentes de las capas de la red.
- ▶ El output se compara con el valor real.
- ▶ Se calcula la derivada de la función de perdida con respecto del output.



Algoritmo Backpropagation - Fase Forward

- ▶ Los inputs ingresan a la red neuronal.
- ▶ Dados los pesos actuales, se calculan las todas las componentes de las capas de la red.
- ▶ El output se compara con el valor real.
- ▶ Se calcula la derivada de la función de perdida con respecto del output.
- ▶ Esta derivada debe ser calculada con respecto a los pesos en todas las capas de la red.



Algoritmo Backpropagation - Fase Backward

- Obtener el gradiente de la función de pérdida con respecto de los pesos.



Algoritmo Backpropagation - Fase Backward

- Obtener el gradiente de la función de pérdida con respecto de los pesos.
- Para esto se necesita usar la **regla de la cadena**.



Algoritmo Backpropagation - Fase Backward

- ▶ Obtener el gradiente de la función de pérdida con respecto de los pesos.
- ▶ Para esto se necesita usar la **regla de la cadena**.
- ▶ Los gradientes se obtienen hacia atrás, empezando por el output.

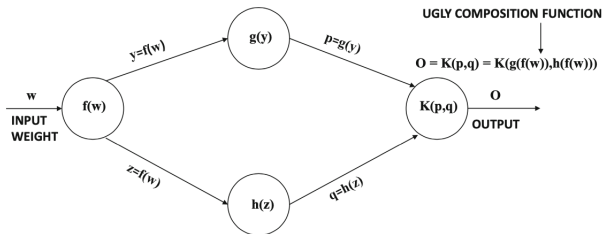


Algoritmo Backpropagation - Fase Backward

- ▶ Obtener el gradiente de la función de pérdida con respecto de los pesos.
- ▶ Para esto se necesita usar la **regla de la cadena**.
- ▶ Los gradientes se obtienen hacia atrás, empezando por el output.
- ▶ Por esto, la fase se llama **Backward**.



Algoritmo Backpropagation



$$\begin{aligned}\frac{\partial o}{\partial w} &= \frac{\partial o}{\partial p} \cdot \frac{\partial p}{\partial w} + \frac{\partial o}{\partial q} \cdot \frac{\partial q}{\partial w} \quad [\text{Multivariable Chain Rule}] \\ &= \frac{\partial o}{\partial p} \cdot \frac{\partial p}{\partial y} \cdot \frac{\partial y}{\partial w} + \frac{\partial o}{\partial q} \cdot \frac{\partial q}{\partial z} \cdot \frac{\partial z}{\partial w} \quad [\text{Univariate Chain Rule}] \\ &= \underbrace{\frac{\partial K(p,q)}{\partial p} \cdot g'(y) \cdot f'(w)}_{\text{First path}} + \underbrace{\frac{\partial K(p,q)}{\partial q} \cdot h'(z) \cdot f'(w)}_{\text{Second path}}\end{aligned}$$

Algoritmo Backpropagation

- Sea \mathcal{M} el conjunto de los r nodos de entrada, \mathcal{J} el conjunto de los t nodos ocultos y \mathcal{K} es el conjunto de los s nodos output.



Algoritmo Backpropagation

- Sea \mathcal{M} el conjunto de los r nodos de entrada, \mathcal{J} el conjunto de los t nodos ocultos y \mathcal{K} es el conjunto de los s nodos output.
- Sea el i -ésimo vector de entrada de la capa oculta

$$V_{i,k} = \sum_{j \in \mathcal{J}} \alpha_{kj} Z_{i,j} = \alpha_{k0} + \mathbf{Z}_i^\top \boldsymbol{\alpha}_k, \quad k \in \mathcal{K},$$

con $\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,t})^\top$, $\boldsymbol{\alpha}_k = (\alpha_{k1}, \dots, \alpha_{kt})^\top$ y $Z_{i,0} = 1$.



Algoritmo Backpropagation

- Sea \mathcal{M} el conjunto de los r nodos de entrada, \mathcal{J} el conjunto de los t nodos ocultos y \mathcal{K} es el conjunto de los s nodos output.

- Sea el i -ésimo vector de entrada de la capa oculta

$$V_{i,k} = \sum_{j \in \mathcal{J}} \alpha_{kj} Z_{i,j} = \alpha_{k0} + \mathbf{Z}_i^\top \boldsymbol{\alpha}_k, \quad k \in \mathcal{K},$$

con $\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,t})^\top$, $\boldsymbol{\alpha}_k = (\alpha_{k1}, \dots, \alpha_{kt})^\top$ y $Z_{i,0} = 1$.

