

# TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2020

**Trainability:**

**Relu, Batch Normalization, Initialization,**

**and Residual Connections (ResNet)**

## Universality Assumption

We often assume DNNs are universally expressive (can model any function) and trainable (the desired function can be found by SGD).

Universal trainability is clearly false but can still usefully guide architecture design.

## Universality Assumption: Expressiveness

DNNs generalize digital circuits.

Consider Boolean Values  $P, Q$  — numbers that are either close to 0 or close to 1.

$$P \wedge Q \approx \sigma(100 * P + 100 * Q - 150)$$

$$P \vee Q \approx \sigma(100 * P + 100 * Q - 50)$$

$$\neg P \approx \sigma(100 * (1 - P) - 50)$$

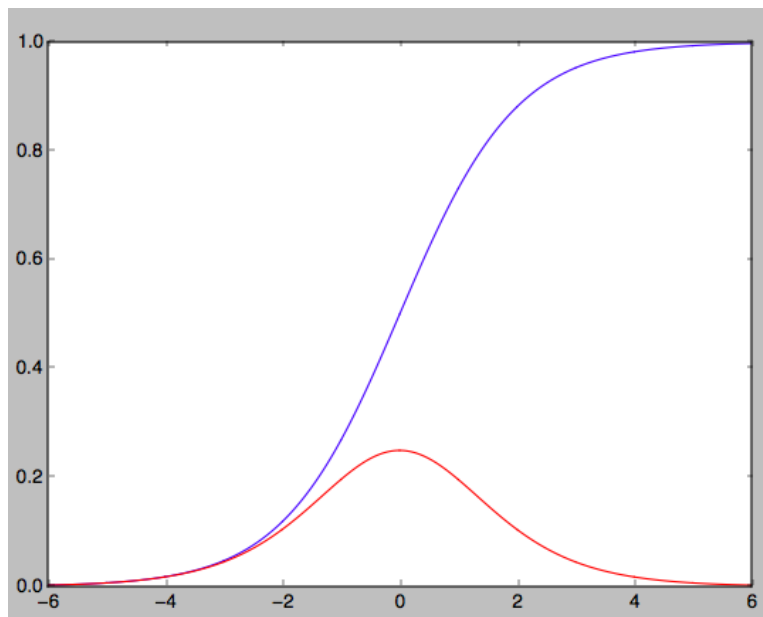
## Universality Assumption: Trainability

The main issue in making deep neural networks trainable is maintaining meaningful gradients.

There are various difficulties.

## Activation Function Saturation

Consider the sigmoid activation function  $1/(1 + e^{-x})$ .

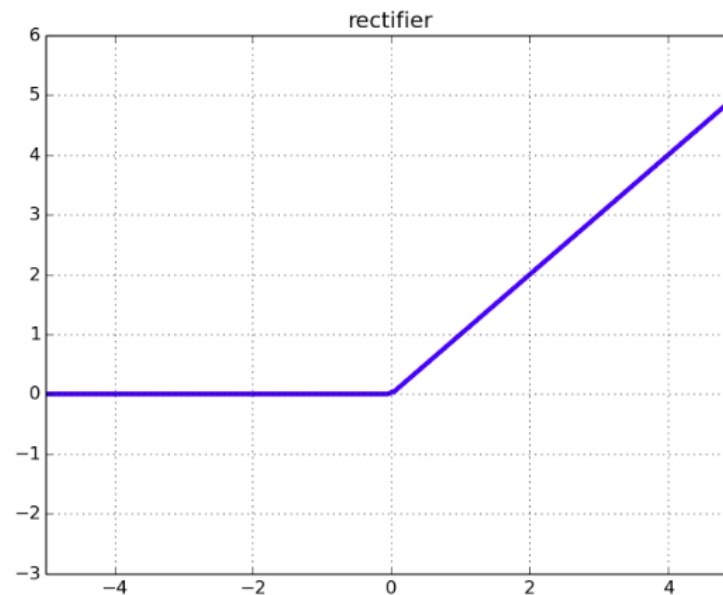


The gradient of this function is quite small for  $|x| > 4$ .

In deep networks backpropagation can go through many sigmoids and the gradient can “vanish”

# The Rectified Linear Unit Activation Function (Relu)

$$\text{Relu}(x) = \max(x, 0)$$



The activation function  $\text{Relu}(x)$  does not saturate for  $x > 0$ .

