



(<http://lib.csdn.net/base/deeplearning>)

深度学习 (<http://lib.csdn.net/base/deeplearning>) -

深度学习训练以及预测优化 (<http://lib.csdn.net/deeplearning/node/750>) -

预测加速&模型压缩 (<http://lib.csdn.net/deeplearning/knowledge/1741>)

👁 440 💬 22

## Deep Compression阅读理解及Caffe源码修改

作者：may0324 (<http://my.csdn.net/may0324>)

最近又转战CNN模型压缩了。。。 (我真是一年换N个坑的节奏)，阅读了HanSong的15年16年几篇比较有名的论文，启发很大，这篇主要讲一下Deep Compression那篇论文，因为需要修改caffe源码，但网上没有人po过，这里做个第一个吃螃蟹的人，记录一下对这篇论文的理解和源码修改过程，方便日后追本溯源，同时如果有什么纰漏也欢迎指正，互相交流学习。

这里就从Why-How-What三方面来讲讲这篇文章。

### Why

首先讲讲为什么CNN模型压缩刻不容缓，我们可以看看这些有名的caffe模型大小：

1. LeNet-5 1.7MB
2. AlexNet 240MB
3. VGG-16 552MB

LeNet-5是一个简单的手写数字识别网络，AlexNet和VGG-16则用于图像分类，刷新了ImageNet竞赛的成绩，但是就其模型尺寸来说，根本无法移植到手机端App或嵌入式芯片当中，就算是想通过网络传输，较高的带宽占用率也让很多用户望尘莫及。另一方面，大尺寸的模型也对设备功耗和运行速度带来了巨大的挑战。随着深度学习的不断普及和caffe, tensorflow, torch等框架的成熟，促使越来越多的学者不用过多地去花费时间在代码开发上，而是可以毫无顾及地不断设计加深网络，不断扩充数据，不断刷新模型精度和尺寸，但这样的模型距离实用却仍是望其项背。

在这样的情形下，模型压缩则成为了亟待解决的问题，其实早期也有学者提出了一些压缩方法，比如weight prune(权值修剪)，权值矩阵SVD分解等，但压缩率也只是冰山一角，远不能令人满意。今年standford的HanSong的ICLR的一篇论文Deep Compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding一经提出，就引起了巨大轰动，在这篇论文工作中，他们采用了3步，在不损失（甚至有提升）原始模型精度的基础上，将VGG和Alexnet等模型压缩到了原来的35~49倍，使得原本上百兆的模型压缩到不到10M，令深度学习模型在移动端等的实用成为可能。

### How

Deep Compression 的实现主要有三步，如下图所示：

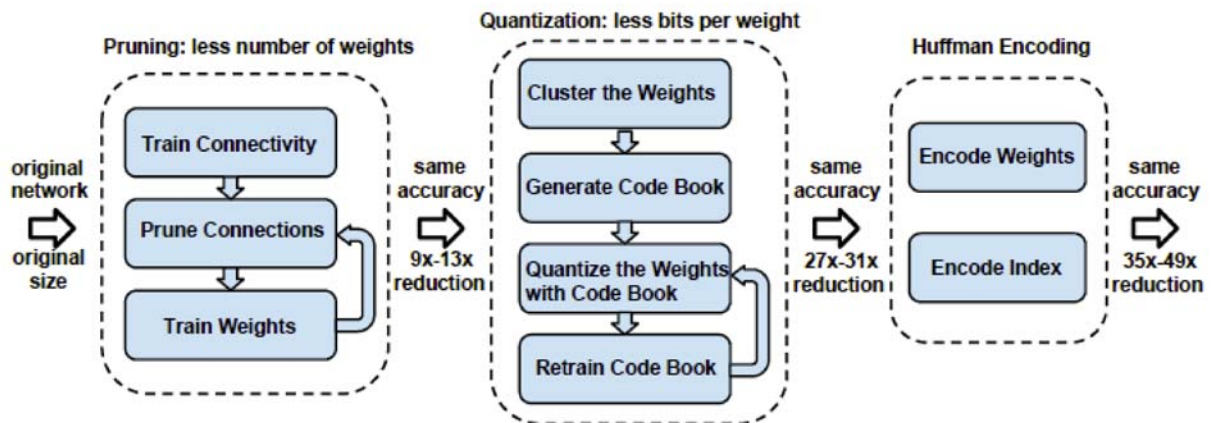


Figure 1: The three stage compression pipeline: pruning, quantization and Huffman coding. Pruning reduces the number of weights by  $10\times$ , while quantization further improves the compression rate: between  $27\times$  and  $31\times$ . Huffman coding gives more compression: between  $35\times$  and  $49\times$ . The compression rate already included the meta-data for sparse representation. The compression scheme doesn't incur any accuracy loss.

