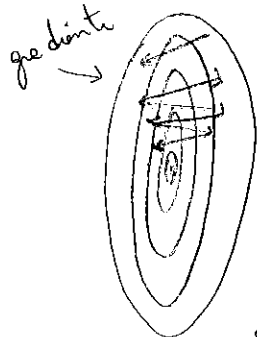


Optimización (Actualización) [empezar con "learning rate decay"]

+ Descenso de gradiente con momentum



pero me gustaría
descender en dirección ↓

Solución posible: no bajar por el gradiente
ni las consideras gradientes anteriores y
"mezclarlos" promediando



promedio \Rightarrow

la dirección horizontal se conserva

Pero no quiero promediar todo la historia \Rightarrow le doy mayor
peso relativo a los últimos valores \Rightarrow promedio exponencial móvil!

En cada iteración mantengo el promedio móvil de $\frac{\partial \mathcal{L}}{\partial W}$ para cada
parámetro W

$\beta = 0$
recupero de-
cay standard

$$\bar{\mathcal{L}} := [\beta] \bar{\mathcal{L}} + (1-\beta) \frac{\partial \mathcal{L}}{\partial W} \quad \leftarrow \text{promedio móvil}$$

$$W := W - \lambda \bar{\mathcal{L}} \quad \leftarrow \text{desciendo por el promedio del gradiente}$$

Interpretación: velocidad + aceleración + fricción
- la regla de actualización la puedo escribir como

$$W := W + \left(-\lambda \beta \bar{\mathcal{L}} - \lambda (1-\beta) \frac{\partial \mathcal{L}}{\partial W} \right)$$

esto induce otro cálculo α para el promedio con λ'

$$V_{\partial W} := \alpha V_{\partial W} - \lambda' \frac{\partial \mathcal{L}}{\partial W} \quad \leftarrow \text{velocidad } V_{\partial W} \text{ (para el cambio en } W)$$

$$W := W + V_{\partial W} \quad \leftarrow \text{parámetro } \alpha \text{ más de fricción}$$

implementación más típica de

momentum \Rightarrow hiperparámetros: α, λ' ($\alpha = 0$ recupera descenso de gradiente standard)

Variante Nesterov: Para cada parámetro W

$$\tilde{W} := W + \alpha V_{\partial W} \quad \leftarrow \text{lo hago por todos los parámetros}$$

$$V_{\partial W} := \alpha V_{\partial W} - \lambda' \frac{\partial \mathcal{L}}{\partial W} \Big|_{\tilde{W}} \quad \leftarrow \text{antes de hacer backpropagation para los parámetros } \tilde{W}$$

$$W := W + V_{\partial W}$$

el gradiente es evaluado en $W + \alpha V_{\partial W} \Rightarrow$ he de
de "corregir" el momentum standard.

Idea: Me anticipo al cambio de posición y calculo el gradiente allí

* Algoritmos con tasa de aprendizaje adaptable

Idea: cambiar la tasa de aprendizaje durante el entrenamiento

- Más simple: no cambiar nada (Opción 0)

- No tan simple: hacer que la tasa vaya "decreciendo" de 0 pesos

¿Por qué? o medida que nos acerquemos al mínimo general de "pesos" más pequeños

Idea 1:

$$\lambda := \frac{1}{1 + d * t} \lambda_0$$

↑
hiperparámetro
 $d \geq 0$

Idea 2:

$$\lambda := |d|^t * \lambda_0$$

↑
hiperparámetro
 $d \leq 1$

Idea 3:

actualiza λ después de cada iteración hasta llegar a un mínimo. Pasa de λ_0 a λ_T

$$\lambda := (1 - \frac{t}{T}) \lambda_0 + \frac{t}{T} \lambda_T$$

En todas las cosas puede llegar a un mínimo luego de lo cual se mantiene constante. Todo se puede hacer por epoch o minibatch (esto define el t)

λ_0 tasa inicial
 λ_T tasa final (mínima)
actualización por T

- Adagrad: adaptar la tasa de aprendizaje de cada parámetro por separado escalándolos de manera inversamente proporcional a la raíz cuadrada de la suma de todos los valores históricos el cuadrado

Para cada parámetro W en cada iteración:

Valor sugerido de $\lambda: 0.01$

$$v_W := v_W + \left(\frac{\partial \mathcal{L}}{\partial W} \right) \left(\frac{\partial \mathcal{L}}{\partial W} \right) \quad \text{multiplicación punto a punto}$$

$$W := W - \left(\lambda \frac{1}{\sqrt{v_W}} \right) \left(\frac{\partial \mathcal{L}}{\partial W} \right) \quad \left(\begin{array}{l} \text{para implementar se usa} \\ \frac{1}{\sqrt{v_W} + \epsilon} \quad \text{para estabilidad} \\ \text{numérica con } \epsilon = 10^{-7} \end{array} \right)$$

↑
sqrt punto a punto

→ La dirección con el mayor gradiente disminuye su tasa de aprendizaje más rápido

→ Con: promedios desde el principio puede penalizar demasiado temprano la tasa

→ Adagrad fue diseñado para superficies convexas en donde tiene muy buenas propiedades (teóricas y prácticas), pero DL es altamente no convexo