## Repeated Multiplication by Network Weights

Consider a deep CNN.

$$L_{i+1} = \text{Relu}(\text{Conv}(\Phi_i, L_i))$$

For  $i$  large,  $L_i$  has been multiplied by many weights.

If the weights are small then the neuron values, and hence the weight gradients, decrease exponentially with depth. **Vanishing Gradients.**

If the weights are large, and the activation functions do not saturate, then the neuron values, and hence the weight gradients, increase exponentially with depth. **Exploding Gradients.**

## Repeated Multiplication by Network Weights

The problem of repeated multiplication by network weights can be addressed with careful initialization.

We want an initialization for which the values stay in the active regions of the activation functions — zero mean and unit variance.



## Initialization

Consider a linear threshold unit

$$y[j] = \sigma(W[j, I]x[I] - B[j])$$

We want the scalar  $y[j]$  to have zero mean and unit variance.

Xavier initialization initializes  $B[j]$  to zero and randomly draws  $W[j, i]$  from a uniform distribution on  $\left(-\sqrt{3/I}, \sqrt{3/I}\right)$ .

Assuming  $x[i]$  has zero mean and unit variance, this gives zero mean and unit variance for  $W[j, I]x[I]$ .

## Batch Normalization

We can also enforce zero mean, unit variance, values dynamically with normalization layers.

In vision networks this is most commonly done with Batch Normalization.

## Batch Normalization

Given a tensor  $x[b, j]$  we define  $\tilde{x}[b, j]$  as follows.

$$\hat{\mu}[j] = \frac{1}{B} \sum_b x[b, j]$$

$$\hat{\sigma}[j] = \sqrt{\frac{1}{B-1} \sum_b (x[b, j] - \hat{\mu}[j])^2}$$

$$\tilde{x}[b, j] = \frac{x[b, j] - \hat{\mu}[j]}{\hat{\sigma}[j]}$$

At test time a single fixed estimate of  $\mu[j]$  and  $\sigma[j]$  is used.

## Spatial Batch Normalization

For CNNs we convert a tensor  $x[b, x, y, j]$  to  $\tilde{x}[b, x, y, j]$  as follows.

$$\hat{\mu}[j] = \frac{1}{BXY} \sum_{b,x,y} x[b, x, y, j]$$

$$\hat{\sigma}[j] = \sqrt{\frac{1}{BXY - 1} \sum_{b,x,y} (x[b, x, y, j] - \hat{\mu}[j])^2}$$

$$\tilde{x}[b, x, y, j] = \frac{x[b, x, y, j] - \hat{\mu}[j]}{\hat{\sigma}[j]}$$

## Adding an Affine Transformation

$$\check{x}[b, x, y, j] = \gamma[j]\tilde{x}[b, x, y, j] + \beta[j]$$

Here  $\gamma[j]$  and  $\beta[j]$  are parameters of the batch normalization.

This allows the batch normalization to learn an arbitrary affine transformation (offset and scaling).

It can even undo the normalization.

## Batch Normalization

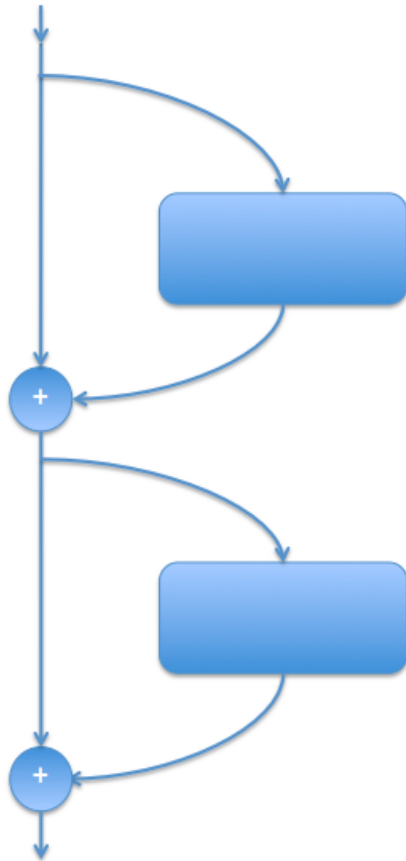
Batch Normalization appears to be generally useful in CNNs but is not always used.

Not so successful in RNNs.

It is typically used just prior to a nonlinear activation function.

