

Práctica 1: Implementación del perceptrón multicapa

Asignatura: Introducción a los modelos computacionales, 4º Grado
Ingeniería Informática (Escuela Politécnica Superior de Córdoba -
Universidad de Córdoba)

Trabajo realizado por:

-Antonio Gómez Giménez (32730338G)
i72gogia@uco.es



Índice:

1. Descripción de los modelos de redes neuronales utilizados:	2
2. Descripción en pseudocódigo de los pasos del algoritmo de retropropagación:	3
3. Experimentos y análisis de resultados:	6
4. Resumen:	21



1. Descripción de los modelos de redes neuronales utilizados:

Principalmente, los modelos que usamos en estas redes neuronales se basan en número de capas y el número de neuronas que nos pide el usuario, es decir, aparte de la capa de entrada y la capa de salida, ambas siempre están presentes, el usuario indica cuantas capas oculta desea, como mínimo una capa oculta (sin incluir la capa de salida).

Aparte del número de capas, el usuario proporciona el número de neuronas por capas, es decir, si dice por ejemplo cinco, entonces en todas las capas del modelo (exceptuando capa de entrada y salida, ya que se especifican las neuronas de esta capa en la propia base de datos proporcionada por el usuario) hay cinco neuronas en total, no puede haber capas con cinco neuronas y otras capas con un número distinto de cinco neuronas.

En resumen, la arquitectura de todo el modelo depende de los datos introducidos por el usuario (número de capas y número de neuronas por cada capa, que en este caso siempre es el mismo).



2. Descripción en pseudocódigo de los pasos del algoritmo de retropropagación:

Voy a obviar la parte de extracción de datos y la ejecución del main y voy a centrarme en la parte donde se realiza el run del algoritmo:

```
runOnlineBackPropagation(DataSet  entrenamiento,  DataSet  test,  número
iteraciones, error iteraciones, error test){
```

```
    if(comprobamos_validacion()){
        dividimos_dataset_entrenamiento();
    }

    hacer{
        trainOnline(entrenamiento);
        test(test);
    }mientras_que(no se cumpla las condiciones de parada);

    if(comprobamos_validacion()){
        hacer{
            test(validacion);
        }mientras_que(no se cumpla las condiciones de parada);
    }

}
```

*Cabe destacar en esta parte del algoritmo que `dividimos_dataset_entrenamiento()` se encuentra ya implementado en esa parte del código (no es una función en sí, está puesta para generalizar lo que habría que realizar en ese punto). Las condiciones de parada son las tres especificadas en el pdf.

Se van a explicar la función **trainOnline** y **test** ya que son bastante interesantes:

trainOnline es una función que contiene un bucle que recorre todos los patrones del dataset y llama a la función **performEpochOnline** al cual se le pasan las entradas, las salidas y el número total de patrones.



test es igual que la función anterior pero en vez de llamar a **performEpochOnline** llama directamente a **feedInputs** y **forwardPropagate** para poder así sacar el error medio de todos los patrones, es decir, el error entre el valor que debería salir y el que nuestra red nos ha dado.

performEpochOnline es una función importante donde se aloja el algoritmo de retropropagación error, por tanto el pseudocódigo sería el siguiente:

```
performEpochOnline(valores entradas, valores salidas, numero_patrones){  
  
    limpiarAcumulador();  
    feedInputs(entradas);  
    forwardPropagate();  
    backPropagateError();  
    accumulateChange();  
    weightAdjustment();  
  
}
```

limpiarAcumulador es una función donde todos los ΔW (de todas las capas menos la primera y para todas las neuronas de esas capas) los ponemos a cero. Esto es útil para la versión online, ya que para cada iteración se resetea a cero aplicando esos pesos en **weightAdjustment**. Cuando esto sea offline se aplica fuera del bucle permitiendo así que se no se aplique cada cambio y que se aplique el cambio acumulándose. En resumen, el acumulador se resetea para cada patrón.

feedInputs es una función donde se alimentan las entradas, es decir, para la primera capa, para todas sus neuronas, se rellenan todos los *out* con los valores de entradas.

forwardPropagate esta función realiza la propagación hacia delante. Para todas las capas menos la primera y para todas las neuronas de esas capas, acumulamos en su *out* el sesgo (almacenado en $w[0]$, vector de pesos), y añadimos para todas los pesos anteriores por sus respectivas salidas (bucle for desde 1 ya que no cogemos el sesgo). Una vez tenemos todo el net (el *out* acumulado, el sumatorio en sí), le aplicamos la sigmoide ($1/(1+\exp(-net))$). Este valor se guarda en *out*.

backPropagateError esta función realiza la propagación hacia atrás. Se divide en dos partes.

La primera parte consiste un un bucle que recorre todas las neuronas de la última capa donde almacenamos en Δ de esa neurona la siguiente fórmula:



$$-2(\text{salida-out}(\text{salida de nuestra neurona de salida}) * \text{out} * (1 - \text{out}));$$

De esta forma calculamos el *delta* de la neuronas de la capa de salida ya que la fórmula no es exactamente igual en las capas ocultas.

En la segunda parte, para cada capa desde la penúltima hasta la segunda capa, para todas las neuronas, más uno por el sesgo, y para todas las neuronas de la capa siguiente, acumulamos el peso de de todas las neurona de la capa siguiente, excepto el peso por los *deltas* de la capa siguiente.

Una vez acumulado todo eso, lo multiplicamos por la salida de la capa actual por 1-la salida actual. Por último, limpiamos el acumulador.

En si lo que hemos hecho en esa acumulación es calcular los pesos de las neuronas por las salidas posteriores a la actual.

accumulateChange esta función acumula los cambios realizados por un patrón, por ello es necesario antes de realizar los cambios guardar los *deltasW* actuales en *lastDeltaW* para almacenarlos.

Para cada capa excepto la primera, para cada neurona y para todas las neuronas de la capa anterior, guardamos en *deltaW* lo que ya había en *deltaW* más el *delta* de la neurona actual por la salida de la neurona de la capa anterior.

El *deltaW* que calcula el sesgo se calcula de una forma un poco diferente ya que no se multiplica por la salida de la neurona de la capa anterior si no que es por 1 ya que es el sesgo.

weightAdjustment esta función ajusta los pesos de toda la red para cada capa excepto la primera y para cada neurona de esas capas. Donde el nuevo peso *W* pasa a ser el peso *W* anterior más el decremento de la tasa de aprendizaje por el *deltaW* menos el factor de momento por el decremento de la tasa de aprendizaje por el *lastDeltaW*.

