

## Un plan de mejoras en 5 etapas parte 2

### De las redes feed forward panditas a las profundas

#### 3. Aprimoramos el descenso por el gradiente Descenso por el gradiente estocástico

##### Introduciendo aleatoriedad

☑ **Minibatch:** Este algoritmo consiste en dividir el conjunto de entrenamiento en mini grupos por lo que en cada época se hace una reasignación aleatoria de todos los ejemplos de manera que  $Qm = p$  (se tienen  $p$  elementos del conjunto de entrenamiento divididos en  $m$  grupos de tamaño  $Q$ ). Con esta técnica se pretende introducir aleatoriedad en el método del descenso por el gradiente al sortear época a época los elementos del conjunto de entrenamiento. En general se encuentra que la cantidad de minibatches debe ser lo suficientemente alta para lograr introducir la aleatoriedad deseada, sin embargo si la cantidad de particiones es demasiado alta también lo será el costo computacional pues serán muchas veces que se realice el descenso por el gradiente. Con respecto a esto Keskar, Mudigere, Nocedal, Smelyanskiy, & Tang (2016) han comentado que cuando los batch son muy amplios (con muchos elementos de entrenamiento) la función costo adquiere una forma tal que el método del descenso por el gradiente tiende a converger hacia mínimos más agudos (pendientes altas). Por su parte, dividir los ejemplos en conjuntos de elementos más pequeños para analizar produce una convergencia del sistema hacia mínimos más planos (pendientes menores). Esto quiere decir que los batch más largos tienden hacia mínimos de los cuales es difícil despegarse mientras que los batch más pequeños, al ser atraídos hacia mínimos con formas más planas tienen la capacidad de escapar de las regiones de mínimos locales.

☑ **Dropout:** Para reducir el overfitting y permitir que una red neuronal artificial responda correctamente cuando es testeada, se aplica el dropout. Este algoritmo implica suprimir un cierto número de neuronas de manera aleatoria por lo que cuando se calcula la salida no se considera toda la información. Desde este punto de vista, cada neurona presente en las capas ocultas son omitidas al azar con una probabilidad  $p = 0,5$  tal que estas neuronas no dependen del resto. Con esto se pretende evitar el establecimiento de complejas adaptaciones donde cada característica a detectar es útil y posible sólo en el contexto planteado por todas las demás. Es decir, cada neurona aprende a detectar una característica que permite producir la respuesta correcta en una multiplicidad de contextos internos posibles y variados en el cual opera la red. En sentido gráfico, el método del descenso por el gradiente no es capaz de salir de mínimos locales y requiere tiempos muy largos para salir de un punto de inflexión, tiempo que crece exponencialmente a medida que aumenta el número de componentes en la cual es plano. Este problema se facilita con la aplicación del algoritmo dropout al suprimir una de las dimensiones donde podría estar ubicado un mínimo local.

Es posible encontrar un paralelismo entre este funcionamiento de la red artificial (dropout) con el funcionamiento del cerebro dado que las neuronas requieren de un tiempo posterior a la emisión de señales para recuperar la energía liberada cuando las transmiten a neuronas vecinas de manera que, por un tiempo, se podría decir que permanecen apagadas.