It is intuitively justified in terms of “internal covariate shift”: as the inputs to a layer change the zero mean unit variance property underlying Xavier initialization are maintained.

## Residual Connections (ResNet)

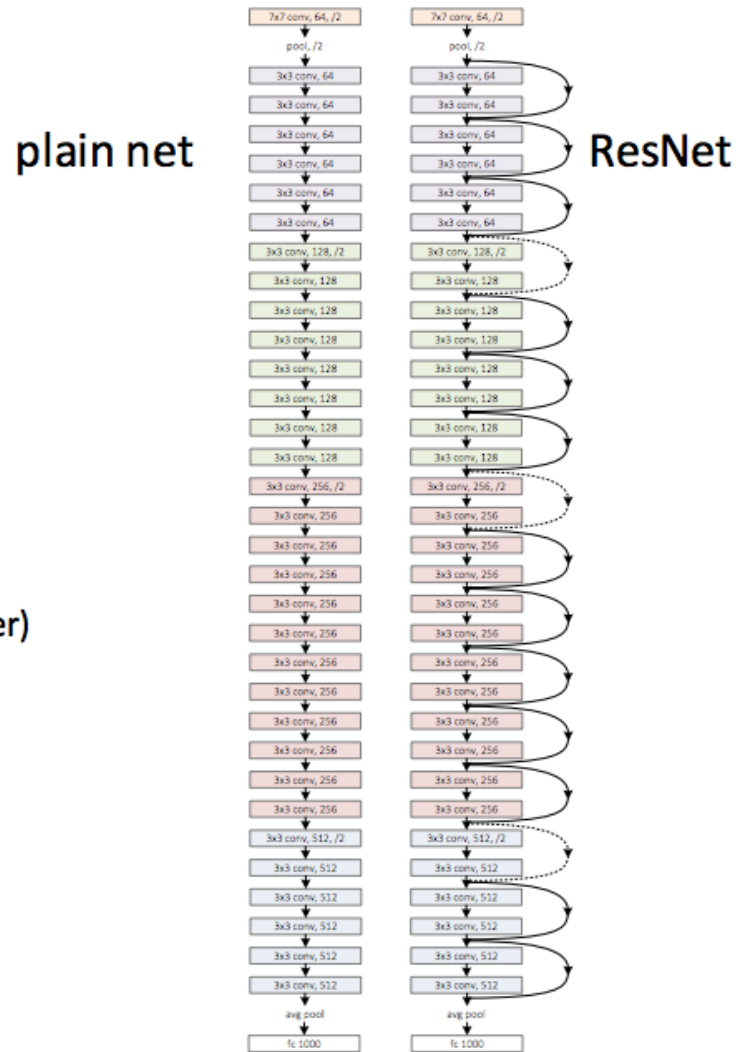


A residual connection produces the sum of the previous layer and the new layer.

The residual connection connects input to output directly and hence preserves gradients.

ResNets were introduced in late 2015 (Kaiming He et al.) and revolutionized computer vision.

# ResNet32



[Kaiming He]



## Simple Residual Skip Connections in CNNs (stride 1)

$$R_{\ell+1}[B, X, Y, J] = \text{Conv}(W_{\ell+1}[X, Y, J, J], B_{\ell+1}[J], L_{\ell}[B, X, Y, J])$$

for  $b, x, y, j$

$$L_{\ell+1}[b, x, y, j] = L_{\ell}[b, x, y, j] + R_{\ell+1}[b, x, y, j]$$

(Recall that we use capital letter indices to denote entire tensors and lower case letters for particular indices.)

## Simple Residual Skip Connections in CNNs (stride 1)

$$R_{\ell+1}[B, X, Y, J] = \text{Conv}(W_{\ell+1}[X, Y, J, J], B_{\ell+1}[J], L_{\ell}[B, X, Y, J])$$

for  $b, x, y, j$

$$L_{\ell+1}[b, x, y, j] = L_{\ell}[b, x, y, j] + R_{\ell+1}[b, x, y, j]$$

Note that in the above equations  $L_{\ell}[B, X, Y, J]$  and  $R_{\ell+1}[B, X, Y, J]$  are the same shape.

In the actual ResNet  $R_{\ell+1}$  is computed by two or three convolution layers.