Es decir, calculamos el momento del peso de esa neurona. le aplicamos el decremento de la tasa de aprendizaje por capas y se lo restamos a el peso actual, cabe destacar que el peso también hay que ajustarlo para el sesgo (*W[0]*).



3. Experimentos y análisis de resultados:

En el **primer experimento** se ha utilizado la base de datos del problema XOR. Esta base de datos representa el problema de clasificación no lineal del XOR. Se utiliza el mismo fichero para entrenamiento y para test. Consta de dos entradas, una salida y cuatro patrones.

En esta primera parte, utilizaremos el factor de momento, no utilizaremos el conjunto de validación, no emplearemos decremento de la tasa de aprendizaje ($F=1$) y el número de iteraciones es de 1000.

Momento

1 capa oculta

-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.1142 62	0.0181 597	0.0073 1603	0.0058 0858	0.00392 239	0.0522 612	0.10172 9
	0.0073 279	2.1715 8e-05	1.0635 e-06	1.0094 3e-06	1.24574 e-07	0.0098 346	0.01474 76
Test(media y desviación típica)	0.1142 62	0.0181 597	0.0073 1603	0.0058 0858	0.00392 239	0.0522 612	0.10172 9
	0.0073 279	2.1715 8e-05	1.0635 e-06	1.0094 3e-06	1.24574 e-07	0.0098 346	0.01474 76

1 capa oculta

-momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.1142 62	0.0181 597	0.0073 1603	0.0058 0858	0.00392 239	0.0522 612	0.10172 9
	0.0073 279	2.1715 8e-05	1.0635 e-06	1.0094 3e-06	1.24574 e-07	0.0098 346	0.01474 76
Test(media y desviación típica)	0.1142 62	0.0181 597	0.0073 1603	0.0058 0858	0.00392 239	0.0522 612	0.10172 9
	0.0073 279	2.1715 8e-05	1.0635 e-06	1.0094 3e-06	1.24574 e-07	0.0098 346	0.01474 76



1 capa oculta
-momento 1,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.1142 62	0.0181 597	0.0073 1603	0.0058 0858	0.00392 239	0.0522 612	0.10172 9
	0.0073 279	2.1715 8e-05	1.0635 e-06	1.0094 3e-06	1.24574 e-07	0.0098 346	0.01474 76
Test(media y desviación típica)	0.1142 62	0.0181 597	0.0073 1603	0.0058 0858	0.00392 239	0.0522 612	0.10172 9
	0.0073 279	2.1715 8e-05	1.0635 e-06	1.0094 3e-06	1.24574 e-07	0.0098 346	0.01474 76

2 capas ocultas
-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.2367 95	0.1920 81	0.0095 0605	0.0055 9973	0.00218 362	0.0011 3643	0.00072 0157
	0.0003 04678	0.0062 5068	2.3834 e-06	1.2208 6e-06	3.55609 e-08	6.2241 8e-10	1.37982 e-09
Test(media y desviación típica)	0.2367 95	0.1920 81	0.0095 0605	0.0055 9973	0.00218 362	0.0011 3643	0.00072 0157
	0.0003 04678	0.0062 5068	2.3834 e-06	1.2208 6e-06	3.55609 e-08	6.2241 8e-10	1.37982 e-09

2 capas ocultas
-momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.2367 95	0.1920 81	0.0095 0605	0.0055 9973	0.00218 362	0.0011 3643	0.00072 0157
	0.0003 04678	0.0062 5068	2.3834 e-06	1.2208 6e-06	3.55609 e-08	6.2241 8e-10	1.37982 e-09



Test(media y variandesviación típica)	0.236795	0.192081	0.00950605	0.00559973	0.00218362	0.00113643	0.000720157
	0.000304678	0.00625068	2.3834e-06	1.22086e-06	3.55609e-08	6.22418e-10	1.37982e-09

2 capas ocultas

-momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.236795	0.192081	0.00950605	0.00559973	0.00218362	0.00113643	0.000720157
	0.000304678	0.00625068	2.3834e-06	1.22086e-06	3.55609e-08	6.22418e-10	1.37982e-09
Test(media y desviación típica)	0.236795	0.192081	0.00950605	0.00559973	0.00218362	0.00113643	0.000720157
	0.000304678	0.00625068	2.3834e-06	1.22086e-06	3.55609e-08	6.22418e-10	1.37982e-09

Tras ver estos resultados, podemos obviar la variable momento ya que sea cual sea el momento introducido no afecta a la arquitectura (probablemente no haya mínimos locales). La mejor arquitectura es 100 nodos para 2 capas ocultas, ofreciendo un error de 0.000720157 ya que da igual el sobreentrenamiento porque los datos de test y entrenamientos son los mismos. No se penaliza entonces el sobreentrenamiento. Comprobamos la mejor arquitectura con $F=2$ ya que no hay validación en este caso, no tendría sentido.

2 capas ocultas con 100 neuronas para $F=1$ y $F=2$

Nodos	$F=1$	$F=2$
Entrenamiento (media y desviación típica)	0.000720157	0.00379216
	1.37982e-09	6.83007e-08
Test(media y desviación típica)	0.000720157	0.00379216
	1.37982e-09	6.83007e-08



En conclusión, con todos los datos obtenidos anteriormente podemos ver que tanto el momento en esta base de datos no tienen importancia ya que al ser un problema tan sencillo es probable que no haya mínimos locales donde en esos casos es donde más útil es el momento. El decremento de la tasa de aprendizaje empeora este modelo en concreto ya que con este número de iteraciones probablemente no le da tiempo a ajustar los pesos iniciales.

El test y el entrenamiento son iguales ya que los datos son los mismos, por tanto, lo que nos interesa es sobreentrenar de tal forma que a más iteraciones mejor ya que se reduce el error ajustándose a la solución (por ejemplo con 5000 iteraciones para el entrenamiento como el test tienen un error medio y varianza de $0.00014554 \pm 4.34502e-11$).

En el **segundo experimento** se ha utilizado la base de datos de la función seno donde esta base de datos está compuesta por 120 patrones de entrenamiento y 41 patrones de test. Ha sido obtenida añadiendo cierto ruido aleatorio a la función seno. Esta base de datos consta de una entrada y una salida.

En esta base de datos ocurre igual que con el momento ya que para $m=0,5$, $0,9$ y $1,5$, los resultados son los mismos, aun así se van a calcular por si se obtuviera algún cambio.