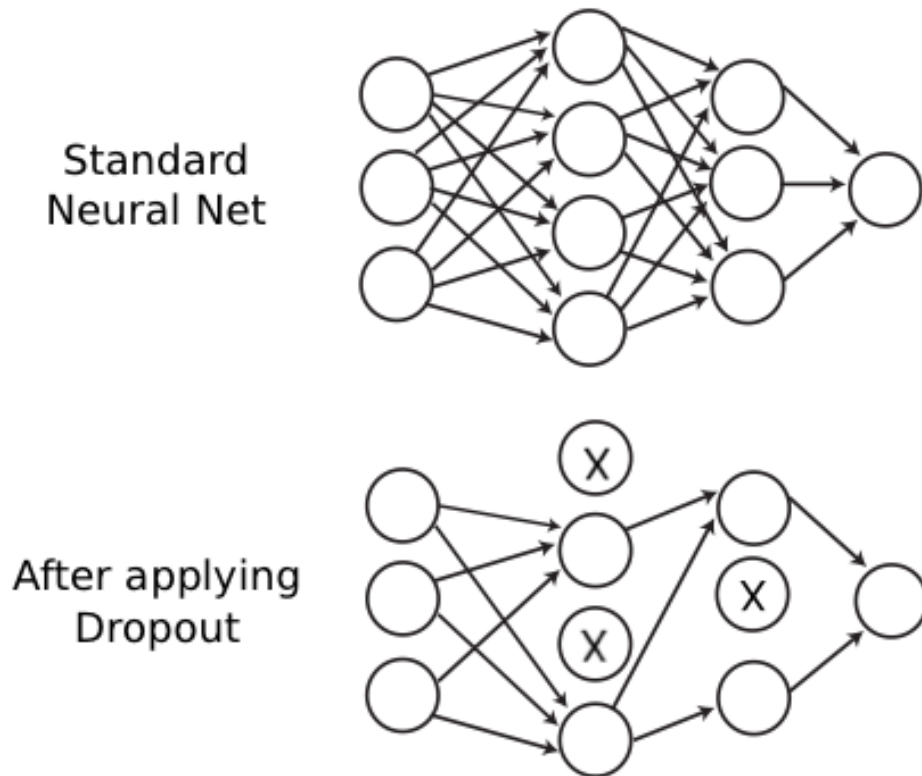


FIG. 39 **Dropout** During the training procedure neurons are randomly “dropped out” of the neural network with some probability  $p$  giving rise to a thinned network. This prevents overfitting by reducing correlations among neurons and reducing the variance in a method similar in spirit to ensemble methods.

\*Improving neural networks by preventing co-adaptation of feature detectors by Geoffrey Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov (2012).

### 3. Aprimoramos el descenso por el gradiente Descenso por el gradiente estocástico

## 4) Ponemos memoria y adaptación al descenso

### 4. Agregamos memoria y adaptación al descenso por el gradiente

Agregamos Momento

$$W_{ij}^{(t+1)} = W_{ij}^{(t)} + \Delta W_{ij}^{(t)}$$

$$\Delta W_{ij}^{(t)} = \underbrace{\gamma \Delta W_{ij}^{(t-1)}}_{\text{termino de momento}} - \underbrace{\eta \frac{\partial E}{\partial W_{ij}}}$$

termino de momento

$$0 \leq \gamma \leq 1$$

Esto le permite al método ganar fuerza cuando el incremento es persistente pero pequeño

Si andamos por una zona "plana" de la función error, y el parámetro de momento es cercano a 1

$$\Delta w_{ij}^{(t)} \approx \Delta w_{ij}^{(t-1)}$$

$$(1 - \gamma) \Delta w_{ij}^{(t)} = -\eta \frac{\partial E}{\partial w_{ij}}$$

$$\Delta w_{ij} = -\frac{\eta}{(1 - \gamma)} \frac{\partial E}{\partial w_{ij}}$$



avanzamos con un paso más largo

Podemos comenzar con  $\gamma = \frac{1}{2}$  y aumentarlo a medida que descendemos hasta llegar a

$$\gamma = 0.99$$

# Nesterov Accelerated Gradient (NAG)

Ilya Sutskever 2012

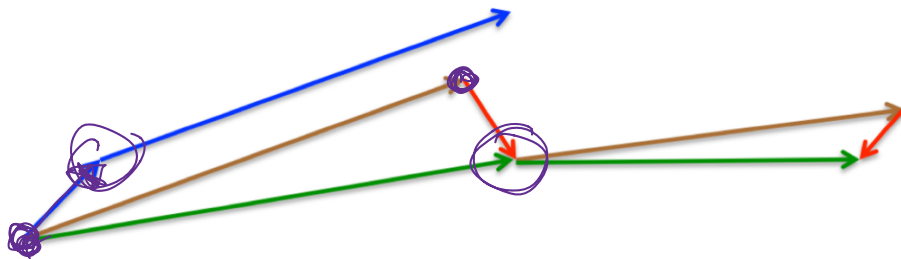
$$W_{ij}^{(t+1)} = W_{ij}^{(t)} + \Delta W_{ij}^{(t)}$$

$$\Delta W_{ij}^{(t)} = \gamma \Delta W_{ij}^{(t-1)} - \gamma \frac{\partial E}{\partial W_{ij}} (W_{ij}^{(t-1)} + \gamma \Delta W_{ij}^{(t-1)})$$

Primero damos un salto grande en la dirección del gradiente acumulado previo

Desde ese punto calculamos el gradiente

Sumamos ambos



Los vectores azules son los saltos usuales.

