

Squeeze-Excitation Networks

An easy to implement attention mechanism

Jacques Nel

September 20, 2020

1 Introduction

The convolutional operator, used in CNNs, has an essentially local receptive field. Recently, there have been some spectacular results achieved by neural network architectures that rely on attention mechanisms. The transformer, for example, consists entirely of an attention mechanism, and manages to achieve state-of-the-art results on NLP tasks.¹

The *Squeeze-and-Excitation (SE)* block allows the model to model complex inter-channel dynamics, and allows a global flow of information, similar to an attention mechanism.² SE blocks are simple to implement, and can simply be added to any pre-existing model architecture. SE blocks can significantly improve a model's performance, at a marginal additional computational cost.

2 SE architecture

The SE block is build around a convolutional operator. Let $\mathbf{X} \in \mathbb{R}^{H' \times W' \times C'}$ denote the input signal, with height, width, and channel depth being H' , W' , and C' respectively.

$$\mathbf{F}_{tr} : \mathbf{X} \in \mathbb{R}^{H' \times W' \times C'} \rightarrow \mathbf{U} \in \mathbb{R}^{H \times W \times C}$$

Take \mathbf{F}_{tr} to be a convolutional operator with learned filter kernels $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_C]$. In other words

$$\mathbf{u}_c = \mathbf{v}_c * \mathbf{X} = \sum_{s=1}^{C'} \mathbf{v}_c^s * \mathbf{x}^s.$$

2.1 Squeeze operation

Each filter map operates with a local receptive field and can not utilize information outside of its region. First, a global channel statistic is generated by a *squeeze* operation, consisting of an average pooling layer.

¹Vaswani et al., "Attention Is All You Need", NIPS 2017.

²Hu et al., "Squeeze-and-Excitation Networks".

$$z_c = \mathbf{F}_{sq}(\mathbf{u}_c) = \frac{1}{HC} \sum_{i=1}^H \sum_{j=1}^W u_c(i, j).$$

2.1.1 Excitation

Next, an *excitation* layer allows channel-wise dependence to be learned on the output of the *squeeze* operation. This simply consists of a simple gating mechanism with two fully connected layers, a sigmoid, and ReLu activation function.

$$\mathbf{s} = \mathbf{F}_{ex}(\mathbf{z}, \mathbf{W}_1, \mathbf{W}_2) = \sigma(\mathbf{W}_2 \text{ReLu}(\mathbf{W}_1 \mathbf{z})),$$

where $\mathbf{W}_1 \in \mathbb{R}^{\frac{C}{\rho} \times C}$ and $\mathbf{W}_2 \in \mathbb{R}^{C \times \frac{C}{\rho}}$. The reduction ratio ρ , preferably a divisor of C , allows one to control the number of parameters. The *excitation* operator acts as a bottleneck. In effect, the SE block acts as a self-attention mechanism, acting on a wider domain than the receptive field of the individual convolution filters.

Finally, the output of the SE block is produced by scaling \mathbf{U} with the signal \mathbf{s} .

$$\tilde{\mathbf{x}}_c = \mathbf{F}_{scale}(\mathbf{u}_c, s_c) = s_c \mathbf{u}_c.$$

We also use a skip connection, that bypasses the SE block entirely.

3 Implementation

The idea of an SE block is extremely simple to implement.

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from torch.optim import Adam
5 from torch.utils.data.dataloader import DataLoader
6
7
8 class SeBlock2d(nn.Module):
9
10     def __init__(self,
11                 channels,
12                 ratio=2):
13
14         super().__init__()
15
16         self.conv = nn.Conv2d(in_channels=channels,
17                               out_channels=channels,
18                               kernel_size=3,
19                               padding=1)
20
21         self.gate_1 = nn.Linear(in_features=channels,
22                                 out_features=channels // ratio)
23         self.gate_2 = nn.Linear(in_features=channels // ratio,
24                                 out_features=channels)
25
26     def forward(self, x):
27
28         # Generate feature maps U.
29         u = self.conv(x)
30         h, w = x.shape[-2:]
31
32         # Apply global average pooling to calculate channel descriptor z.
33         z = 1. / (h + w) * torch.sum(u, dim=(-2, -1))
34
35         # Pass signal through two fully connected layers of excitation operator.
36         s = torch.sigmoid(self.gate_2(F.relu(self.gate_1(z))))
37         s = s.unsqueeze(-1).unsqueeze(-1)
38
39         # Scale by activation s.
40         x_out = s * u
41
42     return x_out
```

Figure 1: Squeeze-and-Excitation architecture

