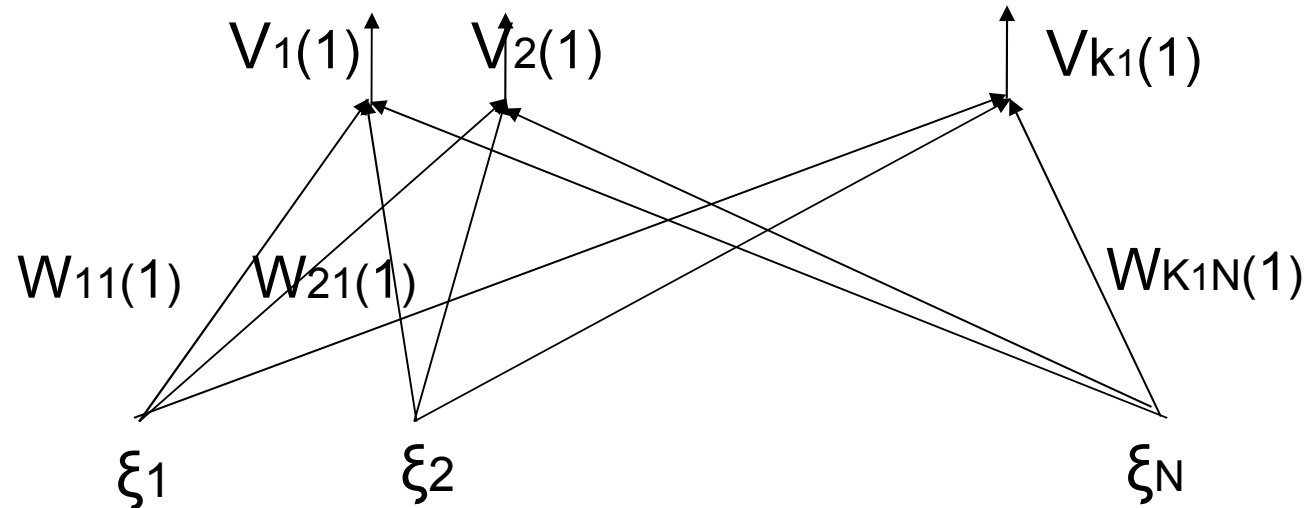
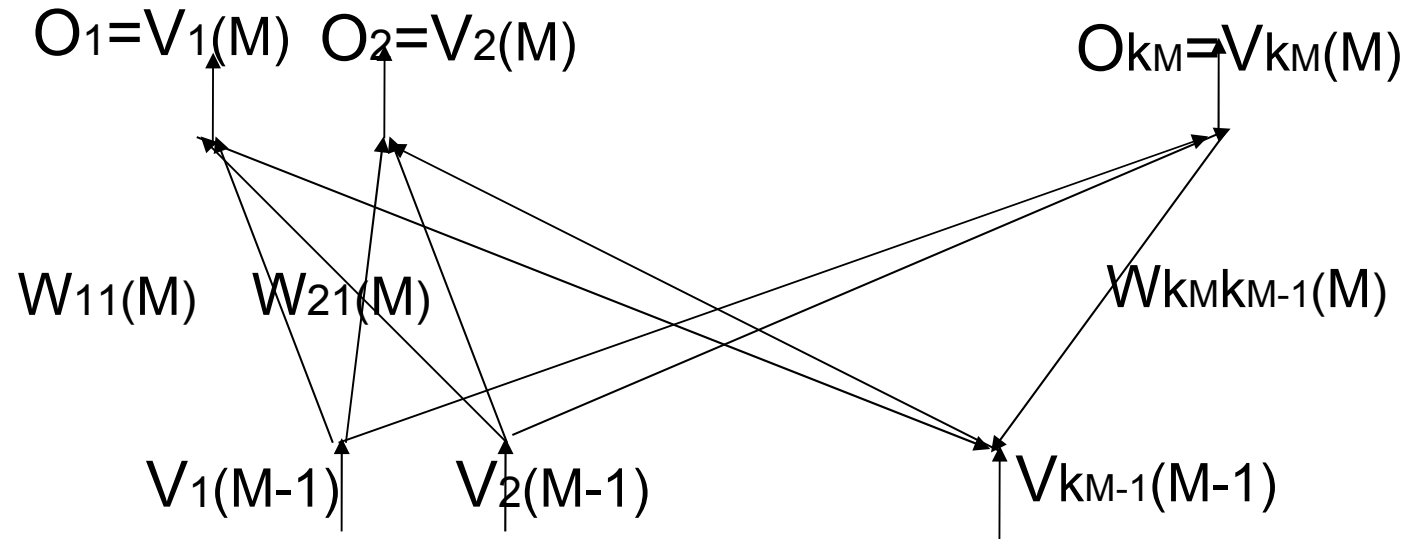