[2012] - RMS Prop. Idea: imitar a AdaGrad pero ponderando la suma de los cuadrados con un promedio exponencial móvil (root mean squared propagation) punto a punto

β propuesto:
0.9
 $\lambda: 0.001$

$$S_{\partial W} = \beta S_{\partial W} + (1-\beta) \frac{\partial L}{\partial W} * \frac{\partial L}{\partial W}$$

$$W := W - \lambda \frac{1}{\sqrt{S_{\partial W}}} * \frac{\partial L}{\partial W}$$

(igual que antes para implementarlo)
se usa $\frac{1}{\sqrt{S_{\partial W} + \epsilon}}$ con $\epsilon \approx 10^{-7} + 10^{-8}$

se puede combinar con momentum tb (aunque en ese caso queda con equivalente al siguiente)

- Adam [2015] (Adaptive Moments) RMSProp + Momentum

$$\overline{\partial W} := \beta_1 \overline{\partial W} + (1-\beta_1) \frac{\partial L}{\partial W} \quad \leftarrow \text{momentum}$$

$$S_{\partial W} := \beta_2 S_{\partial W} + (1-\beta_2) \frac{\partial L}{\partial W} * \frac{\partial L}{\partial W} \quad \leftarrow \text{RMS Prop.}$$

$$W := W - \lambda \frac{1}{\sqrt{S_{\partial W}}} * \overline{\partial W}$$

Típicamente se implementa con "bias correction" \Rightarrow al inicio no tenemos con qué promediar (en cálculos de $\overline{\partial W}$ y $S_{\partial W}$) lo que implica que enteremos considerando como si los primeros valores fueran todos 0 \Rightarrow esto inducirá un bias hacia 0 al inicio del entrenamiento. Corregiremos este hecho dividiendo por un número menor que 1 (en amplificando el promedio inicial) pero que tenderá rápidamente a 1 para volver al promedio que calculamos. Usamos $(1-\beta_1^t)$ para $\overline{\partial W}$ y $(1-\beta_2^t)$ para $S_{\partial W}$ con t la iteración del descenso de gradiente.

Formulación final

$$\overline{\partial W} := \beta_1 \overline{\partial W} + (1-\beta_1) \frac{\partial L}{\partial W}$$

$$\hat{\partial W} := \frac{\overline{\partial W}}{1-\beta_1^t}$$

$$S_{\partial W} := \beta_2 S_{\partial W} + (1-\beta_2) \frac{\partial L}{\partial W} * \frac{\partial L}{\partial W}$$

$$\hat{S}_{\partial W} := \frac{S_{\partial W}}{1-\beta_2^t}$$

$\epsilon = 10^{-8}$ para estabilidad

$$W := W - \lambda \frac{1}{\sqrt{\hat{S}_{\partial W}}} * \hat{\partial W}$$

Ten hiperparámetros $\lambda, \beta_1, \beta_2$. Valores un-
guidos: $\lambda = 0.001, \beta_1 = 0.9, \beta_2 = 0.999$
- cambia solo este!!! si fuera necesario

Eligiendo hiperparámetros: hay de muchos hiperparámetros :(
no se puede probar que funcione la mejor combinación

Tip 1: para architecture usar alguna que al estado del arte diga que tiene sentido (aunque se deben cambiar tamaño de capas y cantidad si hay overfitting)

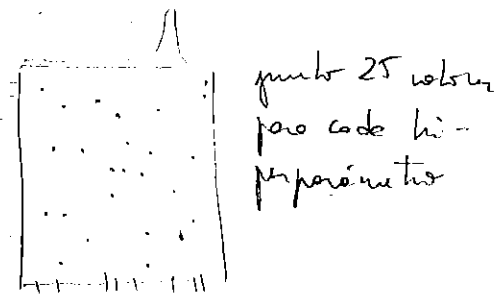
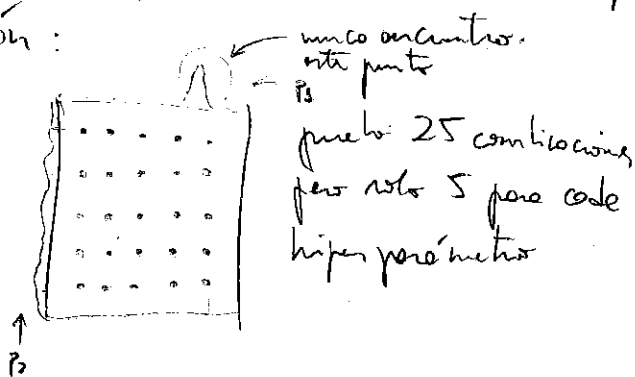
Tip 2: para learning rate o parámetro de decay, buscar en escala exponencial 0.1, 0.01, 0.001, etc

si quisiera buscar entre 0.1 y 0.001 \rightarrow elegir un $t \in [1, 3]$ y probar 10^{-t}
para capas o neuronas escala lineal está bien

Tip 3: tasa de aprendizaje debe ser el hiperparámetro más importante

Tip 4: hacer búsqueda random de los hiperparámetros (NO GRID SEARCH!!!)

versión:



Tip 5: Zoom-in en regiones donde los hiperparámetros dieron buenos resultados y continuar ahí