包括Pruning（权值修剪），Quantization（权值共享和量化），Huffman Coding（Huffman编码）。

### 1. Pruning

如果你调试过caffe模型，观察里面的权值，会发现大部分权值都集中在-1~1之间，即非常小，另一方面，神经网络的提出就是模仿人脑中的神经元突触之间的信息传导，因此这数量庞大的权值中，存在着不可忽视的冗余性，这就为权值修剪提供了根据。pruning可以分为三步：

- step1. 正常训练模型得到网络权值；
- step2. 将所有低于一定阈值的权值设为0；
- step3. 重新训练网络中剩下的非零权值。

经过权值修剪后的稀疏网络，就可以用一种紧凑的存储方式CSC或CSR（compressed sparse column or compressed sparse row）来表示。这里举个栗子来解释下什么是CSR

假设有一个原始稀疏矩阵A

$$A = \begin{bmatrix} 4.0 & 1.0 & 0.0 & 0.0 & 2.5 \\ 0.0 & 4.0 & 1.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 4.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 & 4.0 & 0.0 \\ 2.5 & 0.0 & 0.0 & 0.5 & 4.0 \end{bmatrix}$$

CSR可以将原始矩阵表达为三部分，即AA,JA,IC

$$\begin{aligned} AA &= 4.0 & 1.0 & 2.5 & 4.0 & 1.0 & 1.0 & 4.0 & 1.0 & 4.0 & 2.5 & 0.5 & 4.0 \\ JA &= 1 & 4 & 6 & 9 & 11 & 14 \\ IC &= 1 & 2 & 5 & 2 & 3 & 2 & 3 & 5 & 3 & 4 & 1 & 4 & 5 \end{aligned}$$

其中，AA是矩阵A中所有非零元素，长度为a，即非零元素个数；

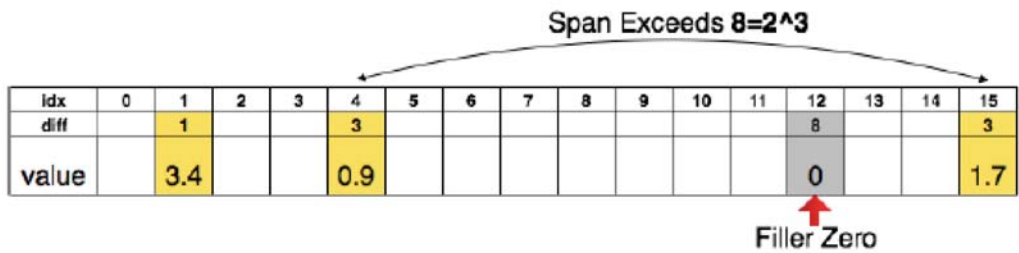
JA是矩阵A中每行第一个非零元素在AA中的位置，最后一个元素是非零元素数加1，长度为n+1，n是矩阵A的行

数；

IC是AA中每个元素对应的列号，长度为a。

所以将一个稀疏矩阵转为CSR表示，需要的空间为 $2 \cdot a + n + 1$ 个，同理CSC也是类似。

可以看出，为了达到压缩原始模型的目的，不仅需要在保持模型精度的同时，prune掉尽可能多的权值，也需要减少存储元素位置index所带来的额外存储开销，故论文中采用了存储index difference而非绝对index来进一步压缩模型，如下图所示：



其中，第一个非零元素的存储的是他的绝对位置，后面的元素依次存储的是与前一个非零元素的索引差值。在论文中，采用固定bit来存储这一差值，以图中表述为例，如果采用3bit，则最大能表述的差值为8，当一个非零元素距其前一个非零元素位置超过8，则将该元素值置零。（这一点其实也很好理解，如果两个非零元素位置差很多，也即中间有很多零元素，那么将这一元素置零，对最终的结果影响也不会很大）

