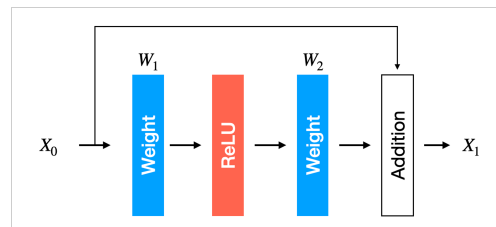


## 1. Residual Connection

Neural networks learn by optimization with gradients. The deeper a neural network grows, the more likely it will suffer from either vanishing gradients or exploding gradients problem since downstream layers (ie. earlier layers) will be optimized with either exponentially small or large gradients due to chain rule. Exploding gradients can be resolved using **gradient clipping**, whereas tackling vanishing gradients is more tricky. Here we study how skip connections can help.

- (a) One common way to tackle vanishing gradients is by adding residual connections. Fig 1 shows a simplified example of a residual block. Assuming that the weights are affine layer with zero biases, **what's the expression of  $X_1$  in terms of  $W_1, W_2$  and  $X_0$ ?**  $W_i \in \mathbb{R}^{n \times n}$ ,  $X_i \in \mathbb{R}^n$ .



**Figure 1:** Diagram for Residual Block

- (b) **Compute  $\frac{\partial X_1}{\partial X_0}$ .** Based on what you see numerically, **why does residual connection preserves gradient norms better?**
- (c) In practice, when transitioning between different layers or blocks within a ResNet architecture, the dimensions of the tensors may change. One common scenario is when the number of channels is doubled and the spatial dimensions are halved, often through a strided convolution. In such cases, the identity shortcut connection (the skip connection that bypasses one or more layers) also needs to change dimensions to match. **Describe how we should modify the expression for  $X_1$  in a residual block to incorporate a skip connection from  $X_0$ , when  $X_0$  does not match the shape of  $X_1$ ?**

## 2. Understanding Upsampling

Many CNN architectures rely on upsampling (particularly for tasks such as image segmentation). We will show how upsampling via linear interpolation can be done using convolutions and transpose convolutions. Let  $\mathbf{x} = [x_1, x_2, x_3]^T$  be a 1D discrete signal. Consider an upsampled signal  $\mathbf{y} = [y_1, y_2, y_3, y_4, y_5]^T$  that is a piecewise linear interpolation given by:

$$\mathbf{y} = \left[ x_1, \frac{x_1 + x_2}{2}, x_2, \frac{x_2 + x_3}{2}, x_3 \right]^T$$

- First, define an upsampled signal  $\tilde{\mathbf{x}} = [x_1, 0, x_2, 0, x_3]^T$  by inserting zeros between samples. **What is the convolutional kernel  $\mathbf{k}$  such that  $\mathbf{y}$  is the result of a convolution of  $\tilde{\mathbf{x}}$  with  $\mathbf{k}$  using stride 1 and zero-padding 1?**
- Now, rather than in two steps, we want to derive  $\mathbf{y}$  directly from  $\mathbf{x}$  in a single transpose convolution. **What is the kernel  $\mathbf{k}$  and stride  $s$  and padding  $p$  such that  $\mathbf{y}$  is the result of a transpose convolution of  $\mathbf{x}$  with  $\mathbf{k}$  using stride  $s$  and zero-padding  $p$ ?**

## 3. Regularization and Dropout

This question shows that a form of dropout is equivalent to ridge regularization.

Recall that linear regression optimizes the following learning objective:

$$\mathcal{L}(\mathbf{w}) = \|\mathbf{y} - X\mathbf{w}\|_2^2 \quad (1)$$

One way of using **dropout** during SGD on the  $d$ -dimensional input features  $\mathbf{f}_i$  is to keep each feature of each sample at random, with distribution  $\sim_{i.i.d} \text{Bernoulli}(p)$  (and zeroing it out if not kept), and then performing a SGD step.

It turns out that such dropout makes our learning objective effectively become

$$\mathcal{L}(\tilde{\mathbf{w}}) = E_{R \sim \text{Bernoulli}(p)} \left[ \|\mathbf{y} - (R \odot X)\tilde{\mathbf{w}}\|_2^2 \right] \quad (2)$$

where  $\odot$  is the element-wise product and the random binary matrix  $R \in \{0, 1\}^{n \times d}$  is such that  $R_{i,j} \sim_{i.i.d} \text{Bernoulli}(p)$ . We use  $\tilde{\mathbf{w}}$  to remind you that this is learned by dropout.

Recalling how Tikhonov-regularized (generalized ridge-regression) least-squares problems involve solving:

$$\mathcal{L}(\mathbf{w}) = \|\mathbf{y} - X\mathbf{w}\|_2^2 + \|\Gamma\mathbf{w}\|_2^2 \quad (3)$$

for some suitable matrix  $\Gamma$ , it turns out we can manipulate (2) to eliminate the expectations and get:

$$\mathcal{L}(\tilde{\mathbf{w}}) = \|\mathbf{y} - pX\tilde{\mathbf{w}}\|_2^2 + p(1-p)\|\tilde{\Gamma}\tilde{\mathbf{w}}\|_2^2 \quad (4)$$

with  $\tilde{\Gamma}$  being a diagonal matrix with  $\Gamma_{ii} = \sqrt{\sum_k X_{ki}^2} = \|\mathbf{f}_i\|_2$ .

- Suppose we have learned a weight vector  $\tilde{\mathbf{w}}$  using (4) (i.e. with dropout). How do we transform it into some  $\mathbf{w}$ , such that  $\mathbf{w}$  is a solution to the traditionally regularized problem (3)? And what would be the  $\Gamma$  matrix for this  $\mathbf{w}$ ? How does the choice of  $p$  affect  $\Gamma$ ?**

*(Hint: This is related to how we adjust weights learned using dropout training for using them at inference time. PyTorch by default does this adjustment during training itself, but here, we are doing dropout slightly differently with no adjustments during training.)*

*The question sounds weird, so here's an imaginary scenario to help make it easier to understand.*

*Suppose after hours of training on a dataset, we have learned a weight vector  $\tilde{\mathbf{w}}$  using dropout training, and then the boss tells us that we need to actually use Tikhonov-regularized training (but we are allowed to design whatever  $\Gamma$  we choose). How do we find  $\mathbf{w}$  without doing training anew? And how do we choose  $\Gamma$  so that the resulting  $\mathbf{w}$  actually solves (3)?*

- (b) Since  $\Gamma$  in (3) is an invertible matrix, we can define  $\tilde{X}, \tilde{w}$ , such that Equation (3)  $\mathcal{L}(\tilde{\mathbf{w}})$  has the same form as classical ridge regression:

$$\mathcal{L}(\tilde{\mathbf{w}}) = \|\mathbf{y} - \tilde{X}\tilde{\mathbf{w}}\|_2^2 + \lambda \|\tilde{\mathbf{w}}\|_2^2 \quad (5)$$

**What are  $\tilde{X}, \tilde{w}$  in terms of the original data matrix  $X$  and  $\Gamma$  and  $w$ ?**

- (c) Using the fact that  $\Gamma = \text{diag}(\|\mathbf{f}_1\|, \|\mathbf{f}_2\|, \dots, \|\mathbf{f}_d\|)$ , **what can you say about the norms of the columns of the effective training matrix  $\tilde{X}$ ? Comment briefly on the similarity between dropout and batch-normalization.**

### Contributors:

- Kevin Li .
- Naman Jain.
- Joey Hong.
- Saagar Sanghavi.
- Jerome Quenum.
- Anant Sahai.