

Kfac assumption

Yaroslav Bulatov

April 10, 2017

Suppose we have a left matmul operation, $y = Wx$. The gradient with respect to W evaluated on example i can be written as

$$G_i = b_i a'_i$$

Quantities b and a are column vectors representing backprops and activations. Note that this has the same shape as W .

For purposes of computing covariance, we flatten this gradient

$$g_i = \text{vec}(G_i) = a_i \otimes b_i$$

Covariance can now be expressed as

$$\text{cov} = \frac{1}{n} \sum_i g_i g'_i = \frac{1}{n} \sum_i (a_i a'_i) \otimes (b_i b'_i)$$

To obtain natural gradient preconditioner we seek a transformation that renders new covariance matrix diagonal.

Let $\hat{a}_i = Ua$ and $\hat{b}_i = Vb$, the corresponding covariance is

$$\text{cov} = \frac{1}{n} \sum_i U a_i a'_i U' \otimes V b_i b'_i V'$$

It's easy to find U and V if we assume the quantity above is equalent to

$$\text{cov} \approx \frac{1}{n} \left(\sum_i U a_i a'_i U' \right) \otimes \left(\sum_i V b_i b'_i V' \right) = \frac{1}{n} U A A' U' \otimes V B B' V$$

From this approximation it follows that we should pick $U = A^{-1}$ and $V = B^{-1}$