做完权值修剪这一步后，AlexNet和VGG-16模型分别压缩了9倍和13倍，表明模型中存在着较大的冗余。

## 2.Weight Shared & Quantization

为了进一步压缩网络，考虑让若干个权值共享同一个权值，这一需要存储的数据量也大大减少。在论文中，采用kmeans算法来将权值进行聚类，在每一个类中，所有的权值共享该类的聚类质心，因此最终存储的结果就是一个码书和索引表。

### 1.对权值聚类

论文中采用kmeans聚类算法，通过优化所有类内元素到聚类中心的差距（within-cluster sum of squares）来确定最终的聚类结果：

$$\arg \min_C \sum_{i=1}^k \sum_{w \in c_i} |w - c_i|^2$$

式中， $W = \{w_1, w_2, \dots, w_n\}$ 是n个原始权值， $C = \{c_1, c_2, \dots, c_k\}$ 是k个聚类。

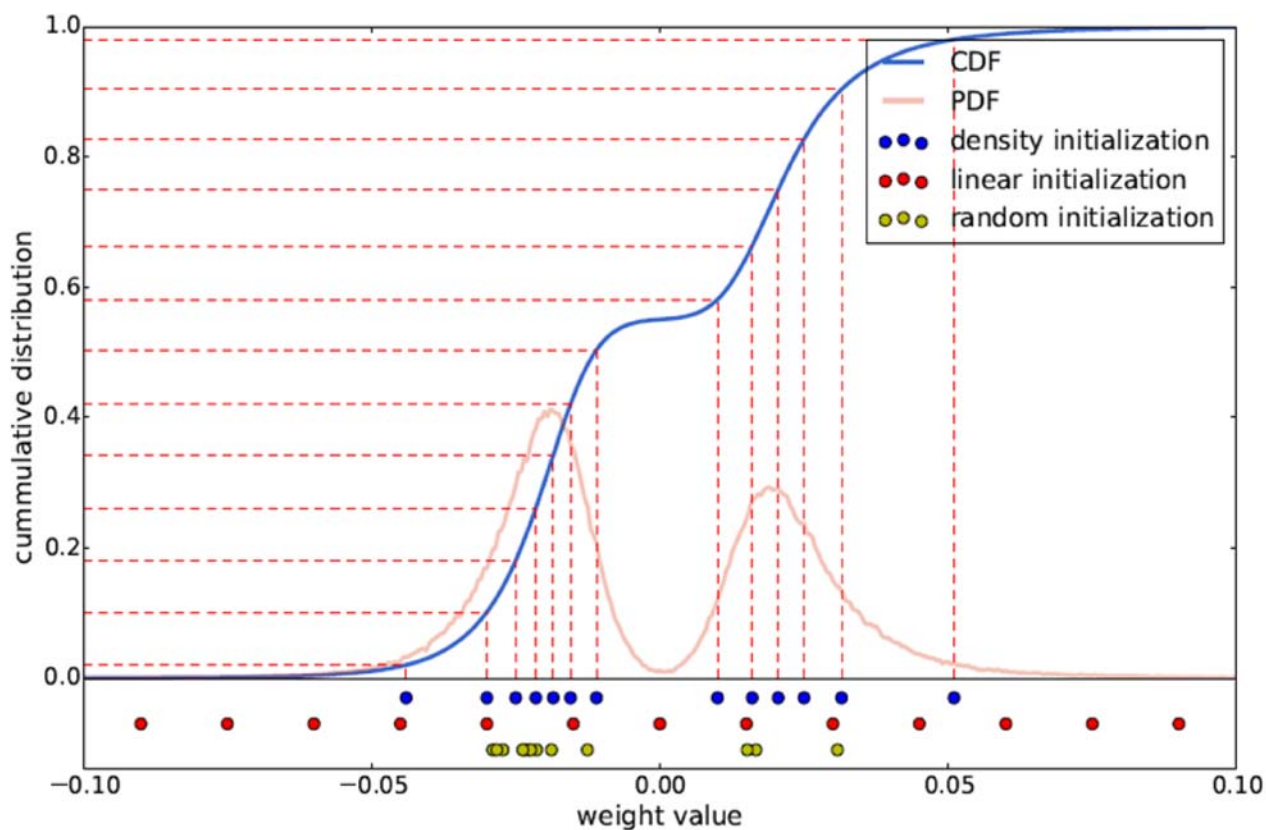
需要注意的是聚类是在网络训练完毕后做的，因此聚类结果能够最大程度地接近原始网络权值分布。