Calculamos todas las medias y todas las varianzas para ver la mejor estructura posible, para, sobre ella, realizar las pruebas de validación y de decremento de la tasa de aprendizaje.

1 capa oculta

-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0296 618	0.0293 966	0.0293 94	0.0288 663	0.02794 22	0.0285 211	0.02925 13
	2.2463 7e-07	5.7004 8e-07	7.5066 1e-07	5.6495 e-07	2.04824 e-08	1.6314 8e-06	7.31577 e-07
Test(media y desviación típica)	0.0361 994	0.0358 996	0.0355 651	0.0349 588	0.03429 04	0.0365 18	0.03880 19
	3.3325 1e-08	2.2320 4e-07	4.8431 4e-07	4.5818 9e-07	2.76407 e-08	3.3373 9e-06	4.00061 e-06



1 capa oculta
-momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0296 618	0.0293 966	0.0293 94	0.0288 663	0.02794 22	0.0285 211	0.02925 13
	2.2463 7e-07	5.7004 8e-07	7.5066 1e-07	5.6495 e-07	2.04824 e-08	1.6314 8e-06	7.31577 e-07
Test(media y desviación típica)	0.0361 994	0.0358 996	0.0355 651	0.0349 588	0.03429 04	0.0365 18	0.03880 19
	3.3325 1e-08	2.2320 4e-07	4.8431 4e-07	4.5818 9e-07	2.76407 e-08	3.3373 9e-06	4.00061 e-06

1 capa oculta
-momento 1,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0296 618	0.0293 966	0.0293 94	0.0288 663	0.02794 22	0.0285 211	0.02925 13
	2.2463 7e-07	5.7004 8e-07	7.5066 1e-07	5.6495 e-07	2.04824 e-08	1.6314 8e-06	7.31577 e-07
Test(media y desviación típica)	0.0361 994	0.0358 996	0.0355 651	0.0349 588	0.03429 04	0.0365 18	0.03880 19
	3.3325 1e-08	2.2320 4e-07	4.8431 4e-07	4.5818 9e-07	2.76407 e-08	3.3373 9e-06	4.00061 e-06

2 capas ocultas
-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0298 902	0.0298 567	0.0296 298	0.0300 067	0.02842 19	0.0199 247	0.01386 09
	1.0632 9e-09	1.0298 3e-08	6.7480 7e-08	3.4230 3e-08	1.80613 e-06	3.4479 9e-05	1.10485 e-07



Test(media y desviación típica)	0.036006	0.0361708	0.0366069	0.0364279	0.0354835	0.0263005	0.019979
	9.33555e-09	1.469e-08	6.19859e-08	2.17063e-09	2.74781e-06	4.22101e-05	3.61838e-07

2 capas ocultas

momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0298902	0.0298567	0.0302206	0.0300067	0.0284219	0.0199247	0.0138609
	1.06329e-09	1.02983e-08	5.36034e-08	3.42303e-08	1.80613e-06	3.44799e-05	1.10485e-07
Test(media y desviación típica)	0.036006	0.0361708	0.0361824	0.0364279	0.0354835	0.0263005	0.019979
	9.33555e-09	1.469e-08	4.94313e-09	2.17063e-09	2.74781e-06	4.22101e-05	3.61838e-07

2 capas ocultas

momento 1,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0298902	0.0298567	0.0302206	0.0300067	0.0284219	0.0199247	0.0138609
	1.06329e-09	1.02983e-08	5.36034e-08	3.42303e-08	1.80613e-06	3.44799e-05	1.10485e-07
Test(media y desviación típica)	0.036006	0.0361708	0.0361824	0.0364279	0.0354835	0.0263005	0.019979
	9.33555e-09	1.469e-08	4.94313e-09	2.17063e-09	2.74781e-06	4.22101e-05	3.61838e-07

Tras analizar estas tablas podemos observar que la mejor configuración es 2 capas con 100 neuronas por capa oculta y que el factor de momento es irrelevante en este caso, ocurre lo mismo que con la base de datos de la XOR. Usaremos la mejor arquitectura de la red (dos capas ocultas con cien neuronas por capa) para comprobar si se realiza una mejora modificando los parámetros de validación y de F.



Cabe destacar que en este caso sí tiene sentido hacer validación ya que el disponemos de varios patrones y el fichero de entrenamiento y test es diferente ya que en el caso anterior al ser igual nos interesaba sobreentrenar ya que buscamos únicamente las soluciones de la XOR y son pocos patrones.

Realizamos las pruebas con distintas validaciones y distintos valores de F.

Validación 0 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.0138609	0.0273926
	1.10485e-07	1.97879e-07
Test(media y desviación típica)	0.019979	0.0331969
	3.61838e-07	0,0006671

Validación 0.15 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.0276061	0.0399456
	6.52888e-05	6.75407e-05
Test(media y desviación típica)	0.031705	0.0387452
	2.49477e-05	2.27884e-05

Validación 0.25 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.0287871	0.0386645
	8.14369e-05	3.66387e-05
Test(media y desviación típica)	0.0494352	0.0583314
	0.000854701	0.00124021



Tras realizar las pruebas podemos observar que la mejor estructura es aquella en la que no se realiza validación y F es 1. Esto se debe a que no le da tiempo a aprender a las neuronas de la capa 1 lo suficiente, por ello con $F=1$ si le permite un mejor entrenamiento.

Respecto a la validación, al eliminar patrones para su uso a la validación, estamos eliminando patrones de entrenamiento. Esto hace que si no hay muchos patrones de entrenamiento, nuestra red aprenda peor. Esto se aprecia, ya que a medida que aumentamos el ratio de patrones para validación, aumenta el error resultante que nos da nuestra red.

En el **tercer experimento** se ha utilizado la base de datos quake donde esta base de datos está compuesta por 1633 patrones de entrenamiento y 546 patrones de test. Se corresponde con una base de datos en la que el objetivo es averiguar la fuerza de un terremoto (medida en escala sismológica de Richter). Como variables de entrada, utilizamos la profundidad focal, la latitud en la que se produce y la longitud, por tanto, consta de tres entradas y una salida.

En esta base de datos ocurre igual que con el momento ya que para $m=0,5, 0,9$ y $1,5$, los resultados son los mismos, aun así se van a calcular por si se obtuviera algún cambio.

