

MOTIVATIONS → REDUCE INTERNAL COVARIATE SHIFT

2 INPUTS TO A NETWORK VARY OVER TIME RESULTING IN CHANGES TO THE RANGE OF NETWORK ACTIVATION OUTPUTS

THIS IS INTERNAL COVARIATE SHIFT

ANSWER IS BATCH NORM

NORMALIZATION GENERALLY IS DONE TO INPUTS TAKE IT A STEP FURTHER AND DO IT TO EACH LAYER'S INPUTS

NORMALIZATION VIA MINI-BATCH STOCHASTIC GRADIENT DESCENT

GET MEAN AND VARIANCE FOR ACTIVATIONS ONLY THROUGH THE ACTIVATIONS OF THAT MINIBATCH

↳ AVOIDS GETTING COMPUTATIONALLY EXPENSIVE AND NON-CONVEX

LEARNING PARAMETERS γ & β TO MAINTAIN A SPECIFIC REPRESENTATION FOR THE NETWORK

Input: Values of x over a mini-batch: $B = \{x_1, \dots, x_m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

SHOULD HAVE DISTRIBUTION OF 0 AND VARIANCE OF 1

$$\rightarrow \gamma = \text{BN}_{\gamma, \beta}(x)$$

DEPENDS ON EXAMPLE x BUT ALSO ALL

EXAMPLES IN THE MINIBATCH

$$\sum_{i=1}^m \hat{x}_i = 0, \frac{1}{m} \sum_{i=1}^m \hat{x}_i^2 = 1$$

TRAINING AND INFERENCE

NO BATCH NORM DURING INFERENCE

↓ CHANGES DISTRIBUTION OF INPUT WHICH IS BAD DURING INFERENCE

BN TRANSFORM GOES IMMEDIATELY BEFORE THE NONLINEARITY

IN CONVOLUTIONS γ & β ARE LEARNED PER FEATURE
MAP RATHER THAN ACTIVATION

↓ ACTIVATION AS IN NODE