

Batch Normalization

Related to input transformation. So that loss landscape looks smooth

answered below

Question
should we use moving
avgs instead of current
avgs for normalize

① mean = $\frac{1}{n} \sum A_i^p$ all activations

Std = $\sqrt{\frac{1}{n} \sum (A_i^p - \mu)^2}$

innovation
each batch
norm layer
allows for
different
mean and
variance

② $BN_i^p = \gamma \hat{A}_i^p + \beta$
learnable parameters

② normalize

$$\hat{A}_i^p = \frac{A_i^p - \mu^p}{\sigma^p}$$

④ Exponential Moving Average

$$\mu_{mov}^p = \alpha \mu_{mov}^p + (1-\alpha) \mu_i^p$$

(EMA) $\sigma_{mov}^p = \alpha \sigma_{mov}^p + (1-\alpha) \sigma_i^p$

hyperparameter

"momentum"

Why EMA?

Ideally

If I want averages of all activations I can keep track of entire data for every minibatch. But that's expensive. So an exponential moving average is a good approximation

these parameters are not used during training they are only used for inference

Batch Normalization

Related to input transformation. So that loss landscape looks smooth

answered below

Question
Should we use moving
avgs instead of current
avgs for normalize

① mean = $\frac{1}{M} \sum A_i^p$ all activations

Std = $\sqrt{\frac{1}{M} \sum (A_i^p - \mu)^2}$

? innovation
each batch
norm layer
allows for
different
mean and
variance

② $BN_i^p = \underbrace{\gamma}_\text{learnable parameters} \hat{A}_i^p + \underbrace{\beta}_\text{learnable parameters}$

② normalize

$$\hat{A}_i^p = \frac{A_i^p - \mu^p}{\sigma^p}$$

③ Exponential Moving Average

$$\mu_{movi} = \alpha \mu_{movi} + (1-\alpha) \mu_i$$

(EMA) $\sigma_{movi} = \alpha \sigma_{movi} + (1-\alpha) \sigma_i$

hyperparameter
"momentum"

Why EMA?

Ideally

If I want averages of all activations I can keep track of entire data for every minibatch. But that's expensive. So an exponential moving average is a good approximation

these parameters are not used during
they are only used
for inference