Calculamos todas las medias y todas las varianzas para ver la mejor estructura posible, para, sobre ella, realizar las pruebas de validación y de decremento de la tasa de aprendizaje.

1 capa oculta

-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0300 477	0.0299 628	0.0297 4	0.0296 278	0.02957 24	0.0295 996	0.02955 11
	5.0806 e-09	1.0048 8e-08	8.8982 4e-10	2.8053 2e-09	4.17144 e-09	3.0772 9e-09	2.17945 e-09
Test(media y desviación típica)	0.0271 501	0.0270 997	0.0270 088	0.0269 324	0.02703 57	0.0269 945	0.02708 75
	6.9578 5e-09	5.3492 9e-09	3.8485 2e-09	9.9295 3e-09	1.1647e -09	2.5834 e-09	1.48474 e-09



1 capa oculta

-momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0300 477	0.0299 628	0.0297 4	0.0296 278	0.02957 24	0.0295 996	0.02955 11
	5.0806 e-09	1.0048 8e-08	8.8982 4e-10	2.8053 2e-09	4.17144 e-09	3.0772 9e-09	2.17945 e-09
Test(media y desviación típica)	0.0271 501	0.0270 997	0.0270 088	0.0269 324	0.02703 57	0.0269 945	0.02708 75
	6.9578 5e-09	5.3492 9e-09	3.8485 2e-09	9.9295 3e-09	1.1647e -09	2.5834 e-09	1.48474 e-09

1 capa oculta

-momento 1,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0300 477	0.0299 628	0.0297 4	0.0296 278	0.02957 24	0.0295 996	0.02955 11
	5.0806 e-09	1.0048 8e-08	8.8982 4e-10	2.8053 2e-09	4.17144 e-09	3.0772 9e-09	2.17945 e-09
Test(media y desviación típica)	0.0271 501	0.0270 997	0.0270 088	0.0269 324	0.02703 57	0.0269 945	0.02708 75
	6.9578 5e-09	5.3492 9e-09	3.8485 2e-09	9.9295 3e-09	1.1647e -09	2.5834 e-09	1.48474 e-09

2 capas ocultas

-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0299 749	0.0298 693	0.0297 66	0.0296 445	0.02954 78	0.0293 941	0.02917 82
	5.4327 1e-09	9.3783 1e-09	6.6343 6e-09	1.0411 4e-08	8.50589 e-09	2.6057 5e-08	1.82489 e-08



Test(media y desviación típica)	0.0270 621	0.0269 948	0.0269 043	0.0269 292	0.02695 46	0.0270 462	0.02704 38
	5.0165 4e-09	2.5131 4e-09	3.2552 e-09	1.3731 6e-09	0,00000 6074	1.9591 3e-08	3.16951 e-09

2 capas ocultas

momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0299 749	0.0298 693	0.0297 66	0.0296 445	0.02954 78	0.0293 941	0.02917 82
	5.4327 1e-09	9.3783 1e-09	6.6343 6e-09	1.0411 4e-08	8.50589 e-09	2.6057 5e-08	1.82489 e-08
Test(media y desviación típica)	0.0270 621	0.0269 948	0.0269 043	0.0269 292	0.02695 46	0.0270 462	0.02704 38
	5.0165 4e-09	2.5131 4e-09	3.2552 e-09	1.3731 6e-09	0,00000 6074	1.9591 3e-08	3.16951 e-09

2 capas ocultas

momento 1,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0299 749	0.0298 693	0.0297 66	0.0296 445	0.02954 78	0.0293 941	0.02917 82
	5.4327 1e-09	9.3783 1e-09	6.6343 6e-09	1.0411 4e-08	8.50589 e-09	2.6057 5e-08	1.82489 e-08
Test(media y desviación típica)	0.0270 621	0.0269 948	0.0269 043	0.0269 292	0.02695 46	0.0270 462	0.02704 38
	5.0165 4e-09	2.5131 4e-09	3.2552 e-09	1.3731 6e-09	0,00000 6074	1.9591 3e-08	3.16951 e-09

En este caso, sorprendentemente, tras analizar las tablas, la mejor estructura no es de dos capas con cien neuronas sino que la mejor estructura es de dos capas con nueve neuronas, esto ocurre ya que el entrenamiento siempre mejora su error pero si nos fijamos, el error del test empieza a incrementar, esto ocurre por que se está empezando a sobreentrenar, es decir, nuestra red se está empezando a ajustar



demasiado a nuestro modelo y por consiguiente, nuestro modelo empezará a generalizar peor.

Por tanto, la elección correcta como mejor estructura es aquella donde no se empieza a sobreentrenar y esta es con dos capas ocultas con nueve neuronas por capa.

Usaremos esta arquitectura de la red para comprobar si se realiza una mejora modificando los parámetros de validación y de F.

Validación 0 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.029766	0.0301729
	6.63436e-09	3.47915e-11
Test(media y desviación típica)	0.0269043	0.0272865
	3.2552e-09	3.71298e-11

Validación 0.15 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.0226564	0.0229613
	3.11424e-05	3.0239e-05
Test(media y desviación típica)	0.0282778	0.0285661
	3.59646e-07	5.34941e-07

Validación 0.25 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.0178831	0.0183702
	6.33076e-05	6.27883e-05
Test(media y desviación típica)	0.0287387	0.0289998
	5.22033e-07	7.48867e-07



Tras realizar las pruebas podemos observar que la mejor estructura es aquella en la que no se realiza validación y F es 1. De hecho, el error de entrenamiento es bastante peor que en validación 0.25 con $F=1$ pero el error de test es mucho menor en validación cero con $F=1$, esto ocurre por que con la validación se generaliza mucho peor, lo que me sorprende es que en estos casos con menos patrones esta realizando sobreentrenamiento, supongo que no puede generalizar mejor por la poca variación de los patrones recibidos al limitar el número de patrones por culpa de la validación. Respecto a el parámetro F ocurre lo mismo que en el experimento anterior, donde las neuronas de la primera capa no aprenden lo suficiente.

En el **cuarto experimento** se ha utilizado la base de datos parkinsons donde esta base de datos está compuesta por 4406 patrones de entrenamiento y 1469 patrones de test. Contiene, como entradas o variables independientes, una serie de datos clínicos de pacientes con la enfermedad de Parkinson y datos de medidas biométricas de la voz, y, como salidas o variables dependientes, el valor motor y total del UPDRS. Por tanto tenemos dos entradas y dos salidas.