PERCEPTRON MULTICAPA



Si $M=1$ (perceptrón simple) $\begin{cases} \bullet \text{fácil de entrenar} \\ \bullet \text{Soluciona familia de problemas restringida} \end{cases}$

Si $M \geq 2$ (perceptrón multicapa) $\begin{cases} \bullet \text{mucho más difícil de entrenar} \\ \quad (\text{durante mucho tiempo no se supo cómo}) \\ \bullet \text{puede representar cualquier función booleana y, más aun, continua en general} \end{cases}$

Teorema (Funahashi 1989):

Sean: g no constante, acotada y monótona creciente; $C \subset \mathbb{R}^N$ compacto

$$f: C \rightarrow \mathbb{R}^L$$

Dado $\varepsilon > 0$, existe K entero y constantes reales W_{ij}, θ_j ($j=1 \dots K, i=1 \dots L$)

w_{ik} ($i=1 \dots K, k=1 \dots N$) tales que si

$$\tilde{f}_i(\xi_1, \dots, \xi_N) := \sum_{j=1}^K W_{ij} g\left(\sum_{k=1}^N w_{jk} \xi_k - \theta_j\right)$$

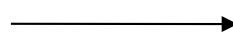
entonces

$$\|f - \tilde{f}\|_{L^\infty(C)} = \max_{\xi \in C} |f(\xi_1, \dots, \xi_N) - \tilde{f}(\xi_1, \dots, \xi_N)| < \varepsilon$$