- El output aquí se define como $\tilde{Y}_{i,k} = g_k(V_{i,k})$, $k \in \mathcal{K}$.



Algoritmo Backpropagation



- Sea \mathcal{M} el conjunto de los r nodos de entrada, \mathcal{J} el conjunto de los t nodos ocultos y \mathcal{K} es el conjunto de los s nodos output.

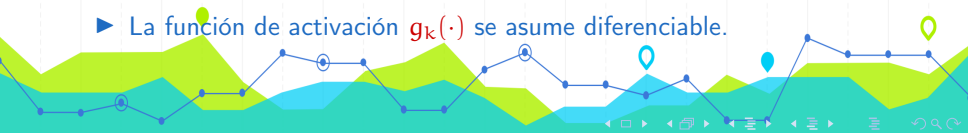
- Sea el i -ésimo vector de entrada de la capa oculta

$$V_{i,k} = \sum_{j \in \mathcal{J}} \alpha_{kj} Z_{i,j} = \alpha_{k0} + \mathbf{Z}_i^\top \boldsymbol{\alpha}_k, \quad k \in \mathcal{K},$$

con $\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,t})^\top$, $\boldsymbol{\alpha}_k = (\alpha_{k1}, \dots, \alpha_{kt})^\top$ y $Z_{i,0} = 1$.

- El output aquí se define como $\tilde{Y}_{i,k} = g_k(V_{i,k})$, $k \in \mathcal{K}$.

- La función de activación $g_k(\cdot)$ se asume diferenciable.



Algoritmo Backpropagation

- Sea $e_{i,k} = Y_{i,k} - \tilde{Y}_{i,k}$, $k \in \mathcal{K}$ y la función de error

$$E_i = \frac{1}{2} \sum_{k \in \mathcal{K}} e_{i,k}^2 = \frac{1}{2} \sum_{k \in \mathcal{K}} (Y_{i,k} - \tilde{Y}_{i,k})^2,$$

con $i = 1, \dots, n$.



- Sea $e_{i,k} = Y_{i,k} - \tilde{Y}_{i,k}$, $k \in \mathcal{K}$ y la función de error

$$E_i = \frac{1}{2} \sum_{k \in \mathcal{K}} e_{i,k}^2 = \frac{1}{2} \sum_{k \in \mathcal{K}} (Y_{i,k} - \tilde{Y}_{i,k})^2,$$

con $i = 1, \dots, n$.

- Buscamos minimizar $ESS = n^{-1} \sum_{i=1}^n E_i$.



- Sea $e_{i,k} = Y_{i,k} - \tilde{Y}_{i,k}$, $k \in \mathcal{K}$ y la función de error

$$E_i = \frac{1}{2} \sum_{k \in \mathcal{K}} e_{i,k}^2 = \frac{1}{2} \sum_{k \in \mathcal{K}} (Y_{i,k} - \tilde{Y}_{i,k})^2,$$

con $i = 1, \dots, n$.

- Buscamos minimizar $ESS = n^{-1} \sum_{i=1}^n E_i$.
- Actualizamos α_i de acuerdo a: $\alpha_{i+1} = \alpha_i + \Delta \alpha_i$, donde $\Delta \alpha_i = -\eta \frac{\partial E_i}{\partial \alpha_i} = \left(-\eta \frac{\partial E_i}{\partial \alpha_{i,j,h}} \right) = (\Delta \alpha_{i,k,j})$ y $(\alpha_{i,k,j})$ denota el vector de dimensión ts de los pesos de la capa oculta.



Algoritmo Backpropagation

- Usando la regla de la cadena,

$$\begin{aligned}\frac{\partial E_i}{\partial \alpha_{i,kj}} &= \frac{\partial E_i}{\partial e_{i,k}} \cdot \frac{\partial e_{i,k}}{\partial \tilde{Y}_{i,k}} \cdot \frac{\partial \tilde{Y}_{i,k}}{\partial V_{i,k}} \cdot \frac{\partial V_{i,k}}{\partial \alpha_{i,kj}} \\ &= e_{i,k} \cdot (-1) \cdot g'_k(V_{i,k}) \cdot Z_{i,j} \\ &= -e_{i,k} g'_k(\alpha_{i,k0} + \mathbf{Z}_i^\tau \alpha_{i,k}) Z_{i,j}\end{aligned}$$