En esta base de datos ocurre igual que con el momento ya que para $m=0,5$, $0,9$ y $1,5$, los resultados son los mismos, aun así se van a calcular por si se obtuviera algún cambio.

Calculamos todas las medias y todas las varianzas para ver la mejor estructura posible, para, sobre ella, realizar las pruebas de validación y de decremento de la tasa de aprendizaje.

1 capa oculta

-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0347 872	0.0289 687	0.0187 598	0.0169 087	0.01500 69	0.0132 07	0.01273 31
	2.1924 5e-07	4.0179 3e-06	4.5077 7e-07	2.9039 e-06	6.92032 e-07	1.1736 6e-06	1.70371 e-06
Test(media y desviación típica)	0.0366 313	0.0293 154	0.0197 846	0.0182 164	0.01670 45	0.0158 95	0.01478 93
	2.9707 9e-07	4.2305 9e-06	8.9171 6e-07	2.8522 3e-06	5.87734 e-07	1.7521 4e-06	2.18482 e-06



1 capa oculta
-momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0347 872	0.0289 687	0.0187 598	0.0169 087	0.01500 69	0.0132 07	0.01273 31
	2.1924 5e-07	4.0179 3e-06	4.5077 7e-07	2.9039 e-06	6.92032 e-07	1.1736 6e-06	1.70371 e-06
Test(media y desviación típica)	0.0366 313	0.0293 154	0.0197 846	0.0182 164	0.01670 45	0.0158 95	0.01478 93
	2.9707 9e-07	4.2305 9e-06	8.9171 6e-07	2.8522 3e-06	5.87734 e-07	1.7521 4e-06	2.18482 e-06

1 capa oculta
-momento 1,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0347 872	0.0289 687	0.0187 598	0.0169 087	0.01500 69	0.0132 07	0.01273 31
	2.1924 5e-07	4.0179 3e-06	4.5077 7e-07	2.9039 e-06	6.92032 e-07	1.1736 6e-06	1.70371 e-06
Test(media y desviación típica)	0.0366 313	0.0293 154	0.0197 846	0.0182 164	0.01670 45	0.0158 95	0.01478 93
	2.9707 9e-07	4.2305 9e-06	8.9171 6e-07	2.8522 3e-06	5.87734 e-07	1.7521 4e-06	2.18482 e-06

2 capas ocultas
-momento 0,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0322 799	0.0228 121	0.0134 504	0.0091 3982	0.00824 027	0.0038 0575	0.00360 742
	6.0737 2e-07	2.0718 7e-06	3.6605 9e-06	1.7714 3e-06	8.46237 e-06	1.2751 9e-06	7.2039e -07



Test(media y desviación típica)	0.0335 933	0.0234 906	0.0147 677	0.0110 509	0.01102 35	0.0062 1885	0.00615 301
	1.0736 e-06	3.2121 e-06	5.5669 6e-06	2.8541 8e-06	8.4375e -06	1.5291 4e-06	2.72816 e-06

2 capas ocultas

momento 0,9

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0322 799	0.0228 121	0.0134 504	0.0091 3982	0.00824 027	0.0038 0575	0.00360 742
	6.0737 2e-07	2.0718 7e-06	3.6605 9e-06	1.7714 3e-06	8.46237 e-06	1.2751 9e-06	7.2039e -07
Test(media y desviación típica)	0.0335 933	0.0234 906	0.0147 677	0.0110 509	0.01102 35	0.0062 1885	0.00615 301
	1.0736 e-06	3.2121 e-06	5.5669 6e-06	2.8541 8e-06	8.4375e -06	1.5291 4e-06	2.72816 e-06

2 capas ocultas

momento 1,5

Nodos	2	4	9	16	32	64	100
Entrenamiento (media y desviación típica)	0.0322 799	0.0228 121	0.0134 504	0.0091 3982	0.00824 027	0.0038 0575	0.00360 742
	6.0737 2e-07	2.0718 7e-06	3.6605 9e-06	1.7714 3e-06	8.46237 e-06	1.2751 9e-06	7.2039e -07
Test(media y desviación típica)	0.0335 933	0.0234 906	0.0147 677	0.0110 509	0.01102 35	0.0062 1885	0.00615 301
	1.0736 e-06	3.2121 e-06	5.5669 6e-06	2.8541 8e-06	8.4375e -06	1.5291 4e-06	2.72816 e-06

Tras ver estos resultados, este caso es similar al de la base de datos seno donde la mejor estructura es dos capas ocultas con cien neuronas por capa ocultas.

Como era de esperarse el momento es irrelevante ya que las tablas para cada momento son idénticas (aunque se me hace un poco raro que para esta base de datos no le afecte).



Cabe destacar que para esta base de datos al aumentar su número de iteraciones es muy probable que siga mejorando su error pero el tiempo de cómputo es muy elevado, solo en una ejecución de dos capas con cien neuronas por capa puede llegar a tardar cerca de cuarenta minutos. Posteriormente comprobaremos si la validación soluciona este problema.

Usaremos esta arquitectura de la red para comprobar si se realiza una mejora modificando los parámetros de validación y de F.

Validación 0 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.00360742	0.0110653
	7.2039e-07	5.31489e-06
Test(media y desviación típica)	0.00615301	0.0131777
	2.72816e-06	6.84293e-06

Validación 0.15 probando F1 y F2

Nodos	F=1	F=2
Entrenamiento(media y desviación típica)	0.0053372	0.00902755
	3.24185e-06	5.11616e-06
Test(media y desviación típica)	0.0105457	0.015275
	1.15317e-05	4.94121e-06

Validación 0.25 probando F1 y F2

No he podido sacar estos resultados por problemas de ejecución probablemente por problemas de espacio al crear el conjunto de validación. El programa ejecuta tres seed pero a partir de la cuarta ya no funciona.



4. Resumen:

Tras analizar todo lo visto en las tablas podemos sacar en conclusión que :

-El **momento** en las redes que se entrenan de forma online es inutil. Esto se debe a que el momento multiplica el estado anterior de variación de los pesos y este es cero ya que de forma online los cambios se aplican en cada iteración. Para ser más exactos, el momento multiplica ΔW en nuestro programa y en la versión online cada vez que realizamos una iteración esta ΔW la reseteamos a cero. Por consiguiente, el momento solo tiene sentido cuando se tenga en cuenta el momento anterior y esto sucederá en el algoritmo offline.

-Para **F** nos hemos dado cuenta de que sucede y lo hemos explicado.

