

## e-Capítulo 7

# Redes neuronales

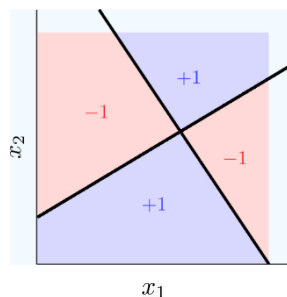
Las redes neuronales son un modelo de "inspiración biológica" que ha tenido un éxito considerable en aplicaciones que van desde la predicción de series temporales hasta la visión. Las redes neuronales son una generalización del perceptrón que utiliza una transformación de características que se *aprende a* partir de los datos. Como tales, son un modelo muy potente y flexible.

Las redes neuronales son un buen modelo candidato para aprender de los datos porque pueden aproximar eficazmente funciones objetivo complejas y vienen con buenos algoritmos para ajustar los datos. Comenzaremos con las propiedades básicas de las redes neuronales y cómo entrenarlas a partir de datos. Presentaremos una serie de técnicas útiles para ajustar los datos minimizando el error en la muestra. Dado que las redes neuronales son un modelo muy flexible, con un gran poder de aproximación, es fácil sobreajustar los datos; estudiaremos una serie de técnicas para controlar el sobreajuste específicas de las redes neuronales.

### 7.1 El perceptrón multicapa (MLP)

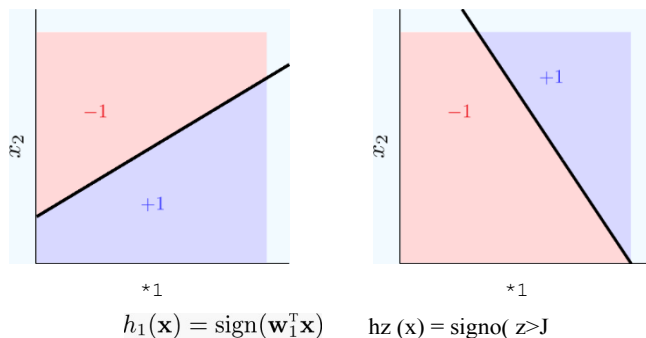
El perceptrón no puede implementar funciones de clasificación sencillas. Para ilustrarlo, utilizamos el target de la derecha, que está relacionado con la función booleana XoR. En este ejemplo, / no puede escribirse como  $\text{sign}(w \cdot x)$ . Sin embargo, / se compone de dos partes lineales. De hecho, como pronto veremos, podemos descomponer J en dos perceptrones simples, correspondientes a las líneas de la figura, y luego combinar las salidas de estos dos perceptrones de forma sencilla para recuperar J.

Los dos perceptrones se muestran a continuación.



La analogía con las neuronas biológicas, aunque inspiradora, no debe llevarse demasiado lejos; al fin y al cabo, construimos aviones con alas que no aletean. Del mismo modo, las redes neuronales, cuando se aplican al aprendizaje a partir de datos, no se parecen demasiado a sus homólogas biológicas.

Ya ser Abu M osta fa Malik M a gdon lsm ail, Hsuan Tien Lin: Enero de 2015.  
Todos los derechos reservados. Prohibido el uso comercial o la redistribución de  
cualquier forma.

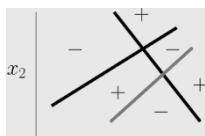


El objetivo  $J$  es igual a  $-1-1$  cuando exactamente uno de  $f_1$ ,  $f_2$  es igual a  $-1-1$ . Esta es la función booleana XOR:  $f = \text{XOR}(f_1, f_2)$  donde  $-1-1$  representa "VERDADERO" y  $-1$  representa "FALSO". Podemos reescribir  $J$  utilizando las operaciones más sencillas OR y AND:  $f_1 = +1$  si al menos una de  $f_1$ ,  $f_2$  es igual a  $-1-1$  y  $\text{AND}(f_1, f_2) = -1-1$  si ambas  $f_1$ ,  $f_2$  son iguales a  $-1-1$ . Utilizando la notación booleana estándar (multiplicación por AND, suma para OR y barra para negación),

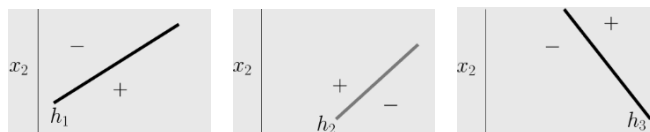
$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$

### Ejercicio 7.1

Consideremos una función objetivo / cuyas regiones "+" y "-" se ilustran a continuación.



El objetivo / tiene tres componentes de perceptrón  $h_1, h_2, h_3$ .



Demuestra que

$$f = \overline{h_1} h_2 h_3 + h_1 \overline{h_2} h_3 + h_1 h_2 \overline{h_3}$$

¿Existe una forma sistemática de pasar de un objetivo que es una descomposición de perceptrones a una fórmula booleana como ésta? [Pista: considere sólo las regiones de / que son '-' y utilice la forma normal disyuntiva (o de ANDS).]

El ejercicio 7.1 muestra que un objetivo complicado, que se compone de perceptrones,

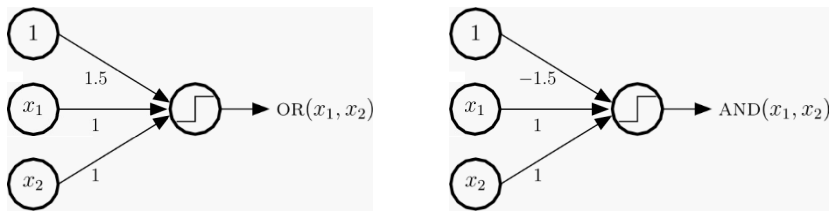
---

es una disyunción de conjunciones (OR Of ANDs) aplicadas al componente

perceptrones. Se trata de una idea útil, porque el perceptrón puede aplicar oR y AND:

$$\text{OR}(z_1, z_2) = \text{signo}(z_1 + z_2 - 3)$$
$$\text{AND}(x_1, x_2) = \text{sign}(x_1 + x_2 - 1.5).$$

Esto implica que estos objetivos más complicados son, en última instancia, meras combinaciones de perceptrones. Para ver cómo combinar los perceptrones para obtener  $f$ , introducimos una representación gráfica de perceptrones, empezando por OR y AND:

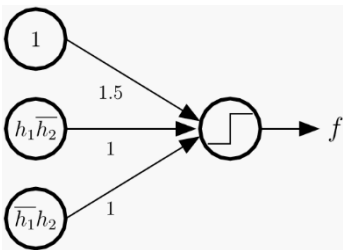


Un nodo envía un valor a una flecha. El peso de una flecha multiplica esta salida y pasa el resultado al siguiente nodo. Todo lo que llega a este siguiente nodo se suma y luego se transforma mediante  $\text{sign}()$  para obtener la salida final.

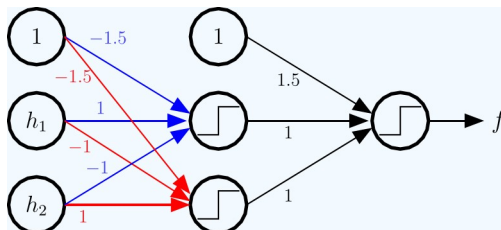
Ejercicio 7.2

- (a) Los booleanos or y AND de dos entradas pueden ampliarse a más de dos entradas:  $\text{OR}(z_1, \dots, z_M)$  es 1 si alguna de las  $M$  entradas es 1;  $\text{AND}(z_1, \dots, z_M)$  es 1 si todas las entradas son iguales a 1. Dar gráfico representaciones de  $\text{OR}(z_1, \dots, z_M)$  y  $\text{AND}(z_1, \dots, z_M)$ .
- (b) Dar la representación gráfica del perceptrón:  $f(x) = \text{sign}(w \cdot x)$ .
- (c) Dar la representación gráfica de  $\text{OR}(z_1, \dots, z_M)$ .

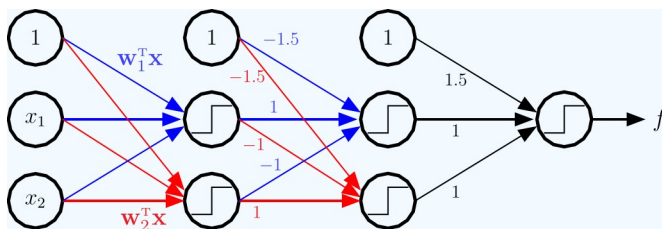
En un MLP para un objetivo complejo. Dado que  $J = b_1 z_1 + b_2 z_2$ , que es un OR de las dos entradas  $z_1$  y  $z_2$ , utilizamos primero el perceptrón OR, para obtener:



Las dos entradas  $\text{fit}_2$  y  $\text{fit}_2$  son AND. Como tales, pueden simularse mediante la salida de dos perceptrones AND. Para tratar la negación de las entradas del AND, negamos los pesos multiplicando las entradas negadas (como se ha hecho en el Ejercicio 7.1(c)). La representación gráfica resultante de  $J$  es:



Los pesos azul y rojo simulan los dos AND requeridos. Por último, como  $\text{fi} = \text{sign}(w^T x)$  y  $\text{fi}_2 = \text{sign}(w_2^T x)$  son perceptrones, ampliamos los nodos  $\text{fi}$  y  $\text{fi}_2$  para obtener la representación gráfica de  $J$ .



El siguiente ejercicio te pide que calcules una fórmula algebraica explícita para  $J$ . La representación gráfica visual es mucho más ordenada y fácil de generalizar.

### Ejercicio 7.3

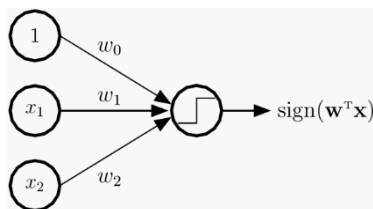
Utiliza la representación gráfica para obtener una fórmula explícita para  $J$  y demuéstalo:

$$J(x) = \text{sign}(\text{sign}(\#1(x) \# z(x)) \text{sign}(\#1(x) \# z(x) + j)),$$

donde  $\#1(x) = \text{sign}(w_1 x)$  y  $\#z(x) =$

$\text{sign}(z)$

Comparemos la forma gráfica de  $J$  con la forma gráfica del perceptrón simple, que se muestra a la derecha. Se utilizan más capas de nodos entre la entrada y la salida para implementar  $J$ , en comparación con el perceptrón simple, por lo que lo llamamos perceptrón multicapa (MLP). Las capas adicionales



---

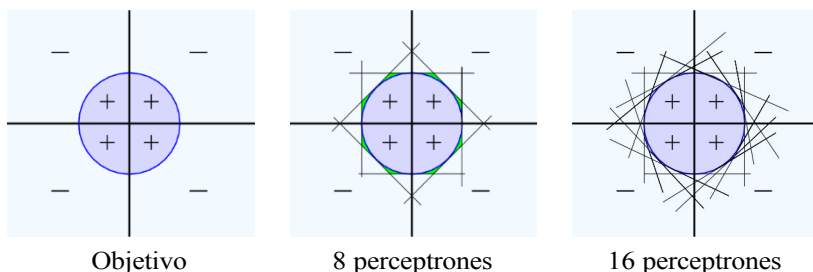
se denominan *capas ocultas*.

Observe que las capas sólo alimentan hacia adelante a la capa siguiente (no hay flechas que apunten hacia atrás ni saltos a otras capas). La capa de entrada (más a la izquierda) no se cuenta como una capa, por lo que en este ejemplo hay 3 capas (2 capas ocultas con 3 nodos cada una y una capa de salida con 1 nodo). El perceptrón simple no tiene capas ocultas, sólo entrada y salida. La adición de capas ocultas es lo que nos permitió implementar el objetivo más complicado.

### Ejercicio 7.4

Para la función objetivo del Ejercicio 7.1, dé el MLP en forma gráfica, así como la forma algebraica explícita.

Si  $J$  puede descomponerse en perceptrones utilizando un oR de ANDS, entonces puede implementarse mediante un perceptrón de 3 capas. Si  $J$  no es estrictamente descomponible en perceptrones, pero la frontera de decisión es suave, entonces un perceptrón de 3 capas puede acercarse arbitrariamente a la implementación de  $J$ . A continuación se presenta una ilustración de "prueba por imágenes" para una función objetivo de disco:



La demostración formal es análoga al teorema del cálculo según el cual cualquier función continua en un conjunto compacto puede aproximarse de forma arbitraria mediante funciones escalonadas. El perceptrón es el análogo de la función escalón. Así pues, hemos encontrado una generalización del perceptrón simple que se parece mucho al propio perceptrón simple, salvo por la adición de más capas. Hemos adquirido la capacidad de modelar funciones objetivo más complejas añadiendo más nodos (unidades *ocultas*) en las capas ocultas esto corresponde a permitir más perceptrones en la descomposición de  $J$ . De hecho, un MLP de 3 capas adecuadamente grande puede aproximar casi cualquier función objetivo y ajustarse a cualquier conjunto de datos, por lo que es un modelo de aprendizaje muy potente. Utilícelo con cuidado. Si su MLP es demasiado grande puede perder capacidad de generalización.

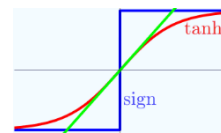
Una vez fijado el tamaño del MLP (número de capas ocultas y número de unidades ocultas en cada capa), se aprenden los pesos de cada enlace (flecha) ajustando los datos. Consideremos el perceptrón simple,



//  $(X_t - @ (W "X) .$

Cuando  $\delta(s) = \text{sign}(s)$ , aprender los pesos ya era un problema combinatorio difícil y contaba con una variedad de algoritmos, incluido el algoritmo de bolsillo, para ajustar los datos (Capítulo 3). El problema de optimización combinatoria es aún más difícil con el MLP, por la misma razón, a saber, que la función  $\text{sign}(\cdot)$  no es suave; una aproximación suave y diferenciable a  $\text{sign}(\cdot)$  nos permitirá utilizar métodos analíticos, en lugar de métodos puramente combinatorios, para encontrar los pesos óptimos. Por lo tanto, aproximamos o "suavizamos" la función  $\text{sign}(\cdot)$  utilizando la función  $\tanh(\cdot)$ . La MLP se denomina a veces red neuronal de umbral (duro) porque la función de transformación es un umbral duro en cero. Aquí, elegimos  $\delta(z) = \tanh(z)$  que es in-

entre lineal y el umbral duro: casi lineal para  $z \approx 0$  y casi  $-1$  para  $|z|$  grande. La función  $\tanh(\cdot)$  es otro ejemplo de *sigmoide* (porque su forma se parece a una "s" aplanada), relacionada con la sigmoide que utilizamos para la regresión logística.<sup>2</sup> Estas redes se denominan redes neurales sigmoidales. Del mismo modo que podríamos utilizar los pesos

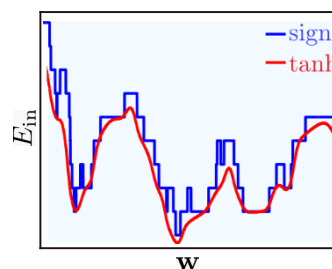


aprendidas de la regresión lineal para la clasificación, podríamos utilizar pesos aprendidos utilizando la red neuronal sigmoideal con función de activación  $\tanh(\cdot)$  para la clasificación sustituyendo la función de activación de salida por  $\text{sign}(\cdot)$ .

### Ejercicio 7.5

Dados  $w_1$  y  $\epsilon > 0$ , hallar  $w_2$  tal que  $|\text{sign}(w_1 x_n) - \tanh(w_2 x_n)| < \epsilon$  para  $x_n \in \mathcal{X}$ .  
 [Pista: Para  $\epsilon$  pequeño,  $\text{sign}(z) \approx \tanh(oz)$  para  $o$  grande.]

El ejemplo anterior muestra que la función  $\text{sign}(\cdot)$  puede aproximarse mucho a la función  $\tanh(\cdot)$ . Una ilustración concreta de esto se muestra en la figura de la derecha. La figura muestra cómo varía el error en la muestra  $E_{in}$  con uno de los pesos en  $w$  en un problema de ejemplo para el perceptrón (curva azul) en comparación con la versión sigmoideal (curva roja). La aproximación sigmoideal captura el general forma del error, de modo que si minimizamos el error sigmoideal dentro de la muestra, obtenemos una buena aproximación a la minimización del error de clasificación dentro de la muestra.

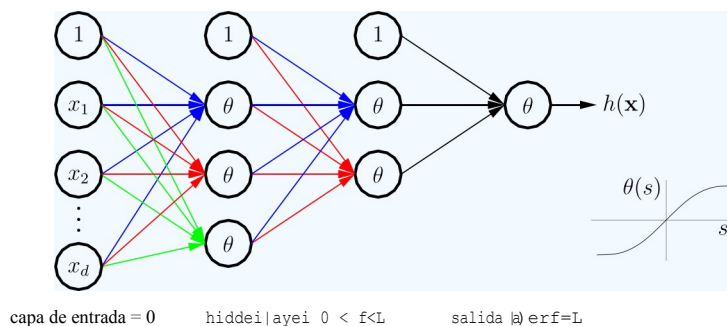


---

2En la regresión logística, utilizamos la sigmoidea porque queríamos una probabilidad como salida. Aquí, utilizamos el  $\tanh(-)$  'suave' porque queremos una función objetivo amigable para optimizar.

## 7.2 Redes neuronales

La red neuronal es nuestro MLP "suavizado". Comencemos con una representación gráfica de una *red neuronal de avance* (el único tipo que consideraremos).



La representación gráfica muestra una función en nuestro conjunto de hipótesis. Aunque esta visión gráfica es estética e intuitiva, con la información "fluyendo" desde las entradas en el extremo izquierdo, a lo largo de los enlaces y a través de los nodos ocultos, en última instancia a la salida  $f(\mathbf{x})$  en el extremo derecho, será necesario describir algorítmicamente la función que se está calculando. Las cosas se van a complicar y esto requiere una notación muy sistemática.

### 7.2.1 Notación

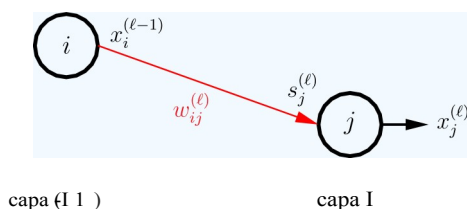
Hay capas etiquetadas por  $f_i = 0, 1, 2, \dots$ . En nuestro ejemplo anterior,  $L = 3$ , es decir, tenemos tres capas (la capa de entrada  $f = 0$  no suele considerarse una capa y sirve para alimentar las entradas). La capa  $f = L$  es la capa de salida, que determina el valor de la función. Las capas intermedias,  $0 < f < L$ , son las capas ocultas. Utilizaremos el superíndice "para referirnos a una capa concreta. Cada capa  $f$  tiene una "dimensión"  $d_f$  lo que significa que tiene  $d_f - 1$  nodos, etiquetados  $0, 1, \dots, d_f - 1$ . Cada capa tiene un nodo especial, que se llama nodo  $\theta$  (etiquetado 0). Este nodo de sesgo está configurado para tener una salida 1, que es análoga a la convención ficticia  $z_0 = 1$  que teníamos para los modelos lineales.

