

ECC: Platform-Independent Energy-Constrained Deep Neural Network Compression via a Bilinear Regression Model



Haichuan Yang¹, Yuhao Zhu¹, Ji Liu^{2,1} ¹University of Rochester, ²Kwai Al Lab at Seattle

Energy Estimation Model

Bilinear Model: approximate the energy consumption. Assumption: For each layer *j*, the energy estimation can be modeled by the interaction between s_i and s_{i+1} : $|\mathcal{U}|$

$$\hat{\mathcal{E}}(\mathbf{s}) = a_0 + \sum_{j=1}^{j-1} a_j$$

Estimate the energy model via a data-driven approach: $-\mathcal{E}(\mathbf{s}))^{2}$] a_0, a_0

$$\min_{a,\ldots,a_{|\mathcal{U}|}\geq 0} \mathrm{E}_{\mathbf{s}}[(\hat{\mathcal{E}}(\mathbf{s})$$

Optimization Algorithm

ADMM reformulation:

 $\min_{\mathcal{W},\mathbf{s}} \max_{z \ge 0, \mathbf{y} \ge \mathbf{0}} \ell(\mathcal{W}) + \mathcal{L}_1(\mathcal{W}, \mathbf{x})$

 $\phi(\mathbf{w}^{(j)}) \leq$

$$\mathcal{L}_1 := \frac{\rho_1}{2} \sum_j \left[\phi(w^{(j)}) - s_j \right]_+^2 + \sum_j y_j(\phi(w^{(j)}) - s_j)$$
$$\mathcal{L}_2 := \frac{\rho_2}{2} \left[\hat{\mathcal{E}}(\mathbf{s}) - E_{\text{budget}} \right]_+^2 + z(\hat{\mathcal{E}}(\mathbf{s}) - E_{\text{budget}})$$

Update
$$\mathcal{W}$$
:
 $\mathcal{W} \leftarrow \operatorname*{argmin}_{\mathcal{W}} \ell(\mathcal{W}) + \mathcal{W}$
Update s:
 $\mathbf{s} \leftarrow \mathbf{s} - \beta \nabla_{\mathbf{s}} (\mathcal{L}_1(\mathcal{W}, \mathbf{s}, \mathbf{y}))$
Update \mathbf{y}, \mathbf{z} :
 $y_j \leftarrow [y_j + \rho_1(\phi(w^{(j)}))]$

 $y_{j} \leftarrow \left[y_{j} + \rho_{1}(\phi(w^{(j)}) - s_{j}) \right]_{+}$ $z \leftarrow \left[z + \rho_{2}(\hat{\mathcal{E}}(\mathbf{s}) - E_{\text{budget}}) \right]_{+}$

 $l_j S_j S_{j+1}$

$$\mathbf{s}, \mathbf{y} + \mathcal{L}_2(\mathbf{s}, z)$$

$$\leq s_j \qquad \hat{\mathcal{E}}(\mathbf{s}) \leq E_{\text{budge}}$$

$\mathcal{L}_1(\mathcal{W}, \mathbf{s}, \mathbf{y})$

$\mathbf{y}) + \mathcal{L}_2(\mathbf{s}, z))$



- inference;



Conclusion

Proposed a bilinear energy consumption model of DNN

Leveraged the energy estimation model to formulate DNN compression as constrained optimization; Outperformed state of the arts.