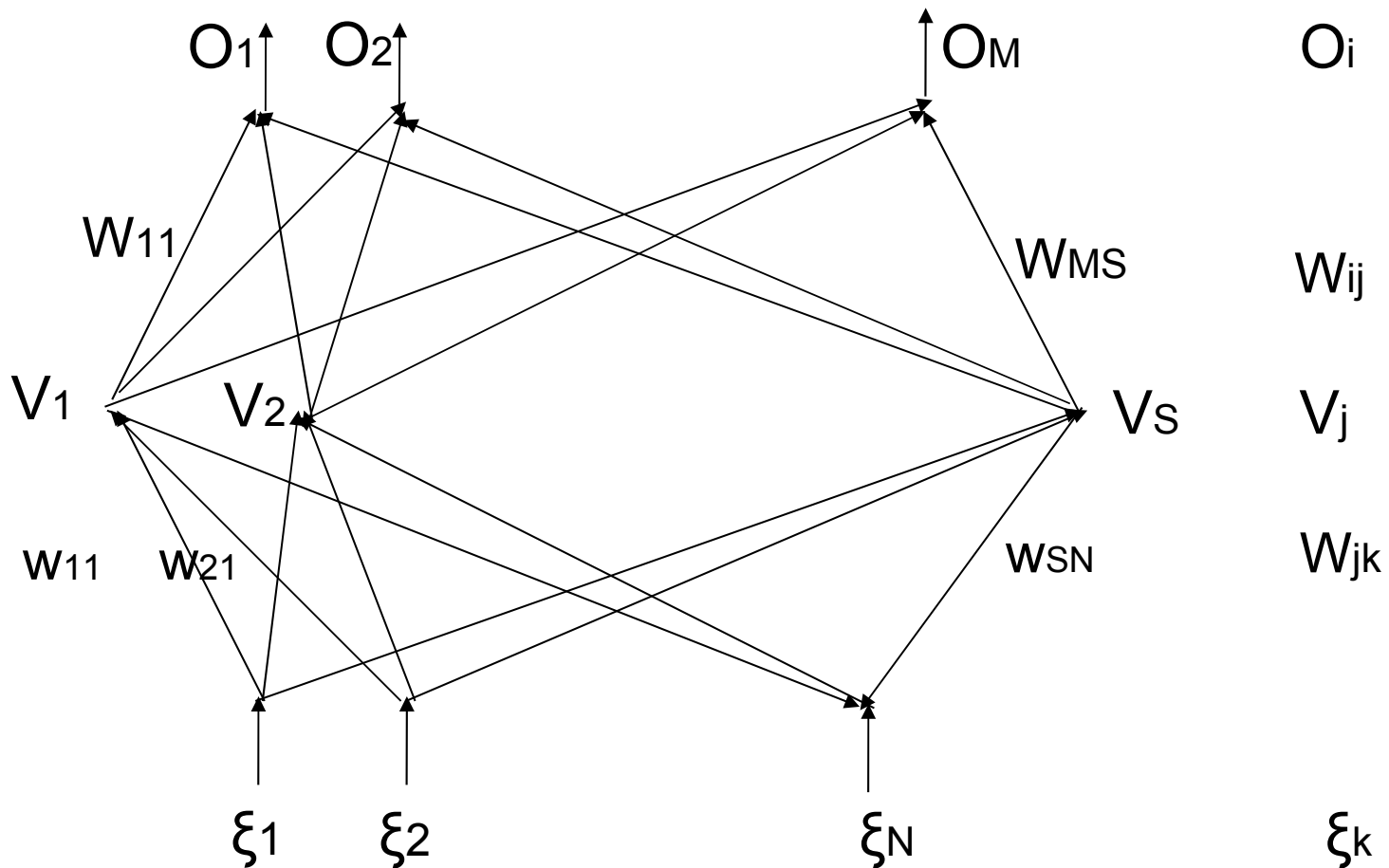
(vale para cualquier métrica L^p , en particular para la cuadrática)

NOTACION Y DEFINICIONES

Tomemos $M=2$



Perceptrón bicapa (una capa oculta)



Patrones: pares $(\xi(\mu), \zeta(\mu))$ $\mu=1, \dots, p$

$\xi(\mu) = (\xi_1(\mu), \xi_2(\mu), \dots, \xi_N(\mu))$
 $\zeta(\mu) = (\zeta_1(\mu), \zeta_2(\mu), \dots, \zeta_M(\mu))$

Para el patrón μ como entrada, la salida será:

$$Q_i^\mu = g(h_i^\mu) = g\left(\sum_{j=1}^S W_{ij} v_j^\mu\right) = g\left(\sum_{j=1}^S W_{ij} g\left(\sum_{k=1}^N w_{jk} \xi_k^\mu\right)\right) \quad c=1, \dots, L$$

En notación matricial:

$$\vec{Q}^\mu = g(W \vec{V}^\mu) = g(W g(w \vec{\xi}^\mu))$$

Función de costo o error:

$$E(w, W) = \frac{1}{2} \sum_{\mu} (\xi_i^\mu - Q_i^\mu)^2 = \frac{1}{2} \sum_{\mu} \left(\xi_i^\mu - g\left(\sum_j W_{ij} g\left(\sum_k w_{jk} \xi_k^\mu\right)\right) \right)^2$$

cuya continuidad y diferenciabilidad dependerán de g .

Pediremos g al menos diferenciable.