Cada flecha representa un peso o fuerza de conexión de un nodo en una capa a un nodo en la capa  $f+1$ . Nótese que los nodos de sesgo no tienen pesos entrantes. No hay otros pesos de conexión. Un nodo con un peso entrante indica que alguna señal se introduce en este nodo. Cada nodo de este tipo con una entrada tiene un  $w_{f+1,f}$ . Si  $s_f = \text{señal}(s)$ , entonces tenemos el MLP para la clasificación. Como hemos mencionado antes, utilizaremos una versión suave del MLP con  $B(z) = \tanh(z)$  para aproximar la función  $\text{sign}(\cdot)$ . El  $\tanh(\cdot)$  es un umbral suave o sigmoide, y ya vimos un sigmoide relacionado

\*En una configuración más general, los pesos pueden conectar dos nodos cualesquiera, además de ir hacia atrás (es decir, puede haber ciclos). Este tipo de redes se denominan redes neuronales recurrentes, y no las consideramos aquí.

cuando hablamos de la regresión logística en el capítulo 3. En última instancia, cuando hacemos clasificación, sustituimos la sigmoide de salida por el umbral duro  $\text{sign}()$ . Como comentario, si estuviéramos haciendo regresión en su lugar, toda nuestra discusión pasa con la transformación de salida que se sustituye por la función de identidad (sin transformación) para que la salida sea un número real. Si estuviéramos haciendo regresión logística, reemplazaríamos el sigmoide de salida  $\tanh()$  por el sigmoide de regresión logística.

El modelo de red neuronal Hpq se especifica una vez que se determina el *arquitectura* de la red neuronal, es decir, la dimensión de cada capa  $d^{(1)}, d^{(2)}, \dots, d^{(L)}$  (donde  $L$  es el número de capas). Una hipótesis finita se especifica seleccionando pesos para los enlaces. Acerquémonos a un nodo de la capa oculta capa  $l$ , para ver qué pesos hay que especificar.

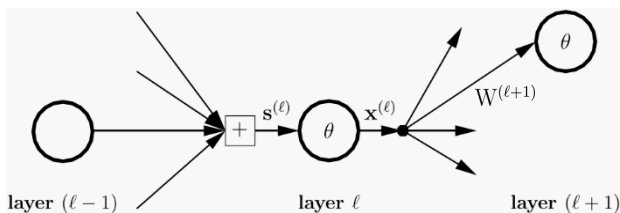


Un nodo tiene una señal de entrada  $s$  y una salida  $z$ . Los pesos de los enlaces al nodo desde la capa anterior están en  $w$ , por lo que los pesos están indexados por la capa a la que van. Así, la salida de los nodos de la capa  $l-1$  se multiplica por los pesos  $w$ . Utilizamos subíndices para indexar los nodos de una capa.

Así, es el peso que entra en el nodo  $j$  de la capa  $l$  desde el nodo  $i$  de la capa anterior, la señal que entra en el nodo  $j$  de la capa  $l$  es  $s_j^{(l)}$ , y la salida de este nodo

es  $z_j^{(l)}$ . Hay algunos nodos especiales en la red. Los nodos cero de cada capa son nodos constantes, con salida 1. No tienen peso entrante, pero sí saliente. Los nodos de la capa de entrada  $l=0$  son para los valores de entrada, y no tienen peso entrante ni función de transformación.

En la mayoría de los casos, sólo necesitamos tratar la red capa por capa, por lo que introducimos la notación vectorial y matricial para ello. Recogemos todas las señales de entrada a los nodos  $1, \dots, d^{(l)}$  en la capa  $l$  en el vector  $s^{(l)}$ . Del mismo modo, recogemos la salida de los nodos  $0, \dots, d^{(l)}$  en el vector  $x^{(l)}$ ; nótese que  $x^{(l)}$   $(1 \times d^{(l)})$  debido al nodo de polarización 0. Hay enlaces que conectan las salidas de todos los nodos de la capa anterior con las entradas de la capa  $l$ . Así, en la capa  $l$ , tenemos a  $(d^{(l-1)} + 1) \times d^{(l)}$  matriz de pesos  $W^{(l)}$ . La entrada  $(i, j)$  de  $W^{(l)}$  va del nodo  $i$  de la capa anterior al nodo  $j$  de la capa  $l$ .



parámetros layer  $\ell$

señales de entrada	$\mathbf{s}^{(\ell)}$	$d^{(\ell)}$ vector de entrada de dimensión $d^{(\ell)}$
salidas	$\mathbf{x}^{(\ell)}$	$d^{(\ell)}$ - 1 vector de salida de dimensión $d^{(\ell)}$
ponderaciones	$W^{(\ell)}$	$(d^{(\ell-1)} - 1) \times d^{(\ell)}$ matriz de dimensión $d^{(\ell)}$
ponderaciones	$W^{(\ell+1)}$	$(d^{(\ell)} - 1) \times d^{(\ell+1)}$ matriz de dimensión $d^{(\ell+1)}$

Después de fijar los pesos  $W^{(\ell)}$  para  $\ell = 1, \dots, \tilde{n}$ , habrá especificado una hipótesis de red neuronal concreta  $f_{\mathbf{C}}^H$ . Recopilamos todas estas matrices de pesos en un único parámetro de peso  $\mathbf{w}$  ( $W^{(1)}, W^{(2)}, \dots, W^{(\tilde{n})}$ ), y a veces escribiremos  $f_{\mathbf{C}}(\mathbf{x}; \mathbf{w})$  para indicar explícitamente la dependencia de la hipótesis respecto a  $\mathbf{w}$ .

## 7.2.2 Propagación hacia delante

La hipótesis  $f_{\mathbf{C}}(\mathbf{x})$  de la red neuronal se calcula mediante el algoritmo *de propagación hacia delante*. En primer lugar, observe que las entradas y salidas de una capa están relacionadas por la función de transformación,

$$\mathbf{x}^{(\ell)} = \frac{1}{\theta(s^{(\ell)})} \cdot \quad (7.1)$$

donde  $\theta(s^{(\ell)})$  es un vector cuyos componentes son  $\theta(s^{(\ell)})$ . Para obtener el vector de entrada en la capa  $\ell$ , calculamos la suma ponderada de las salidas de la capa anterior

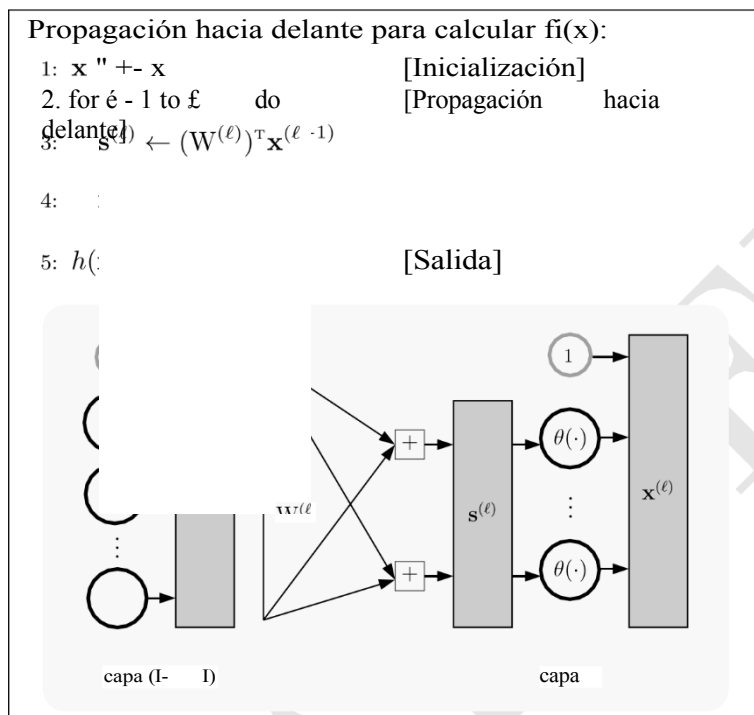
capa, con pesos especificados en  $W^{(\ell)}$ :  $s^{(\ell)} = \sum_{p=1}^{d^{(\ell-1)}} w_{p\ell} x_p^{(\ell-1)}$ . Este proceso se representa de forma compacta mediante la ecuación matricial

$$\mathbf{s}^{(\ell)} = (W^{(\ell)})^T \mathbf{x}^{(\ell-1)}. \quad (7.2)$$

Todo lo que queda es inicializar la capa de entrada a  $\mathbf{x}^{(0)} = \mathbf{x}$  (por lo que  $d^{(0)} = d$ , la dimensión de entrada)<sup>4</sup> y utilizar las ecuaciones (7.2) y (7.1) en la siguiente cadena,

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{W^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{W^{(2)}} \mathbf{s}^{(2)} \xrightarrow{\theta} \mathbf{x}^{(2)} \dots \rightarrow \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x}).$$

'Recordemos que los vectores de entrada también se aumentan con  $z = 1$ .



Tras la propagación hacia delante, el vector de salida  $\mathbf{x}^{(\ell)}$  en cada capa  $\ell = 0, \dots, L$  se ha calculado.

### Ejercicio 7.6

Sean  $V$  y  $@$  el número de nodos y pesos de la red neuronal,

$$d''', \quad @ = d''' (d''' + 1).$$

$\ell=0$    $\ell=1$

En términos de  $V$  y  $@$ , cuántos cálculos se realizan en la propagación hacia delante (sumas, multiplicaciones y evaluaciones de  $\theta$ ).

{Respuesta.  $0(@)$  multiplicaciones y sumas, y  $0(V)$  evaluaciones}.

Si queremos calcular  $A_{ip}$ , sólo necesitamos  $f_i(\mathbf{x}_q)$  y  $p_q$ . Para la suma de cuadrados,

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n; \mathbf{w}) - y_n)^2$$

$$= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n^{(L)} - y_n)^2.$$



Ahora discutiremos cómo minimizar  $E_p$  para obtener los pesos aprendidos. Será una aplicación directa del descenso de gradiente, con un algoritmo especial que calcula el gradiente eficientemente.

### 7.2.3 Algoritmo de retropropagación

En el capítulo 3 estudiamos un algoritmo para llegar a un mínimo local de una superficie de error suave dentro de la muestra, a saber, el descenso de gradiente: inicializar los pesos a  $w(0)$  y para  $t = 1, 2, \dots$  actualizar los pesos dando un paso en la dirección del gradiente negativo,

lo llamamos (por lotes) descenso *por gradiente*. Para implementar el descenso por gradiente, necesitamos el gradiente.

#### Ejercicio 7.7

Para el perceptrón sigmoideal,  $f(x) = \tanh(w \cdot x)$ , dejemos que el error en la muestra sea  $E_{in}(w) = \frac{1}{2} \sum_{n=1}^N (\tanh(w \cdot x_n) - y_n)^2$ . Demuestre que

$$\frac{\partial E_{in}}{\partial w} = - \sum_{n=1}^N (\tanh(w \cdot x_n) - y_n) (1 - \tanh^2(w \cdot x_n)) x_n.$$

Si  $w$  es grande, ¿qué ocurre con el gradiente; cómo se relaciona esto con por qué es difícil optimizar el perceptrón.

Consideremos ahora la red neuronal multicapa sigmoideal con  $\phi(z) = \tanh(z)$ . Como  $f(x)$  es suave, podemos aplicar el descenso de gradiente a la función de error resultante. Para ello, necesitamos el gradiente  $\nabla E_{in}(w)$ . Recordemos que el vector de pesos  $w$  contiene todas las matrices de pesos  $W^{(1)}, \dots, W^{(L)}$ , y necesitamos las derivadas con respecto a todos estos pesos. A diferencia del perceptrón sigmoideal del Ejercicio 7.7, para la red sigmoideal multicapa no existe una expresión simple de forma cerrada para el gradiente. Consideremos un error en la muestra que es la suma de los errores puntuales sobre los puntos de datos (como lo es el error en la muestra al cuadrado),

$$E_{in}(w) = \frac{1}{2} \sum_{n=1}^N e_n^2.$$

donde  $e_n = \phi(f(x_n; w)) - y_n$ . Para el error al cuadrado,  $e = (\phi(f(x; w)) - y)^2$ . Para calcular el gradiente de  $E_{in}$ , necesitamos su derivada con respecto a cada matriz de pesos:

$$\frac{\partial E_{in}}{\partial W^{(\ell)}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial e_n}{\partial W^{(\ell)}}. \quad (7.3)$$

El elemento básico de (7.3) es la derivada parcial del error en el punto de datos  $e_n$ , con respecto al  $W^{(\ell)}$ . Una forma rápida y sucia de obtener  $\frac{\partial e_n}{\partial W^{(\ell)}}$  es

utilizar el enfoque numérico de diferencias finitas. La complejidad de obtener las derivadas parciales con respecto a cada peso es  $O(Q^2)$ , donde  $Q$  es el número de pesos (véase el problema 7.6). A partir de (7.3), tenemos que calcular estas derivadas para cada punto de datos, por lo que el enfoque numérico es prohibitivo desde el punto de vista computacional. Ahora derivamos un elegante algoritmo de programación dinámica conocido como *retropropagación*. La *retropropagación* nos permite calcular las derivadas parciales con respecto a cada peso de manera eficiente, utilizando  $O(Q)$  computación. Describimos la retropropagación para obtener la derivada parcial del error  $e$ , pero el algoritmo es general y se puede utilizar para obtener la derivada parcial de cualquier función de la salida  $f_i(x)$  con respecto a los pesos.

La retropropagación se basa en varias aplicaciones de la regla de la cadena para escribir derivadas parciales en la capa  $\ell$  utilizando derivadas parciales en la capa  $(\ell - 1)$ . Para describir el algoritmo, definimos el *vector de sensibilidad* para la capa  $\ell$ , que es la sensibilidad (gradiente) del error  $e$  con respecto a la señal de entrada  $s^\ell$  que entra en la capa  $\ell$ . Denotamos la sensibilidad por  $d^\ell$ ,

$$d^\ell = \frac{\partial e}{\partial s^\ell}.$$

La sensibilidad cuantifica cómo cambia  $e$  con  $s^\ell$ . Utilizando la sensibilidad, podemos escribir las derivadas parciales con respecto a los pesos  $W^\ell$  como

$$\frac{\partial e}{\partial W^\ell} = x^{(\ell-1)} (d^\ell)^\top. \quad (74)$$

Derivaremos esta fórmula más adelante, pero por ahora examinémosla de cerca. Las derivadas parciales de la izquierda forman una matriz de dimensiones  $(d^{(\ell-1)} \times d^\ell)$  y el "producto exterior" de los dos vectores de la derecha da exactamente una matriz. Las derivadas parciales tienen contribuciones de dos componentes. (i) El vector de salida de la capa  $\ell$  de la que proceden los pesos; cuanto mayor sea la salida, más sensible será  $e$  a los pesos de la capa. (ii) El vector de sensibilidad de la capa  $\ell$  a la que van los pesos; cuanto mayor sea el vector de sensibilidad, más sensible será  $e$  a los pesos de esa capa.

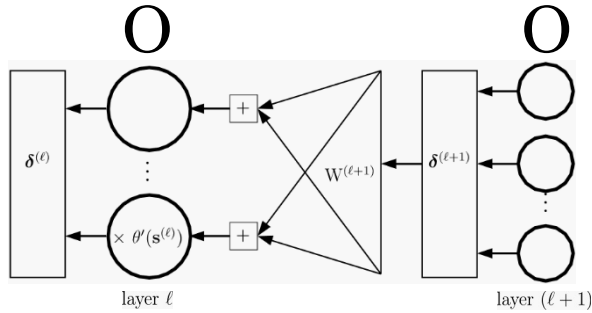
Las salidas  $x^\ell$  para cada capa  $\ell = 1, \dots, L$  pueden calcularse mediante una propagación hacia delante. Por tanto, para obtener las derivadas parciales, basta con obtener los vectores de sensibilidad  $d^\ell$  para cada capa  $\ell = 1$  (recuerde que no hay señal de entrada a la capa  $L = 0$ ). Resulta que los vectores de sensibilidad pueden obtenerse ejecutando una versión ligeramente modificada de la red neuronal *hacia atrás*, y de ahí el nombre de retropropagación. En la propagación hacia delante, cada capa produce el vector  $x^\ell$  y en la retropropagación, cada capa produce (hacia atrás) el vector  $d^\ell$ . En la propagación hacia delante, calculamos  $x^\ell$  a partir de  $x^{(\ell-1)}$  y en la propagación hacia atrás, calculamos  $d^\ell$  a partir de  $d^{(\ell+1)}$ . La idea básica se ilustra en la siguiente figura.

---

La programación dinámica es una elegante técnica algorítmica en la que se construye una  
 @ Abu-Mostafa, Magdon-Ismael, Lin: Ene-2015 e-Chap:7-12

---

solución a un problema complejo utilizando las soluciones a problemas relacionados pero más sencillos.



Como se puede ver en la figura, la red neuronal está ligeramente modificada sólo en que hemos cambiado la función de transformación de los nodos. En la propagación hacia delante, la transformación era la sigmoidea  $\delta(\cdot)$ . En la retropropagación, la transformación es la *multiplicación por  $\delta'$*  ( $s'$ ), donde  $s'$  es la entrada al nodo. Así que la función de transformación es ahora diferente para cada nodo, y depende de la entrada al nodo, que depende de  $x$ . Esta entrada se calculó en la propagación hacia adelante. Para la función de transformación  $\tanh(\cdot)$ ,  $\tanh'(s'') = 1 - \tanh^2(s'') = 1 - x'' x'''$ , donde significa multiplicación por componentes.

