

Adaptive detrending to accelerate convolutional gated recurrent unit training for contextual video recognition

Minju Jung^{a,c}, Haanvid Lee^b, Jun Tani^{c,*}

^a School of Electrical Engineering, Korea Advanced Institute of Science and Technology, Daejeon, Republic of Korea

^b School of Computing, Korea Advanced Institute of Science and Technology, Daejeon, Republic of Korea

^c Cognitive Neurorobotics Research Unit, Okinawa Institute of Science and Technology Graduate University, Okinawa, Japan



ARTICLE INFO

Article history:

Received 23 October 2017

Received in revised form 21 March 2018

Accepted 14 May 2018

Available online 22 May 2018

Keywords:

Detrending

Normalization

Internal covariate shift

Convolutional neural networks (CNNs)

Recurrent neural networks (RNNs)

Convolutional recurrent neural networks (ConvRNNs)

ABSTRACT

Video image recognition has been extensively studied with rapid progress recently. However, most methods focus on short-term rather than long-term (contextual) video recognition. Convolutional recurrent neural networks (ConvRNNs) provide robust spatio-temporal information processing capabilities for contextual video recognition, but require extensive computation that slows down training. Inspired by normalization and detrending methods, in this paper we propose “adaptive detrending” (AD) for temporal normalization in order to accelerate the training of ConvRNNs, especially of convolutional gated recurrent unit (ConvGRU). For each neuron in a recurrent neural network (RNN), AD identifies the trending change within a sequence and subtracts it, removing the internal covariate shift. In experiments testing for contextual video recognition with ConvGRU, results show that (1) ConvGRU clearly outperforms feed-forward neural networks, (2) AD consistently and significantly accelerates training and improves generalization, (3) performance is further improved when AD is coupled with other normalization methods, and most importantly, (4) the more long-term contextual information is required, the more AD outperforms existing methods.

© 2018 Elsevier Ltd. All rights reserved.

1. Introduction

Convolutional neural networks (CNNs) (Lecun, Bottou, Bengio, & Haffner, 1998) show remarkable performance on the ImageNet challenge dataset, consisting of 1000 classes and 1.2 million training images (Krizhevsky, Sutskever, & Hinton, 2012). Encouraged by this success, several approaches exploit the spatial processing capability of CNNs in video recognition tasks (Simonyan & Zisserman, 2014; Tran, Bourdev, Fergus, Torresani, & Paluri, 2015). Two-stream CNNs (Simonyan & Zisserman, 2014) and convolutional 3D (C3D) networks (Tran et al., 2015) are the most commonly used networks. Two-stream CNNs combine classification abilities of spatial- and temporal-stream networks, being composed of a spatial-stream network that processes individual RGB frames and a temporal-stream network that processes stacked optical flow over several frames. C3D networks extend 2D convolution to 3D convolution by adding time as a third dimension, processing stacked consecutive RGB frames. However, both networks employ a stacking strategy that utilizes only a limited number of temporal correlations between stacked frames in order to recognize videos. Once the temporal window advances to the next position,

information from the previous stack is completely dropped. This creates a problem of contextual recognition that requires the extraction of long-range temporal correlations (Jung, Hwang, & Tani, 2015).

In this paper, we attempt to overcome this limitation using recently introduced convolutional recurrent neural networks (ConvRNNs) that replace the weight multiplication of RNNs with convolution in order to exploit spatial and temporal information processing capabilities of CNNs and recurrent neural networks (RNNs), respectively (Ballas, Yao, Pal, & Courville, 2015; Kalchbrenner et al., 2016; Shi et al., 2015). By extracting spatio-temporal features hierarchically, ConvRNNs handle complex problems in the space-time domain, such as precipitation nowcasting (Shi et al., 2015), video recognition (Ballas et al., 2015), and video prediction (Kalchbrenner et al., 2016). Also, problems restricted to the spatial domain can be handled by ConvRNNs in an iterative manner (Romera-Paredes & Torr, 2016). For example, in instance segmentation, ConvRNNs sequentially segment one instance of an image at a time (Romera-Paredes & Torr, 2016). However, training ConvRNNs is painfully slower than training feed-forward CNNs, which receive a single frame or stacked multiple frames for video recognition, because recurrent connections require additional computation. Moreover, it is hard to parallelize computation of ConvRNNs due to the sequential nature of RNNs, which require

* Corresponding author.

E-mail address: jun.tani@oist.jp (J. Tani).

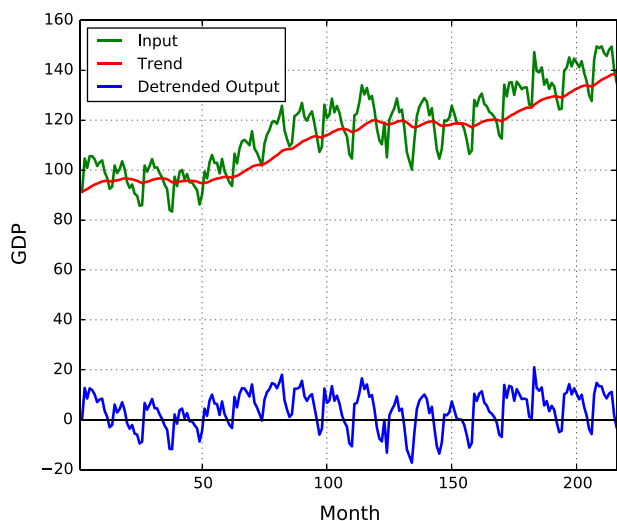


Fig. 1. Example of conventional detrending with Brazilian GDP. The detrended output is obtained by subtracting the trend from the original input. In this example, we use an exponential moving average (EMA) with a fixed decay factor of 0.95 to define the trend.

computations from previous time steps in advance for computing the current time step. Thus, finding a way to achieve faster learning convergence has been a barrier to practical development of ConvRNNs.

