Visualizing and Understanding Deep Convolutional Neural **Network in Image-based Structural Health Monitoring**

Yuqing Gao and Khalid M. Mosalam

Department of Civil and Environmental Engineering, University of California Berkeley

Abstract

In this work, we emphasize the importance of understanding the trained deep convolutional neural network in the vision-based structural health monitoring (SHM) and investigate the applications of the recent widely-used visual interpretation methods, namely Guided backpropagation (GBP), Gradient-based class activation map (Grad-CAM) and Guided Grad-CAM, in the basic vision task of SHM, i.e., spalling detection.

Introduction

In the past decade, deep convolutional neural network (CNN) is increasingly attracting wide attention from researchers in many fields. It achieves the state-of-the-art performance and is proved to be far better than traditional methods. However, deep CNN is treated as a "black box," whose internal working principle is hard to understand by a human. This drawback significantly limits its usage in the real-world applications, e.g., vision-based structural health monitoring (SHM).

Visual Interpretation Methodologies

(4)

(5)

Guided Backpropagation (GBP)

$$R_{i,j}^{l} = \mathbb{1}(F_{i,j}^{l} > 0) \cdot \mathbb{1}(G_{i,j}^{l} > 0) \cdot G_{i,j}^{l}$$
(3)

Gradient-based Class Activation Map (Grad-CAM)



Figure 2 Formulation of restoring



As an early work in vision-based SHM, we introduce Guided backpropagation (GBP), Gradient-based class activation map (Grad-CAM) and Guided Grad-CAM, to understand the network behavior through visualizing the backward signal or activated features. We conduct a comprehensive exploration towards the performance of these interpretation approaches in basic SHM vision tasks.

Signal Flow in Deep CNN

In the forward pass of CNN, F^{l} denotes the output from the *l*-th layer before activation function. A^{l} represents the activated feature maps from F^{l} , e.g., $A^{l} = \text{ReLU}(F^{l})$, where only positive entries in F^{l} are passed to A^{l} . S is the class score vector computed by the deep CNN, and S^{c} is one entry of it for a particular class $c, c \in \{1, 2, ..., C\}$. In the backward pass, G^{l} denotes gradient of the output with respect to A^{l} in the *l*-th layer, i.e., $G^{l} = \frac{\partial S}{\partial A^{l}}$. When the signal further passes back through the activation function, the restoring gradient to bottom layers is denoted as R^{l} , where only positive gradient can be passed back. W_n^c represents the weight of the *n*-th feature map of A^l with respect to class c, as a scalar; $A_{i,j;n}^l$ is the entry with spatial location (i, j) in the *n*-th feature map of A^l . $\sum_{i,j} A^l_{i,j;n}$ is a shorthand notation of $\sum_{i} \sum_{j} A_{i,j;n}^{l}$ representing the global average pooling over a single feature map ignoring some constant multipliers.

$$P^l = \mathbb{1}(C^l > 0) \cdot C^l \tag{1}$$

$$\widetilde{\mathcal{W}_n^c} = \frac{1}{M} \sum_i \sum_j \frac{\partial S^c}{\partial A_{i,j:n}^l}$$

$$\mathcal{L}_{Grad-CAM}^{c} = \sum_{n} \widetilde{\mathcal{W}_{n}^{c}} \cdot A_{n}^{l}$$

Guided Grad-CAM

Up-sampling $L_{Grad-CAM}^{c}$ to the image size, and further elementwise multiplying with the GBP saliency map.





Acknowledgement

The author would like to acknowledge the funding support of TBSI.