Ahora queremos introducir un método "adaptativo" que considere conexiones "locales" para cada sinapsis

$$\Delta w_{ij}(t) = -\gamma g_{ij}(t) \cdot \frac{\partial E}{\partial w_{ij}}$$

↑  
universal

↑  
específico

Hacemos inicialmente

$$g_{ij}(0) = 1 \quad \forall i, j$$

$$\text{IF} \left[ \frac{\partial E(t-1)}{\partial w_{ij}} \cdot \frac{\partial E(t)}{\partial w_{ij}} \right] > 0 \quad \text{THEN}$$

$$g_{ij}(t+1) = g_{ij}(t) + 0.5$$

ELSE

$$g_{ij}(t+1) = 0.95 * g_{ij}(t)$$

END

• Esto lo hacemos por mini batch!

- limitamos los valores de  $g_{ij}(t)$

$$0.1 \leq g_{ij}(t) \leq 10 \quad \forall t, i, j$$

$$0.01 \leq g_{ij}(t) \leq 100 \quad \forall t, i, j$$

- usamos minibatches NO muy pequeños.  
Esto evita a que no haya cambios de signo en el gradiente solo porque tenemos "muestras" pequeñas. Recuerden que el error es una "función aleatoria" (de muchas variables aleatorias).
- Mezclamos momentos con perss adaptativos.
- lo adaptativo se hace "por eje"



## RPROP

Tenemos presente que el gradiente puede ser muy diferente para pesos diferentes y esto puede cambiar a lo largo del descenso por el gradiente

Combinamos la idea de adaptar solo por el signo con la idea de adaptar pesos por referado

$$\text{IF } \left[ \frac{\partial E(t-1)}{\partial w_{ij}} \cdot \frac{\partial E(t)}{\partial w_{ij}} \right] > 0$$

$$\text{THEN } g(t+1) = g(t) * 1.2$$

$$g(t+1) = g(t) * 0.5$$

END

Este método anda mal para minibatches

# RMS PROP

Buscamos un método tipo RPROP pero que nos permite usar mini batches.

$$\text{Mean Square } (w_{ij}, t) = 0.9 \text{ Mean Square } (w_{ij}, t-1) + 0.1 \left( \frac{\partial E}{\partial w_{ij}} \right)^2$$

✓  
 $s_t^{ij}$

$$g_t = \frac{\partial E}{\partial w_{ij}^{(t)}}$$

$$s_t^{ij} = \beta s_{t-1}^{ij} + (1-\beta) (g_t^{ij})^2$$

$$\Delta w_{ij}^{(t)} = -\eta \frac{g_t}{\sqrt{s_t^{ij} + \epsilon}}$$

$$w_{ij}^{(t+1)} = w_{ij}^{(t)} + \Delta w_{ij}^{(t)}$$

$$\beta \approx 0.9$$

$$\epsilon \approx 10^{-8}$$

$$\eta \approx 10^{-3}$$

$\Delta\Delta\Delta M$

(adaptive moment estimation)

$$g_t = \frac{\partial E}{\partial w_{ij}}$$

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

$$\hat{m}_t = \frac{m_t}{1-\beta_1}$$

$$\Delta_t = \beta_2 \Delta_{t-1} + (1-\beta_2) g_t^2$$

$$\hat{\Delta}_t = \frac{\Delta_t}{(1-\beta_2)}$$

$$E[m_t] = E[g_t]$$

$$E[\Delta_t] = E[g_t^2]$$

$$w_{ij}^{t+1} = w_{ij}^t + \Delta w_{ij}^t$$

$$\Delta w_{ij}^t = -\eta_t \frac{\hat{m}_t}{\sqrt{\hat{\Delta}_t + \epsilon}}$$

$$\beta_1 \approx 0.9$$

$$\beta_2 \approx 0.99$$

$$\epsilon \approx 10^{-8}$$

$$\eta \approx 10^{-3}$$

# Regularizaciones

Agregamos términos adicionales a la función error (o costo, o loss) que "limiten" los valores de los pesos según nuestros deseos

$$E(\bar{w}, \bar{y}) = \sum_{k=1}^P \sum_i e(\bar{w}, y_i^k) + \lambda R(f)$$

hyperparameters

función que tenemos que aprender

Regularización  $L_1$  o Lasso (least shrinkage and selection operator)