EL ALGORITMO

g diferenciable $\Rightarrow E(W)$ diferenciable \Rightarrow puede aplicarse descenso por gradiente

$$\Rightarrow \Delta W = -\eta \nabla E \quad (\eta \text{ velocidad})$$

Conexiones capa oculta - capa de salida:

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}} = \eta \sum_{\mu} (\xi_i^{\mu} - o_i^{\mu}) g'(h_i^{\mu}) v_j^{\mu} = \eta \sum_{\mu} \delta_i^{\mu} v_j^{\mu}$$

$$\text{donde } \delta_i^{\mu} := (\xi_i^{\mu} - o_i^{\mu}) g'(h_i^{\mu}) \quad (1)$$

\Rightarrow regla δ , igual a un perceptrón simple, como si los v_j^{μ} fueran las entradas.

Conexiones entrada - capa oculta:

$$\Delta W_{jk} = -\eta \frac{\partial E}{\partial W_{jk}} = -\eta \sum_{\mu=1}^P \frac{\partial E}{\partial v_j^{\mu}} \frac{\partial v_j^{\mu}}{\partial W_{jk}} =$$

$$= \eta \sum_{\mu} \left(\sum_i \underbrace{(\xi_i^{\mu} - o_i^{\mu}) g'(h_i^{\mu}) W_{ij}}_{\delta_i^{\mu}} g'(h_j^{\mu}) \xi_k^{\mu} \right) = \eta \sum_{\mu} \delta_i^{\mu} W_{ij} g'(h_j^{\mu}) \xi_k^{\mu} = \eta \sum_{\mu} \delta_j^{\mu} \xi_k^{\mu}$$

siendo ahora

$$\delta_j^\mu := g'(u_j^\mu) \sum_i W_{ij} \delta_i^\mu \quad (2)$$

En general, para cualquier número de capas vale:

$$\Delta w_{pj} = \eta \sum_\mu \delta_{out} \cdot V_{in}$$

V_{in} entradas de la capa anterior o entradas reales

δ_{out} como en (1) o en (2) dependiendo de si es la última capa de conexiones o una anterior.

Observación: los δ de una capa oculta se calculan a partir de los de las unidades que esa capa alimenta (de ahí el nombre de **error backpropagation**).

IMPLEMENTACIÓN

Aprendizaje $\left\{ \begin{array}{l} \text{Sincrónico: primero se calculan todas las salidas } (v_i) \\ \text{y luego todos los } \delta \text{ } (\Rightarrow \text{batch}) \\ \text{Asincrónico: se entrena con un patrón por vez} \end{array} \right.$

g habituales $\left\{ \begin{array}{l} g_1(h) = f_{\beta}(h) = \frac{1}{1 + \exp(-2\beta h)} \Rightarrow \text{rango } (0, 1) \\ g_2(h) = \tanh \beta h \Rightarrow \text{rango } (-1, 1) \end{array} \right.$

Se cumple entonces: $\begin{array}{l} g_1'(h) = 2\beta g_1(1-g_1) \\ g_2'(h) = \beta(1-g_2^2) \end{array} \Rightarrow \text{facilitan el cálculo de los } \delta_i$

Pasos a seguir (versión asincrónica o secuencial)

(mantenemos la notación: M número de capas

V_i^m salida de la i -ésima neurona de la m -ésima capa

$$V_i^0 = \xi_i$$

w_{ij}^m : conexión de V_j^{m-1} a V_i^m

1. Inicializar los w (pequeños y al azar)

2. Elegir un patrón (μ) ; $V_i^0 = \xi_i^\mu \quad \forall i$

3. Etapa forward: $V_i^m = g(h_i^m) = g\left(\sum_j w_{ij}^m V_j^{m-1}\right) \quad \forall i, m$ hasta los V_i^M finales

4. $\delta_i^M = g'(h_i^M) (\xi_i^\mu - V_i^M)$ deltas de la capa de salida para el patrón considerado

