

## Adam optimizer

Saturday, October 9, 2021 4:22 PM

Adaptive moment estimation  $\Rightarrow$  Combination of momentum + RMSprop + exponential decay optimization

$$\left. \begin{array}{l} \textcircled{1} \text{ Momentum} \Rightarrow m = \beta_1 m + (1 - \beta_1) \frac{\partial L}{\partial W} \Big|_W \\ \textcircled{2} \quad \quad \quad \beta_t = \gamma \cdot \beta_{t-1} + (1 - \gamma) \cdot \frac{\partial L}{\partial W} \Big|_W \\ \textcircled{3} \text{ Bias correction} \Rightarrow \hat{m} = \frac{m}{1 - \beta_1^t} \Rightarrow t = \text{iteration} \\ \textcircled{4} \quad \quad \quad \hat{\beta}_t = \frac{\beta_{t-1}}{1 - \beta_{t-1}^t} \\ \textcircled{5} \quad W = W - \eta \frac{\hat{m}}{\sqrt{\hat{\beta}_t + \epsilon}} \end{array} \right\} \rightarrow \text{Combination of everything}$$

- If the data is sparse, use the self-applicable methods, namely Adagrad, Adadelata, RMSprop, Adam.
- RMSprop, Adadelata, Adam have similar effects in many cases.
- Adam just added bias-correction and momentum on the basis of RMSprop,
- As the gradient becomes sparse, Adam will perform better than RMSprop.

Adam is the most popular one