Ioffe and Szegedy (2015) argue that internal covariate shift is responsible for the increased training time in feed-forward neural networks, including multi-layer perceptrons (MLPs) and CNNs, and they suggest batch normalization (BN) to normalize the input distribution of a neuron for each mini-batch, as a way to reduce training time. BN successfully removes internal covariate shift, thereby significantly accelerating training with improved generalization, and this technique has become standard for training feed-forward neural networks. Some studies use BN with RNNs because unrolled RNNs over time can be seen as deep neural networks in terms of time as well as depth (Cooijmans, Ballas, Laurent, & Courville, 2017; Laurent, Pereyra, Brakel, Zhang, & Bengio, 2016). However, BN is incompatible with RNNs, regardless of computing global statistics along the time domain (Laurent et al. 2016) or local statistics at each time step (Cooijmans et al., 2017). Use of global statistics ignores statistics at each time step, but use of local statistics does not accommodate training sequences of variable lengths. As an alternative, layer normalization (LN) (Ba, Kiros, & Hinton, 2016) eliminates dependencies between mini-batch samples that obviate the use of BN with RNNs. LN computes statistics over all neurons in each layer and accelerates training of RNNs and MLPs, but not CNNs. Neither BN nor LN is generally applied to ConvRNNs.

The current paper focuses on the time domain in order to accelerate training of ConvRNNs. Much of time series analysis and many forecasting methods can be applied only to stationary time series. Detrending transforms non-stationary time series to stationary series by identifying the change as a trend and removing it. This method is straightforward, and is illustrated in the context of the Brazilian gross domestic product¹ in Fig. 1. The current research applies this method to normalize sequences of neurons in RNNs. Our key insight here is that the hidden state of a gated recurrent unit (GRU) (Cho et al., 2014) can be considered as a trend that can be approximated by the form of an exponential moving average with an adaptively changing decay factor. Based on this insight, we propose a novel temporal normalization method,

“adaptive detrending” (AD), for use with GRU and convolutional gated recurrent unit (ConvGRU), which is a variant of ConvRNNs extended from GRU. The implications of AD are fourfold:

- AD is easy to implement, reducing computational cost and consuming less memory than competing methods.
- AD eliminates temporal internal covariate shift.
- AD controls the degree of detrending (or normalization) through decay factor adaptability.
- AD is fully compatible with existing normalization methods.

2. Background

2.1. Batch normalization

Internal covariate shift slows training of deep neural networks, because the distribution of layer inputs changes continuously as lower layer parameters are updated. Batch normalization (BN) (Ioffe & Szegedy, 2015) has recently been proposed to reduce internal covariate shift by normalizing network activation as follows:

$$\mu = \frac{1}{m} \sum_{i=1}^m x_i \quad (1)$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2 \quad (2)$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad (3)$$

$$y_i = \gamma \hat{x}_i + \beta \quad (4)$$

where x is the activations of a neuron in a mini-batch of size m , μ and σ^2 are the mean and variance of a mini-batch, respectively, \hat{x} is normalized input, ϵ is an infinitesimal constant for numerical stability, and y is an affine transformation of normalized inputs \hat{x} . During training, the input distribution to a layer is transformed to a fixed distribution with a zero mean and unit variance, regardless of the change in parameters of lower layers. Additionally, an affine transformation with two learnable parameters γ and β follows normalization in order to recover the original activation when required. BN accelerates training and improves generalization of CNNs on ImageNet classification tasks.

Due to its success in feed-forward neural networks, BN has been applied to RNNs to speed training and improve generalization (Cooijmans et al., 2017; Laurent et al., 2016). In Laurent et al. (2016), BN is applied only to vertical (input-to-hidden) and not to horizontal (hidden-to-hidden) connections because the repeated rescaling of horizontal connections induces vanishing and exploding gradient problems. Also, the mean and variance for BN are computed by averaging along not only the mini-batch axis but also the time axis, which is called “sequence-wise normalization”. On the other hand, Cooijmans et al. (2017) develop “step-wise normalization” and show that (1) applying BN to horizontal as well as vertical connections is possible by properly initializing γ of an affine transformation and beneficial for reducing temporal internal covariate shift, and (2) using statistics for each time step separately preserves initial transient phase information. However, with this method, estimation of statistics at each time step degrades along the time axis due to variation in length of training and test sequences. During training, mini-batch configuration involves the use of zero, or last frame padding for shorter sequences. Furthermore, statistics for each time step are estimated only up to the length of the longest training sequence T_{max} . After training, accurate statistics for test sequences longer than the longest training sequence T_{max} cannot be generated. Due to these factors, performance suffers.

¹ http://www2.stat.duke.edu/~mw/data-sets/ts_data/brazil_econ.

Download English Version:

<https://daneshyari.com/en/article/6862887>

Download Persian Version:

<https://daneshyari.com/article/6862887>

[Daneshyari.com](https://daneshyari.com)