En la figura, la capa  $(\ell + 1)$  produce (hacia atrás) el vector de sensibilidad  $\delta^{(\ell+1)}$ , que se multiplica por los pesos  $W^{(\ell+1)}$ , se suma y se pasa a los nodos de la capa  $\ell$ . Los nodos de la capa  $\ell$  se multiplican por  $\theta'(s^{(\ell)})$  para obtener  $\delta^{(\ell)}$ . Usando  $\otimes$ , una notación abreviada para este paso de retropropagación es:

$$\delta^{(\ell)} = \theta'(s^{(\ell)}) \otimes [W^{(\ell+1)} \delta^{(\ell+1)}]_1^{d^{(\ell)}}, \quad (7.5)$$

donde el vector  $W^{(\ell+1)} \delta^{(\ell+1)}$  contiene los componentes  $1, \dots, d^{(\ell)}$  del vector  $W^{(\ell+1)} \delta^{(\ell+1)}$  (excluyendo el componente de sesgo que tiene índice 0). Esta fórmula no es sorprendente. La sensibilidad de  $e$  a las entradas de la capa  $\ell$  es proporcional

a la pendiente de la función de activación en la capa  $\ell$  (una mayor pendiente significa una menor

un cambio en  $s''$  tendrá un mayor efecto sobre  $x''$ ), el tamaño de los pesos que salen de la capa (pesos más grandes significan que un pequeño cambio en  $s''$  tendrá más impacto en  $s^{(\ell+1)}$ ) y la sensibilidad en la capa siguiente (un cambio en la capa  $\ell$  afecta a las entradas de la capa  $\ell + 1$ , por lo que si  $e$  es más sensible a la capa  $\ell + 1$ , entonces también será más sensible a la capa  $\ell$ ).

Derivaremos esta recursión hacia atrás más adelante. Por ahora, observe que si sabemos  $\delta^{(\ell+1)}$ , entonces se puede obtener  $\delta^{(\ell)}$ . Utilizamos  $\delta^{(\ell)}$  para sembrar el

proceso hacia atrás, y podemos obtener que explícitamente porque e  $(x'' - p)^2$   
 $(B(s^*)) - p)^2$ .

Por lo tanto,

$$\begin{aligned}\delta^{(L)} &= \frac{\partial e}{\partial \mathbf{s}^{(L)}} \\ &= \frac{\partial}{\partial \mathbf{s}^{(L)}} (\mathbf{x}^{(L)} - y)^2 \\ &= 2(\mathbf{x}^{(L)} - y) \frac{\mathbf{x}^{(L)}}{\mathbf{s}^{(L)}} \\ &= 2(\mathbf{x}^{(L)} - y) \theta'(\mathbf{s}^{(L)}).\end{aligned}$$

Cuando la transformación de salida es  $\tanh(\cdot)$ ,  $\theta'(\mathbf{s}^{(L)}) = 1 - (\mathbf{z}^{(L)})^2$  (clasificación); cuando la transformación de salida es la identidad (regresión),  $\theta'(\mathbf{s}^{(L)}) = 1$ . Ahora, utilizando (7.5), podemos calcular todas las sensibilidades:

$$\delta^{(1)} \leftarrow \delta^{(2)} \dots \leftarrow \delta^{(L-1)} \leftarrow \delta^{(L)}.$$

Tenga en cuenta que como sólo hay un nodo de salida,  $\mathbf{s}^*$  es un escalar, y también lo es  $\delta^{(L)}$ . El siguiente cuadro de algoritmos resume la retropropagación.

**Retropropagación para calcular sensibilidades  $\delta^{(L)}$ .**

Entrada: un punto de datos  $(\mathbf{x}, p)$ .

0: Ejecutar propagación hacia adelante en  $\mathbf{x}$  para calcular  $y$  y guardar:

$$\begin{aligned}\mathbf{s}^{(I)} & \text{ para } I = 1, \dots, \tilde{n}; \\ \mathbf{x}^{(I)} & \text{ para } I = 0, \dots, \tilde{n}. \\ \mathbf{z}^{(L)} &= 2(\mathbf{x}^{(L)} - y) \theta'(\mathbf{s}^{(L)}) \quad [\text{Inicialización}] \\ \theta'(\mathbf{s}^{(L)}) &= \begin{cases} 1 - (\mathbf{x}^{(L)})^2 & \theta(s) = \tanh(s); \\ 1 & \theta(s) = s. \end{cases}\end{aligned}$$

2: para  $\ell = \tilde{n} - 1$  a 1 hacer [Back-Propagation]

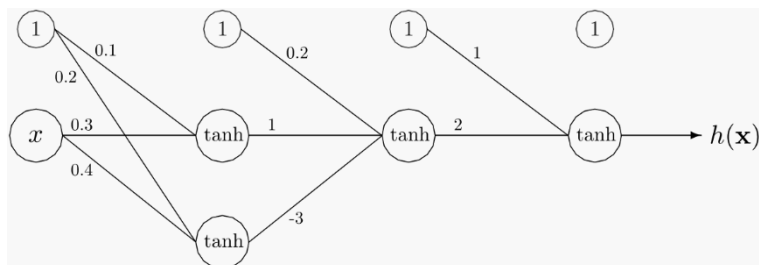
3: Let  $\theta'(\mathbf{s}^{(\ell)}) = [1 - \mathbf{x}^{(\ell)} \otimes \mathbf{x}^{(\ell)}]_1^{d^{(\ell)}}$ .

4: Calcular la sensibilidad  $\delta^{(\ell)}$  a partir de  $\delta^{(\ell+1)}$ :

$$\delta^{(\ell)} \leftarrow \theta'(\mathbf{s}^{(\ell)}) \otimes [\mathbf{W}^{(\ell+1)} \delta^{(\ell+1)}]_1^{d^{(\ell)}}$$

En el paso 3, asumimos transformaciones de nodo oculto  $\tanh$ . Si las funciones de transformación de la unidad oculta no son  $\tanh(\cdot)$ , la derivada del paso 3 debe actualizarse en consecuencia. Usando la propagación hacia adelante, calculamos  $\mathbf{x}^{(I)}$  para  $I = 0, \dots, \tilde{n}$  y utilizando la retropropagación, calculamos  $\delta^{(I)}$  para  $I = 1, \dots, \tilde{n}$ . Por último, obtenemos la derivada parcial del error en un punto de datos simple mediante la ecuación (7.4). Nada ilumina mejor las partes móviles que trabajar un ejemplo de principio a fin.

Ejemplo 7.1. Considere la siguiente red neuronal.



Hay una sola entrada, y las matrices de pesos son:

$$W^{(1)} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}, \quad p(z) = \begin{bmatrix} 0.2 \\ 1 \\ -3 \end{bmatrix}, \quad p() = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Para el punto de datos  $z = 2$ ,  $p = 1$ , la propagación hacia delante da:

$$\begin{array}{c|c|c|c|c|c|c} \mathbf{x}^{(0)} & \mathbf{s}^{(1)} & \mathbf{x}^{(1)} & \mathbf{s}^{(2)} & \mathbf{x}^{(2)} & \mathbf{s}^{(3)} & \mathbf{x}^{(3)} \\ \hline \begin{bmatrix} 1 \\ 2 \end{bmatrix} & \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0.7 \\ 1 \end{bmatrix} & \begin{bmatrix} 1 \\ 0.60 \\ 0.76 \end{bmatrix} & \begin{bmatrix} 1 \\ -1.48 \\ 8 \end{bmatrix} & \begin{bmatrix} 1 \\ -0.99 \\ 0 \end{bmatrix} & \begin{bmatrix} -0.8 \\ 8 \end{bmatrix} & -0.666 \end{array}$$

Más arriba mostramos cómo se calcula  $\mathbf{s}^{(i)}$  ( $\mathbf{W}^{(i)} \mathbf{x}^{(i-1)}$ ). La retropropagación da

$$\begin{array}{c|c|c} \delta^{(3)} & \delta^{(2)} & \delta^{(1)} \\ \hline [-1.855] & [(1 - 0.9^2) \cdot 2 \cdot -1.855] = [-0.69] & \begin{bmatrix} -0.44 \\ 0.88 \end{bmatrix} \end{array}$$

Hemos mostrado explícitamente cómo se obtiene  $\delta^{(2)}$  a partir de  $\delta^{(3)}$ . Ahora es sencillo combinar los vectores de salida  $\mathbf{x}^{(3)}$  con los vectores de sensibilidad  $\delta^{(3)}$  utilizando (7.4) para obtener las derivadas parciales que se necesitan para el gradiente:

$$\frac{de}{d\mathbf{v}^{(1)}} = {}^{(0)}(\delta^{(1)})^T = \begin{bmatrix} -0.44 & 0.88 \\ -0.88 & 1.75 \end{bmatrix}, \quad \frac{de}{d\mathbf{m}^{(2)}} = \begin{bmatrix} -0.69 \\ -0.42 \\ -0.53 \end{bmatrix}, \quad \frac{de}{d\mathbf{v}^{(1)}} = \begin{bmatrix} -0.44 \\ 0.88 \end{bmatrix}, \quad \frac{de}{d\mathbf{v}^{(1)}} = \begin{bmatrix} -1.85 \\ 1.67 \end{bmatrix}$$

### Ejercicio 7.8

Repita los cálculos del ejemplo 7.1 para el caso en que la transformación de salida sea la identidad. Debe calcular  $\mathbf{s}^{(i)}$ ,  $\mathbf{x}^{(i)}$ ,  $\mathbf{d}^{(i)}$  y  $de/d\mathbf{W}^{(i)}$

Derivemos (7.4) y (7.5), que son las ecuaciones centrales de la retropropagación. No hay nada más que la aplicación repetida de la regla de la cadena. Si quieres confiar en nuestras matemáticas, no te perderás gran cosa por seguir adelante.

Comienza el salto seguro: Si confías en nuestras matemáticas, puedes saltarte esta parte sin comprometer la secuencia lógica. Un **recuadro verde** similar te indicará cuándo volver a unirte.

Para empezar, veamos más de cerca la derivada parcial,  $\partial e / \partial W^{(\ell)}$ . La situación se ilustra en la figura 7.1.

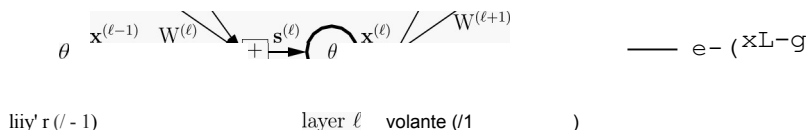


Figura 7.1: Cadena de dependencias de  $W^{(\ell)}$  a  $x^{(\ell)}$ .

Podemos identificar la siguiente cadena de dependencias por las que  $W^{(\ell)}$  influye en la salida  $x^{(\ell)}$ , y por tanto en el error  $e$ .

$$W^{(\ell)} \rightarrow s^{(\ell)} \rightarrow x^{(\ell)} \rightarrow s^{(\ell+1)} \rightarrow \dots \rightarrow x^{(L)} = h.$$

Para derivar (7.4), nos limitamos a un único peso y utilizamos la regla de la cadena. Para un solo peso en  $W^{(\ell)}$ , un cambio en  $W^{(\ell)}$  sólo afecta a  $s^{(\ell)}$  y así por la regla de la cadena,

$$\frac{\partial e}{\partial w_{ij}^{(\ell)}} = \frac{\partial s_j^{(\ell)}}{\partial w_{ij}^{(\ell)}} \frac{\partial e}{\partial s_j^{(\ell)}} = x_i^{(\ell-1)} \cdot \delta_j^{(\ell)},$$

donde la última igualdad se deduce porque  $s_j^{(\ell)} = \sum_i w_{ij}^{(\ell)} x_i^{(\ell-1)}$  y por definición de  $\delta_j^{(\ell)}$ . Hemos derivado la forma componente de (7.4).

Derivamos ahora la forma componente de (7.5). Como  $e$  depende de  $s^{(L)}$  sólo a través de  $x^{(L)}$  (véase la figura 7.1), por la regla de la cadena, tenemos:

$$\delta_j^{(\ell)} = \frac{\partial e}{\partial s_j^{(\ell)}} = \frac{\partial e}{\partial x_j^{(L)}} \frac{\partial x_j^{(L)}}{\partial s_j^{(\ell)}} = \delta_j^{(L)} \cdot g'(x_j^{(\ell)}).$$



Para obtener la derivada parcial  $\frac{f_i}{f_{i,x'}}$ , necesitamos entender cómo cambia  $e$  debido a cambios en  $x'$ ). De nuevo, a partir de la figura 7.1, un cambio en  $x'$ ) sólo afecta a

$s^{(+*)}$  y, por tanto, e. Dado que un componente concreto de  $x^{(i)}$  puede afectar a todos los componentes de  $s^{(+)}$ , tenemos que sumar estas dependencias utilizando la regla de la cadena:

$$\frac{de}{g^{(l)}} \sum_{j \in J}^{d^{(\ell+1)}} \frac{\partial s_k^{(\ell+1)}}{\partial \mathbf{x}_j^{(\ell)}} \cdot \frac{\partial e}{\partial s_k^{(\ell+1)}} - \sum_{j \in J}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}.$$

Juntando todo esto, hemos llegado a la versión componente de (7.5)

$$\delta_j^{(\ell)} = \theta'(\mathbf{s}_j^{(\ell)}) \sum_{k=1}^{d^{(\ell+1)}} w_{jk}^{(\ell+1)} \delta_k^{(\ell+1)}, \quad (7.6)$$

Intuitivamente, el primer término proviene del impacto de  $s^+$  sobre  $x^+$ ; el sumatorio es el impacto de  $x^+$  sobre  $s^+$ , y el impacto de  $s^+$  sobre  $f$  es lo que nos devuelve las sensibilidades en la capa (I-I-1), dando lugar a la recursión hacia atrás.

Fin del salto seguro: Los que saltaron se reúnen ahora con nosotros para discutir cómo la retropropagación nos da Wfiq.

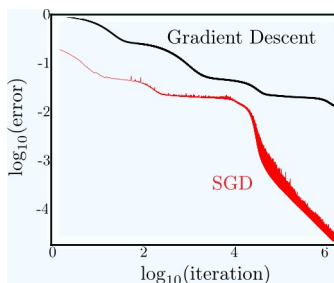
La retropropagación funciona con un punto de datos  $(x, p)$  y pesos  $w$  ( $W'$ , ...,  $W^*$ ). Puesto que ejecutamos una propagación hacia delante y hacia atrás para calcular las salidas  $x''$  y las sensibilidades  $d'''$ , el tiempo de ejecución es del orden del número de pesos de la red. Calculamos una vez para cada punto de datos  $(x, p)$  para obtener  $WA, q(x, p)$  y, utilizando la suma en (7.3), agregamos estos gradientes de punto único para obtener el gradiente  $ho/cf$  completo  $WAip$ . A continuación resumimos el algoritmo.

**Algoritmo para calcular**  $A_{iq}(w)$  y  $W A_{iq}(w)$ .  
**Entrada:**  $w \quad (W^*), \dots, W^*)\}; P \quad (x, y) \dots (xy, pq)$ .  
**Salida:** errorin () y gradiente  $g \quad (G'', \dots, G'')$ .  
 1: Inicializar:  $Aip = 0$  y  $G' = 0$   $W'$ ) para  $\epsilon = 1, \dots, L$ : Para cada punto de datos  $(xp, pp)$ ,  $n = 1, \dots, N$ , do  
 2:     Calcular  $x''$  para  $fi = 0, \dots, L$                                 propagación hacia atrás  
 3:     Calcula  $d''$ ) para  $fi = L, \dots, 1$                                 delante  
 4:     #  $i_n$     #  $i_n - / - (\times'')$  )<sup>2</sup>                                retropropagaciónJ  
 5:     para  $I = 1, \dots, \tilde{n}$  do  
 6:          $G^{(\ell)}(x_n) = [x^{(\ell-1)}(\delta^{(\ell)})^T]$   
 7:          $yC) \quad GC) \quad y \quad GJ(x, y)$

( G ' ) ( x p ) es el gradiente en el punto de datos x p ). La actualización del peso para una sola iteración del descenso de gradiente de tasa de aprendizaje fija es  $W'' + W''$

-  $\nabla G$  ", para  $\epsilon = 1, \dots, L$ . Hacemos todo esto para una iteración del descenso de gradiente, un cálculo costoso para sólo un pequeño paso. ,  $L$ . Hacemos todo esto para una iteración del descenso de gradiente, un cálculo costoso para sólo un pequeño paso.

En el Capítulo 3, discutimos el *descenso de gradiente s/ofcios- lie (S 'D)* como una alternativa más eficiente al modo por lotes. En lugar de esperar al gradiente agregado  $G'$  al final de la iteración, se actualizan inmediatamente los pesos a medida que se procesa secuencialmente cada punto de datos utilizando el gradiente de punto único en el paso 7 del algoritmo:  $W'' = W' - \eta \nabla J(x_p)$ . En esta versión secuencial, todavía se ejecuta una propagación hacia delante y hacia atrás para cada punto de datos.



punto de datos, pero realiza  $N$  actualizaciones de los pesos. A la derecha se muestra una comparación del descenso gradiente por lotes con SGD. Utilizamos 500 ejemplos de entrenamiento de los datos de dígitos y una red neuronal de 2 capas con 5 unidades ocultas y una tasa de aprendizaje  $\eta = 0,01$ . La curva SGD es errática porque no se está minimizando el error total en cada punto de datos, sino que se realizan  $N$  actualizaciones de los pesos. La curva SGD es errática porque no se minimiza el error total en cada iteración, sino el error en un punto de datos específico. Un método para amortiguar este comportamiento errático consiste en disminuir la tasa de aprendizaje a medida que avanza la minimización.

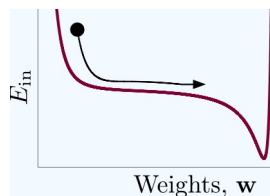
La velocidad a la que se minimiza  $J$  depende en gran medida del algoritmo de optimización que se utilice. SGD parece significativamente mejor que el simple descenso de gradiente, pero podemos hacerlo mucho mejor, incluso SGD no es muy eficiente. En la Sección 7.5, discutimos algunos otros métodos potentes (por ejemplo, conjugate gradients) que pueden mejorar significativamente el descenso de gradiente y el SGD, al hacer un uso más eficaz del gradiente.

**Inicialización y terminación.** Elegir los pesos iniciales y decidir cuándo detener el descenso de gradiente puede ser complicado, en comparación con la regresión logística, porque el error en la muestra ya no es convexo. Según el Ejercicio 7.7, si los pesos se inicializan demasiado grandes para que  $\tanh(w \cdot x_p) \in [-1, 1]$ , entonces el gradiente será cercano a cero y el algoritmo no llegará a ninguna parte. Esto es especialmente problemático si se inicializan los pesos con el signo equivocado. Normalmente es mejor inicializar los pesos a valores  $\text{randn}()$  donde  $\tanh(w \cdot x_p) \approx 0$  para que el algoritmo tenga la flexibilidad de mover los pesos fácilmente para ajustarse a los datos. Una buena opción es inicializar usando pesos aleatorios gaussianos, en  $N(0, \sigma^2)$  donde  $\sigma$  es pequeño. Pero, ¿cómo de pequeño debe ser  $\sigma$ ? Una heurística sencilla es que queremos que  $|w \cdot x_p|^2$  sea pequeño. Dado que  $E_p |w \cdot x_p|^2 = \sigma^2 \|x_p\|^2$ , deberíamos elegir  $\sigma$  de modo que  $\sigma \max_p \|x_p\| \ll 1$ .