A mi parecer **F** tiene más sentido en redes con más número de iteraciones de esta forma es mejor modificar los pesos de forma más importante aquellos que están en las entradas, ya que con 1000 iteraciones no le da tiempo prácticamente a entrenar correctamente los pesos más cercanos a las entradas en nuestro caso cuando hay dos capas nos referimos a la primera capa.

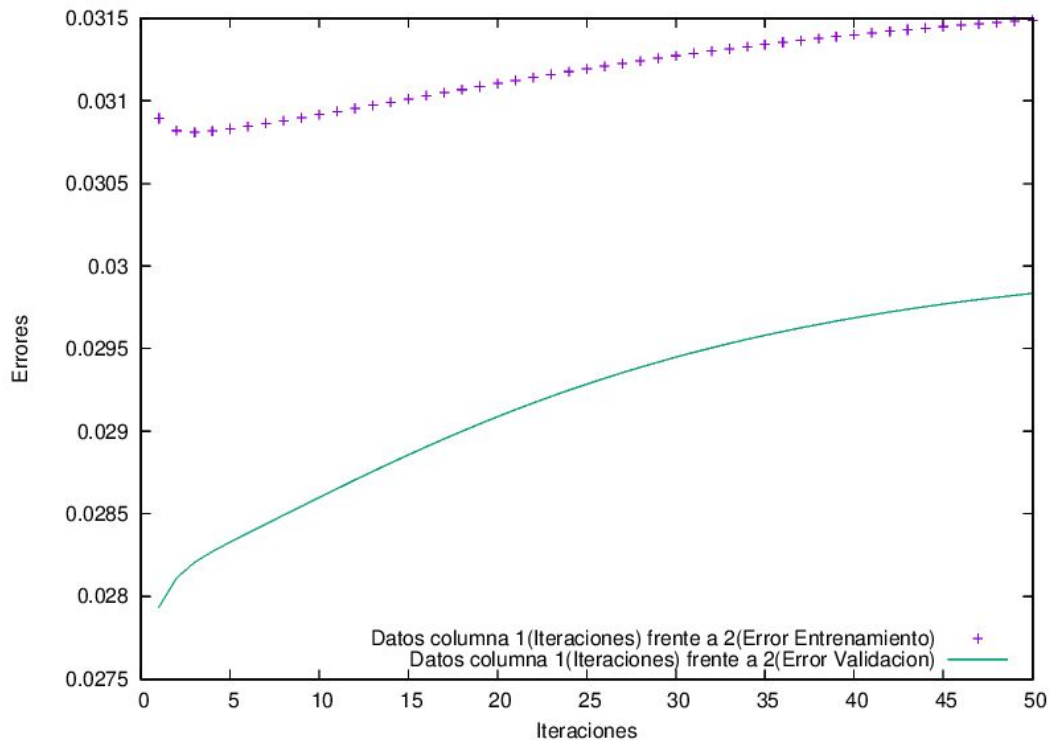
-Respecto a la **validación**, todos los resultado donde se realiza la validación son peores pero la utilidad de la validación se encuentra en evitar el sobreentrenamiento(en nuestros casos no se aplica por el pequeño número de iteraciones y por la poca cantidad de patrones, ya que si hay pocos patrones y nos excedemos con los patrones de validación, en el entrenamiento no será capaz de generalizar por falta de patrones) y reducir el número de iteraciones.

Esto último se nota mucho sobre todo en aquellas bases de datos donde hay muchos patrones ya que por cada iteración se tarda bastante y al reducir el número de iteraciones mejoramos enormemente la velocidad de cómputo respecto a aumentar el error obtenido.

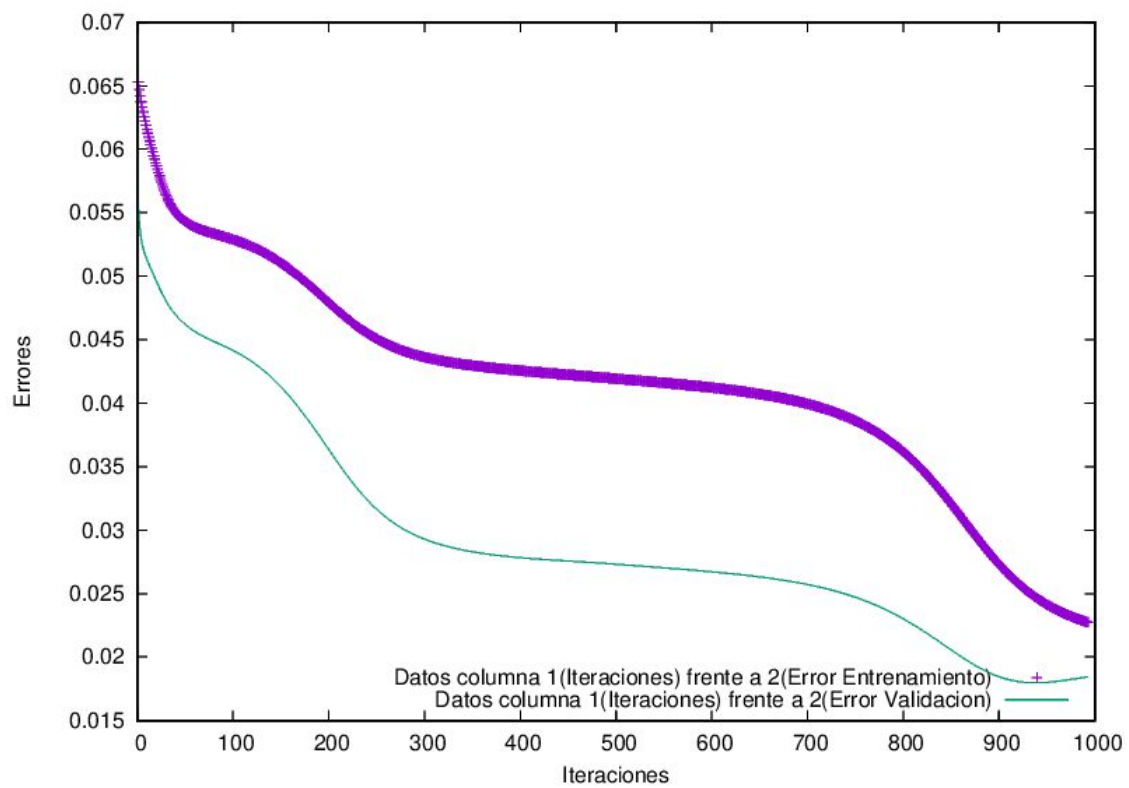
Para apreciar cómo funciona validación vamos a observar las siguientes gráficas, estas gráficas representan la base de datos seno con dos capas y cien nodos por capas donde se ha escogido la mejor solución (menor error) y hemos observado que valores da de error de entrenamiento y validación con respecto a la iteración para cada una de las semillas.



