

03/10/23

Deep Learning

- Recap: CNN, VGG
- Challenges in training deep networks
- Introduction to batch normalization
- The simplest variant of the VGG model has 182M params.
- Modern DL models have in excess of 500M params!
- This leads to challenges in training such models
- We will discuss popular methods to deal with such challenges.
- Batch Normalization

- Recall the SGD method

- It helped when in the # training samples is large

- However, a drawback is noisy estimates of gradients

- We also have to deal with internal covariate shifts due to the minibatches

- Batch-norm addresses this problem by "whitening" the inputs to a layer in the net

- Let $B = \{\underline{x}^{(1)}, \dots, \underline{x}^{(m)}\}$ be the minibatch of samples in SGD : with $\underline{x} \in \mathbb{R}^d$

- For each dim $k=1 \dots d$ do the following

$$\mu_k = \frac{1}{m} \sum_{i=1}^m x_k^{(i)}$$

$$\sigma_k^2 = \frac{1}{m-1} \sum_{i=1}^m (x_k^{(i)} - \mu_k)^2$$

For each sample $i=1 \dots m$ do the following

$$\hat{z}_k^{(i)} = \frac{x_k^{(i)} - \mu_k}{\sqrt{\sigma_k^2 + \epsilon}} ; \quad \epsilon \text{ is a small constant}$$

$$y_k^{(i)} = \gamma_k \cdot \hat{z}_k^{(i)} + \alpha_k$$

γ_k, α_k are learnable params
 Scaling & bias used to ensure the expressivity of the model

- Note that the effective input of minibatch samples to the layer is now

$$\mathbf{y}_{\text{BN}} = \{ \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)} \}$$

Also note that the dim of \mathbf{y} is the same as \mathbf{z}

- During training, the sample mean and variance of the entire dataset is computed and saved.
- During inference, an input sample is effectively whitened using the "global" mean and variance.