- Usando la regla de la cadena,

$$\begin{aligned}\frac{\partial E_i}{\partial \alpha_{i,kj}} &= \frac{\partial E_i}{\partial e_{i,k}} \cdot \frac{\partial e_{i,k}}{\partial \tilde{Y}_{i,k}} \cdot \frac{\partial \tilde{Y}_{i,k}}{\partial V_{i,k}} \cdot \frac{\partial V_{i,k}}{\partial \alpha_{i,kj}} \\ &= e_{i,k} \cdot (-1) \cdot g'_k(V_{i,k}) \cdot Z_{i,j} \\ &= -e_{i,k} g'_k(\alpha_{i,k0} + \mathbf{Z}_i^\tau \alpha_{i,k}) Z_{i,j}\end{aligned}$$

- La actualización de $\alpha_{i,kj}$ usando el **descenso por gradiente** se puede expresar como

$$\alpha_{i+1,kj} = \alpha_{i,kj} - \eta \frac{\partial E_i}{\partial \alpha_{i,kj}} = \alpha_{i,kj} + \eta \delta_{i,k} Z_{i,j},$$

donde $\delta_{i,k} = -\frac{\partial E_i}{\partial \tilde{Y}_{i,k}} \cdot \frac{\partial \tilde{Y}_{i,k}}{\partial V_{i,k}} = e_{i,k} g'_k(V_{i,k})$.

Algoritmo Backpropagation

- En nuestra notación δ_{ik} es el producto entre e_{ik} y la derivada de $g_k(\cdot)$.



Algoritmo Backpropagation

- ▶ En nuestra notación δ_{ik} es el producto entre e_{ik} y la derivada de $g_k(\cdot)$.
- ▶ δ_{ik} se conoce también como **sensibilidad** o **gradiente local**.



Algoritmo Backpropagation

- ▶ En nuestra notación δ_{ik} es el producto entre e_{ik} y la derivada de $g_k(\cdot)$.
- ▶ δ_{ik} se conoce también como **sensibilidad** o **gradiente local**.
- ▶ η es el parámetro de aprendizaje del algoritmo.



Algoritmo Backpropagation

- ▶ En nuestra notación δ_{ik} es el producto entre e_{ik} y la derivada de $g_k(\cdot)$.
- ▶ δ_{ik} se conoce también como **sensibilidad** o **gradiente local**.
- ▶ η es el parámetro de aprendizaje del algoritmo.
- ▶ La siguiente parte del algoritmo busca la actualización de los pesos que conectan la capa oculta con la capa de entrada.



Algoritmo Backpropagation

- Sea $u_{i,j} = \sum_{m \in \mathcal{M}} \beta_{i,jm} x_{i,m} = \beta_{i,j0} + \mathbf{X}_i^\top \boldsymbol{\beta}_{i,j}$, $j \in \mathcal{J}$, con $\mathbf{X}_i = (x_{i,1}, \dots, x_{i,r})^\top$, $\boldsymbol{\beta}_{i,j} = (\beta_{i,j1}, \dots, \beta_{i,jr})^\top$, y $x_{i,0} = 1$.



Algoritmo Backpropagation

- ▶ Sea $U_{i,j} = \sum_{m \in \mathcal{M}} \beta_{i,jm} X_{i,m} = \beta_{i,j0} + \mathbf{X}_i^\top \boldsymbol{\beta}_{i,j}$, $j \in \mathcal{J}$, con $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,r})^\top$, $\boldsymbol{\beta}_{i,j} = (\beta_{i,j1}, \dots, \beta_{i,jr})^\top$, y $X_{i,0} = 1$.
- ▶ El output es $Z_{ij} = f_j(U_{ij})$.



Algoritmo Backpropagation

- ▶ Sea $U_{i,j} = \sum_{m \in \mathcal{M}} \beta_{i,jm} X_{i,m} = \beta_{i,j0} + \mathbf{X}_i^\top \boldsymbol{\beta}_{i,j}$, $j \in \mathcal{J}$, con $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,r})^\top$, $\boldsymbol{\beta}_{i,j} = (\beta_{i,j1}, \dots, \beta_{i,jr})^\top$, y $X_{i,0} = 1$.
- ▶ El output es $Z_{ij} = f_j(U_{ij})$.
- ▶ Sea $\boldsymbol{\beta}_i = (\boldsymbol{\beta}_{i,1}^\top, \dots, \boldsymbol{\beta}_{i,t}^\top)^\top = (\beta_{i,jm})$ el vector $(r+1)t$ de pesos entre la capa de entrada y la oculta.