5. Etapa backward: retropropagación de errores

$$\delta_i^{m-1} = g'(h_i^{m-1}) \sum_j w_{ji}^m \delta_j^m \quad m = M, M-1, \dots, 2$$

$$\Delta w_{ij}^m = \eta \delta_i^m V_j^{m-1}$$

$$w_{ij}^{nue} = w_{ij}^{vie} + \Delta w_{ij}$$

7. Ir a 2 (seleccionar patrón)

EXTENSIONES Y VARIANTES

Dos defectos principales de BP: *lento*

Mínimos locales

Algunas mejoras:

- *Momento*:

$$\Delta w_{ij}(n+1) = -\eta \partial E / \partial w_{ij} + \alpha \Delta w_{ij}(n)$$

--> Si superficie de costo plana, acelera en un factor $1/(1-\alpha)$
Si hay oscilaciones, las fluctuaciones son escaladas por η

- *Parámetros adaptivos*:

$$\Delta \eta = \begin{cases} +a & \text{si } \Delta E < 0 \text{ en los últimos pasos} \rightarrow \text{crece aritméticamente} \\ -\beta \eta & \text{si } \Delta E > 0 \rightarrow \text{decrece geométricamente} \\ 0 & \text{en otro caso} \end{cases}$$

Si $\Delta E > 0 \rightarrow \eta$ decrece \rightarrow se anula la modificación

$a = 0$ hasta un paso exitoso (si se estaba usando momento)

- *Otras técnicas determinísticas*: steepest descent

Gradientes conjugados

Quasi-Newton

- *Técnicas estocásticas*

- Gradientes conjugados: $d^{(n+1)} = -\nabla E^{(n+1)} + \beta^{(n+1)} d^{(n)}$

Cada nueva dirección es un compromiso entre la del gradiente y la anterior.

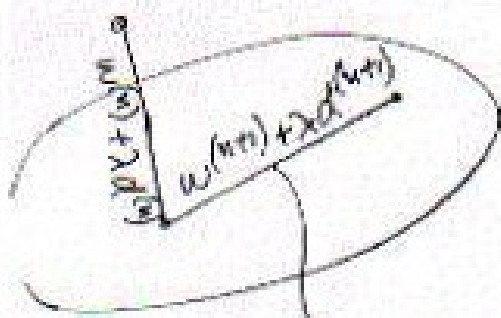
Condición: $\beta^{(n+1)}$ tal que $d^{(n+1)}$ no cambie la componente del gradiente a lo largo de la dirección previa (componente que ya era nula):

$$d^{(n)} \cdot \nabla E(w^{(n+1)} + \lambda d^{(n)}) = 0$$

$$d^{(n)} \cdot (\nabla E(w^{(n+1)}) + H \lambda d^{(n)}) = 0$$

$$\underbrace{d^{(n)} \cdot \nabla E(w^{(n+1)})}_0 + \lambda d^{(n)} H d^{(n)} = 0$$

$$\Leftrightarrow d^{(n)} H d^{(n+1)} = 0 \quad d^{(n)} \text{ y } d^{(n+1)} \text{ vectores conjugados}$$



Parado en un punto de ese segmento, la dirección del gradiente debe ser \perp a la del paso anterior

Posibles β :

- Fletcher-Reeves ('64, la original):

$$\beta^{(n+1)} = \frac{(\nabla E^{(n+1)})^2}{(\nabla E^{(n)})^2}$$

- Polak-Ribière:

$$\beta^{(n+1)} = \frac{(\nabla E^{(n+1)} - \nabla E^{(n)})^T \cdot \nabla E^{(n+1)}}{(\nabla E^{(n)})^2}$$

- Las últimas N (dimension) direcciones son mutuamente conjugadas.
- No requiere conocer H .
- Superior a BP y SD en general. Encuentra el mínimo de una sup. cuadrática en N pasos.
- Computacionalmente más complejo. Sensible al λ de cada paso.