
GroupNorm

A Preprint

Abstract

Authors suggest GroupNorm as an alternative to BatchNorm as it has significant problems when it comes to computer vision such as small batch size which is used for memory purposes but using BatchNorm error increases rapidly.

1 Introduction

A family of feature normalization methods, including BN, LN, IN, and GN, perform the following computation:

$$\hat{x}_i = \frac{1}{\sigma_i}(x_i - \mu_i)$$

where x - feature, in case of 2D images $i = (i_N, i_C, i_H, i_W)$ where N - batch axis, C - channels, H, W - height and width respectively

μ_i, σ_i defined as

$$\mu_i = \frac{1}{m} \sum_{k \in S_i} x_k \quad \sigma_i = \sqrt{\frac{1}{m} \sum_{k \in S_i} (x_k - \mu_i)^2 + \epsilon} \quad (1)$$

ϵ - small constant, S_i is the set of pixels in which the mean and std are computed (m is the size of the set).

$$\text{In BatchNorm} \quad S_i = \{k | k_C = i_C\} \quad (2)$$

$$\text{In LayerNorm} \quad S_i = \{k | k_N = i_N\} \quad (3)$$

$$\text{In InstanceNorm} \quad S_i = \{k | k_N = i_N, k_C = i_C\} \quad (4)$$

In BatchNorm μ and σ are computed along (N, H, W) axis which basically means

$$\mu_i = \frac{1}{N \cdot H \cdot W} \cdot \sum_{n=1}^N \sum_{h=1}^H \sum_{w=1}^W x_{n,c,h,w}$$

In LayerNorm μ and σ are computed along (C, H, W)

In InstanceNorm μ and σ are computed along (H, W)

And also! all these methods learn linear transformation to compensate for the possible loss of representational ability:

$$y_i = \gamma \hat{x}_i + \beta_i$$

2 GroupNorm

Formally, GroupNorm is defined as

$$\mathcal{S}_i = \{k | k_N = i_N, [\frac{k_C}{C \cdot G}] = [\frac{i_C}{C \cdot G}]\} \quad \text{where } G - \text{number of groups (predefined hyperparameter = 32)}$$

The main point in GroupNorm is that When the batch size is small in different normalizations, these statistics become noisy and unreliable estimates of the true underlying data distribution. This instability leads to increased error, but in GroupNorm it computes its normalization statistics independently of the batch size. For each sample in a batch, GN divides the channels into a predefined number of groups (G). Then, for each group within that single sample, it calculates the mean and variance using only the features belonging to that group and that sample, across their spatial dimensions (H, W). Because the calculation is confined to a single sample (and its groups of channels), the batch size (N) doesn't influence the statistics themselves. This makes GN's performance stable even with very small batch sizes