### 2. 聚类中心初始化

常用的初始化方式包括3种：

- a) 随机初始化。即从原始数据中随机产生k个观察值作为聚类中心。
- b) 密度分布初始化。现将累计概率密度CDF的y值分布线性划分，然后根据每个划分点的y值找到与CDF曲线的交点，再找到该交点对应的x轴坐标，将其作为初始聚类中心。
- c) 线性初始化。将原始数据的最小值到最大值之间的线性划分作为初始聚类中心。

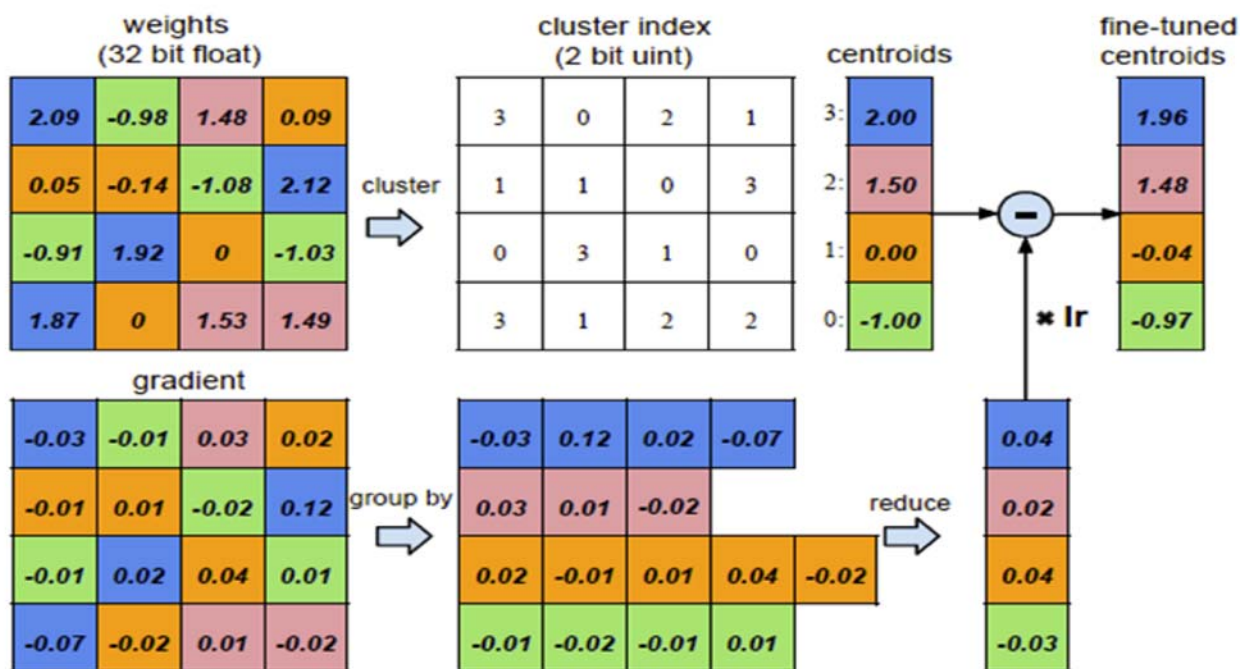
三种初始化方式的示意图如下所示：



由于大权值比小权值更重要（参加HanSong15年论文），而线性初始化方式则能更好地保留大权值中心，因此文中采用这一方式，后面的实验结果也验证了这个结论。

### 3. 前向反馈和后项传播

前向时需要将每个权值用其对应的聚类中心代替，后向计算每个类内的权值梯度，然后将其梯度和反传，用来更新聚类中心，如图：



共享权值后，就可以用一个码书和对应的index来表征。假设原始权值用32bit浮点型表示，量化区间为256，即



8bit，共有n个权值，量化后需要存储n个8bit索引和256个聚类中心值，则可以计算出压缩率compression ratio:  
 $r = 32 * n / (8 * n + 256 * 32) \approx 4$   
可以看出，如果采用8bit编码，则至少能达到4倍压缩率。