**Ejercicio 7.9**

¿Qué puede salir mal si se inicializan todos los pesos exactamente a cero?

¿Cuándo nos detenemos? Es arriesgado basarse únicamente en el tamaño del gradiente para detenerse. Como se ilustra a la derecha, podría detenerse prematuramente cuando la iteración alcanza una región relativamente plana (lo que es más común de lo que podría sospecharse). En la práctica, lo mejor es una combinación de criterios de parada, por ejemplo, detenerse sólo cuando se produce una mejora marginal del error junto con un gradiente pequeño. más un límite superior en el número de iteraciones.



### 7.2.4 Regresión para la clasificación

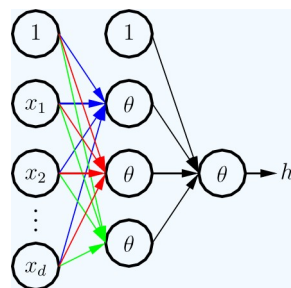
En el capítulo 3, mencionamos que podía utilizar los pesos resultantes de la regresión lineal como pesos de perceptrón para la clasificación, y puede hacer lo mismo con las redes neuronales. Específicamente, ajuste los datos de clasificación (pq - -1-1) como si fuera un problema de regresión. Esto significa que utiliza la función de identidad como transformación del nodo de salida, en lugar de  $\tanh()$ . Esto puede ser de gran ayuda debido a las "regiones planas" a las que es susceptible la red cuando se utiliza el descenso de gradiente, lo que ocurre a menudo en el entrenamiento. La razón de estos periodos planos en la optimización es la naturaleza excepcionalmente plana de la función  $\tanh$  cuando su argumento se hace grande. Si por alguna razón los pesos se hacen grandes hacia el principio del entrenamiento, entonces la superficie de error empieza a parecer plana, porque la  $\tanh$  se ha saturado. Ahora, el descenso de gradiente no puede hacer ningún progreso y usted podría pensar que está en un mínimo, cuando en realidad está lejos de un mínimo. El problema de una superficie de error plana se mitiga considerablemente cuando la transformación de salida es la identidad, porque puedes recuperarte de un mal movimiento inicial si resulta que te lleva a pesos grandes (la salida lineal nunca se satura). Para un ejemplo concreto de un error prematuramente plano en la muestra, vea la figura del Ejemplo 7.2 en la página 25.

## 7.3 Aproximación frente a generalización

Un MLP suficientemente grande con 2 capas ocultas puede aproximar funciones de decisión suaves arbitrariamente bien. Resulta que basta con una sola capa oculta'. Una red neuronal con una sola capa oculta que tiene  $m$  unidades ocultas ( $d!$   $m$ ) implementa una función de la forma

$$h(\mathbf{x}) = \theta \left( w_{01}^{(2)} + \sum_{j=1}^m w_{j1}^{(2)} \theta \left( \sum_{i=0}^d w_{ij}^{(1)} x_i \right) \right)$$

@ Abu-Mostafa, Magdon-Ismael, Lin: Ene-2015



e-Chap:7-19

"Aunque una capa oculta es suficiente, no es necesariamente la forma más eficiente de ajustarse a los datos; por ejemplo, puede existir una red de 2 capas ocultas mucho más pequeña.

Se trata de una representación engorrosa para una red tan sencilla. Una notación simplificada para este caso especial es mucho más conveniente. Para los pesos de la segunda capa, utilizaremos simplemente  $w_j$ , ... y utilizaremos  $v$  para denotar la columna  $j$  de la matriz de pesos de la primera capa  $W(*)$ , para  $j = 1 \dots m$ . Con esta notación más sencilla, la hipótesis tiene un aspecto mucho más agradable:

$$h(\mathbf{x}) = \theta \left( w_0 + \sum_{j=1} w_j \theta(\mathbf{v}_j^T \mathbf{x}) \right)$$

Red neuronal frente a transformadas no lineales. Recordemos el modelo lineal del capítulo 3, con la transformada no lineal  $\Phi(\mathbf{x})$  que transforma  $\mathbf{x}$  en  $\mathbf{z}$ :

$$\mathbf{x} \rightarrow \mathbf{z} = \Phi(\mathbf{x}) = [1, \phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_m(\mathbf{x})]^T.$$

El modelo lineal con transformada no lineal es una hipótesis de la forma

$$h(\mathbf{x}) = \theta \left( w_0 + \sum_{j=1}^M w_j \phi_j(\mathbf{x}) \right).$$

Las  $\phi_j(\cdot)$  se denominan funciones de base. A primera vista, la red neuronal y el modelo lineal parecen casi idénticos, estableciendo  $g(\mathbf{v}; \mathbf{x}) = d(\mathbf{x})$ . Sin embargo, hay una sutil diferencia que tiene un gran impacto práctico. Con la transformada no lineal, las funciones de base  $\phi_j(\cdot)$  se fijan de antemano antes de observar los datos. Con la red neuronal, la "función base"  $g(\mathbf{v}; \mathbf{x})$  tiene un parámetro  $\mathbf{v}$  en su interior, y podemos ajustar  $\mathbf{v}$  sin ver los datos. En primer lugar, esto tiene un impacto computacional porque el parámetro  $\mathbf{v}$  aparece *dentro de* la no linealidad  $g(\cdot)$ ; el modelo ya no es lineal en sus parámetros. Vimos un efecto similar con los centros de la red de función de base radial en el Capítulo 6. Los modelos que no son lineales en sus parámetros plantean un reto computacional importante cuando se trata de ajustarlos a los datos. En segundo lugar, significa que podemos tener verdaderos *hubs* /unc/iones a /fie *datos*. Las *funciones de base* sintonizable, aunque computacionalmente más difíciles de ajustar a los datos, nos dan mucha más flexibilidad para ajustar los datos que las funciones de base fija. Con  $m$  funciones de base sintonizables se tiene aproximadamente el mismo poder de aproximación para ajustar los datos que con *md* funciones de base fijas. Para  $d$  grandes, las funciones de base sintonizables tienen considerablemente más potencia.

### Ejercicio 7\10

No es de extrañar que añadir nodos en la capa oculta proporcione a la red neuronal una mayor capacidad de aproximación, ya que se añaden más parámetros.



¿Cuántos parámetros de peso hay en una red neuronal con arquitectura especificada por  $d$  ( $d'''$ ,  $d''$ ,  $d'$ ,  $d$ ), un vector que da el número de nodos en cada capa? Evalúe su fórmula para una red de 2 capas ocultas con 10 nodos ocultos en cada capa oculta.

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2 < \frac{(2RC_f)^2}{m}$$

donde  $\rho = \max_p \|\mathbf{x}_p\|$  es el "radio" de los datos. El error en la muestra disminuye inversamente con el número de unidades ocultas. Para la clasificación, existe un resultado similar con una dependencia ligeramente peor de  $m$ . Con alta probabilidad,

$$E_{in} \leq E_{out}^* + O(C_f/\sqrt{m}),$$

donde  $E_{out}$  es el error fuera de muestra del clasificado óptimo que discutimos en el Capítulo 6. El mensaje es que  $E_{out}$  puede hacerse pequeño eligiendo una capa oculta suficientemente grande.

**Generalización y dimensión VC.** Para  $m$  suficientemente grande, podemos conseguir que  $E_{out}$  sea pequeño, así que lo que queda es asegurar  $E_{in} \leq E_{out}$ . Debemos fijarnos en la dimensión VC. Para la red neuronal de dos capas de umbral duro (MLP) donde  $\delta(z) = \text{sign}(z)$ , mostramos un simple límite en la dimensión VC:

$$VC \leq (const) m d \log(md). \quad (7.7)$$

Para una red neuronal sigmoide general,  $d$  puede ser infinito. Para la sigmoide  $\tanh(\cdot)$ , con nodo de salida  $\text{sign}(\cdot)$ , se puede demostrar que  $VC = O(U \log Q)$  donde  $U$  es el número de nodos ocultos y  $Q$  es el número de pesos; para el caso de dos capas

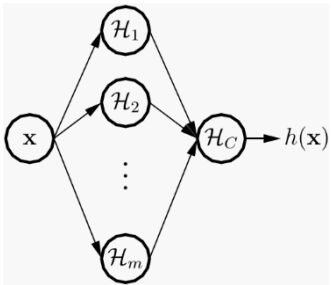
$$VC = O(md(m + d)).$$

La red  $\tanh(\cdot)$  tiene mayor dimensión VC que el MLP de 2 capas, lo que no es sorprendente porque  $\tanh(\cdot)$  puede aproximarse a  $\text{sign}(\cdot)$  eligiendo pesos suficientemente grandes. Por lo tanto, todas las dicotomías que pueden ser implementadas por el MLP también pueden ser implementadas por la red neuronal  $\tanh(\cdot)$ .

<sup>7</sup>No describimos los detalles de cómo se puede medir la complejidad de un objetivo. Una @ Abu-Mostafa, Magdon-Ismael, Lin: Ene-2015 e-Chap:7-21

medida es el tamaño de los componentes de "alta frecuencia" de  $J$  en su transformada de Fourier. Otra medida más restrictiva es el número de derivadas acotadas que tiene /.

Para derivar (7.7), en realidad mostraremos un resultado más general. Consideremos el conjunto de hipótesis ilustrado por la red de la derecha. El nodo oculto  $i$  de la capa oculta implementa una función  $f_i$ ;  $C_{7-i}$  que mapea  $\mathcal{X} \rightarrow \{-1, 1\}$ ; el nodo de salida implementa una función  $f_C$  que mapea  $\{-1, 1\}^m$  a  $\mathcal{Y}$ . Este nodo de salida combina las salidas de los nodos de la capa oculta para implementar la hipótesis



$$h(\mathbf{x}) = h_C(h_1(\mathbf{x}), \dots, h_m(\mathbf{x})).$$

(Para el MLP de 2 capas, todos los conjuntos de hipótesis son perceptrones).

Supongamos que la dimensión VC de  $H$ , es  $d$  y la dimensión VC de  $H_y$  es  $d'$ . Fijemos  $\mathbf{x}_1, \dots, \mathbf{x}_N$  y las hipótesis  $f_1, \dots, f_m$  implementadas por los nodos ocultos. Las hipótesis  $f_1, \dots, f_m$  son ahora funciones de base fija que definen una transformada a  $\mathbb{R}^m$ ,

$$\begin{matrix} h_1(\mathbf{x}_1) & \dots & h_1(\mathbf{x}_N) \\ \mathbf{x}_1 \rightarrow \mathbf{z}_1 = & \dots & \mathbf{x}_N \rightarrow \mathbf{z}_N = \\ h_m(\mathbf{x}_1) & \dots & h_m(\mathbf{x}_N) \end{matrix}$$

Los puntos transformados son vectores binarios en  $\mathbb{R}^m$ . Dados  $f_1, \dots, f_m$ , los puntos  $\mathbf{x}_1, \dots, \mathbf{x}_N$  se transforman en una disposición de puntos  $\mathbf{z}_1, \dots, \mathbf{z}_N$  en  $\mathbb{R}^m$ . Utilizando nuestra flexibilidad para elegir  $f_1, \dots, f_m$ , ahora acotamos el número de posibles arreglos  $\mathbf{z}_1, \dots, \mathbf{z}_N$  que podemos obtener.

Los primeros componentes de todos los  $\mathbf{z}_i$  vienen dados por  $f_1(\mathbf{x}_1), \dots, f_1(\mathbf{x}_N)$ , que es una dicotomía de  $\mathbf{x}_1, \dots, \mathbf{x}_N$  implementada por  $f_1$ . Como la dimensión VC de  $H$  es  $d$ , hay como máximo  $N^d$  dicotomías de este tipo.<sup>8</sup> Es decir, hay como máximo  $N^d$  formas diferentes de elegir asignaciones a *todos los* primeros componentes de la  $\mathbf{z}_i$ . Del mismo modo, una asignación a todos los componentes  $\mathbf{z}_i$  puede elegirse como máximo de  $N^d$  maneras. Así pues, el número total de disposiciones posibles para  $\mathbf{z}_1, \dots, \mathbf{z}_N$  es

$$\prod N^{d_i} = N^{\sum_{i=1}^m d_i}.$$

Cada una de estas disposiciones puede dicotomizarse como máximo de  $N$  maneras<sup>d'</sup>, ya que la dimensión VC de  $H_y$  es  $d'$ . Cada una de estas dicotomías para un arreglo particular da una dicotomía de los datos  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . Por lo tanto, el número máximo de dicotomías diferentes que podemos implementar en  $\mathbf{x}_1, \dots, \mathbf{x}_N$  está limitado por el producto: el número de disposiciones posibles por el número de maneras

<sup>8</sup>Recordemos que para cualquier conjunto de hipótesis con dimensión VC  $d$  y cualquier  $N > d$ ,  $\text{rn}(N)$  (el número máximo de dicotomías implementables) está limitado por  $(ed/d)^d < N^d$  (para simplificar suponemos que  $d > 2$ ).

de dicotomizar un arreglo particular. Hemos demostrado que

$$m(N) \leq N^{d_c} \cdot N^{\sum_{i=1}^m d_i} = N^{d_c + \sum_{i=1}^m d_i}.$$

Sea  $D = d - 1 - d_i$ . Después de un poco de álgebra (dejada al lector), si  $N \geq 2D \log_2 D$ , entonces (N) de lo que concluimos que  $d \geq 2D \log_2 D$ . Para el MLP de 2 capas,  $d = d - 1 - 1$  y  $d_i = rn - 1 - 1$ , por lo que tenemos que  $D = d - d_i = m(d - 2) - 1 - 1 - O(md)$ . Por tanto,  $d \geq O(md \log(md))$ . Nuestro análisis parece muy burdo, pero es casi ajustado: es posible romper  $n(md)$  puntos con  $m$  unidades ocultas (véase el problema 7.16), por lo que el límite superior puede ser flojo como mucho en un factor logarítmico. Utilizando la dimensión VC, la barra de error de generalización del capítulo 2 es  $O(\sqrt{N}/N)$  que para el org(d MLP de 2 capas es  $O(\sqrt{d} \log(d)/N)$ .

Obtendremos una buena generalización si  $m$  no es demasiado grande y podemos ajustar los datos

si  $m$  es suficientemente grande. Se impone un equilibrio. Por ejemplo, elegir  $m$  como  $N^{-1/2}$  on tout in  $\mathbb{R}^n$  d Aip  $-1$  fi q. Es decir,  $A_{n,t}$  q, (el rendimiento óptimo) a medida que crece  $N$ , y  $m$  crece de forma sublineal con  $N$ . En la práctica, el régimen "asintótico" es un lujo y no basta con fijar  $m$ . Estos resultados teóricos son una buena orientación, pero el mejor rendimiento fuera de la muestra suele obtenerse cuando se controla el sobreajuste utilizando la validación (para seleccionar el número de unidades ocultas) y la regularización para evitar el sobreajuste.

Concluimos con una nota sobre dónde se sitúan las redes neuronales en el debate paramétrico-no paramétrico. Hay parámetros explícitos que aprender, por lo que para- métrico parece correcto. Pero también destacan los rasgos distintivos de los modelos no paramétricos: la red neuronal es genérica y flexible y puede alcanzar un rendimiento óptimo cuando  $N$  crece. Ni los modelos paramétricos ni los no paramétricos abarcan toda la historia. Optamos por etiquetar las redes neuronales como *semi-param e Iric*.

## 7.4 Regularización y validación

La red neuronal multicapa es potente y, unida al descenso por gradiente (un buen algoritmo para minimizar Aip), tenemos una receta para el sobreajuste. Analizamos algunas técnicas prácticas de ayuda.

### 7.4.1 Penalizaciones de complejidad basadas en el peso

Al igual que con los modelos lineales, se puede regularizar el aprendizaje utilizando una penalización de complejidad minimizando un error aumentado (error penalizado en la muestra). El regularizador de caída de peso al cuadrado es popular, ya que tiene un error aumentado:

$$E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \sum_{i,j} (w_{ij}^{(\ell)})^2$$

El parámetro de regularización  $\lambda$  se selecciona mediante validación, como se explica en el Capítulo 4. Para aplicar el descenso de gradiente, necesitamos  $\nabla_{\mathbf{w}} J(\mathbf{w})$ . El término de penalización añade

al gradiente un término proporcional a los pesos,

$$\frac{\partial E_{\text{aug}}(\mathbf{w})}{\partial \mathbf{W}^{(\ell)}} = \frac{\partial E_{\text{in}}(\mathbf{w})}{\partial \mathbf{W}^{(\ell)}} + \frac{2\lambda}{N} \mathbf{W}^{(\ell)}$$

Sabemos cómo obtener  $\partial E_{\text{aug}} / \partial \mathbf{W}^{(\ell)}$  mediante retropropagación. El término de penalización añade un componente a la actualización del peso que está en la dirección negativa de  $\mathbf{w}$ ,

es decir, hacia pesos cero — de ahí el término decaimiento del peso.

Otro regularizador similar es la eliminación del ceipfit, con error aumentado:

$$L_{\text{aug}}(\mathbf{w}, \mathbf{A}) = \text{Min}(\cdot) + \frac{\lambda}{N} \sum_{i,j} \frac{(w_{ij}^{(\ell)})^2}{1 + (w_{ij}^{(\ell)})^2}$$

Para un peso pequeño, el término de penalización es muy parecido al decaimiento del peso, y lo reducirá a cero. Para un peso grande, el término de penalización es aproximadamente una constante y contribuye poco al gradiente. Los pesos pequeños decaen más rápido que los pesos grandes, y el efecto es "eliminar" esos pesos más pequeños.

Ejercicio 7.11

Para la eliminación del peso, demuestre:

$$\frac{\partial E_{\text{aug}}}{\partial w_{ij}^{(\ell)}} = \frac{\partial E_{\text{in}}}{\partial w_{ij}^{(\ell)}} + \frac{2\lambda}{N} \cdot \frac{w_{ij}^{(\ell)}}{(1 + (w_{ij}^{(\ell)})^2)^2}$$

Argumente que la eliminación de pesos reduce más rápidamente los pesos pequeños que los grandes.

## 7.4.2 Parada anticipada

Otro método de regularización, que a primera vista no se parece a la regularización, es la parada temprana. Un método iterativo como el descenso de gradiente no explora todo el conjunto de hipótesis de una vez. Con más iteraciones, se explora más del conjunto de hipótesis. Esto significa que utilizando menos iteraciones, se explora un conjunto de hipótesis más pequeño y se obtiene una mejor generalización".

