

OPTIMIZER

Stochastic Gradient Descent (with Momentum)

for (x, y) in dataset:

$$J = \nabla l(\theta | x, y)$$

$$m = J + \text{momentum} * m$$

$$\theta = \theta - \epsilon * m$$

→ init momentum = 0

- Pros: Works in most cases
- Cons: Need to tune your learning rate

What happen if learning rate too high?

→ Loss spikes, because:

- Magnitude of weights are much larger

$$\|\theta\| = \|w\| \nearrow$$

- Magnitude of gradient are much larger

$$\|\nabla l_{\theta}\| = \|\nabla w l\| \nearrow$$

⇒ And since $w = w - lr * \text{gradient}$. Larger learning rate leads to larger "fluctuation" in weight, which leads to loss spike.

How to prevent loss spike?

1. RProp

Basic idea: scale gradient by its magnitude $\|\nabla_{\theta} l\|$

$$\begin{matrix} \dots \\ m = J / J.\text{norm}() + \text{momentum} * m \\ \dots \end{matrix}$$

2. RMS Prop:

Compute v : a running average of $\|\nabla_{\theta} l\|^2$
scale by $1/\sqrt{v}$

$$\begin{matrix} \dots \\ v = \beta_2 * v + (1 - \beta_2) * J.\text{square}() \\ m = J / v.\text{sqr}t() + \text{momentum} * m \\ \dots \end{matrix}$$