3.Huffman Coding

Huffman 编码是最后一步，主要用于解决编码长短不一带来的冗余问题。因为在论文中，作者针对卷积层统一采用8bit编码，而全连接层采用5bit，所以采用这种熵编码能够更好地使编码bit均衡，减少冗余。

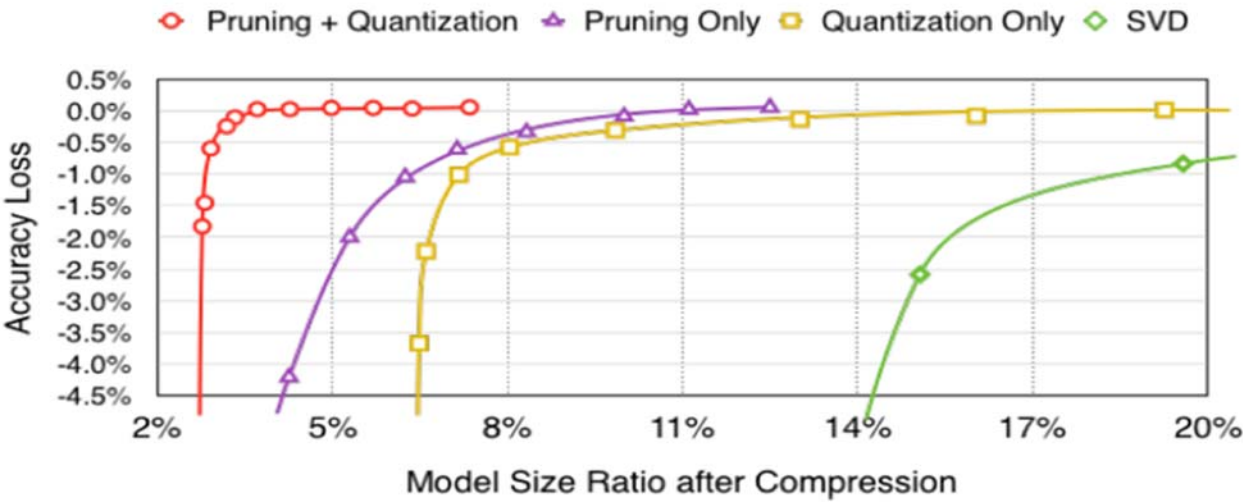
4.Evaluation

实验结果就是能在保持精度不变（甚至提高）的前提下，将模型压缩到前所未有的小。直接上图有用数据说话。

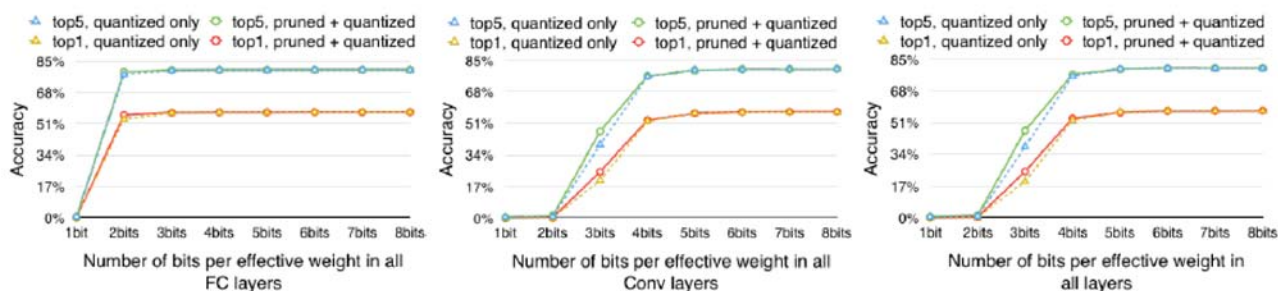
Network	Top-1 Error	Top-5 Error	Parameters	Compress Rate
LeNet-300-100 Ref	1.64%	-	1070 KB	
LeNet-300-100 Compressed	1.58%	-	<b>27 KB</b>	<b>40×</b>
LeNet-5 Ref	0.80%	-	1720 KB	
LeNet-5 Compressed	0.74%	-	<b>44 KB</b>	<b>39×</b>
AlexNet Ref	42.78%	19.73%	240 MB	
AlexNet Compressed	42.78%	19.70%	<b>6.9 MB</b>	<b>35×</b>
VGG-16 Ref	31.50%	11.32%	552 MB	
VGG-16 Compressed	31.17%	10.91%	<b>11.3 MB</b>	<b>49×</b>

