

**UNIVERSIDAD DE SANTIAGO DE CHILE
FACULTAD DE INGENIERIA
DEPARTAMENTO DE INGENIERIA INFORMÁTICA**



**ANÁLISIS DE LA EFICIENCIA DEL ENTRENAMIENTO DE REDES
NEURONALES PROFUNDAS BASADO EN SIMULATED ANNEALING**

Felipe Alberto Reyes González

Profesor Guía: Victor Parada

Tesis de grado presentada en conformidad a los
requisitos para obtener el grado de Magíster en
Ingeniería Informática

Santiago, Chile

2017

© Felipe Alberto Reyes González- 2017



• Algunos derechos reservados. Esta obra está bajo una Licencia Creative Commons Atribución-Chile 3.0. Sus condiciones de uso pueden ser revisadas en:

<http://creativecommons.org/licenses/by/3.0/cl/>.

TABLA DE CONTENIDO

1	Introducción	1
1.1	Antecedentes y motivación	1
1.2	Descripción del problema	3
1.3	Solución propuesta	4
1.3.1	Características de la solución	4
1.3.2	Propósito de la solución	4
1.4	Objetivos y alcances del proyecto	4
1.4.1	Objetivo general	4
1.4.2	Objetivos específicos	4
1.4.3	Alcances	5
1.5	Metodología y herramientas utilizadas	5
1.5.1	Metodología de trabajo	5
1.5.2	Herramientas de desarrollo	6
2	Aspectos teóricos y revisión de la literatura	7
2.1	Aspectos teóricos	7
2.1.1	El perceptrón multicapa y su arquitectura	7
2.1.2	Propagación de la entrada y el algoritmo de retropropagación	8
2.1.3	El desvanecimiento del gradiente	12
2.2	Revisión de la literatura	14
2.2.1	Retropropagación del error	14
2.2.2	Levenberg-Marquardt	15
2.2.3	Regularización bayesiana	16
2.2.4	Retropropagación resiliente	17
	REFERENCIAS BIBLIOGRÁFICAS	18

ÍNDICE DE TABLAS

Tabla 1.1	Especificaciones del equipo	6
Tabla 2.1	Algunas funciones de activaciones.	10

ÍNDICE DE ILUSTRACIONES

Figura 2.1	Perceptrón multicapa	8
Figura 2.2	Funciones de activación mas utilizadas.	10
Figura 2.3	Gradiente descendente	14

INTRODUCCIÓN

ANTECEDENTES Y MOTIVACIÓN

Las redes neuronales (*Neural Networks*, NN) son sistemas de procesamiento de información que basan su estructura en una analogía de las redes neuronales biológicas. Consisten en un conjunto de elementos de procesamiento simple llamados nodos, estos nodos están dispuestos en una estructura jerárquica y conectadas entre sí por un valor numérico llamado peso que, mediante un proceso de entrenamiento, varía su valor.

La actividad que una neurona realiza en una NN consiste en ponderar las entradas de la neurona por los pesos de las conexiones de la neurona para luego ser sumadas y entregadas a la función de activación asociada (McCulloch y Pitts, 1943). La salida corresponderá a la respuesta que la neurona genera a la entrada.

Las neuronas con una función de activación umbral, fueron estudiadas por Rosenblatt (1962) quién las denominó *perceptrón*. El modelo básico se compone de 3 capas conectadas consecutivamente. La primera capa corresponde a la capa de entrada, que recibe el patrón de entrada a clasificar. La segunda capa contiene las neuronas de asociación o detección de características. Y la tercera capa es la capa de salida, que contiene las neuronas que reconocen los patrones. La limitación de los perceptrones se debe a la capacidad de clasificación basado en el umbral, lo que no permite clasificar patrones más complejos.

El conjunto de n neuronas se llamará capa, y una NN puede estar compuesta de una o más capas. Cada capa estará compuesta por una cantidad de neuronas que no necesariamente será la misma para todas las capas, y estarán dispuestas en forma consecutiva de tal manera que las capas se conecten unas con otras y siempre hacia adelante. La primera capa, la de entrada, recibirá un patrón que será entregado a las distintas neuronas que la capa posea. Cada neurona de la capa de entrada procesará los datos y generará una salida que servirá de entrada para la capa siguiente, repitiendo el proceso para cada una de las capas de la NN hasta llegar a la capa de salida, en cuyo caso la salida representará la respuesta de la red, concretando así el ciclo.

Las NN han sido utilizadas para la clasificación de entradas, y han sido diseñados diversos métodos para entrenar la red y que los pesos se adapten, de tal manera que la salida de la red sea representativa de la salida esperada; a este método de entrenamiento se le llama *supervisado*. Dentro de los métodos clásicos de entrenamiento se encuentra el método del gradiente descendente, el método de Newton, el gradiente conjugado, el método quasi-Newton, o el algoritmo Levenberg-Marquardt. El más utilizado es el método del gradiente, que consiste en actualizar los

pesos de las distintas neuronas en función de la dirección contraria al gradiente de la función de activación, logrando minimizar el error.

