Today :. Review

    • Adam

    • Batch Normalization

    • Regularization: norm penalty, early stopping, dropout

• Adam: Adaptive Moments - a combination of RMS prop & momentum

    •    $\underline{s} \leftarrow \rho_1 \underline{s} + (1-\rho_1) \underline{g}$      $\rho_1, \rho_2$ : hyper parameters

    •    $\underline{r} \leftarrow \rho_2 \cdot \underline{r} + (1-\rho_2) \cdot \underline{g} \odot \underline{g}$

    •    $\hat{\underline{s}} = \dfrac{\underline{s}}{1-\rho_1^t}$    (bias correction)      $t$: iteration count

    •    $\hat{\underline{r}} = \dfrac{\underline{r}}{1-\rho_2^t}$    (bias correction).

    •    $\Delta\theta = -\dfrac{\varepsilon}{\sqrt{\underline{s}}+\hat{\underline{r}}} \odot \hat{\underline{s}}.$

    •    $\theta \leftarrow \theta + \Delta\theta$

<u>Batch normalization:</u>    $\hat{y}_i = w_1 \cdots w_\ell \, x$

            $\underline{w} = [w_1 \cdots w_\ell]^T$

            $\underline{w} \leftarrow \underline{w} - \varepsilon \underline{g}$      $\underline{g} = [g_1 \cdots g_\ell]^T$

            $\hat{y}_{i+1} = (w_1 - \varepsilon g_1) \cdots (w_\ell - \varepsilon g_\ell) \cdot x$

$$= \cdots + \varepsilon^2 g_1 g_2 \underbrace{\prod_{i=3}^{l} w_i}_{} + $$
$$\varepsilon^l g_1 \cdots g_l$$ ✓

observation: The choice of $\varepsilon$ should be such that the contribution of the second order gradient relation is reduced, while also reducing the contribution of the other powers of $\varepsilon$.

This is a _hard_ problem.  $\left( h_i = \underline{h_{i-1}} \cdot w_i \right.$, $h_i$: output at $i$th layer,
How can batch normalization help?  $w_i$: weight of $i$th layer $\left. \right)$

Batch normalization: If



$$H = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \end{bmatrix} D$$

$$H' = \frac{H - \underline{M}}{\underline{\sigma}}$$