

BackProp específico para clasificación

- Se suele utilizar la función de activación *softmax* en la capa de salida.
- Dado un conjunto de entrenamiento $S = \{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$, con $\mathbf{x}_n \in \mathbb{R}^{M_0}$, $\mathbf{t}_n \in \{0, 1\}^{M_2}$, ($M_2 \equiv C$), esto permite establecer como criterio de optimización, alternativo al error cuadrático, la **entropía cruzada**:

$$q_S(\Theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{l=1}^{M_2} t_{n,l} \log s_l^2(\mathbf{x}_n; \Theta)$$

- *Problema de entrenamiento*: encontrar Θ tal que la **entropía cruzada** sea mínima.
Solución: **DESCENSO POR GRADIENTE**:

$$\Delta \theta_{ij}^l = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^l} \quad 1 \leq l \leq 2, \quad 1 \leq i \leq M_l, \quad 0 \leq j \leq M_{l-1}$$

Derivación del BackProp para clasificación

- For $N = 1$, $q_S(\Theta) = - \sum_{l=1}^{M_2} t_l \log s_l^2(\mathbf{x}; \Theta)$
- Actualización de los pesos de la capa de salida $\theta_{i,j}^2$ para una muestra genérica $(\mathbf{x}, \mathbf{t}) \equiv (\mathbf{x}_n, \mathbf{t}_n)$

$$q(\Theta) \equiv q_n(\Theta) = - \sum_{l=1}^{M_2} \left(t_l \log s_l^2 \right); \quad s_l^2 = g \left(\phi_l^2 \right); \quad \phi_l^2 = \sum_{m=0}^{M_1} \theta_{lm}^2 s_m^1$$

$$\begin{aligned} \frac{\partial q}{\partial \theta_{ij}^2} &= \frac{\partial q}{\partial s_i^2} \frac{\partial s_i^2}{\partial \theta_{ij}^2} = \frac{\partial q}{\partial s_i^2} \frac{d s_i^2}{d \phi_i^2} \frac{\partial \phi_i^2}{\partial \theta_{ij}^2} \\ &\quad \downarrow \quad \downarrow \quad \downarrow \\ &= - \left(\frac{t_i}{s_i^2} g'(\phi_i^2) \right) s_j^1 \stackrel{\text{def}}{=} -\delta_i^2 s_j^1 \end{aligned}$$

$$\frac{\partial q}{\partial \theta_{ij}^2} = -\delta_i^2 s_j^1, \quad \delta_i^2 \stackrel{\text{def}}{=} \frac{t_i}{s_i^2} g'(\phi_i^2)$$

$$\Delta \theta_{i,j}^2 = -\rho \frac{\partial q_S}{\partial \theta_{i,j}^2} = \rho \frac{t_i}{s_i^2} f'(z_i^2) s_j^1 = \rho \delta_i^2 s_j^1$$

- Actualización de los pesos de la capa oculta $\theta_{i,j}^1$: idéntica al BackProp para regresión.

Algoritmo BACKPROP para clasificación

Entrada: Una topología, datos de entrenamiento S , un factor de aprendizaje ρ , pesos iniciales θ_{ij}^l

$1 \leq l \leq L, 1 \leq i \leq M_l, 1 \leq j \leq M_{l-1}$ y condiciones de convergencia

Salidas: Pesos de las conexiones que minimizan el error cuadrático medio de S

Método:

Mientras no converja

Para $1 \leq l \leq L, 1 \leq i \leq M_l, 0 \leq j \leq M_{l-1}$, inicializar $\Delta\theta_{ij}^l = 0$ Fin-para

Para toda muestra de aprendizaje $(\mathbf{x}, \mathbf{t}) \in S$

Desde la capa de entrada a la capa de salida ($l = 1, \dots, L$):

Para $1 \leq i \leq M_l$: Calcular $\phi_i^l(\mathbf{x})$ y $s_i^l(\mathbf{x}) = g(\phi_i^l(\mathbf{x}))$ Fin-para

Desde la salida a la entrada ($l = L, \dots, 1$),

Para cada nodo ($1 \leq i \leq M_l$)

$$\text{Calcular } \delta_i^l(\mathbf{x}) = \begin{cases} g'(\phi_i^l(\mathbf{x})) \frac{t_i}{s_i^L(\mathbf{x})} & \text{si } l == L \\ g'(\phi_i^l(\mathbf{x})) (\sum_r \delta_r^{l+1}(\mathbf{x}) \theta_{ri}^{l+1}) & \text{en otro caso} \end{cases}$$

Para $0 \leq j \leq M_{l-1}$: Calcular: $\Delta\theta_{ij}^l += \delta_i^l(\mathbf{x}) s_j^{l-1}(\mathbf{x})$ Fin-para

Fin-para cada nodo

Fin-para cada muestra

Para $1 \leq l \leq L, 1 \leq i \leq M_l, 0 \leq j \leq M_{l-1}$, Actualizar pesos: $\theta_{ij}^l += \rho \Delta\theta_{ij}^l$ Fin-para

Fin-mientras

Algoritmo BACKPROP (incremental) para clasificación

Entrada: Una topología, datos de entrenamiento S , un factor de aprendizaje ρ , pesos iniciales θ_{ij}^l $1 \leq l \leq L$, $1 \leq i \leq M_l$, $1 \leq j \leq M_{l-1}$ y condiciones de convergencia

Salidas: Pesos de las conexiones que minimizan el error cuadrático medio de S

Método:

Mientras no converja

Para toda muestra de aprendizaje $(x, t) \in S$

Desde la capa de entrada a la capa de salida ($l = 1, \dots, L$):

Para $1 \leq i \leq M_l$, calcular $\phi_i^l(x)$ y $s_i^l(x) = g(\phi_i^l(x))$ Fin-para

Desde la salida a la entrada ($l = L, \dots, 1$),

Para cada nodo ($1 \leq i \leq M_l$)

Calcular $\delta_i^l(x) = \begin{cases} g'(\phi_i^l(x)) \frac{t_i}{s_i^L(x)} & \text{si } l == L \\ g'(\phi_i^l(x)) (\sum_r \delta_r^{l+1}(x) \theta_{ri}^{l+1}) & \text{en otro caso} \end{cases}$

Para $0 \leq j \leq M_{l-1}$, calcular: $\Delta\theta_{ij}^l = \delta_i^l(x) s_j^{l-1}(x)$ Fin-para

Fin-para

Para $1 \leq l \leq L$, $1 \leq i \leq M_l$, $0 \leq j \leq M_{l-1}$: actualizar $\theta_{ij}^l += \rho \Delta\theta_{ij}^l$ Fin-para

Fin-para

Fin-mientras