La técnica de minimización más simple, cuando el gradiente está disponible, es seleccionar la dirección de descenso más empinada y aplicar una búsqueda unidireccional a lo largo de esta dirección. Esta técnica de optimización de descenso más pronunciada se define por

$$\Delta(n)W = -\lambda(n)G(n) \quad (1.1)$$

Donde $\Delta W(n)$ es la variación de los pesos para una iteración y $\lambda(n)$ es un coeficiente que minimiza la función en la dirección de descenso. El algoritmo de retropropagación propuesto por Werbos (1974) y popularizado por Rumelhart et al (1986) se basa en una variación de la técnica de pendiente más pronunciada. No se utiliza búsqueda unidireccional sino un paso de descenso fijo, η , llamado tasa de aprendizaje, que se añade a una fracción de la última variación, α , llamada momentum. El último término introduce algunos elementos del método de gradiente conjugado. Así la actualización se realiza como en la ecuación 1.2

$$W(n) = -\eta G(n) + \alpha \Delta W(n-1) \quad (1.2)$$

El método quasi-Newton utiliza el cálculo de las segundas derivadas de la función objetivo, se obtiene una mejor comprensión de la topología de la función, que conduce a su vez a elegir una dirección descendente más eficiente. Dejar:

$$\Delta W(n) = \lambda(n)S(n) \quad (1.3)$$

Donde la dirección de descenso $S(n)$ está definida por:

$$S(n) = -[H(n)]^{-1}G(n) \quad (1.4)$$

Y donde $H(n)$ es la matriz Hessiana. La principal dificultad con este enfoque es que encontrar la solución de este sistema en cada iteración es una tarea muy tediosa. Los métodos métricos variable, también llamados métodos cuasi-Newton, eluden esta dificultad aproximando directamente el inverso de la matriz de Hessiana, $[H(n)]^{-1}$, de la primera derivada, $G(n)$. Estos métodos son las técnicas de optimización sin restricciones más populares y, entre ellas, BGFS es el método más utilizado.

El método de Newton es un método de segundo orden, pues hace uso de la matriz Hessiana. El

objetivo de este método es encontrar una mejor dirección de entrenamiento mediante el uso de la segunda derivada de la función de pérdida. Considera una aproximación cuadrática de la función en la solución inicial utilizando una expansión de la serie de Taylor. De esta manera, el método de Newton permite mover la solución inicial en función de la inversa de la matriz Hessiana, que pasará a ser la dirección de entrenamiento de Newton.

El método del gradiente del conjugado fue propuesto por Fletcher y Reeves (1964), utiliza sucesivas direcciones conjugadas basadas en el gradiente y el residuo. Si la función objetivo es cuadrática y la dirección de búsqueda es minimizada exactamente en cada iteración, el método converge de forma cuadrática. Leonard y Kramer (1990) han utilizado el método para entrenar redes neuronales como alternativa a la retropropagación. Los pesos de una red se actualizan de acuerdo con una búsqueda unidireccional en la dirección de descenso $S(n)$ de la siguiente forma:

$$\Delta W(n) = \lambda(n)S(n) \quad (1.5)$$

Y la dirección de descenso $S(n)$ se calcula a partir del gradiente de iteraciones pasadas y presentes como muestra la ecuación 1.6

$$S(n) = G(n) + \frac{\|G(n)\|}{\|G(n-1)\|} S((n-1)) \quad (1.6)$$

Para que las direcciones de descenso permanezcan conjugadas, la búsqueda unidireccional, en cada iteración, tiene que ser llevada con alta precisión. Sin embargo, después de muchas iteraciones, las direcciones podrían llegar a ser casi paralelas. Para superar esta dificultad, Fletcher y Reeves (1964) sugieren volver a iniciar el procedimiento igualando la dirección de descenso al gradiente en cada iteración $(t+1)$, donde (t) es el número total de pesos en la red.

DESCRIPCIÓN DEL PROBLEMA

La retropropagación basa su funcionamiento en multiplicaciones sucesivas basadas en el error para poder calcular los gradientes, y a medida que el error se propaga hacia la capa de entrada de la red el gradiente comienza a disminuir su valor por cada capa que atraviesa. Esto significa que el gradiente disminuirá de manera exponencial, lo que representa un problema para redes profundas, ya que las capas mas cercanas a la capa de entrada necesitarán más tiempo para ser entrenadas.

SOLUCIÓN PROPUESTA

Características de la solución

Mediante el uso de el algoritmo *simulated annealing* se busca analizar la eficiencia que la NN alcanza en una red neuronal profunda frente a otros métodos de aprendizaje tales como: SGD y RMSPROP.

Propósito de la solución

El propósito de la solución es aportar en el campo de las redes neuronales y la clasificación de datos, proporcionando un análisis comparativo de la convergencia de distintas redes.

OBJETIVOS Y ALCANCES DEL PROYECTO

En ésta sección se presenta el objetivo general, los objetivos específicos además del alcance y limitaciones de la presenta investigación.

Objetivo general

Evaluar el desempeño del algoritmo *simulated annealing* y su efecto sobre el entrenamiento de redes neuronales profundas en comparación con otros métodos.