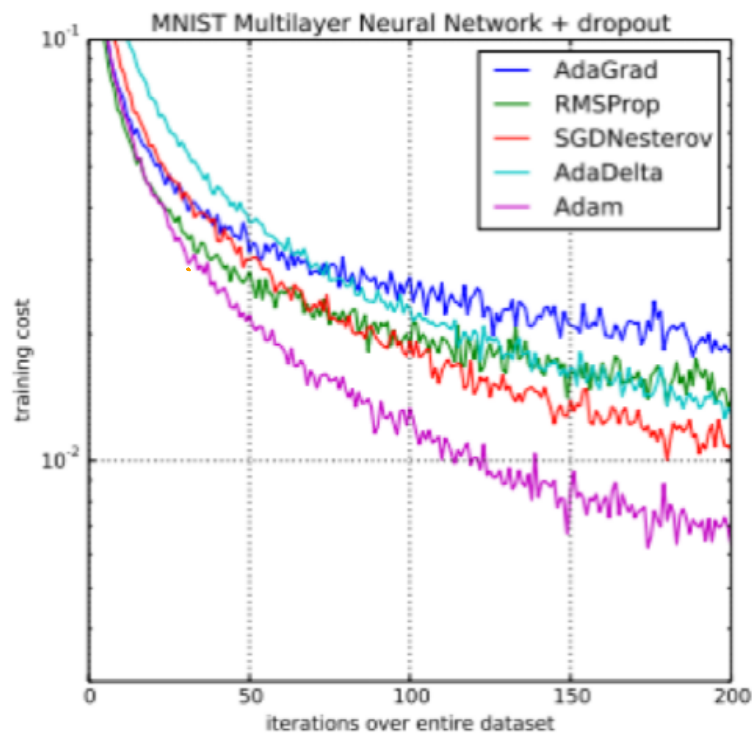
$$\bar{E}_L(\{\bar{w}\}, \{\bar{y}\}) = \sum_{k=1}^n \sum_{i=1}^m (y_i^k - \theta_i^k)^2 + \lambda \sum_{\alpha} |w_{\alpha}|$$