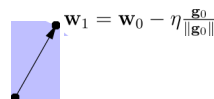
Consideremos el descenso de gradiente de paso fijo con tamaño de paso  $p$ . En el primer paso, empezamos en los pesos  $\mathbf{w}_k$ , y damos un paso de tamaño  $p$  hasta  $\mathbf{w} = \mathbf{w}_k - p \cdot \mathbf{g}_k$ .

Porque hemos dado un paso en la dirección del gradiente negativo, hemos 'mirado' todos los hipótesis en la región sombreada de la derecha.

Esto se debe a que un paso en el gradiente negativo conduce a la disminución más brusca de

$\|\mathbf{w}\|$ , y por tanto  $\mathbf{w}$  minimiza  $\|\mathbf{w}\|$  entre todas las ponderaciones con  $\|\mathbf{w} - \mathbf{w}_k\| = p$ .

Buscamos indirectamente en todo el conjunto de hipótesis

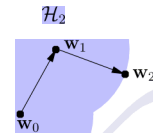


y elegimos la hipótesis  $w_C$  by con el mínimo error dentro de la muestra.

Si queremos ser rigurosos con la corrección, el conjunto de hipótesis explorado podría depender del conjunto de datos, por lo que no podemos aplicar directamente el análisis de la CV, que requiere que el conjunto de hipótesis esté fijado de antemano. Dado que sólo estamos ilustrando la idea principal, pasaremos por alto estos tecnicismos.



Consideremos ahora el segundo paso, ilustrado a la derecha, que pasa a explorar indirectamente el conjunto de hipótesis de pesos con  $\|w - w_1\|_p$ , eligiendo la mejor. Como  $w_1$  ya era el minimizador de  $f_{1,q}$  sobre  $W_1$ , esto significa que  $w_2$  es el minimizador de  $E_{1,p}$  entre todas las hipótesis en  $W_2$ , donde

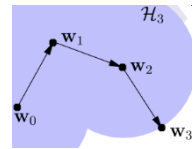


$$\mathcal{H}_2 = \mathcal{H}_1 \cup \{w : \|w - w_1\| \leq \eta\}.$$

Obsérvese que  $W \subset W_2$ . Del mismo modo, definimos el conjunto de hipótesis

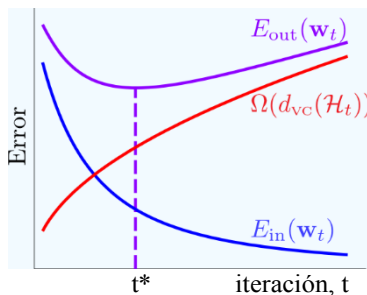
$$\mathcal{H}_3 = \mathcal{H}_2 \cup \{w : \|w - w_2\| \leq \eta\},$$

y en la 3ª iteración, elegimos pesos  $w_3$  que minimizan en  $W_3$ . Podemos continuar este argumento a medida que avanza el descenso gradiente, y definir una secuencia anidada de conjuntos de hipótesis



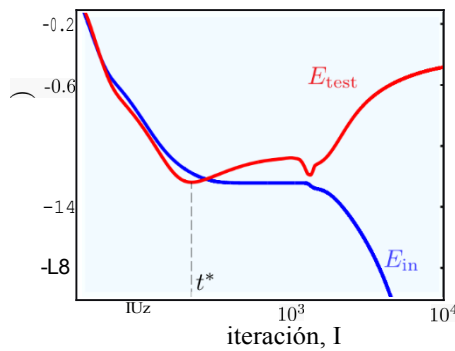
$$\mathcal{H}_1 \subset \mathcal{H}_2 \subset \mathcal{H}_3 \subset \mathcal{H}_4 \subset \dots$$

A medida que  $I$  aumenta,  $\text{in}(t)$  disminuye y  $dpz(F)$  aumenta. Por lo tanto, esperamos ver la compensación aproximación-generalización que se ilustró en la Figura 2.3 (reproducida aquí con la iteración  $I$  como proxy para  $dzz$ ):



La cifra sugiere que puede ser mejor parar pronto en algún  $I^*$ , mucho antes de que alcanzando un mínimo de  $E_{in}$ . De hecho, este panorama se observa en la práctica.

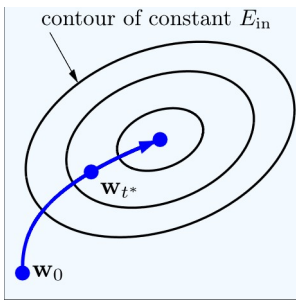
**Ejemplo T.2.** Volvemos a la tarea de dígitos de clasificar "1" frente a todos los demás dígitos, con 70 puntos de datos seleccionados aleatoriamente y una pequeña red neuronal sigmoideal con una sola unidad oculta y un nodo de salida  $\tanh(\cdot)$ . La siguiente figura muestra el error en la muestra y el error de prueba frente al número de iteraciones.



Las curvas refuerzan nuestra discusión teórica: el error de prueba disminuye inicialmente a medida que la ganancia de aproximación supera la peor barra de error de generalización; después, el error de prueba aumenta a medida que la barra de error de generalización empieza a dominar la ganancia de aproximación, y el sobreajuste se convierte en un problema grave.

En el ejemplo anterior, a pesar de utilizar una red neuronal parsimoniosa con un único nodo oculto, el sobreajuste era un problema porque los datos son ruidosos y la función objetivo es compleja, por lo que tanto el ruido estocástico como el determinista son significativos. Tenemos que regularizar.

En el ejemplo, es mejor parar pronto en  $t^*$  y disminuir el aprendizaje al conjunto de hipótesis más pequeño  $H_{t^*}$ . En este sentido, la detención temprana es una forma de regularización. La parada temprana está relacionada con el decaimiento del peso, como se ilustra a la derecha. Se inicializa  $w_0$  cerca de cero; si se detiene pronto en  $w_t$ , se habrá detenido en pesos más cercanos a  $w_k$ , es decir, pesos más pequeños. La parada temprana logra indirectamente pesos más pequeños, que es lo que el decaimiento del peso logra directamente. Para determinar cuándo detener el entrenamiento, utilice un conjunto de validación para controlar la



error de validación en la iteración  $I$  a medida que minimizas el error del conjunto de entrenamiento. Indique los pesos  $w_t$  que tienen el mínimo error de validación cuando termine el entrenamiento.

Después de seleccionar  $I^*$ , es tentador utilizar todos los datos para entrenar para iteraciones  $I^*$ . Desafortunadamente, volver a añadir los datos de validación y entrenar para  $I^*$  iteraciones puede conducir a un conjunto de pesos completamente diferente. La estimación de validación del rendimiento sólo es

válida para  $w_i$ , (los pesos que debe emitir). Esto parece ir en contra de la sabiduría de la curva de aprendizaje decreciente del capítulo 4: si aprendes con más datos, obtienes una hipótesis final mejor.<sup>0</sup>

---

\* Tampoco se recomienda utilizar todos los datos para entrenar hasta un error en la muestra de  $> \frac{1}{n}$  (usual). Además, un error en la muestra de  $< \frac{1}{n}$  puede incluso no ser alcanzable con todos los datos.

**Ejercicio 7\12**

¿Por qué la salida con/- en lugar de entrenar con todos los datos para /" iteraciones no va en contra de la sabiduría de que aprender con más datos es mejor.

*{Pista: "Más **datos es mejor**" se aplica a un modelo fijo ( $J$ -,  $A$ ). La **parada temprana** es la selección del modelo en un conjunto de hipótesis anidadas  $H_1 \subset H_2 \subset \dots$  determinadas por  $O$ " ; . ¿Qué ocurre si se utilizan todos los datos  $D$ )?*

Cuando se utiliza la parada anticipada, existe un equilibrio habitual a la hora de elegir el tamaño del conjunto de validación: si es demasiado grande, habrá pocos datos para entrenar; si es demasiado pequeño, el error de validación no será fiable.

**Ejercicio 7\13**

Supongamos que ejecuta el descenso por gradiente durante 1000 iteraciones. Tiene 500 ejemplos en  $\mathcal{D}$ , y utiliza 450 para  $\mathcal{D}_{\text{in}}$  y 50 para  $\mathcal{D}_{\text{val}}$ . Se obtiene el peso de la iteración 50, con  $A_{(50)} = 0.05$  y  $\text{Strain}(50) = 0.04$ .

- ¿Es  $C_{(50)}$  una estimación insesgada de  $T$ : " $i$  ( $wsD$ )?"
- Utilice el límite de Hoeffding para obtener un límite para  $T$ : " $i$  utilizando  $T$ :" más una barra de error. El límite debe mantenerse con una probabilidad de al menos 0,1.
- ¿Puedes enlazar  $T$ : " $i$  utilizando  $T$ :"; o necesitas más información?

El ejemplo 7.2 también ilustra otro problema común con el nodo de salida sigmoidal: el descenso de gradiente a menudo encuentra una región plana donde  $A_q$  disminuye muy poco." Se puede detener el entrenamiento pensando que se ha encontrado un mínimo local. Esta "parada temprana" por error se denomina a veces la propiedad de "autorregularización" de las redes neuronales sigmoidales. La regularización accidental debida a una convergencia mal interpretada no es fiable. La validación es mucho mejor.

### 7.4.3 Experimentos con datos numéricos

Pongamos en práctica la teoría en la tarea de los dígitos (clasificar "1" frente a todos los demás dígitos). Aprendemos sobre 500 puntos de datos elegidos al azar utilizando una red neuronal sigmoidal con una capa oculta y 10 nodos ocultos. Hay 41 pesos (parámetros sintonizables), por lo que hay más de 10 ejemplos por grado de libertad, lo que es bastante razonable. Utilizamos la transformación de salida de identidad  $8(s)$  -  $s$  para reducir la posibilidad de quedarse atascado en una región plana de la superficie de error. Al final del entrenamiento, utilizamos la transformación de salida  $6(s) = \text{sign}(s)$  para clasificar realmente los datos. Después de más de 2 millones de iteraciones de descenso de gradiente, conseguimos acercarnos a un mínimo local. El resultado

se muestra en la figura 7.2.

No hace falta ser un genio para darse cuenta del sobreajuste. La figura 7.2 demuestra la capacidad de aproximación de una red neuronal de tamaño moderado. Probemos el peso

---

" La función de transformación lineal de la salida ayuda a evitar esas regiones excesivamente planas.

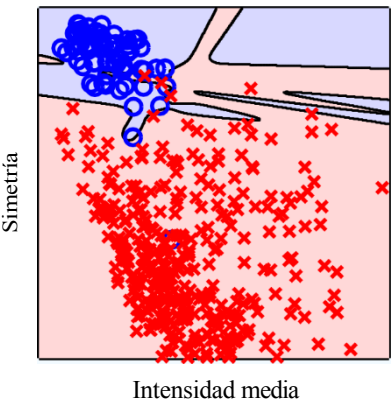
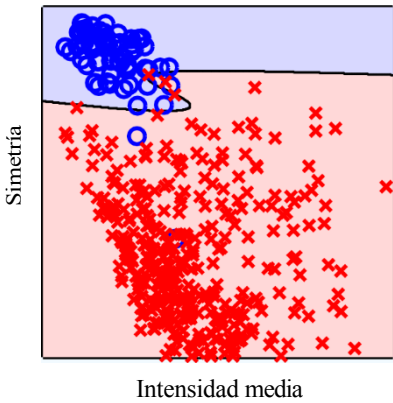


Figura 7.2: Red neuronal de 10 unidades ocultas entrenada con descenso gradiente en 500 ejemplos de los datos de dígitos (sin regularización). Los círculos azules son el dígito "1" y las x rojas son los demás dígitos. El sobreajuste es rampante.

para combatir el sobreajuste. Minimizamos  $A_{pg}(w; \lambda)$  en  $(\lambda)$  con  $I = 0,01$ . Obtenemos un separador mucho más creíble, que se muestra a continuación.



Como ilustración final, probemos la parada anticipada con un conjunto de validación de tamaño 50 (una décima parte de los datos); así, el conjunto de entrenamiento tendrá ahora un tamaño de 450. La dinámica de entrenamiento del descenso de gradiente se muestra en la Figura 7.3(a). La función de transformación de salida lineal ha ayudado a que no haya periodos extremadamente planos en el error de entrenamiento. El límite de clasificación con parada temprana en  $I^*$  se muestra en la Figura 7.3(b). El

---

resultado es similar a la caída de pesos. En ambos casos, el límite de clasificación regularizado es más creíble. En última instancia, lo que importa son las estadísticas cuantitativas, que se resumen a continuación.

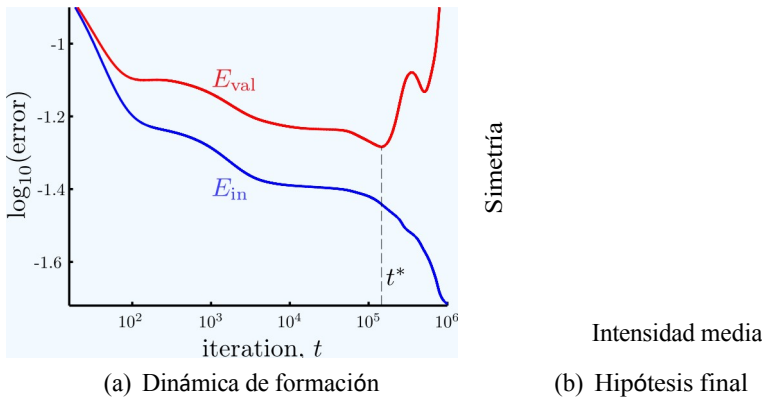


Figura 7.3: Parada temprana con 500 ejemplos de los datos de dígitos. (a) Errores de entrenamiento y validación para el descenso por gradiente con un conjunto de entrenamiento de tamaño 450 y un conjunto de validación de tamaño 50. (b) La hipótesis final "regularizada" obtenida mediante la parada temprana en  $t^*$ , el error de validación mínimo.

	$E_{\text{train}}$	$E_{\text{val}}$	$E_{\text{in}}$	$E_{\text{out}}$
No Regularization	—	—	0.2%	3.1%
Weight Decay	—	—	1.0%	2.1%
Early Stopping	1.1%	2.0%	1.2%	2.0%

## 7.5 Mejorar el descenso por gradiente

El descenso gradiente es un método simple para minimizar  $A_{\text{ip}}$  que tiene problemas de convergencia, especialmente con superficies de error planas. Una solución es minimizar un error más amigable, que es por lo que ayuda el entrenamiento con un nodo de salida lineal. En lugar de cambiar la medida del error, hay mucho margen para mejorar el propio algoritmo. El descenso gradiente da un paso de tamaño  $p$  en la dirección del gradiente negativo. ¿Cómo debemos determinar  $p$  y es el gradiente negativo la mejor dirección en la que moverse?

### Ejercicio 7\14

Consideremos la función de error  $(w - w^*)^T A (w - w^*)$ , donde  $A$  es una matriz definida positiva arbitraria. Fijemos  $w_0$ .

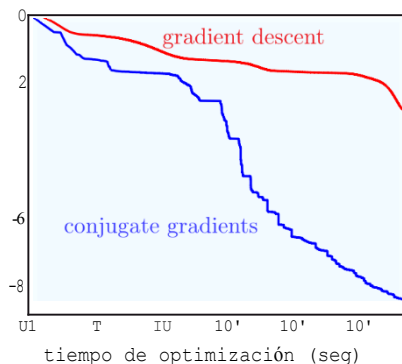
Demuestre que el gradiente  $2A(w - w^*)$  Qué pesos minimizan  $A(w)$ . ¿El descenso de gradiente le mueve en la dirección de estos pesos óptimos?

Concilia tu respuesta con la afirmación del capítulo 3 de que el gradiente es la mejor dirección para dar un paso. {Pista. ¿De qué tamaño era el paso?}



El ejercicio anterior muestra que el gradiente negativo no es necesariamente la mejor dirección para un paso grande, y nos gustaría dar pasos más grandes para aumentar

la eficacia de la optimización. La siguiente figura muestra dos algoritmos: nuestro viejo amigo el descenso por gradiente y nuestro pronto amigo el descenso por gradiente conjugado. Ambos algoritmos minimizan  $A_q$  para una red neuronal de 5 unidades ocultas que ajusta datos de 200 dígitos. La diferencia de rendimiento es espectacular.

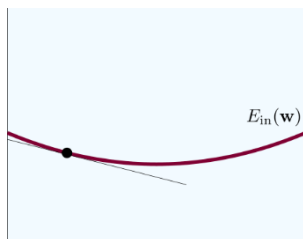


A continuación analizaremos métodos para "reforzar" el descenso de gradiente, pero sólo arañaremos la superficie de este importante tema conocido como optimización numérica. Los dos pasos principales en un procedimiento de optimización iterativo son determinar:

1. ¿En qué dirección hay que buscar un óptimo local?
2. ¿Hasta qué punto hay que dar un paso en esa dirección?

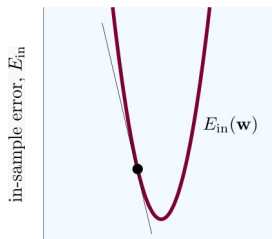
### 7.5.1 Elección de la tasa de aprendizaje $p$

En el descenso por gradiente, la tasa de aprendizaje  $p$  multiplica el gradiente negativo para dar el movimiento  $-p \nabla E_{in}$ . El tamaño del paso dado es proporcional a  $p$ . El tamaño óptimo del paso (y, por tanto, la tasa de aprendizaje  $p$ ) depende de lo ancha o estrecha que sea la superficie de error cerca del mínimo.



pesos,  $w$

ancho: utilizar  $p$  grande.



weights,  $w$

estrecho: utilizar  $p$  minúscula.

Cuando la superficie es más ancha, podemos dar pasos más grandes sin excedernos; como  $\|\nabla E_{in}\|$  es pequeño, necesitamos una  $p$  grande. Como no sabemos de antemano lo ancha que es la superficie, es fácil elegir un valor ineficiente para  $p$ .

Descenso por gradiente con tasa de aprendizaje variable. Una simple heurística que adapta la tasa de aprendizaje a la superficie de error funciona bien en la práctica. Si el error disminuye, se aumenta  $p$ ; si no, el paso ha sido demasiado grande, así que se rechaza la actualización y se disminuye  $p$ . A cambio de un pequeño esfuerzo adicional, obtenemos un aumento significativo del descenso por gradiente.