Objetivos específicos

Los objetivos establecidos para el presente trabajo son descritos a continuación

1. Definir las reglas de aprendizaje a comparar.

2. Construir los conjuntos de datos de entrada y salida a analizar.
3. Establecer los parámetros de las redes neuronales para la experimentación.
4. Establecer los algoritmos de aprendizaje a comparar.
5. Entrenar las redes con los distintos conjuntos de datos.
6. Establecer las conclusiones del trabajo.

Alcances

1. Se analizará la misma arquitectura con diferentes reglas de aprendizaje.
2. Los conjunto de datos para el entrenamiento a utilizar son los propuestos en (Morse y Stanley, 2016).

METODOLOGÍA Y HERRAMIENTAS UTILIZADAS

Metodología de trabajo

Considerando el aspecto investigativo del trabajo, se considera la utilización del método científico. Entre las actividades que componen la metodología, Sampieri (2006) describe los siguientes pasos para desarrollar una investigación:

- **Formulación de la hipótesis:** Las redes neuronales que adolecen del desvanecimiento del gradiente se ven beneficiadas por el uso del algoritmo *simulated annealing* en la convergencia.
- **Marco teórico:** Una revisión de la literatura donde se aborda el problema planteado, para situarse en el contexto actual de los problemas. Se describirán redes neuronales que buscan solucionar el mismo problema.
- **Diseño de la solución:** Se deberá diseñar el experimento para generar los datos que permitan sustentar las comparaciones entre las distintas redes.

- Análisis y verificación de los resultados: Los resultados se analizarán considerando los valores de convergencia de los distintos métodos.
- Presentación de los resultados: Se presentarán tablas que describan los resultados obtenidos y que se consideren pertinentes.
- Conclusiones obtenidas en el desarrollo de la investigación.

Herramientas de desarrollo

Para el desarrollo y ejecución de los experimentos se utilizará un equipo con las siguientes características

Sistema Operativo	Solus 2017.04.18.0 64-bit
Procesador	Intel® Core™i5-2450M CPU @ 2.50GHz x 4
RAM	7.7Gb
Gráficos	Intel® Sandybridge Mobile
Almacenamiento	935.6 GB

Tabla 1.1: Especificaciones del equipo

El software que se utilizará es:

- Lenguaje de programación: Python.
- Sistema de redes neuronales: Keras API (Chollet, 2015).
- Herramienta ofimática: \LaTeX .

ASPECTOS TEÓRICOS Y REVISIÓN DE LA LITERATURA

En esta sección se abarcan los aspectos relacionados al conocimiento general para la comprensión del presente trabajo y la revisión de la literatura asociada al trabajo presentado. Para realizar un análisis de la eficiencia en la convergencia de los algoritmos de aprendizaje es necesario conocer la base teórica de ésta. Para ellos, se describen los conceptos necesarios de las redes neuronales y el desvanecimiento del gradiente. La sección 2.1 se centra en explicar las partes fundamentales para el entendimiento del problema en cuestión.

ASPECTOS TEÓRICOS

El perceptrón multicapa y su arquitectura

Dentro de las redes neuronales, el perceptrón multicapa es una de las arquitecturas más usadas para resolver problemas. Esto es debido a que poseen la capacidad de ser un aproximador universal (Marvin Minsky, 1987). Esto no implica que sea una de las redes más potentes o con mejores resultados, el perceptrón multicapa posee una serie de limitaciones, como el proceso de aprendizaje para problemas que dependan de un gran número de variables, la dificultad para realizar un análisis teórico de la red debido a la presencia de componentes no lineales y a la alta conectividad.

El perceptron multicapa posee una estructura de capas compuestas por neuronas. Cada una de las capas está formada por un conjunto de neuronas y se distinguen tres tipos de capas: la capa de entrada, las capas ocultas y la capa de salida.

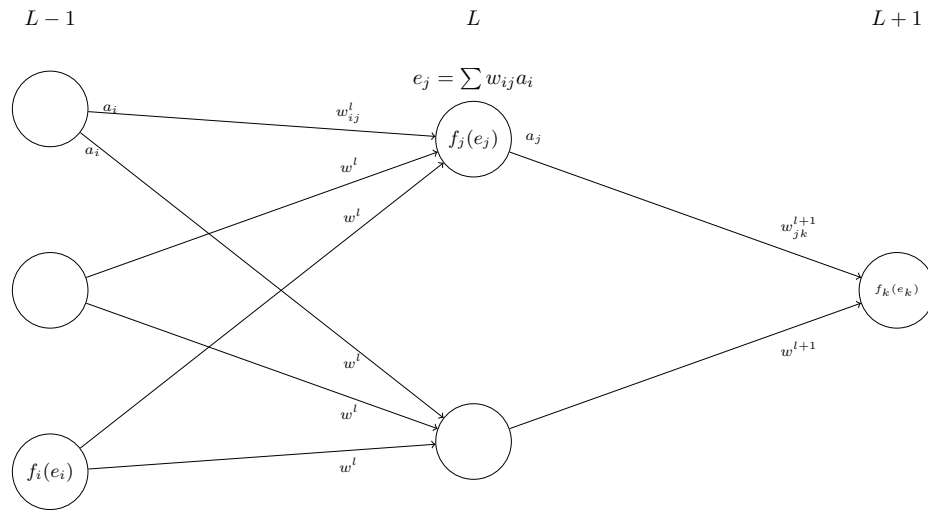


Figura 2.1: Perceptrón multicapa

En la figura 2.1 se observa que las conexiones van siempre hacia adelante. Las neuronas de la capa l se conectan con las neuronas de la capa $l + 1$. Las neuronas de la capa de entrada se encargan de recibir los patrones y propagar dichas señales a las neuronas de la capa siguiente. La última capa, la capa de salida, proporciona la respuesta de la red al patrón presentado. Las neuronas de las capas ocultas realizan el procesamiento de las señales generadas por el patrón de entrada.