## Handling Spacial Reduction

Consider  $L_{\ell}[B, X_{\ell}, Y_{\ell}, J_{\ell}]$  and  $R_{\ell+1}[B, X_{\ell+1}, Y_{\ell+1}, J_{\ell+1}]$

$$X_{\ell+1} = X_{\ell}/s$$

$$Y_{\ell+1} = Y_{\ell}/s$$

$$J_{\ell+1} \geq J_{\ell}$$

In this case we construct  $\tilde{L}_{\ell}[B, X_{\ell+1}, Y_{\ell+1}, J_{\ell+1}]$

$$\text{for } b, x, y, j \quad \tilde{L}_{\ell}[b, x, y, j] = \begin{cases} L_{\ell}[b, s * x, s * y, j] & \text{for } j < J_{\ell} \\ 0 & \text{otherwise} \end{cases}$$

$$L_{\ell+1}[B, X_{\ell+1}, Y_{\ell+1}, J_{\ell+1}] = \tilde{L}_{\ell}[B, X_{\ell+1}, Y_{\ell+1}, J_{\ell+1}] \\ + R_{\ell+1}[B, X_{\ell+1}, Y_{\ell+1}, J_{\ell+1}]$$

## Deeper Versions use Bottleneck Residual Paths

We reduce the number of features to  $K < J$  before doing the convolution.

$$U[B, X, Y, K] = \text{Conv}'(\Phi_{\ell+1}^A[1, 1, J, K], L_\ell[B, X, Y, J])$$

$$V[B, X, Y, K] = \text{Conv}'(\Phi_{\ell+1}^B[3, 3, K, K], U[B, X, Y, K])$$

$$R[B, X, Y, J] = \text{Conv}'(\Phi_{\ell+1}^R[1, 1, K, J], V[B, X, Y, K])$$

$$L_{\ell+1} = L_\ell + R$$

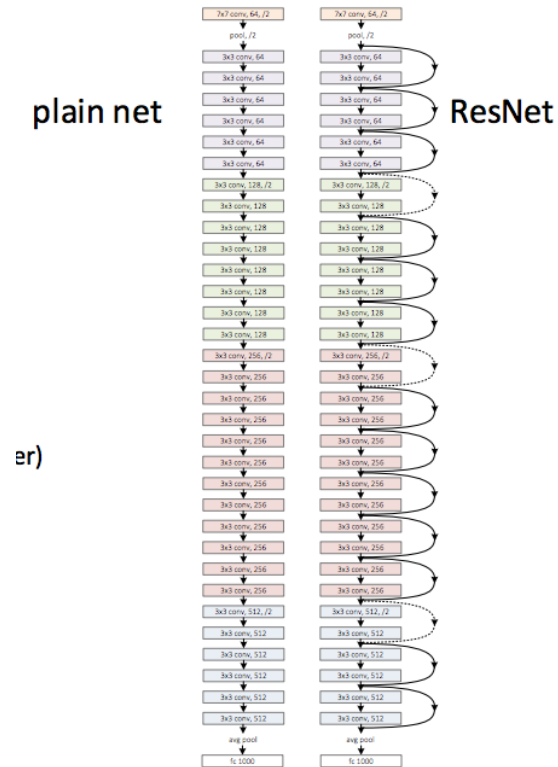
Here  $\text{CONV}'$  may include batch normalization and/or an activation function.

## A General Residual Connection

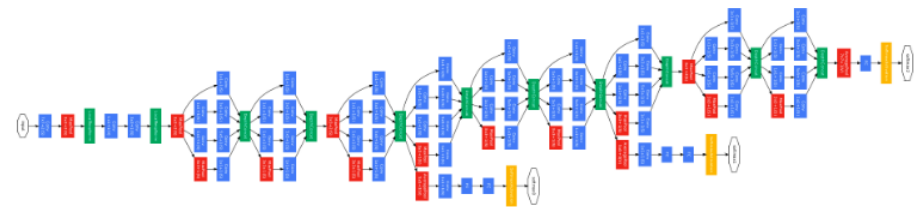
$$y = \tilde{x} + R(x)$$

Where  $\tilde{x}$  is either  $x$  or a version of  $x$  adjusted to match the shape of  $R(x)$ .

# ResNet Simplicity



[Kaiming He]



## **ResNet Power**

ResNet gives powerful image classification.

ResNet is used in folding proteins.

ResNet is the network used in AlphaZero for Go, Chess and Shogi.

Residual connections are now universal in all forms of deep models such as RNNs and Transformers in language processing.

**END**