**Descenso por gradiente con tasa de aprendizaje variable:**

```

1: Inicializar  $w(0)$ , y  $p$  en  $\neq 0$ . Establecer  $\alpha$  y  $\beta$ 
2: mientras (no se cumple el criterio de parada) hacer
3:     Sea  $p = \alpha$  y  $v = \beta \nabla f(w)$ 
4:     si  $A, (w/t) T \quad v/)) < f''(w())$  ento
5:         aceptar:  $(+ \text{aces}() + v())$  "- "
6:     si no
7:         rechazar:  $w(t-t-1) = w(t)$ ;  $p \leftarrow t = \text{fipt}$ .
8:     Iterar hasta el paso siguiente,  $t \leftarrow t-1$ .
```

Por lo general, lo mejor es utilizar un parámetro de incremento conservador, por ejemplo  $\alpha = 0.05$  -  $0.1$ , y un parámetro de decremento un poco más agresivo, por ejemplo  $\beta = 0.5$  -  $0.8$ . Esto se debe a que, si el error no disminuye, entonces se está en una situación inusual y se requiere una acción más drástica.

Después de pensarlo un poco, cabe preguntarse por qué necesitamos una tasa de aprendizaje. Una vez determinada la dirección en la que hay que moverse,  $v(I)$ , ¿por qué no seguir simplemente en esa dirección hasta que el error deje de disminuir? Esto nos lleva al descenso por gradiente *más pronunciado* con *búsqueda de línea*.

Descenso más pronunciado. El descenso gradiente elige una dirección de descenso  $v(I) = -g(I)$  y actualiza los pesos a  $w(I - I - 1) = w(I) - I \cdot v(I)$ . En lugar de elegir  $p$  arbitrariamente, elegiremos el óptimo  $p$  que minimice  $A_{ip}(w(I - I - 1))$ . Una vez que tengas la dirección en la que moverte, hazlo lo mejor posible moviéndote a lo largo de la línea  $w(t) + I v(t)$  y deteniéndote cuando  $A_{ip}$  sea mínimo (de ahí el término búsqueda lineal).

Es decir, elegir un tamaño de paso  $p^*$ , donde

$$p^*(t) = \arg \min_{p \in L} (w(t) + p v(t))^T \nabla f(w(t) + p v(t)).$$

**Descenso más pronunciado (Gradient Descent - Line Search):**

```

1: Inicializar  $w(0)$  y fijar  $I = 0$ ;
2: mientras (no se cumple el criterio de parada) hacer
3:     Sea  $p^* = \arg \min_{p \in L} (w(I) + p \nabla f(w(I)))^T \nabla f(w(I) + p \nabla f(w(I)))$ 
4:     Sea  $p^* = \arg \min_{p \in L} (w(I) + p \nabla f(w(I)))^T \nabla f(w(I) + p \nabla f(w(I)))$ 
5:      $w(t+1) = w(t) + p^* \nabla f(w(t))$ .
6:     Iterar al paso siguiente,  $t \leftarrow t-1$ .
```

La búsqueda lineal del paso 4 es un problema de optimización unidimensional.

La búsqueda lineal es un paso importante en la mayoría de los algoritmos de optimización, por lo que se requiere un algoritmo eficiente. Escriba  $A(p)$  para  $A, q(w(I) - I - l_v(I))$ . El objetivo es encontrar un mínimo de  $A(p)$ . Damos un algoritmo sencillo basado en la búsqueda binaria.

Búsqueda en la línea. La idea es encontrar un intervalo en la línea que contenga un mínimo local. Luego, reducir rápidamente el tamaño de este intervalo manteniendo como invariante el hecho de que contiene un mínimo local.

El invariante básico es un arreglo en  $U$ :

$$\eta_1 < \eta_2 < \eta_3 \text{ with } E(\eta_2) < \min\{E(\eta_1), E(\eta_3)\}.$$

Como  $f$  es continua, debe haber un mínimo local en el intervalo  $p, p_3$ . Consideremos ahora el punto medio del intervalo,

$$\bar{\eta} = \frac{\eta_1 + \eta_3}{2},$$

de ahí el nombre de *algoritmo de bisección*. Supongamos que  $h < c_2$  como se muestra. Si  $f(p) < +(c_2)$  entonces  $(c, c_2)$  es un nuevo arreglo  $U$  más pequeño; y, si  $A(\eta) > +(z)$ , entonces  $(F, 9z, \eta)$  es el nuevo arreglo  $U$  más pequeño. En cualquier caso, el proceso de bisección se puede iterar con el nuevo arreglo en  $U$ . Si  $p$  resulta ser igual a  $c_2$  per-

turb  $p$  ligeramente para resolver la degeneración. Dejamos al lector que determine cómo obtener el nuevo arreglo  $U$  más pequeño para el caso  $p >$

Un algoritmo eficaz para encontrar una disposición inicial en  $U$  es comenzar con  $p_1 - 0$  y  $e$  para algún paso  $e$ . Si  $A(9z)$  (pt), considere la secuencia  $9 = 0, e, 2e, 4e, 8e,$

...

(cada vez que se duplica el paso). En algún momento, el error debe aumentar. Cuando el error aumenta por primera vez, los tres últimos pasos dan una disposición en  $U$ . Si, en cambio,  $A(p) > A(p_2)$ , consideramos la secuencia

$$=^{\wedge}, 0, -c, -2c, -4c, -8^{\wedge},$$

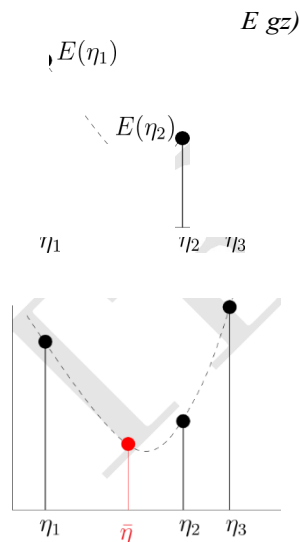
(el paso sigue duplicándose pero en sentido inverso). De nuevo, cuando el error aumenta por primera vez, los tres últimos pasos dan un arreglo en  $U$ .<sup>12</sup>

#### Ejercicio 7\15

Demuestre que  $|f(a) - f(i)|$  decrece exponencialmente en el algoritmo de bisección.

{Pista.' demuestran que dos iteraciones reducen al menos a la mitad el tamaño del intervalo).

El algoritmo de bisección continúa bisecando el intervalo y actualizando a una nueva disposición en  $U$  hasta que la longitud del intervalo  $|p_3 - p|$  es lo suficientemente pequeña, momento en el cual



---

<sup>12</sup> No nos preocupamos por  $\epsilon(pJ) = (02)$ - tales vínculos pueden romperse con pequeñas perturbaciones.

puede devolver el punto medio del intervalo como mínimo local aproximado. Normalmente 20 iteraciones de bisección son suficientes para obtener una solución aceptable. En el problema 7.8 se da un algoritmo de interpolación cuadrática mejor, que en la práctica sólo necesita unas 4 iteraciones.

Ejemplo 7.3. Ilustramos estas tres heurísticas para mejorar el gradiente descenso en nuestro reconocimiento de dígitos (clasificación de "1" frente a otros dígitos). Utilizamos

200 puntos de datos y una red neuronal con 5 unidades ocultas. En la Figura 7.4 mostramos el rendimiento del descenso por gradiente, el descenso por gradiente con tasa de aprendizaje variable y el descenso más pronunciado (búsqueda lineal). La tabla siguiente resume el error en la muestra en varios puntos de la optimización.

Método	Tiempo de optimización		
	10 segundos	1.000 seg.	50.000 seg.
Descenso gradual	0.079	0.0206	0.00144
Descenso Gradiente Estocástico	0.0213	0.00278	0.000022
Tasa de aprendizaje variable	0.039	0.014	0.00010
Descenso más pronunciado	0.043	0.0189	0.000012

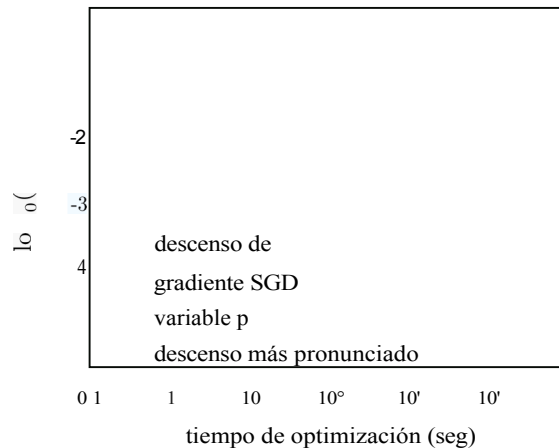


Figura 7.4: Descenso gradual, tasa de aprendizaje variable y descenso más pronunciado utilizando datos de dígitos y una red neuronal de 2 capas y 5 unidades ocultas con salida lineal. Para la tasa de aprendizaje variable,  $\alpha = 1,1$  y  $\beta = 0,8$ .

Obsérvese que el SGD es bastante competitivo. La figura ilustra por qué es difícil saber cuándo dejar de minimizar. Una región plana "atrapó" a todos los métodos, a pesar de que utilizamos una transformación lineal del nodo de salida. Es muy difícil diferenciar entre una región plana (que suele estar causada por un valle muy pronunciado que conduce a un comportamiento en zig-zag

ineficiente) y un verdadero mínimo local.

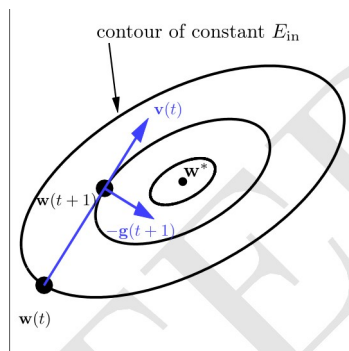
@ Abu-Mostafa, Magdon-Ismail, Lin: Jan-2015

e-Chap:7-33



## 7.5.2 Minimización del gradiente conjugado

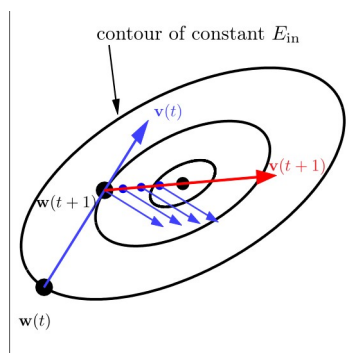
El gradiente conjugado es la reina entre los métodos de optimización porque aprovecha un principio sencillo. No deshaga lo que ya ha conseguido. Cuando finaliza una búsqueda lineal, dado que el error no puede reducirse aún más retrocediendo o avanzando a lo largo de la dirección de búsqueda, debe ser que el nuevo gradiente y la dirección de búsqueda lineal anterior sean ortogonales. Esto significa que se ha conseguido que una de las componentes del gradiente sea cero, es decir, la componente a lo largo de la dirección de búsqueda  $v(t)$  (véase la figura).



Si la siguiente dirección de búsqueda es la negativa del nuevo gradiente, será ortogonal a la dirección de búsqueda anterior.

Se está en un mínimo local cuando el gradiente es cero, y poner un componente a cero es sin duda un paso en la dirección correcta. A medida que se desplaza a lo largo de la siguiente dirección de búsqueda (por ejemplo, el nuevo gradiente negativo), el gradiente cambiará y puede que no permanezca ortogonal a la dirección de búsqueda anterior, una tarea que realizó laboriosamente en la búsqueda de línea anterior. El algoritmo de gradiente conjugado elige la siguiente dirección  $v(t+1)$  de forma que el gradiente a lo largo de esta dirección, *permanecerá perpendicular a la dirección de búsqueda anterior*  $v(t)$ . Esto se llama la dirección conjugada, de ahí el nombre.

Después de una búsqueda lineal a lo largo de esta nueva dirección  $v(t+1)$  para minimizar  $\min'$  se habrán puesto a cero dos componentes del gradiente. En primer lugar, el gradiente permanece perpendicular a la dirección de búsqueda anterior  $v(t)$ .



En segundo lugar, el gradiente será ortogonal a  $v(t+1)$  debido a la búsqueda de la línea (véase la figura). El gradiente a lo largo de la nueva dirección  $v(t+1)$  se muestra con las flechas azules de la figura. Como  $v(t+1)$  es conjugado a  $v(t)$ , observe cómo el gradiente a medida que nos movemos a lo largo de  $v(t+1)$  permanece ortogonal a la dirección anterior  $v(t)$ .

### Ejercicio 7.16

¿Por qué la nueva dirección de búsqueda pasa por los pesos óptimos?

Hemos progresado. Ahora dos componentes del gradiente son cero. En dos

dimensiones, esto significa que el propio gradiente debe ser cero y hemos terminado.

@ Abu-Mostafa, Magdon-Ismail, Lin: Jan-2015

e-Chap:7-34

En una dimensión más alta, si pudiéramos continuar poniendo un componente del gradiente a cero con cada búsqueda de línea, manteniendo todos los componentes anteriores a cero, eventualmente pondríamos cada componente del gradiente a cero y estaríamos en un mínimo local. Nuestra discusión es válida para una función de error cuadrática idealizada. En general, la minimización del gradiente conjugado implementa nuestras expectativas idealizadas aproximadamente. Sin embargo, funciona a las mil maravillas porque la configuración idealizada es una buena aproximación una vez que nos acercamos a un mínimo local, y aquí es donde algoritmos como el descenso de gradiente se vuelven ineficaces.

Pasemos ahora a los detalles. El algoritmo construye la dirección de búsqueda actual como una combinación lineal de la dirección de búsqueda anterior y el gradiente actual,

$$\mathbf{v}(t) = -\mathbf{g}(t) + \mu_t \mathbf{v}(t-1)$$

donde

$$\mu_t = \frac{\mathbf{g}(t+1)^T (\mathbf{g}(t+1) - \mathbf{g}(t))}{\mathbf{g}(t)^T \mathbf{g}(t)}$$

El término  $\mathbf{v}(t-1)$  se denomina término de impulso porque te pide que sigas moviéndote en la misma dirección en la que te estabas moviendo. El multiplicador  $\mu_t$  se denomina parámetro de impulso. El algoritmo completo de descenso de gradiente conjugado se resume en el siguiente cuadro de algoritmo.

**Descenso por gradiente conjugado:**

- 1: Inicializar  $\mathbf{w}(0)$  y fijar  $I = 0$ ; fijar  $\mathbf{v}(-1) = \mathbf{0}$
- 2: mientras no se cumpla el criterio de parada hacer
- 3: Sea  $\mathbf{v}(I) = -\mathbf{g}(I) + \mu_t \mathbf{v}(I-1)$ , donde

$$\mu_t = \frac{\mathbf{g}(t+1)^T (\mathbf{g}(t+1) - \mathbf{g}(t))}{\mathbf{g}(t)^T \mathbf{g}(t)}$$

- 4: Let  $\eta^* = \arg\min_{\eta} E_{\text{in}}(\mathbf{w}(t) + \eta \mathbf{v}(t))$ .
- 5:  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta^* \mathbf{v}(t)$ ;
- 6: Iterar al paso siguiente,  $I \leftarrow I + 1$ ;

La única diferencia entre el descenso por gradiente conjugado y el descenso más pronunciado está en el paso 3, en el que la dirección de búsqueda de la línea es distinta de la del gradiente negativo. Contrariamente a la intuición, la dirección del gradiente negativo no siempre es la mejor dirección para moverse, ya que puede deshacer parte del buen trabajo realizado anteriormente. En la práctica, para las superficies de error que no son exactamente cuadráticas, los  $\mathbf{v}(I)$  son sólo aproximadamente conjugados y se recomienda "reiniciar" el algoritmo poniendo  $\mathbf{v}$  a cero cada cierto tiempo (por ejemplo cada  $d$  iteraciones). Es decir, cada  $d$  iteraciones se lanza una iteración de descenso más pronunciado.

**Ejemplo 7.4.** Siguiendo con el ejemplo de los dígitos, comparamos el gradiente conjugado y el descenso más pronunciado del campeón anterior en la siguiente tabla y en la Figura 7.5.

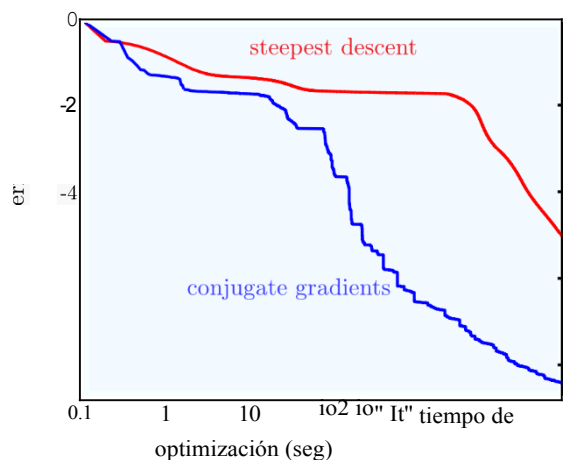


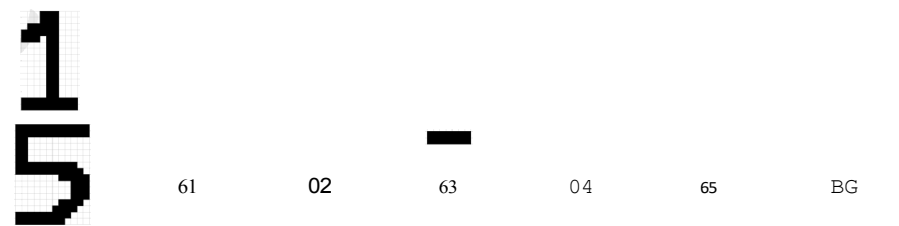
Figura 7.5: Descenso más pronunciado frente a descenso de gradiente conjugado utilizando 200 ejemplos de los datos de dígitos y una red neuronal sigmoidal de 2 capas con 5 unidades ocultas.

Method	Método	Optimization Time
Steepest Descent	0.043	1.000 sec
Conjugate Gradients	0.0200	$1.13 \times 10^{-6}$
Gradien		

La diferencia de rendimiento es espectacular. □

## 7.6 Aprendizaje profundo: Redes con muchas capas

La aproximación universal dice que una sola capa oculta con suficientes unidades ocultas puede aproximar cualquier función objetivo. Pero puede que no sea una forma natural de representar la función objetivo. A menudo, muchas capas imitan mejor el aprendizaje humano. Empecemos con el problema del reconocimiento de dígitos para clasificar "1" frente a "5". Un primer paso natural es descomponer los dos dígitos en componentes básicos, igual que se descompone una cara en dos ojos, una nariz, una boca, dos orejas, etc. He aquí un intento de clasificar un "1" y un "5" prototípicos.





combina  $c$  y  $z_b$  en la salida final. En este punto, es útil rellenar todos los espacios en blanco con un ejercicio.

#### Ejercicio 7\18

Dado que la entrada  $x$  es una imagen, es conveniente representarla como una matriz  $z_j$  de sus píxeles que son negros ( $z_j = 1$ ) o blancos ( $z_j = 0$ ). La forma básica  $\phi_k$  identifica un conjunto de estos píxeles que son negros.