Propagación de la entrada y el algoritmo de retropropagación

El perceptrón multicapa define una relación entre la entrada y la salida. Esta relación se obtiene propagando hacia adelante los valores de las variables de entrada, es por esto que también se les llama redes *feedforward*. Cada neurona de la red procesa la entrada recibida y produce una respuesta que se propaga, mediante las conexiones, hacia las neuronas de la capa siguiente.

Existen dos fases importante dentro del modelo

- Fase de entrenamiento: Se usa un conjunto de datos o patrones de entrenamiento para determinar los pesos que definen el modelo de la NN. Se calculan de manera iterativa, de acuerdo con los valores de entrenamiento, con el objeto de minimizar el error cometido entre la salida obtenida por la NN y la salida deseada.

Los pesos óptimos se obtienen minimizando una función. Uno de los criterios utilizados es la minimización del error cuadrático medio entre el valor de salida y el valor real esperado.

- **Fase de prueba:** Durante el entrenamiento, el modelo se ajusta al conjunto de entrenamiento, perdiendo la habilidad de generalizar su aprendizaje a casos nuevos, a esta situación se le llama sobreajuste.

Para evitar el sobreajuste, se utiliza un segundo grupo de datos diferentes, el conjunto de validación, que permitirá controlar el proceso de aprendizaje.

Si un perceptrón multicapa con C capas y n_c neuronas en la capa c , donde $W_c = (w_{ij}^c)$ es la matriz de pesos, w_{ij}^c representará el peso de la conexión de la neurona i de la capa c hasta la neurona j de la capa siguiente. Denotaremos a_i^c a la activación de la neurona i de la capa c que se calcula de la siguiente manera:

- **Activación de una neurona de la capa de entrada:** Las neuronas se encargan de transmitir la entrada recibida, por lo tanto

$$a_i^1 = x_i, i = 1, 2, \dots, n$$

donde $X = (x_1, x_2, \dots, x_n)$ representa el vector de entrada.

- **Activación de una neurona de la capa oculta:** Las neuronas de una capa oculta procesa la información recibida aplicando la función de activación f a la suma de los productos de la entrada por sus pesos, es decir

$$a_i^c = f \left(\sum_{j=1}^{n_{c-1}} w_{ji}^{c-1} a_j^{c-1} + \theta_i^c \right), i = 1, 2, \dots, n_c; c = 2, 3, \dots, C - 1$$

donde a_j^{c-1} es la salida de la capa anterior a c .

- **Activación de una neurona de la capa de salida:** La activación de una neurona de la capa de salida viene dada por la función de activación f aplicada a la suma de los productos de la entrada por sus pesos, es decir

$$y_i = a_i^C = f \left(\sum_{j=1}^{n_{C-1}} w_{ji}^{C-1} a_j^{C-1} + \theta_i^C \right), i = 1, \dots, n_c$$

donde $Y = (y_1, y_2, \dots, y_{n_c})$ es el vector de salida.

Función	Fórmula	Rango
Identidad	$f(x) = x$	$[-\infty, \infty]$
Lineal por tramos	$f(x) = \begin{cases} -1 & x < -1 \\ a * x & -1 \leq x \leq 1 \\ 1 & x > 1 \end{cases}$	$[-1, 1]$
Sinusoidal	$f(x) = \sin(\omega x + \varphi)$	$[-1, 1]$
Sigmoidal	$f(x) = \frac{1}{1+\exp(-x)}$	$[0, 1]$
Tangente hiperbólica	$\frac{1-\exp(-x)}{1+\exp(-x)}$	$[-1, 1]$

Tabla 2.1: Algunas funciones de activaciones.

La función f es la función de activación de la neurona. Aunque existe gran variedad de funciones de activación (ver tabla 2.1), las funciones de activación más utilizadas son la sigmoidal y la tangente hiperbólica, descritas en las ecuaciones 2.1 y 2.2 respectivamente.

$$f_{sigm}(x) = \frac{1}{1 + \exp(-x)} \quad (2.1)$$

$$f_{tanh}(x) = \frac{1 - \exp(-x)}{1 + \exp(-x)} \quad (2.2)$$

Ambas funciones poseen como imagen un intervalo de valores entre $[0, 1]$ y $[-1, 1]$ como se observa en la figura 2.2.