Seed1:

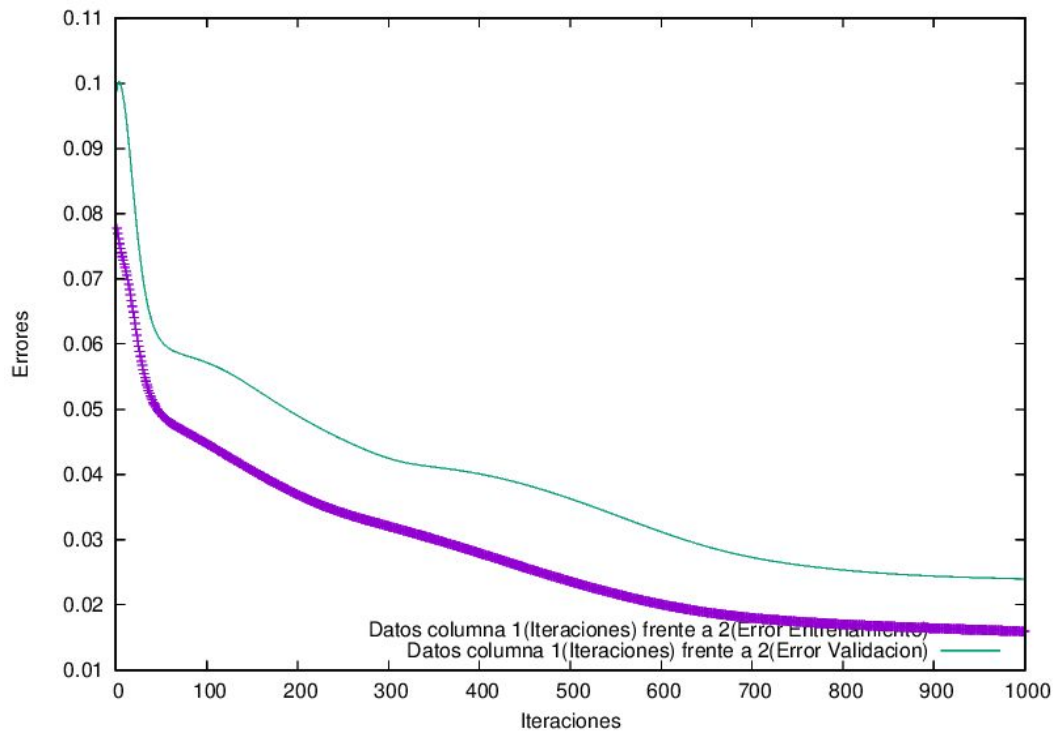


Seed2:

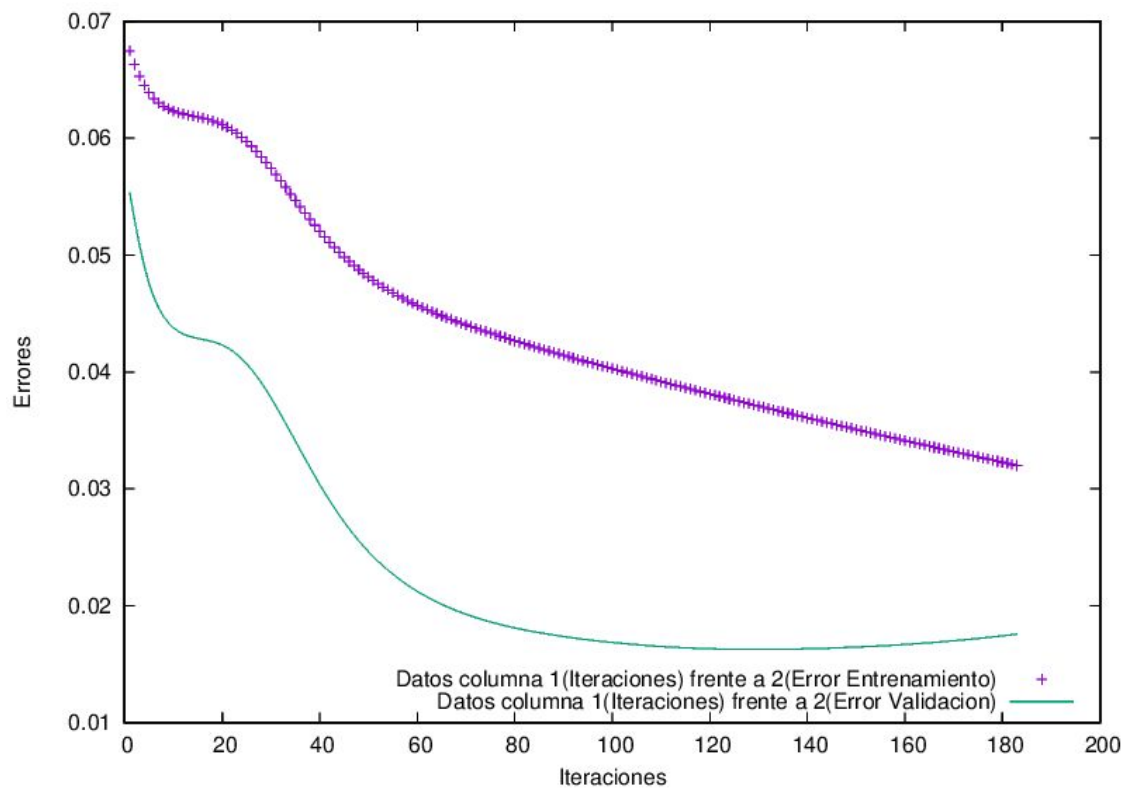




Seed3:

