

Title : Understanding Rectifiers

Explain how the main technique from ``Delving Deep into Rectifiers'' may alleviate a core issue in the training scheme for ``Going Deeper with Convolutions.''

Core issues in ‘Going Deeper with Convolutions’

Improving the performance of deep neural networks = increasing size and depth.

1. Increased size leads to larger number of parameters, making the network inclined to overfitting, creating major bottleneck between labeled examples and limited training set.
2. Increased size leads to dramatically increased use of computational resources ; waste of computation. (for 2 linked convolutional layers, increase in the number of their filters results in quadratic increase of computation.)

Definition of Rectifiers :

Rectified activation units : Rectified Linear activation function (Relu) :

Piecewise linear function that will output the input directly if it is positive, otherwise it will output zero.

Instead of increasing the dimension and depth of the system, applying PReLU (Parametric Rectified Linear Unit) can improve the model fitting with nearly zero extra computation cost and small overfitting risk - alleviating the core issues of increasing dimensions and depth in order to increase the performance of the system.

Explain the core effect of pre-activation from ``Identity Mappings in Deep Residual Networks'' compared to the original residual formulation.

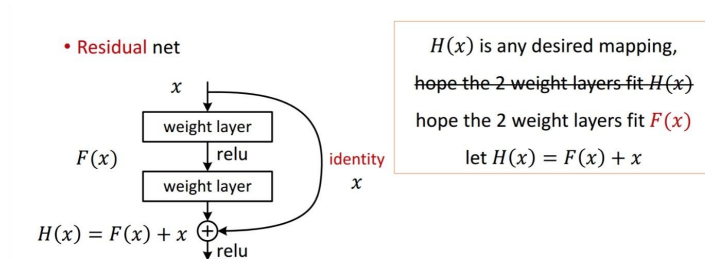


Illustration on Residual set.

Residual network

Optimizing $F(x)$, the residual mapping will result in getting desired $H(x)$.

Problem : Rapid degradation

If the depth of the network deepens, the accuracy gets saturated and degrades rapidly.

Continued

Residual unit: $y_i = h(x_i) + F(x_i, W_i)$

$$x_{L+1} = \text{ply}_i \rightarrow x_{L+1} = x_L + F(x_L, W_L)$$

identity mapping

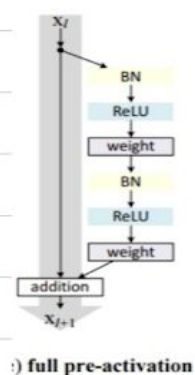
$$\Rightarrow x_L = x_L + \sum_{i=L}^{L-1} F(x_i, W_i)$$

Sum. Residual function

$$\frac{\partial \mathcal{E}}{\partial x_L} = \underbrace{\frac{\partial \mathcal{E}}{\partial x_L}}_{\text{direct propagation}} \underbrace{\frac{\partial x_L}{\partial x_L}}_{\text{propagation through multiple weight layers}} = \frac{\partial \mathcal{E}}{\partial x_L} \left(1 + \frac{\partial}{\partial x_L} \sum_{i=L}^{L-1} F(x_i, W_i) \right) \quad \mathcal{E} \text{ as Loss}$$

Question: Does \hat{f} identity mapping yields the best performance?
 the paper proves that setting \hat{f} as identity map composed of ReLU and BN without intercepting shortcut path.

Full pre activation

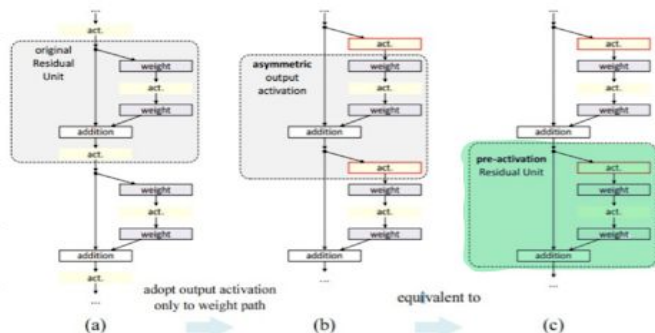


$$\rightarrow y_{L+1} = f(y_L) + F(f(y_L), W_{L+1})$$

impose activation only on F by introducing \hat{f}

$$\rightarrow y_{L+1} = f(y_L) + F(\hat{f}(y_L), W_{L+1})$$

$$= x_{L+1} = x_L + F(\hat{f}(x_L), W_L)$$



The resulting performance of system with full pre activation shows that the accuracy on training set is lower than that of the original model, the overall accuracy on the test was evidently higher.

- optimizing was achieved easily on full pre activation model
- Regularization is done by BN in full pre activation model.