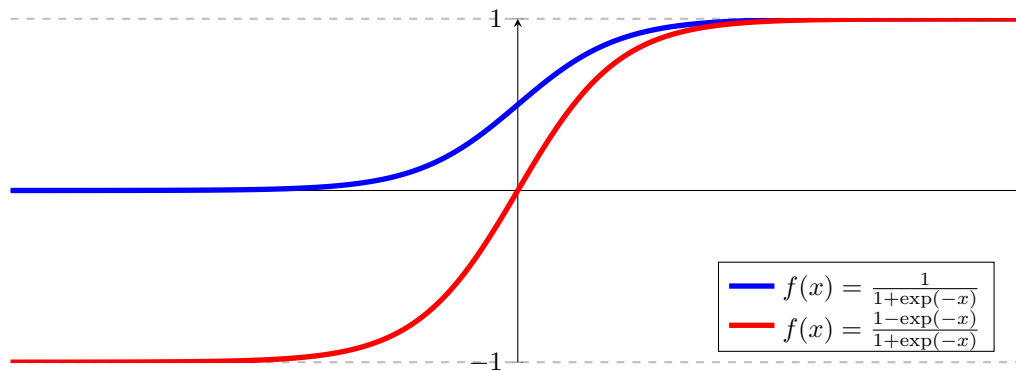


Figura 2.2: Funciones de activación mas utilizadas.

El perceptrón multicapa actualiza sus pesos en función de una regla de aprendizaje, de tal manera que los nuevos pesos permitan reducir el error de salida. Por tanto, para cada patrón de entrada a la red es necesario disponer de un patrón de salida deseada. El objetivo es que la salida de la red sea lo más próxima posible a la salida deseada, debido a esto es que el aprendizaje de la red se describe como un problema de minimización de la siguiente manera

$$\min_W E$$

donde W es el conjunto de parámetros de la red (pesos y umbrales) y E es una función de error que evalúa la diferencia entre las salidas de la red y las salidas deseadas. En la mayor parte de los casos, la función de error se define como:

$$E = \frac{1}{N} \sum_{i=1}^N e(i) \quad (2.3)$$

Donde N es el número de muestras y $e(n)$ es el error cometido por la red para el patrón i , definido de la siguiente manera

$$e(i) = \frac{1}{n_C} \sum_{j=1}^{n_C} (s_j(i) - y^j(n))^2 \quad (2.4)$$

Siendo $Y(i) = (y_1(i), y_2(i), \dots, y_{n_C}(i))$ y $S(i) = (s_1(i), s_2(i), \dots, s_{n_C}(i))$ los vectores de salida y salidas deseadas para el patrón i respectivamente.

De esta manera, si W^* es un mínimo de la función de error E , en dicho punto el error será cercano a cero, y en consecuencia, la salida de la red será próxima a la salida deseada. La presencia de funciones de activación no lineales hace que la respuesta de la red sea no lineal respecto a los parámetros ajustables, por lo que el problema de minimización es un problema no lineal y se hace necesario el uso de técnicas de optimización no lineales para su resolución.

Las técnicas utilizadas suelen basarse en la actualización de los parámetros de la red mediante la determinación de una dirección de búsqueda. En el caso de las redes neuronales multicapa, la dirección de búsqueda más utilizada se basa en la dirección contraria del gradiente de la función de error E , el método de gradiente descendente.

Si bien el aprendizaje de la red busca minimizar el error total de la red, el procedimiento está basado en métodos del gradiente estocástico, que son una sucesión de minimizaciones del error $e(i)$ por cada patrón, en lugar de minimizar el error total E de la red. Aplicando el método del gradiente estocástico, cada parámetro w se modifica para cada patrón de entrada n según la siguiente regla de aprendizaje

$$w(i) = w(n-1) - \alpha \frac{\partial e(i)}{\partial w} \quad (2.5)$$

donde $e(i)$ es el error para el patrón de entrada i dado por la ecuación 2.4, y α es la tasa de aprendizaje, éste último determina el desplazamiento en la superficie del error.

Como las neuronas están ordenadas por capas y en distintos niveles, es posible aplicar el método del gradiente de forma eficiente, resultando en el *algoritmo de retropropagación* (Rumelhart, Hinton, y Williams, 1986) o *regla delta generalizada*. El término retropropagación es utilizado debido

a la forma de implementar el método del gradiente en las redes multicapa, pues el error cometido en la salida de la red es propagado hacia atrás, transformándolo en un error para cada una de las neuronas ocultas de la red.

El algoritmo de retropropagación es el método de entrenamiento más utilizado en redes con conexión hacia adelante. Es un método de aprendizaje supervisado, en el que se distinguen claramente dos fases:

1. Se aplica un patrón de entrada, el cual se propaga por las distintas capas que componen la red hasta producir la salida de la misma. Esta salida se compara con la salida deseada y se calcula el error cometido por cada neurona de salida.
2. Estos errores se transmiten desde la capa de salida, hacia todas las neuronas de las capas anteriores (Fritsch, 1996). Cada neurona recibe un error que es proporcional a su contribución sobre el error total de la red. Basándose en el error recibido, se ajustan los errores de los pesos sinápticos de cada neurona.

El desvanecimiento del gradiente

El problema del gradiente desvaneciente nace en las NN profundas, éstas utilizan funciones cuyo gradiente se encuentran entre 0 y 1. Debido a que estos gradientes pequeños se multiplican durante la retropropagación, tienden a desvanecerse a través de las capas, evitando que la red aprenda.