- (a) Demuestre que  $\phi_k$  puede ser calculada por el nodo de la red neuronal

$$\phi_k(\mathbf{x}) = \tanh\left(w_0 + \sum_i w_{ij} x_{ij}\right)$$

- (b) ¿Cuáles son las entradas del nodo de la red neuronal?
- (c) ¿Qué valores eliges para los pesos? [Sugerencia: considere por separado los pesos de los píxeles  $z_j$  y los  $\phi_k$ ].
- (d) ¿Cómo elegirías  $w_0$ ? (No todos los dígitos se escriben igual, por lo que es posible que una forma básica no siempre se represente exactamente en la imagen).
- (e) Dibuja la red final, rellenando todos los detalles que puedas.

Se habrá dado cuenta de que la salida de  $z_j$  es todo lo que necesitamos para resolver nuestro problema. Este no sería el caso si estuviéramos resolviendo el problema multiclase completo con los nodos  $z_3, \dots, z_g$  correspondientes a los diez dígitos.  $z_g$  correspondientes a los diez dígitos. Además, hemos resuelto nuestro problema con relativa facilidad, ya que nuestra red "profunda" sólo tiene dos capas ocultas. En un problema más complejo, como el reconocimiento de caras, el proceso empezaría igual que aquí, con formas básicas. En el siguiente nivel, constituiríamos formas más complicadas a partir de las formas básicas, pero aún no habríamos llegado a casa. Estas formas más complicadas constituirían formas aún más complicadas hasta que por fin tuviéramos objetos realistas como ojos, boca, orejas, etc. Habría una jerarquía de características "básicas" hasta que resolviéramos nuestro problema al final.

Y ahora, el remate y el quid de la historia. Primero el remate. Enfoque la red que hemos construido y observe lo que hacen las distintas capas. La primera capa construye una representación de bajo nivel de formas básicas; la siguiente capa construye una representación de nivel superior a partir de estas formas básicas. A medida que subimos capas, obtenemos representaciones más complejas a partir de partes más sencillas de la capa anterior: una descomposición "inteligente" del problema, que empieza por lo sencillo y se va haciendo más complejo, hasta que finalmente se resuelve el problema. Esta es la promesa de la red profunda, que proporciona una visión similar a la humana de cómo se está resolviendo el problema basándose en una jerarquía de representaciones más complejas de la entrada. Si bien con una sola capa oculta podríamos alcanzar una solución de precisión similar, no obtendríamos esa percepción. El panorama es halagüeño para nuestro problema de reconocimiento

intuitivo de dígitos, pero aquí está el quid de la cuestión: para un problema de aprendizaje complejo, ¿cómo automatizamos todo esto en un algoritmo informático?

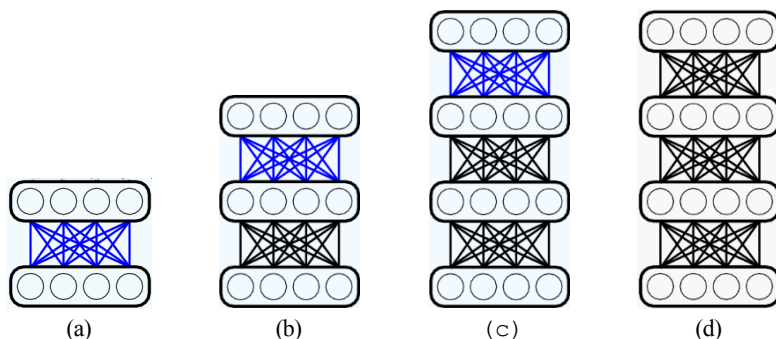


Figura 7.6: Algoritmo de aprendizaje profundo Greedy. (a) Se aprenden los pesos de la primera capa. (b) Se fija la primera capa y se aprenden los pesos de la segunda. (c) Se fijan las dos primeras capas y se aprenden los pesos de la tercera. (d) Los pesos aprendidos pueden utilizarse como punto de partida para ajustar toda la red.

### 7.6.1 Algoritmo codicioso de aprendizaje profundo

Históricamente, las redes neuronales poco profundas (de una sola capa oculta) han prevalecido sobre las redes profundas porque son difíciles de entrenar, presentan muchos mínimos locales y, en relación con el número de parámetros de caída, tienen una gran tendencia a sobreajustarse (la composición de no linealidades suele ser mucho más potente que una combinación lineal de no linealidades). Recientemente, algunas heurísticas sencillas han demostrado empíricamente un buen rendimiento y han devuelto el protagonismo a las redes profundas. De hecho, el mejor algoritmo actual para el reconocimiento de dígitos es una red neuronal profunda entrenada con dichas heurísticas.

La heurística codiciosa tiene una forma general. Se aprenden los pesos de la primera capa  $W^1$  y se fijan.<sup>10</sup> La salida de la primera capa oculta es una transformación no lineal de las entradas  $x^0$  a  $x^1$ . Estas salidas  $x^1$  se utilizan para entrenar los pesos de la segunda capa  $W^2$ , *manteniendo fijos los pesos de la primera capa*. Este es la esencia del algoritmo codicioso, elegir "codiciosamente" los pesos de la primera capa, fijarlos, y luego pasar a los pesos de la segunda capa. Se ignora la posibilidad de que puedan existir mejores ponderaciones de la primera capa si se tiene en cuenta lo que hace la segunda capa. El proceso continúa con las salidas  $x^2$  utilizadas para aprender las ponderaciones  $W^3$ , y así sucesivamente.

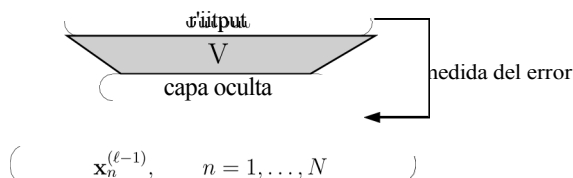


---

"Recordemos que utilizamos el superíndice (-) para denotar la capa fi.

### Algoritmo Greedy de aprendizaje profundo:

1. para  $I = 1, \dots, \tilde{n}$  hacer,  $\tilde{n}$  do
- 2:  $W^{(I)}$  y  $V^{(I)}$  son dadas de iteraciones anteriores. 3:  
Calcular las salidas  $x_n^{(I)}$  de la capa  $I - 1$  para  $n = 1, \dots, N$ .
- 4: Utilizar  $(x_n^{(I)}, y_n)$  para aprender los pesos  $W^{(I)}$  entrenando una red neuronal simple de capa oculta. ( $W^{(I)}$  y  $V^{(I)}$  son fijas).



Tenemos que aclarar el paso 4 del algoritmo. Se aprenden los pesos  $W^{(I)}$  y  $V$ , aunque  $V$  no es necesario en el algoritmo. Para aprender los pesos, minimizamos un error (que dependerá de la salida de la red), y ese error aún no está definido. Para definir el error, primero debemos definir la salida y luego cómo calcular el error a partir de la salida.

Autocodificador no supervisado. Un enfoque consiste en tomar en serio la noción de que la capa oculta ofrece una representación de alto nivel de las entradas. Es decir, deberíamos ser capaces de reconstruir todos los aspectos importantes de la entrada a partir de la salida de la capa oculta. Una prueba natural es reconstruir la propia entrada: la salida será  $x_p$ , una predicción de la entrada  $x$ ; y, el error es la diferencia entre las dos. Por ejemplo, utilizando el error al cuadrado,

$$e_p = \|x_p - x\|^2.$$

Cuando todo está dicho y hecho, obtenemos los pesos sin usar los objetivos  $p$  y la capa oculta da una codificación de las entradas, de ahí el nombre de autocodificador no supervisado. Esto recuerda a la red de funciones de base radial del capítulo 6, en la que utilizamos una técnica no supervisada para aprender los centros de las funciones de base, lo que proporcionó un conjunto representativo de entradas como centros. Aquí, vamos un paso más allá y diseccionamos el propio espacio de entrada en trozos representativos del problema de aprendizaje. Al final, los objetivos tienen que volver a aparecer (normalmente en la capa de salida).

Red profunda supervisada. El enfoque anterior se adhiere al objetivo filosófico de que las capas ocultas proporcionen una representación jerárquica "inteligente" de las entradas. Un enfoque más directo consiste en entrenar la red de dos capas en los objetivos. En este caso, la salida es el objetivo predicho y la

medida de error  $e_p(p, p)$  se calcularía de la forma habitual (por ejemplo, error al cuadrado, error de entropía cruzada, etc.).

En la práctica, no hay un veredicto sobre qué método es mejor. El campo de los autocodificadores no supervisados está ligeramente más concurrido que el de los supervisados. Pruebe los dos y vea qué funciona para su problema, suele ser lo mejor. Una vez que tenga su medida de error, simplemente busque en su caja de herramientas de optimización y minimice el error utilizando su método favorito (descenso de gradiente, descenso de gradiente estocástico, descenso de gradiente conjugado, ...). Una táctica común es utilizar primero el autocodificador no supervisado para establecer los pesos y luego afinar toda la red utilizando el aprendizaje supervisado. La idea es que el paso no supervisado te lleve al mínimo local correcto de la red completa. Pero, independientemente del bando al que pertenezcas, tienes que elegir la arquitectura de la red profunda (el número de capas ocultas y su tamaño), y para eso no hay ninguna poción mágica. Tendrás que recurrir a viejos trucos como la validación, o a una comprensión profunda del problema (nuestra red hecha a mano para la tarea "1" frente a "5" sugiere una red profunda con seis nodos ocultos en la primera capa oculta y dos en la segunda).

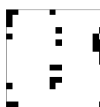
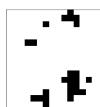
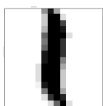
### Ejercicio 7.19

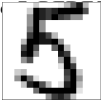
Anteriormente, para nuestro problema de los dígitos, utilizamos la simetría y la intensidad. ¿Cómo se relacionan estas características con las redes profundas? ¿Siguen siendo necesarias?

Ejemplo T.S. Aprendizaje profundo para el reconocimiento de dígitos. Volvamos al problema de clasificación de dígitos '1' frente a '5' utilizando una arquitectura de red profunda

$$d(1), d^{*1}, d^{21}, d^{31}] \quad [256, 6, 2, 1].$$

(La misma arquitectura que construimos antes a mano, con 16 x 16 píxeles de entrada y 1 de salida). Utilizaremos el descenso de gradiente para entrenar las redes de dos capas en el algoritmo codicioso. En el problema 7.7 se da una forma matricial conveniente para el gradiente de la red de dos capas. Para el codificador automático no supervisado, la salida objetivo es la matriz de entrada  $X$ . Para la red profunda supervisada, la salida objetivo es sólo el vector objetivo  $y$ . Utilizamos el enfoque supervisado con 1.000.000 iteraciones de descenso de gradiente para cada paso greedy supervisado utilizando una muestra de 1.500 ejemplos de los datos de dígitos. A continuación se muestra lo que aprendieron las 6 unidades ocultas de la primera capa oculta. Para cada nodo oculto de la primera capa oculta, mostramos los píxeles correspondientes a los 20 pesos entrantes más altos.





61

02

63

04

65

BG

Los datos reales no son tan limpios como nuestros análisis idealizados. No se sorprenda. Sin embargo, podemos discernir que az ha seleccionado los píxeles (formas) en el típico "1" que es poco probable que estén en un típico "5". Las demás características parecen centrarse en el "5" y, en cierta medida, coinciden con las que hemos construido a mano. No vamos a insistir en si la representación refleja la intuición humana.

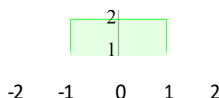
alcance. Lo importante es que este resultado es automático y depende exclusivamente de los datos (aparte de nuestra elección de la arquitectura de red); lo que importa es el rendimiento fuera de la muestra. Para las distintas arquitecturas, realizamos más de 1.000 experimentos de validación seleccionando 500 puntos de entrenamiento aleatorios cada vez y los datos restantes como conjunto de pruebas.

Arquitectura de redes profundas	A;p	prueba
t256,3,2, I]	0	0.170%
t256,6,2, I]	0	0.187%
t256, 12,2, 1/	0	0,187Po
t256,24,2, 1/	0	0,183Po

en US fltways cero porque hay muchos parámetros, incluso con sólo 3 unidades ocultas en la primera capa oculta. Esto huele a sobreajuste. Pero el resultado de la prueba es impresionante, con un 99,8% de precisión, que es lo único que nos importa. Nuestras características de simetría e intensidad construidas a mano eran buenas, pero no tanto. O

## 7.7 Problemas

**Problema 7.1** Implementa la siguiente función de decisión utilizando un perceptrón de 3 capas.



**Problema 7.2** Un conjunto de  $M$  hiperplanos generalmente dividirá el espacio en algún número de regiones. Cada punto de  $\mathbb{R}^d$  puede etiquetarse con un vector de  $M$  dimensiones que determina en qué lado de cada plano se encuentra. Así, por ejemplo, si  $M = 3$ , entonces un punto con un vector  $(-1, -1, -1)$  está en el lado -1 del primer hiperplano, y en el lado -1 del segundo y tercer hiperplanos. Una región se define como el conjunto de puntos con la misma etiqueta.

- Demuestra que las regiones con la misma etiqueta son convexas.
- Demostrar que  $M$  hiperplanos pueden crear a lo sumo  $2^M$  regiones distintas.
- [difícil] ¿Cuál es el número máximo de regiones creadas por  $M$  hiperplanos en  $d$  dimensiones?

{Respuesta:  $\frac{d!}{2} + 1$ ).

{Pista: Utilice la inducción y deje que  $B(M, d)$  sea el número de regiones creadas por  $M$  hiperplanos en el espacio  $d$ . Consideremos ahora la adición del  $(M+1)$  hiperplano. Demostrar que este hiperplano interseca a lo sumo  $B(M, d)$  de las regiones  $B(M, d)$ . Por cada región que cruza, añade exactamente una región, por lo que  $B(M+1, d) = B(M, d) + B(M, d-1)$ . (¿Te suena esta recurrencia?) Evalúa las condiciones de contorno:  $B(M, 0)$  y  $B(0, d)$ , y proceder a partir de ahí. Para ver que el hiperplano  $M+1$  sólo interseca regiones  $B(M, d-1)$ , argumentar de la siguiente manera. Tratar el hiperplano  $M$  como un espacio  $(d-1)$ -dimensional, y proyectar los  $M$  hiperplanos iniciales en este espacio para obtener  $M$  hiperplanos en un espacio  $(d-1)$ -dimensional. Estos  $M$  hiperplanos pueden crear como máximo  $B(M, d-1)$  regiones en este espacio. Argumentemos que esto significa que el hiperplano  $M+1$  sólo interseca a lo sumo  $B(M, d-1)$  de las regiones creadas por los hiperplanos  $M$  en el espacio  $d$ ).

**Problema 7.3** Supongamos que una función objetivo  $f$  (para clasificación) está representada por un número de hiperplanos, donde las diferentes regiones definidas por los hiperplanos (véase el Problema 7.2) podrían clasificarse  $\pm 1$ , como en los ejemplos bidimensionales que hemos considerado en el texto. Sean los hiperplanos  $l_1, l_2, \dots, l_r$ , donde  $l_i(x) = \text{sign}(w_i - x)$ . Consideremos todas las regiones que se clasifican  $\pm 1$ , y que una de tales regiones sea  $r$ . Sea  $c = (c_1, c_2, \dots, c_r)$  la etiqueta de cualquier punto de la región (todos los puntos de una región dada tienen la misma etiqueta); la etiqueta  $c_r$  indica en qué lado  $l_r$  se encuentra el punto. Defina el término AND correspondiente a la región  $r$  mediante

$$f_r = \prod_{i=1}^r l_i^{c_i}, \text{ donde } l_i^{c_i} = \begin{cases} h_i & \text{if } c_i = +1, \\ h_i & \text{if } c_i = -1. \end{cases}$$

Demuestre que  $f = \sum_{r=1}^R f_r$ , donde  $r_1, \dots, r_R$  son todas las regiones positivas. (Usamos multiplicación para los operadores AND y suma para los operadores OR).

**Problema 7.4** Remitiéndonos al problema 7.3, cualquier función objetivo que pueda descomponerse en hiperplanos  $l_1, \dots, l_k$  puede representarse por  $f = \sum_{r=1}^R f_r$ , donde hay  $k$  regiones positivas.

Cuál es la estructura del perceptrón de 3 capas (número de unidades ocultas en cada capa) que implementará esta unión, demostrando el siguiente teorema:

**Teorema.** Cualquier unión de decisión cuyas regiones  $\pm 1$  se definen en términos de las regiones creadas por un conjunto de hiperplanos puede ser implementada por un perceptrón de 3 capas.

**Problema 7.5 [Difícil]** Enuncie y demuestre una versión de un *teorema de Aproximación Universal*:

**Teorema.** Cualquier función objetivo  $f$  (o clasificación) definida en  $\mathbb{R}^d$ , cuyas superficies límite de clasificación son suaves, puede ser aproximada arbitrariamente por un perceptrón de 3 capas.

*{Pista: Descomponer el hipercubo unitario en  $c$ -hipercubos (de ellos), El volumen de estos  $c$ -hipercubos que interseca los límites de clasificación debe tender a cero (¿por qué? - utilizar suavidad). Así, la función que toma el valor de  $f$  en cualquier  $c$ -hipercubo que no interseca la frontera y un valor arbitrario en estos  $c$ -hipercubos frontera se aproximará a  $f$  arbitrariamente cerca, como  $O(\epsilon)$ .}*

**Problema 7.6** La aproximación de diferencia finita para obtener el gradiente se basa en la siguiente fórmula de cálculo:

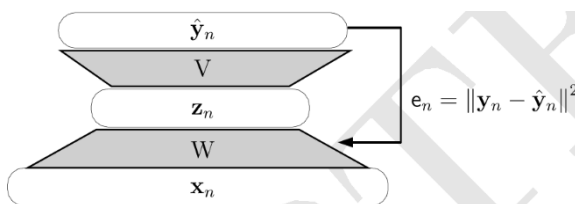
$$\frac{\partial h}{\partial w_{ij}^{(\ell)}} = \frac{h(w_{ij}^{(\ell)} + \epsilon) - h(w_{ij}^{(\ell)} - \epsilon)}{2\epsilon} + O(\epsilon^2),$$



donde  $\frac{\partial E}{\partial w_{ij}}$  denota el valor de la derivada cuando todos los pesos se mantienen en sus valores en  $w$  excepto (o el peso  $w_{ij}$ ), que es perturbado por  $r$ . Para obtener el gradiente, necesitamos la derivada parcial con respecto a cada peso.

