

Problem

Driven by the significance of depth, a question arises: Is learning better networks as easy as stacking more layers? An obstacle to answering this question was the notorious problem of vanishing/exploding gradients [1, 9], which hamper convergence from the beginning

depth, problem of gradient with explicit

Sol1

normalized initialization [23, 9, 37, 13] and intermediate normalization layers

When deeper networks are able to start converging, a degradation problem has been exposed: with the network depth increasing, accuracy gets saturated

#Mianproblem

degradation is not caused by overfitting

There exists a solution by construction to the deeper model: the added layers are identity mapping, and the other layers are copied from the learned shallower model.

#Mainsol

Residual block

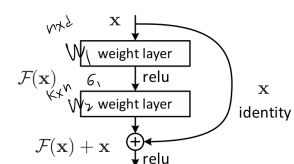


Figure 2. Residual learning: a building block.

$$F(x) = W_2(G_1(W_1(x)))$$

if $x \in \mathbb{R}^d$

$\mathbb{R}^d \xrightarrow{W_1} \mathbb{R}^n \xrightarrow{G_1} \mathbb{R}^n \xrightarrow{W_2} \mathbb{R}^k$

$\mathbb{R}^d \xrightarrow{\text{identity}} \mathbb{R}^k$

element-wise

underlying mapping as $H(x)$, we let the stacked nonlinear layers fit another mapping of $F(x) := H(x) - x$. The original mapping is recast into $F(x) + x$. We hypothesize that it is easier to optimize the residual mapping than to optimize the original, unreferenced mapping. To the extreme, if an identity mapping were optimal, it would be easier to push the residual to zero than to fit an identity mapping by a stack of nonlinear layers.

Formally, denoting the desired underlying mapping as $H(x)$, we let the stacked nonlinear layers fit another mapping of $F(x) := H(x) - x$. The original mapping is recast into $F(x) + x$. We hypothesize that it is easier to optimize the residual mapping than to optimize the original, unreferenced mapping. To the extreme, if an identity mapping were optimal, it would be easier to push the residual to zero than to fit an identity mapping by a stack of nonlinear layers

if identity mapping is nonoptimal.

If one hypothesizes that multiple nonlinear layers can asymptotically approximate complicated functions

they can asymptotically approximate the residual functions, i.e., $H(x) - x$ (assuming that the input and output are of the same dimensions)

The degradation problem suggests that the solvers might have difficulties in approximating identity mappings by multiple nonlinear layers. With the residual learning reformulation, if identity mappings are optimal, the solvers may simply drive the weights of the multiple nonlinear layers toward zero to approach identity mappings

#Mpoint