Si se tiene una NN, la activación de una neurona de una capa intermedia i con función de activación f_i y con entrada

$$net_i(t) = \sum_j w_{ji} y^j(t-1)$$

es

$$y^i(t) = f_i(net_i(t))$$

Además w_{ji} es el peso de la conexión desde la unidad j de la capa anterior hasta la unidad i de la capa actual, $d_k(t)$ será la respuesta esperada de la unidad k de la capa de salida en el tiempo t . Usando el error cuadrático medio (*Mean square error*, MSE), el error de k será

$$E_k(t) = (d_k(t) - y^k(t))^2$$

En un tiempo $\tau \leq t$ cualquiera, el error de una neurona j que no sea una neurona de entrada es

la suma de los errores externos y el error propagado hacia atrás desde la neurona previa será

$$\vartheta_j(\tau) = f'_j(\text{net}_j(\tau)) \left(E_j(\tau) + \sum_i w_{ij} \vartheta_i(\tau + 1) \right)$$

El peso actualizado en el tiempo τ resulta $w_{jl}^{new} = w_{jl}^{old} + \alpha \vartheta_j(\tau) y^l(\tau - 1)$ donde α es la tasa de aprendizaje, y l es una unidad arbitraria conectada a la unidad j .

La propagación hacia atrás de un error que ocurre en una unidad u en un tiempo t hacia una unidad v para q pasos, escala el error de la siguiente manera

$$\frac{\partial \vartheta_v(t - q)}{\partial \vartheta_u(t)} = \begin{cases} f'_v(\text{net}_v(t - 1)) w_{uv} & q = 1 \\ f'_v(\text{net}_v(t - q)) \sum_{l=1}^n \frac{\partial \vartheta_l(t - q + 1)}{\partial \vartheta_u(t)} w_{lv} & q > 1 \end{cases} \quad (2.6)$$

Con $l_q = v$ y $l_0 = u$, el factor de escalamiento es

$$\frac{\partial \vartheta_v(t - q)}{\partial \vartheta_u(t)} = \sum_{l_1=1}^n \cdots \sum_{l_{q-1}=1}^n \prod_{m=1}^q f'_{l_m}(\text{net}_{l_m}(t - m)) w_{l_m l_{m-1}} \quad (2.7)$$

La sumatoria de los n^{q-1} términos $\prod_{m=1}^q f'_{l_m}(\text{net}_{l_m}(t - m)) w_{l_m l_{m-1}}$ escalan el error. Los distintos términos pueden tener signos diferentes, por lo tanto, el aumento del número de unidades n no implica un incremento del error absoluto. Pero con mas unidades se incrementa la expectativa de que el valor absoluto del error aumente. Si $\rho(m, l_m, l_{m-1}) := |f'_{l_m}(\text{net}_{l_m}(t - m)) w_{l_m l_{m-1}}| < 1.0$ para todo m , el producto en (2.7) decrece exponencialmente con q , es decir, el error se desvanece como muestra la figura 2.3. Un error que se desvanece a lo largo del flujo casi no tiene efecto en la actualización de los pesos.