Algoritmo Backpropagation

- ▶ Sea $U_{i,j} = \sum_{m \in \mathcal{M}} \beta_{i,jm} X_{i,m} = \beta_{i,j0} + \mathbf{X}_i^\top \boldsymbol{\beta}_{i,j}$, $j \in \mathcal{J}$, con $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,r})^\top$, $\boldsymbol{\beta}_{i,j} = (\beta_{i,j1}, \dots, \beta_{i,jr})^\top$, y $X_{i,0} = 1$.
- ▶ El output es $Z_{ij} = f_j(U_{ij})$.
- ▶ Sea $\boldsymbol{\beta}_i = (\boldsymbol{\beta}_{i,1}^\top, \dots, \boldsymbol{\beta}_{i,t}^\top)^\top = (\beta_{i,jm})$ el vector $(r+1)t$ de pesos entre la capa de entrada y la oculta.
- ▶ La actualización de estos pesos está dada por $\boldsymbol{\beta}_{i+1} = \boldsymbol{\beta}_i + \Delta \boldsymbol{\beta}_i$, donde $\Delta \boldsymbol{\beta}_i = -\eta \frac{\partial E_i}{\partial \boldsymbol{\beta}_i} = \left(-\eta \frac{\partial E_i}{\partial \beta_{i,jm}} \right) = (\Delta \beta_{i,jm})$.



Algoritmo Backpropagation

- ▶ Sea $U_{i,j} = \sum_{m \in \mathcal{M}} \beta_{i,jm} X_{i,m} = \beta_{i,j0} + \mathbf{X}_i^\top \boldsymbol{\beta}_{i,j}$, $j \in \mathcal{J}$, con $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,r})^\top$, $\boldsymbol{\beta}_{i,j} = (\beta_{i,j1}, \dots, \beta_{i,jr})^\top$, y $X_{i,0} = 1$.
- ▶ El output es $Z_{ij} = f_j(U_{ij})$.
- ▶ Sea $\boldsymbol{\beta}_i = (\boldsymbol{\beta}_{i,1}^\top, \dots, \boldsymbol{\beta}_{i,t}^\top)^\top = (\beta_{i,jm})$ el vector $(r+1)t$ de pesos entre la capa de entrada y la oculta.
- ▶ La actualización de estos pesos está dada por $\boldsymbol{\beta}_{i+1} = \boldsymbol{\beta}_i + \Delta \boldsymbol{\beta}_i$, donde $\Delta \boldsymbol{\beta}_i = -\eta \frac{\partial E_i}{\partial \boldsymbol{\beta}_i} = \left(-\eta \frac{\partial E_i}{\partial \beta_{i,jm}} \right) = (\Delta \beta_{i,jm})$.
- ▶ Al igual que la parte anterior, se utiliza la regla de la cadena para determinar $\frac{\partial E_i}{\partial \boldsymbol{\beta}_i}$.



Algoritmo Backpropagation

- No hay una demostración formal de la convergencia del algoritmo.



Algoritmo Backpropagation

- ▶ No hay una demostración formal de la convergencia del algoritmo.
- ▶ El algoritmo es lento y los estimadores muchas veces inestables.



Algoritmo Backpropagation

- ▶ No hay una demostración formal de la convergencia del algoritmo.
- ▶ El algoritmo es lento y los estimadores muchas veces inestables.
- ▶ Existen problemas de **identificabilidad** en los pesos.



Algoritmo Backpropagation

- ▶ No hay una demostración formal de la convergencia del algoritmo.
- ▶ El algoritmo es lento y los estimadores muchas veces inestables.
- ▶ Existen problemas de **identificabilidad** en los pesos.
- ▶ En la práctica se **deja correr** el algoritmo hasta que se estabilice.



Algoritmo Backpropagation

- ▶ No hay una demostración formal de la convergencia del algoritmo.
- ▶ El algoritmo es lento y los estimadores muchas veces inestables.
- ▶ Existen problemas de **identificabilidad** en los pesos.
- ▶ En la práctica se **deja correr** el algoritmo hasta que se estabilice.
- ▶ Valores grandes de η aceleran la convergencia.



- Modos de aprendizaje.



Consideraciones

- ▶ Modos de aprendizaje.
- ▶ Re-escalamiento.



- ▶ Modos de aprendizaje.
- ▶ Re-escalamiento.
- ▶ Cantidad de capas ocultas y nodos ocultos.



- ▶ Modos de aprendizaje.
- ▶ Re-escalamiento.
- ▶ Cantidad de capas ocultas y nodos ocultos.
- ▶ Inicialización de los pesos.



- ▶ Modos de aprendizaje.
- ▶ Re-escalamiento.
- ▶ Cantidad de capas ocultas y nodos ocultos.
- ▶ Inicialización de los pesos.
- ▶ Poda y sobreajuste.