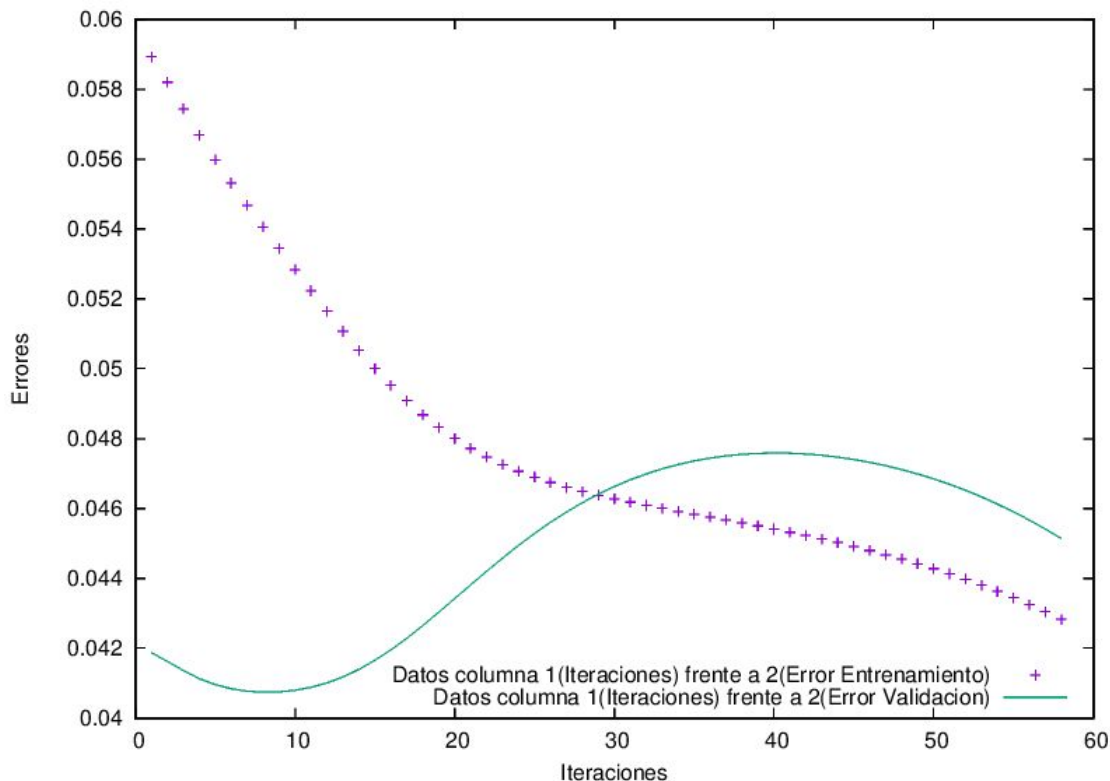


Seed4:





Seed5:



De estas gráficas se puede extraer mucha información.

En la primera gráfica tanto el error en la validación como el error en el entrenamiento sube por tanto se para el número de ejecuciones.

En la segunda gráfica siempre disminuye el error de ambas por tanto actúa como la seed3, es decir, como si no existiera validación realizando así todas las iteraciones posibles.

Respecto a la cuarta ocurre algo interesante, la validación corta el número de iteraciones ya que cuando llega a cierto punto el error de validación aumenta y la diferencia entre el error y el siguiente es notorio, sino lo fuera la validación no pararía las iteraciones.

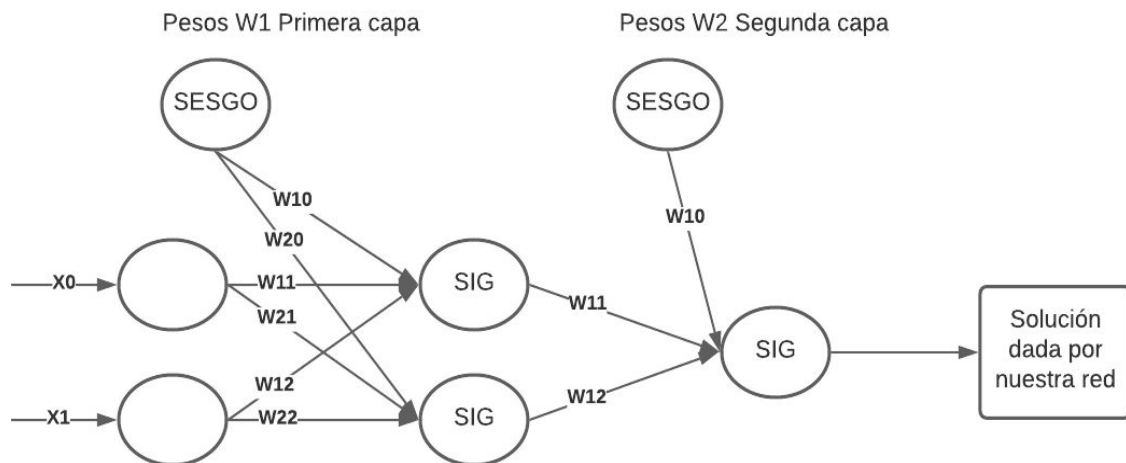
Por último, en la quinta iteración, ocurre algo similar al caso anterior pero tras realizar una subida, posteriormente hace una bajada y esto me da que pensar.

Tras revisar código he comprobado que cuando la validación baje a la validación esta se debe resetear, es decir, en este caso lo que sucede con la validación es que siempre es peor al mínimo aunque al final decrece y nosotros debemos de contemplar que cuando se realice cualquier bajada resetee este contador por que la validación solo se debe parar cuando aumente durante 50 iteraciones seguidas, es decir, no influye el error mínimo de la validación. Por ello, es necesario realiza un



cambio en el código para arreglar este error, ya que en ciertas ocasiones puede darse que termine el número de iteraciones antes de tiempo.

Por último, para que quede todo claro, se va a mostrar gráficamente cómo quedaría la red neural con una estructura de una capa con dos neuronas en la base de datos de XOR.



Las entradas que se le darán a nuestra red son las siguientes:

1	-1
-1	-1
-1	1
1	1

Y los pesos de nuestra red los siguientes:

Capa 1

W10-> 2.79585
W20-> -2.95944
W11-> -2.95888
W21-> -2.72116
W12-> -2.74063
W22-> -2.74222



Capa 2

W10-> -2.05467

W11-> 4.70256

W12-> -5.01248

Las salidas esperadas respecto a las salidas de nuestra red neuronal son las siguientes:

1 -- 0.887675

0 -- 0.112299

1 -- 0.887796

0 -- 0.134979

Como podemos observar, en esta red más o menos se acerca a los valores reales pero dista bastante los valores que deseamos que son aquellos que se acercan más a 1 cuando debe salir 1 y 0 cuando debe ser 0.

Con más capas y más neuronas se acerca a esos valores, ya que entrena mejor la red para alcanzarlos.