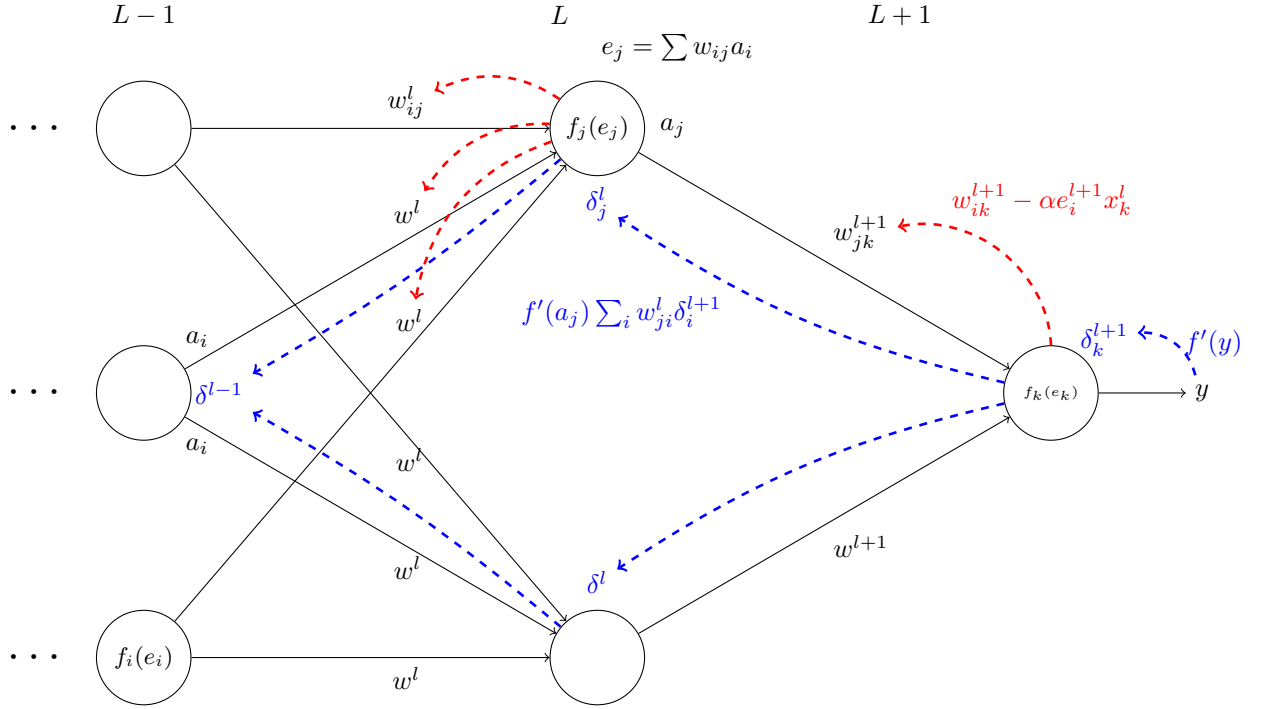


Figura 2.3: Gradiente descendente

REVISIÓN DE LA LITERATURA

El aprendizaje de las NN resulta ser difícil debido a la propagación de la señal a través de unidades no lineales. El primer algoritmo fue la retropropagación del error (Werbos, 1974; Rumelhart y cols., 1986), existiendo distintas alternativas.

Retropropagación del error

El método de retropropagación del error busca minimizar una función de error, como puede ser el error cuadrático medio (*Mean square error*, MSE) entre las respuestas calculadas por la red y las respuestas deseadas, para luego modificar los pesos de la NN en cada época.

El algoritmo consta de dos fases, la primera consiste en propagar la señal desde la capa de entrada hacia la capa de salida a través de las capas ocultas. La segunda fase se realiza la corrección de los pesos de la NN en sentido inverso, es decir, desde la capa de salida hacia la

capa de entrada.

Levenberg-Marquardt

El método de Levenberg-Marquardt (LM) (Levenberg, 1944; Marquardt, 1963) es propuesto originalmente en (Marquardt, 1963), pero es Hagan y Menhaj (1994) quién lo utiliza por primera vez en una NN. El método se basa en el método de Newton, que es un método de segundo orden. El paso básico del método de Newton durante una época n es

$$W(n+1) = W(n) - H(n)^{-1}g(n) \quad (2.8)$$

donde W es el vector de pesos, $g(n)$ es el vector gradiente actual y $H(n)$ es la matriz Hessiana de la función de error respecto de los pesos. El método de Newton converge mas rápido que los métodos de gradiente conjugado.

El algoritmo LM fue diseñado para realizar entrenamientos rápidos de segundo orden sin necesidad de calcular la matriz Hessiana. En el caso que la función de error tiene la forma de una suma de cuadrados, la matriz Hessiana es posible aproximarla como en la ecuación 2.9

$$H = J^T J \quad (2.9)$$

donde J es la matriz Jacobiana que contiene las primeras derivadas de la función de error de la red respecto de los pesos de la NN, mientras que el gradiente se obtendrá como se muestra en la ecuación 2.10

$$g = J^T e \quad (2.10)$$

siendo e el vector de errores de la red. El cálculo de la matriz Jacobiana se reduce a un cálculo del método de retropropagación (Hagan y Menhaj, 1994) mas sencillo que el calculo de la matriz Hessiana.

El método actualizará los pesos de forma similar a como lo hace el método de Newton:

$$W(n+1) = W(n) - [J^T J + \alpha I]^{-1} J^T e \quad (2.11)$$

donde α es la tasa de aprendizaje. Cuando α es igual a cero, el método se convierte en el método

de Newton utilizando una aproximación del calculo de la matriz Hessiana

$$W(n+1) = W(n) - H(n)^{-1}g(n) \quad (2.12)$$

y cuando α es grande, el aporte de la matriz Hessiana es despreciable, obteniendo como resultado el método del gradiente descendente con una tasa de aprendizaje $\frac{1}{\alpha}$

$$W(n+1) = W(n) - \frac{1}{\alpha}g(n) \quad (2.13)$$

Regularización bayesiana

El método de regularización bayesiana (BR) consiste en modificar la función de costo. La función de costo mas utilizada es el MSE, que es la suma de cuadrados de los errores individuales e_i (ver 2.14)

$$MSE_d = \frac{1}{N} \sum_{i=1}^N e_i^2 \quad (2.14)$$

El método de BR consiste en mejorar la capacidad de generalización del modelo. A la función de coste de la ecuación 2.14 se le agrega el término MSE_w , que incluye el efecto de la suma de cuadrados de los pesos de la NN

$$MSE_w = \frac{1}{N} \sum_{i=1}^N w_i^2 \quad (2.15)$$

Y la función de coste se obtendría de la siguiente manera

$$MSE = \beta MSE_d + \mu MSE_w \quad (2.16)$$

donde β y μ son parámetros que deben ser ajustados según la metodología Bayesiana de MacKay (MacKay, 1992a, 1992b). La metodología asume que los pesos son variables aleatorias que son descritas por distribuciones específicas, comunmente Gaussianas. Los parámetros de regularización están relacionados con las varianzas desconocidas asociadas a estas distribuciones, y es posible estimarlas mediante técnicas estadísticas. Una de las particularidades del algoritmo es que proporciona una medida de cuántos parámetros de la NN están siendo utilizados.