suma sobre todos los coeficientes y umbrales

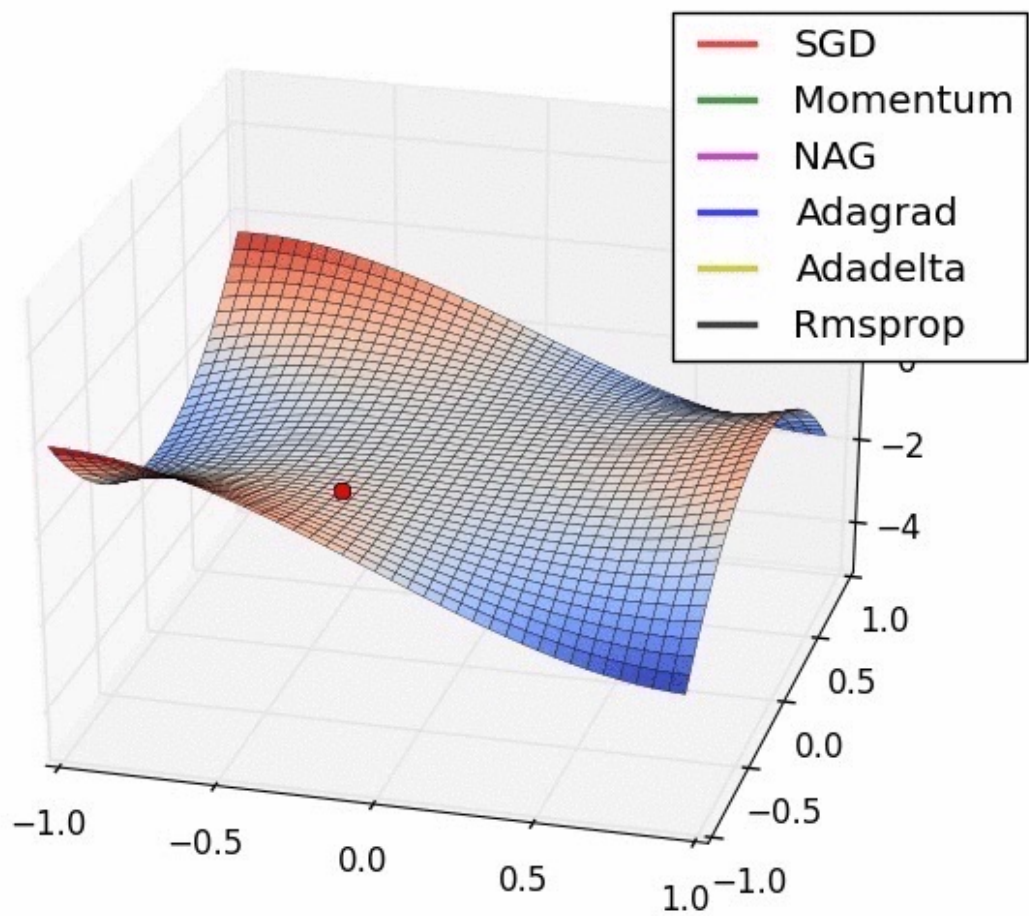
Regularización  $L_2$  o Ridge

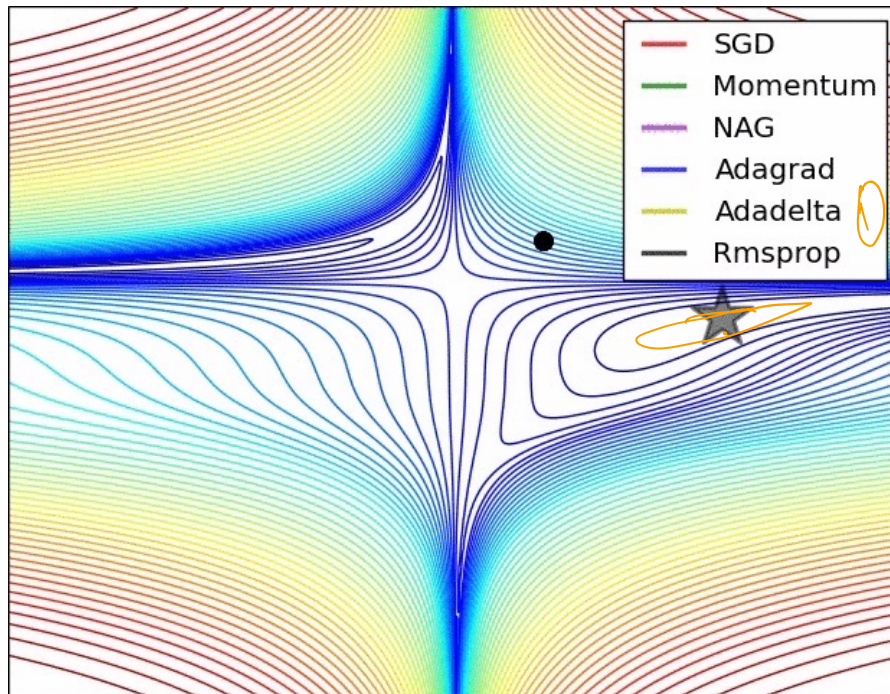
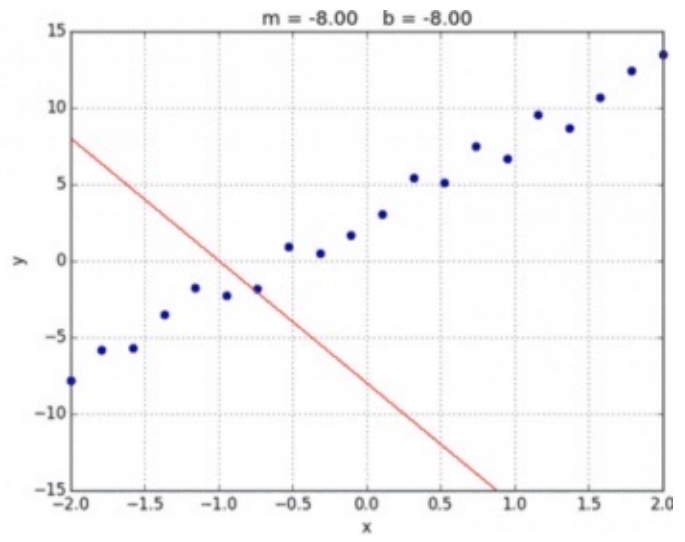
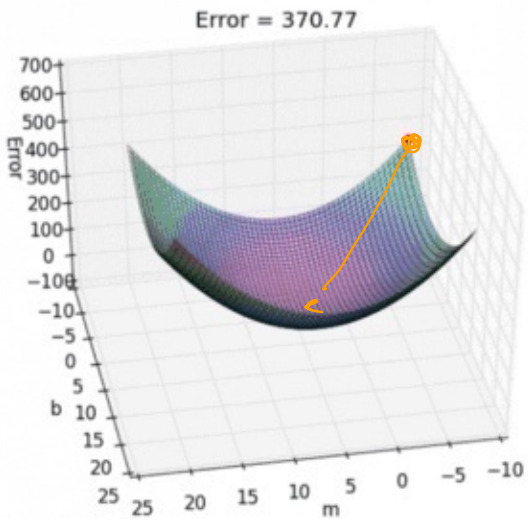
$$\bar{E}_L(\{\bar{w}\}, \{\bar{y}\}) = \sum_{k=1}^n \sum_{i=1}^m (y_i^k - \theta_i^k)^2 + \lambda \sum_{\alpha} w_{\alpha}^2$$

