Explaining the Unexplained: A **CL**ass-Enhanced Attentive **R**esponse (CLEAR) Approach to Understanding Deep Neural Networks

Devinder Kumar*, Alexander Wong & Graham W. Taylor



Current Approaches – Heatmap/Attention based!







Deconvolution: Zeiler et.al. ECCV 14

Guided backpropagation: ICLR 2015





Saliency: Simonyan et.al. CVPR 2013



Deep Taylor Decomp. Montavon et. al. PR journal 2017



Prediction Difference: Zintgraf et. al. ICLR 2017

Interpretations



Class Enhanced Attentive Response (CLEAR) Map



Binary Heatmap



CLEAR Map



Class Enhanced Attentive Response (CLEAR)

$$\hat{h}_l = \sum_{k=1}^{K} z_{k,l} * w_{k,l}.$$

Output response at layer l

$$R_l = G_1 U_1 G_2 U_2 G_{l-1} U_{l-1} G_l z_l$$

Output response of layer l

$$R(\underline{x}|c) = G_1 U_1 G_2 U_2 \dots G_{L-1} U_{L-1} G_L^c z_L. \qquad \text{given} \quad c \ (1 \le c \le N)$$

$$\hat{C}(\underline{x}) = \operatorname*{argmax}_{c} R(\underline{x}|c).$$

Dominant Class Response

 $D_{\hat{C}}(\underline{x}) = R(\underline{x}|\hat{C}).$

Dominant Response

Class Enhanced Attentive Response (CLEAR) Map





MNIST RESULTS



Correctly Classified



Wrongly Classified



SVHN RESULTS



Correctly Classified







-

Accuracy(%)	MNIST	SVHN
Full image	99.26	92.60
with only strong features	79.89	69.12
without strong features	43.45	54.46

Stanford Dog Dataset Results



Stanford Dog Dataset Results



Stanford Dog Dataset Results

Chihuahua Jap. Spaniel Maltese Pekinese Shih-Tzu B. Spaniel Papillon Toy Terrier R. Ridgeback Afghan Hound







Conclusion

- Sparsity in the individual response maps from the last layer kernels : same pattern for all datasets considered.
- Evidence for classes tend to come from very specific localized regions.
- **CLEAR maps** enable the visualization of not only the areas of interest that predominantly influence the decision-making process, but also the **degree of influence** as well as the **dominant class** of influence in these areas.
- Showed efficacy of CLEAR maps both quantitatively and qualitatively.

Thank You!

devinder.kumar@uwaterloo.ca

http://devinderkumar.com