Retropropagación resiliente

En las capas ocultas de las NN se utilizan con frecuencia funciones de activación sigmoideas, y estas reducen la entrada a un rango finito de salida. Se caracterizan por sus pendientes próximas a cero para entradas muy grandes, y esto representa un problema cuando se utiliza el gradiente descendente para el entrenamiento, pues se acentúa el problema del desvanecimiento del gradiente descrito en la sección 2.1.3.

Riedmiller y Braun (1993) presenta el algoritmo de retropropagación resiliente (*Resilient retropropagation*, RPROP), con el que busca eliminar los efectos del desvanecimiento del gradiente, esto mediante el uso del signo de la derivada para determinar la dirección del ajuste de los pesos, haciendo que la magnitud de la derivada no tenga efecto sobre la actualización de los pesos de la NN. Los pesos de la NN se incrementarán por un factor Δ_i cuando la derivada de la función respecto a dicho peso tenga el mismo signo que las dos iteraciones anteriores. Mientras que si la derivada de la función respecto a dicho peso cambia de signo respecto de la iteración anterior los pesos se decrementarán por un factor Δ_d . En caso de que la derivada de la función respecto de dicho peso sea igual a cero, el valor de actualización no varía.

Tieleman y Hinton (2012) presentan el algoritmo de la raíz cuadrática media de retropropagación resiliente (*Root mean square resilient retropropagation*, RMSPROP), una variante del algoritmo RPROP. Se propone mantener una media móvil del cuadrado del gradiente para cada peso

$$MeanSquare(w, t) = 0.9MeanSquare(w, t - 1) + 0.1 \frac{\partial E}{\partial w^{(t)}}^2 \quad (2.17)$$

Y dividiendo el gradiente por $\sqrt{MeanSquare(w, t)}$ hace que el aprendizaje funcione de mejor manera.

REFERENCIAS BIBLIOGRÁFICAS

- Chollet, F. (2015). *Keras*. <https://github.com/fchollet/keras>. GitHub.
- Fletcher, R., y Reeves, C. M. (1964, 1 de febrero). Function minimization by conjugate gradients. *The Computer Journal*, 7(2), 149–154. Descargado de <http://dx.doi.org/10.1093/comjnl/7.2.149> doi: 10.1093/comjnl/7.2.149
- Fritsch, J. (1996). *Modular neural networks for speech recognition* (Masters Thesis). KIT.
- Hagan, M. T., y Menhaj, M. B. (1994, Nov). Training feedforward networks with the marquardt algorithm. *IEEE Transactions on Neural Networks*, 5(6), 989-993. doi: 10.1109/72.329697
- Leonard, J., y Kramer, M. (1990). Improvement of the backpropagation algorithm for training neural networks. *Computers & Chemical Engineering*, 14(3), 337 - 341. Descargado de <http://www.sciencedirect.com/science/article/pii/0098135490870706> doi: [http://dx.doi.org/10.1016/0098-1354\(90\)87070-6](http://dx.doi.org/10.1016/0098-1354(90)87070-6)
- Levenberg, K. (1944). A method for the solution of certain non-linear problems in least squares. *Quarterly Journal of Applied Mathematics*, 11(2), 164–168.
- MacKay, D. J. C. (1992a, mayo). Bayesian interpolation. *Neural Comput.*, 4(3), 415–447. Descargado de <http://dx.doi.org/10.1162/neco.1992.4.3.415> doi: 10.1162/neco.1992.4.3.415
- MacKay, D. J. C. (1992b, May). A practical bayesian framework for backpropagation networks. *Neural Computation*, 4(3), 448-472. doi: 10.1162/neco.1992.4.3.448
- Marquardt, D. W. (1963). An algorithm for least-squares estimation of nonlinear parameters. *Journal of the Society for Industrial and Applied Mathematics*, 11(2), 431-441. Descargado de <http://dx.doi.org/10.1137/0111030> doi: 10.1137/0111030
- Marvin Minsky, S. A. P. (1987). *Perceptrons: An introduction to computational geometry* (Expanded ed.). The MIT Press.
- McCulloch, W. S., y Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4), 115–133. doi: 10.1007/BF02478259
- Morse, G., y Stanley, K. O. (2016). Simple evolutionary optimization can rival stochastic gradient descent in neural networks. En *Proceedings of the genetic and evolutionary computation conference 2016* (pp. 477–484). New York, NY, USA: ACM. Descargado de <http://doi.acm.org/10.1145/2908812.2908916> doi: 10.1145/2908812.2908916

Riedmiller, M., y Braun, H. (1993). A direct adaptive method for faster backpropagation learning: the rprop algorithm. En *IEEE international conference on neural networks* (p. 586-591 vol.1). doi: 10.1109/ICNN.1993.298623

Rosenblatt, F. (1962). *Principles of neurodynamics: perceptrons and the theory of brain mechanisms*. Spartan Books.

Rumelhart, D. E., Hinton, G. E., y Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323, 533–536.

Sampieri, R. (2006). *Metodología de la investigación*. México: McGraw Hill.

Tieleman, T., y Hinton, G. (2012). *Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude*. COURSERA: Neural Networks for Machine Learning.

Werbos, P. J. (1974). *Beyond regression: New tools for prediction and analysis in the behavioral sciences* (Tesis Doctoral no publicada). Harvard University.