suma sobre todos los coeficientes y umbrales



Diferentes métodos de optimizar el descenso por el gradiente (problema MNIST)





## 4- Cambiamos la función error o costo

☑ **Cross Entropy:** La función error se ha definido como la suma de los cuadrados de la diferencia entre los valores deseados y la salida obtenida por la red. Por la naturaleza de esta función y dado que el método del descenso por el gradiente se definió proporcional al  $\frac{\partial E}{\partial w}$ , cuanto más próximos estemos al mínimo donde sea  $E = C = 0$  menor será la velocidad (pendiente) con la que nos acerquemos. Si esta función costo fuera reemplazada por una función logarítmica, el descenso por el gradiente iría incrementando su velocidad de acercamiento al mínimo dado que el gradiente sería cada vez mayor. Esto se suele usar como una métrica para saber cuán cerca estamos de la solución correcta.

La cross entropy o entropía cruzada viene de la teoría de la información. La información es una medida de la cantidad de bits requeridos para codificar y transmitir un evento.

Eventos poco probables tienen mucha información

Eventos muy probables tienen poca información

La entropía es una medida del número de bits requeridos para codificar y enviar un evento obtenido de una distribución de probabilidades.

En teoría de información cuantifican la sorpresa.

Si un evento tiene poca probabilidad, es más "sorprendente"

Si la distribución de probabilidad de los bits tiene un pico nítido  $\Rightarrow$  baja posibilidad de sorpresa

Si la distribución de probabilidad de los bits es balanceada  $\Rightarrow$  alta probabilidad de sorpresa.

Para una variable aleatoria  $x$  obtenida de una distribución  $p(x)$ , la cantidad

$$H(x) = - \sum_{x \in P(x)} p(x) \ln(p(x))$$

La entropía cruzada se construye sobre la idea de entropía y calcula el número de bits requeridos para codificar o transmitir un evento construido a partir de una distribución comparada con otra.



En Inteligencia Artificial es el no de bits necesarios para codificar datos que vienen de una fuente con una distribución  $p$  con un modelo  $q$ .

$$H = - \sum_{x \in \mathcal{P}(x)} P(x) \log(Q(x))$$

Kullback 1959

Hopfield 1987

Sara Solla 1988

Para el perceptron para aprendizaje supervisado

Si  $O_i^k = \frac{1}{1 + e^{-h_i^k}}$ , entonces

$$E(\{w\}, \{J\}) = - \sum_{i, \mu} J_i^\mu \log_2(O_i^k) + (1 - J_i^\mu) \log_2(1 - O_i^k)$$