Demuestre que la complejidad computacional de obtener todas estas derivadas parciales es  $O(W^2)$ . [Pista: hay que hacer dos propagaciones hacia delante para cada peso].

**Problema 7.7** Considere la siguiente red de 2 capas, con el vector de salida  $y$ . Esta es la red de dos capas utilizada para el algoritmo de red profunda codiciosa.



Recoge los vectores de entrada  $x_n$  (junto con una columna de unos) como filas de la matriz de datos de entrada  $X$ , y de forma similar forma  $Z$  a partir de  $z$ . Las matrices objetivo  $Y$  y  $\hat{Y}$  se (ordenan) de forma similar. Supongamos que el nodo de salida es lineal y que la activación de la capa oculta es  $\sigma(\cdot)$ .

(a) Demuestre que el error dentro de la muestra es

$$E = \frac{1}{2} \text{tr}((Y - \hat{Y})(Y - \hat{Y})^T),$$

donde

$$\begin{aligned} X & \text{ es } N \times (d+1) \\ W & \text{ es } (d+1) \times d' \\ Z & \text{ es } N \times d' \\ Y & \text{ es } N \times d \\ V & \text{ es } (d' + 1) \times \dim(y) \\ \hat{Y} & \text{ es } N \times \dim(y) \end{aligned}$$

(Es conveniente descomponer  $V$  en su primera fila  $V_0$  correspondiente a los sesgos y sus restantes filas  $V_1$ :  $1$  es el vector  $N \times 1$  de unos).

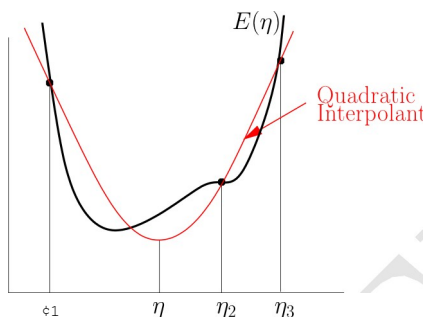
(b) derivar las matrices de gradiente:

$$\frac{\partial E}{\partial V} = Z^T (Y - \hat{Y})$$

$$\frac{\partial E}{\partial W} = X^T \odot (XW \odot (q(XW)V + 1)V - Y)$$

donde  $\odot$  significa multiplicación por elementos. Algunas de las derivadas matriciales de funciones que implican la traza del apéndice pueden ser útiles.

**Problema 7.8** Interpolación cuadrática para la búsqueda de líneas Supongamos que se ha encontrado un arreglo en  $U$ , como se ilustra a continuación.



En lugar de utilizar la bisección para construir el punto  $z$ , la interpolación cuadrática ajusta una curva cuadrática  $C(q)$  o  $q^2 + bq + c$  a los tres puntos y utiliza el mínimo de esta interpolación cuadrática como  $q$ .

- (a) Demuestre que el mínimo  $q$  de la interpolante cuadrática para un arreglo en  $U$  está dentro del intervalo  $tq_1, tq_2$ .
- (b) Sea  $e_i = A(\eta_i)$ ,  $e_3 = C(\eta_3)$ . Obtener la función cuadrática que interpola los tres puntos  $(\eta_1, e_1)$ ,  $(\eta_2, e_2)$ ,  $(\eta_3, e_3)$ . Demostrar que el mínimo de esta interpolante cuadrática viene dado por:

$$q = \frac{1}{2} \frac{(e_1 - e_2)(\eta_2^2 - \eta_3^2) - (e_1 - e_3)(\eta_1^2 - \eta_2^2)}{(e_1 - e_2)(\eta_1 - \eta_3) - (e_1 - e_3)(\eta_1 - \eta_2)}$$

{Pista:  $e_1 = a\eta_1^2 + b\eta_1 + c$ ,  $e_2 = a\eta_2^2 + b\eta_2 + c$ ,  $e_3 = a\eta_3^2 + b\eta_3 + c$ . Solve para  $a, b, c$  y el mínimo de la cuadrática viene dado por  $q = -b/2a$ .}

- (c) Dependiendo de si  $C(z)$  es menor que  $A(\eta_2)$ , y de si  $z$  está a la izquierda o a la derecha de  $\eta_2$ , hay 4 casos. En cada caso, ¿cuál es la disposición en  $U$  más pequeña?
- (d) ¿Y si  $z = \eta_2$  un caso degenerado?

Nota: en general, las interpolaciones cuadráticas convergen muy rápidamente a una  $q$  localmente óptima. En la práctica, 4 iteraciones son más que suficientes.

**Problema T.9** [Convergencia de la minimización de Monte-Carlo] Supongamos que el mínimo global  $w^*$  está en el cubo unitario y la superficie de error es cuadrática cerca de  $w^*$ . Entonces, cerca de  $w^*$ ,

$$E(w) = E(w^*) + \frac{1}{2}(w - w^*)^T H (w - w^*)$$

donde el hessiano  $H$  es definido positivo y simétrico.

- (a) Si se muestrea uniformemente  $\mathbf{w}$  en el cubo unitario, demuestre que

$$P[E \leq E(\mathbf{w}^*) + \epsilon] = \frac{\int_{\mathbf{x}^* H \mathbf{x} < 2\epsilon} d^d \mathbf{x}}{\text{m\u00f3l}} = \frac{d^d}{2^d} \frac{V_d(r)}{V_d(1)}$$

donde  $V_d(r)$  es el volumen de la esfera  $d$ -dimensional de radio  $r$ ,

$$V_d(r) = \pi^{d/2} r^d / \Gamma(\frac{d}{2} + 1).$$

[Pistas:  $P(\mathbf{c}(\mathbf{w}'') + \mathbf{c} = P(\mathbf{w} - \mathbf{w}'')^T H(\mathbf{w} - \mathbf{w}''))$ ,  
 Supongamos que la ortogonal  $L$  motriz  $A$  diagonaliza  $H$ :  
 $A^T H A = \text{diag}(\lambda_1, \dots, \lambda_d)$ . Cambiar variables  $\mathbf{u} = L^T (\mathbf{w} - \mathbf{w}'')$  y  $\mathbf{u} = L^T (\mathbf{w} - \mathbf{w}'')$ ]

- (b) Suponga que muestrea  $M$  veces y elige las ponderaciones con error mínimo,  $\mathbf{w}_M$ . Demuestre que

$$P[E(\mathbf{w}_M) > E(\mathbf{w}^*) + \epsilon] \approx \left(1 - \frac{1}{\pi d} \left(\mu \frac{\epsilon}{\sqrt{d}}\right)^d\right)^M,$$

donde  $\mu = \sqrt{\frac{1}{d} \sum \lambda_i^2}$  y  $A$  es la media geométrica de los valores propios de  $H$ .

(Puede utilizar  $\Gamma(d/2) = \frac{\sqrt{\pi}}{2} \Gamma(d/2 - 1)$ .)

- (c) Demuestre que  $\frac{1}{M} \sum_{i=1}^M \log \frac{1}{p_i} \geq \log \frac{1}{p}$ , entonces con probabilidad al menos  $1 - \epsilon$ ,

(Puede utilizar  $\log(1 - a) \geq -a/\epsilon$  para  $a$  pequeña y  $(rd)^{1/d} \geq 1$ .)

**Problema 7.10** Para una red neuronal con al menos 1 capa oculta y transformaciones  $\tanh$  en cada nodo que no sea de entrada, cuál es el gradiente (con respecto a los pesos) si todos los pesos se fijan en cero.

¿Es buena idea inicializar los pesos a cero?

**Problema 7.11 [Tasa óptima de aprendizaje]** Supongamos que estamos en las proximidades de un mínimo local,  $\mathbf{w}^*$ , de la superficie de error, o que la superficie de error es cuadrática. La expresión de la función de error viene dada por

$$E(\mathbf{w}_t) = E(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w}_t - \mathbf{w}^*)^T H (\mathbf{w}_t - \mathbf{w}^*) \quad (7.8)$$

de donde es fácil ver que el gradiente viene dado por  $\mathbf{g}_t = H(\mathbf{w}_t - \mathbf{w}^*)$ . Las actualizaciones de peso vienen dadas entonces por  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_t$ , y restando  $\mathbf{w}^*$  de ambos lados, vemos que

$$\mathbf{e}_{t+1} = (I - \eta H) \mathbf{e}_t \quad (7.9)$$

Como  $H$  es simétrico, se puede formar una base ortonormal con sus vectores propios. Proyectando  $\mathbf{e}_t$  y  $\mathbf{e}_{t+1}$  sobre esta base, vemos que en esta base, cada



y dejando que  $c(a)$  sea el componente  $a^{**}$  en esta base, vemos que

$$- \frac{1}{\lambda_{min}} - m, \frac{1}{\lambda_{max}} \quad (7.10)$$

por lo que vemos que cada componente exhibe convergencia lineal con su propia co-eficiencia de convergencia  $k_g \rightarrow 1 - q$ . El peor componente dominará la convergencia por lo que nos interesa elegir  $q$  de forma que se minimice el  $k_u$  de mayor magnitud. Dado que  $H$  es definida positiva, todos los  $\lambda_j$  son positivos, por lo que es fácil ver que uno debe elegir  $q$  para que  $1 - q \lambda_1$ ,  $1 - A$  y  $1 - A$ , o uno debe elegir. Resolviendo para la  $q$  óptima, se encuentra que

$$q_{opt} = \frac{2}{\lambda_{min} + \lambda_{max}} \quad k_{opt} = \frac{1}{1 - c} \quad (7.11)$$

donde  $c = n/J$  es el número de condición de  $H$ , y es una medida importante de la estabilidad de  $H$ . Cuando  $c \rightarrow 0$ , se suele decir que  $H$  está mal condicionado. Entre otras cosas, esto afecta a la capacidad de calcular numéricamente la inversa de  $H$ .

**Problema 7.12** [Con una tasa de aprendizaje variable, supongamos que  $\eta_k \rightarrow 0$  que satisface  $1/\eta_k \rightarrow \infty$  con  $\eta_k^{1/2} \rightarrow 0$ , por ejemplo se podría

elegir  $q_k = 1/(1 + k)$ . Demostrar que el descenso de gradiente convergerá a un mínimo local.

**Problema 7.13** [Aproximación de diferencia finita al hessiano]

- (a) Consideremos la función  $C(v, z)$ . Demostrar que la aproximación por diferencias finitas a las derivadas parciales de segundo orden vienen dadas por

$$\frac{\partial^2 E}{\partial w_1^2} = \frac{1}{4h^2} (C(v, z+2h) + C(v, z-2h) - 2C(v, z))$$

$$\frac{\partial^2 E}{\partial w_2^2} = \frac{1}{4h^2} (C(v, z+2h) + C(v, z-2h) - 2C(v, z))$$

$$\frac{\partial^2 E}{\partial w_2^2} = \frac{1}{4h^2} (E(w_1+h, w_2+h) + E(w_1-h, w_2-h) - E(w_1+h, w_2-h) - E(w_1-h, w_2+h))$$

- (b) Proporcione un algoritmo para calcular la aproximación por diferencias finitas a la matriz hessiana para  $\text{Min}(w)$ , el error en la muestra para una red neuronal multicapa con pesos  $w = [w^1, \dots, w^L]$ .
- (c) Calcule el tiempo de ejecución asintótica de su algoritmo en función del número de pesos de su red y del número de puntos de datos.

**Problema 7.14** Supongamos que damos un paso fijo en alguna dirección, nos preguntamos cuál es la dirección óptima (o este paso fijo suponiendo que el modelo cuadrático para la superficie de error es exacto:

$$E_{in}(w_t + \delta w) = Z_{in}(w_t) + g_t^T \Delta w + \frac{1}{2} \Delta w^T H_t \Delta w.$$

Así que queremos minimizar  $\mathcal{A}_m(\mathbf{A}_w)$  con respecto a  $\mathbf{A}_w$  sujeto a la restricción de que el tamaño del paso es  $q$ , es decir, que  $\mathbf{A}_w \mathbf{A}_w^T = q^2 \mathbf{I}$ .

- (a) Demostrar que el Lagrangiano para este problema de minimización restringida es:

$$\mathcal{L} = \mathcal{E}_{\text{in}}(\mathbf{w}_t) + \mathbf{g}_t^T \Delta \mathbf{w} + \frac{1}{2} \Delta \mathbf{w}^T (\mathbf{H} - 2\alpha \mathbf{I}) \Delta \mathbf{w}, \quad (7.12)$$

donde  $\alpha$  es el multiplicador de Lagrange.

- (b) Resuelve (o  $\mathbf{A}_w$  anda y demuestra que satisfacen las dos ecuaciones:

$$\mathbf{A}_w = (\mathbf{H} + 2\alpha \mathbf{I})^{-1} \mathbf{g}_t$$

$$\Delta \mathbf{w}^T \Delta \mathbf{w} = q^2.$$

- (c) Demuestre que  $\mathbf{A}_w$  satisface la ecuación implícita:

$$\frac{1}{2} (\mathbf{A}_w^T \mathbf{g}_t + \mathbf{A}_w^T \mathbf{H} \mathbf{A}_w) = 0.$$

Argumentar que el segundo término es  $O(1)$  y el primero es  $O(|\mathbf{g}_t|/q)$ . Por tanto,  $\mathbf{A}_w$  es grande para un tamaño de paso  $q$  pequeño.

- (d) Supongamos que  $q$  es grande. Demostrar que, a orden principal en  $1/q$ ,

$$\mathbf{A}_w \approx \frac{1}{q} \mathbf{g}_t$$

Por lo tanto  $\mathbf{A}_w$  es grande, coherente con la expansión de  $\mathbf{A}_w$  a orden principal en  $1/q$ . {Pista: expandir  $\mathbf{A}_w$  al orden principal en  $1/q$ .}

- (e) Utilizando (d), demuestre que  $\mathbf{A}_w \approx \frac{1}{q} \mathbf{g}_t$ .

### Problema 7.15

Sea  $\mathbf{H}$  una matriz simétrica positiva definida. Sea  $\mathbf{g}_k$  un vector en  $\mathbb{R}^n$ . Sea  $\mathbf{H}_k$  la aproximación hessiana del producto exterior es  $\mathbf{H}$  en  $\mathbf{g}_k$ . Sea  $\mathbf{H}_k$  la suma parcial a  $k$ , y sea  $\mathbf{H}_k$  su inversa.

La aproximación hessiana del producto exterior es  $\mathbf{H}$  en  $\mathbf{g}_k$ . Sea  $\mathbf{H}_k$  la suma parcial a  $k$ , y sea  $\mathbf{H}_k$  su inversa.

- (a) Demuestra que  $\mathbf{H}_k^{-1} \mathbf{g}_k = \mathbf{H}^{-1} \mathbf{g}_k$ .
- $$\mathbf{H}_k^{-1} \mathbf{g}_k = \mathbf{H}^{-1} \mathbf{g}_k - \frac{\mathbf{H}^{-1} \mathbf{g}_k \mathbf{g}_k^T \mathbf{H}^{-1} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{H}^{-1} \mathbf{g}_k} \mathbf{H}_k^{-1} \mathbf{g}_k$$
- {Pistas:  $\mathbf{H}_k^{-1} \mathbf{g}_k = \mathbf{H}^{-1} \mathbf{g}_k - \frac{\mathbf{H}^{-1} \mathbf{g}_k \mathbf{g}_k^T \mathbf{H}^{-1} \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{H}^{-1} \mathbf{g}_k} \mathbf{H}_k^{-1} \mathbf{g}_k$ }

- (b) Utilice la parte (a) para dar un algoritmo  $O(N^2)$  para calcular  $\mathbf{H}_k^{-1} \mathbf{g}_k$ , el mismo tiempo que se tarda en calcular  $\mathbf{H}$ . ( $N$  es el número de dimensiones en  $\mathbf{g}$ ).

Nota: normalmente, este algoritmo se inicializa con  $\mathbf{H}_0 = c \mathbf{I}$  para algún  $c$  pequeño. Así que el algoritmo realmente calcula  $(\mathbf{H} + c \mathbf{I})^{-1} \mathbf{g}$ ; los resultados no son muy sensibles a la elección de  $c$ , siempre que  $c$  sea pequeño.

**Problema 7.16** En el texto, hemos calculado un límite superior en la dimensión VC o (el perceptrón de 2 capas es  $dv \leq O(md \log md)$ ) donde  $la$  es el número de unidades ocultas en la capa oculta. Demostrar que este límite es esencialmente apretado mostrando que  $d \leq U(md)$ . Para ello, mostrar que es posible encontrar  $md$  puntos que pueden ser destrozados cuando  $la$  es incluso como follows.

Consideremos cualquier conjunto de  $N$  puntos  $x_1, \dots, x_N$  en posición general con  $N \geq md$ .  $N$  puntos en  $d$  dimensiones están en posición general si ningún subconjunto de  $d+1$  puntos yace en un hiperplano de  $d-1$  dimensiones. Ahora, considere cualquier dicotomía en estos puntos con  $r$  o los puntos clasificados 41. Sin pérdida de generalidad, vuelva a etiquetar los puntos de modo que  $x_1, \dots, x_r$  sean 41.

- Demuestre que, sin pérdida de generalidad, puede suponer que  $r \geq N/2$ . Por lo tanto, para el resto del problema puede suponer que  $r \geq N/2$ .
- Dividir los puntos positivos en grupos de tamaño  $d$ . El último grupo puede tener menos de  $d$  puntos. Demostrar que el número de grupos es como máximo  $\lceil N/d \rceil$ . Etiquetar estos grupos  $\theta_i$  para  $i = 1 \dots \lceil N/d \rceil$ .
- Demostrar que (o cualquier subconjunto de  $k$  puntos con  $k \leq d$ , existe un hiperplano que contiene esos puntos y no otros.
- Por la parte anterior, sea  $w_i^T x$  es el hiperplano que pasa por los puntos del grupo  $\theta_i$ , y no contiene ningún otro. Entonces

$$w_i^T x_n + b_i = 0$$

si y sólo si  $x \in \theta_i$ . Demuestre que es posible encontrar  $h$  suficientemente pequeño para que para  $x \in \theta_i$

$$|w_i^T x + b_i| < h,$$

y para  $x \notin \theta_i$

$$|w_i^T x + b_i| > h.$$

- Demuestre que para  $x \in \theta_i$ ,

$$\text{sign}(w_i^T x + b_i + h) + \text{sign}(-w_i^T x - b_i + h) = 2,$$

y para  $x \notin \theta_i$

$$\text{sign}(w_i^T x + b_i + h) - \text{sign}(-w_i^T x - b_i + h) = 0$$

- Utiliza los resultados obtenidos hasta ahora para construir un MLP de 2 capas con  $2r$  unidades ocultas que implemente la dicotomía (que era arbitraria). Completa el argumento para demostrar que  $d \leq md$ .