$E \geq 0$  para cualquier  $\{w\}$  o  $\{J\}$

$$E = 0 \Rightarrow O_i^k = J_i^\mu \quad \forall i, \mu$$

Si  $O_i^k = \tanh(h_i^k)$ , entonces

$$E(\{w\}, \{J\}) = - \sum_{i, \mu} \left\{ (1 + J_i^\mu) \log_2(1 + O_i^k) + (1 - J_i^\mu) \log_2(1 - O_i^k) \right\}$$

$E \geq 0$  para cualquier  $\{w\}$  o  $\{J\}$

$$E = 0 \Rightarrow O_i^k = J_i^\mu$$



# Observaciones Generales

- Si el conjunto de entrenamientos (data set en la jerga de IA) es muy redundante (tiene datos parecidos), el gradiente que obtenemos con una mitad de los datos es igual al gradiente que obtenemos con la otra mitad. Podemos entonces aprender con una mitad, y vamos cambiando "la mitad", eligiéndola al azar en cada época.

No tiene que ser mitad. Podemos sufrir en cada época  $p$  neuronas y todas sus sinapsis.

- Si usamos todo el gradiente, podemos usar métodos de optimización global, sobre todo heurísticos, pero si el no de sinapsis es muy alto y los datos son redundantes, es mejor usar dropout y minibatches.

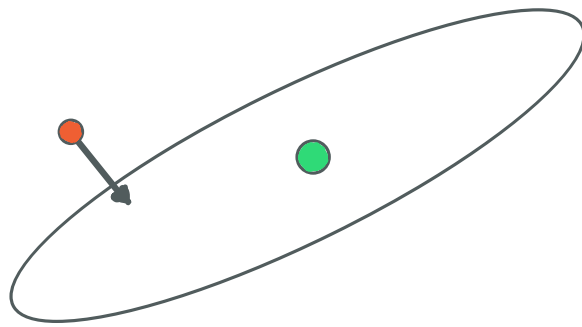
- Cuanto mayores son los minibatches, más eficiente es el método computacionalmente

- Si dos capas ocultas de neuronas tienen los mismos umbrales y las mismas sinapsis entrando y saliendo, tendrán el mismo gradiente y seguirán siempre iguales. Por esto es importante evitar estas decisiones simétricas.
- Si una neurona oculta tiene muchas entradas (large fan-in en la jerga), pequeños cambios en los pesos pueden producir un resaca en la actualización (overshooting). Otra, que no sigue bien la topología del gradiente. Por esto es importante normalizar los valores iniciales de forma tal que sean proporcionales a  $\frac{1}{\sqrt{N}}$ , donde  $N$  es el número de entradas

En el caso de neuronas lineales vemos que la función error es una función cuadrática en muchas direcciones y plana en otras

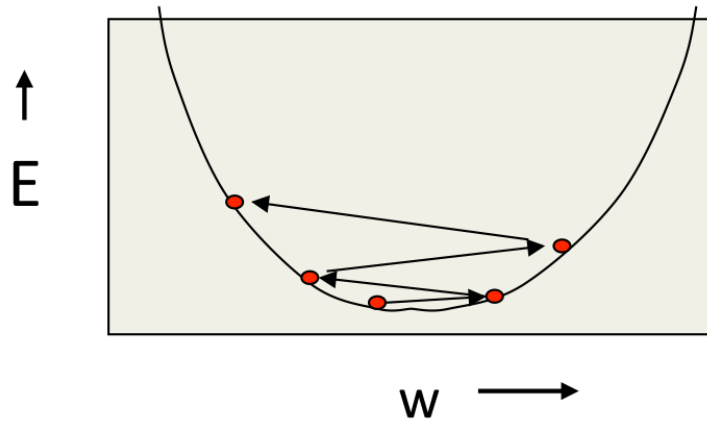
En el caso de funciones no lineales todo se complica pero localmente podemos pensar alrededor de cada mínimo como un bowl cuadrático con diferentes características.

Bajar por el gradiente en un batch completo garantiza que el error disminuye, pero no apunta nunca hacia el mínimo



El gradiente es grande en la dirección en la cual debemos "avanzar" poco, pero es pequeño en la dirección en la cual debemos "avanzar" mucho.

Si el learning rate es demasiado grande el descenso por el gradiente se queda "atrapado en un barranco", o "salten el barranco"



Lo que deseamos es:

- o movernos rápidamente en las direcciones con pequeños gradientes
- o movernos lentamente en las direcciones con gradientes grandes.