Network	Top-1 Error	Top-5 Error	Parameters	Compress Rate
Baseline Caffemodel (BVLC)	42.78%	19.73%	240MB	1×
Fastfood-32-AD (Yang et al., 2014)	41.93%	-	131MB	2×
Fastfood-16-AD (Yang et al., 2014)	42.90%	-	64MB	3.7×
Collins & Kohli (Collins & Kohli, 2014)	44.40%	-	61MB	4×
SVD (Denton et al., 2014)	44.02%	20.56%	47.6MB	5×
Pruning (Han et al., 2015)	42.77%	19.67%	27MB	9×
Pruning+Quantization	42.78%	19.70%	8.9MB	27×
<b>Pruning+Quantization+Huffman</b>	<b>42.78%</b>	<b>19.70%</b>	<b>6.9MB</b>	<b>35×</b>

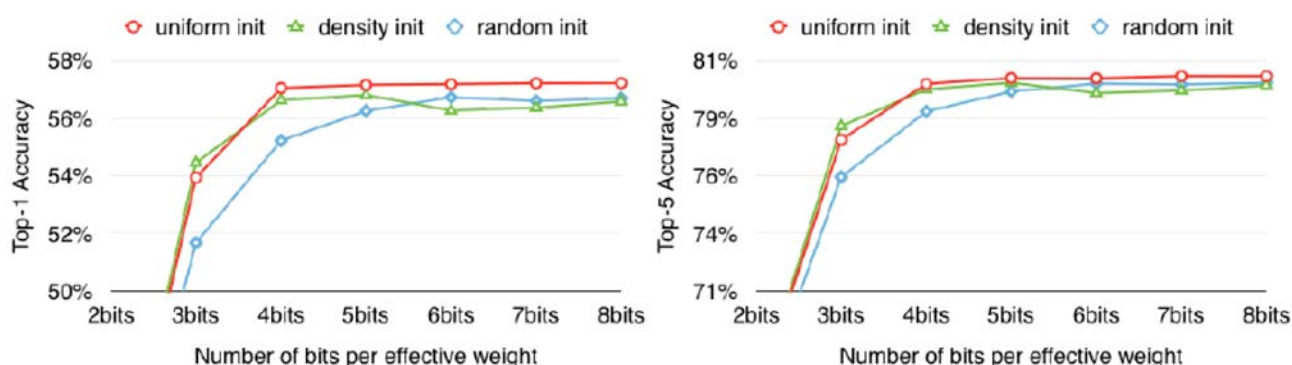
5.Discussion



不同模型压缩比和精度的对比，验证了pruning和quantization一块做效果最好。



不同压缩bit对精度的影响，同时表明conv层比fc层更敏感，因此需要更多的bit表示。



不同初始化方式对精度的影响，线性初始化效果最好。

#CONV bits / #FC bits	Top-1 Error	Top-5 Error	Top-1 Error Increase	Top-5 Error Increase
32bits / 32bits	42.78%	19.73%	-	-
8 bits / 5 bits	42.78%	19.70%	0.00%	-0.03%
8 bits / 4 bits	42.79%	19.73%	0.01%	0.00%
4 bits / 2 bits	44.77%	22.33%	1.99%	2.60%

卷积层采用8bit，全连接层采用5bit效果最好。

## What

此部分讲一讲修改caffe源码的过程。其实只要读懂了文章原理，修改起来很容易。

对pruning过程来说，可以定义一个mask来“屏蔽”修剪掉的权值，对于quantization过程来说，需定义一个indice来存储索引号，以及一个centroid结构来存放聚类中心。

在include/caffe/layer.hpp中为Layer类添加以下成员变量：

```
vector<int> masks_;
vector<int> indices_;
vector<Dtype> centroids_;
```

以及成员函数：

```
virtual void ComputeBlobMask(float ratio) {}
/**
```

由于只对卷积层和全连接层做压缩，因此，只需修改这两个层的对应函数即可。

在include/caffe/layers/base\_conv\_layer.hpp添加成员函数

```
virtual void ComputeBlobMask(float ratio) {}
```

这两处定义的函数都是基类的虚函数，不需要具体实现。

在include/caffe/layers/conv\_layer.hpp中添加成员函数声明：

```
virtual void ComputeBlobMask(float ratio);  
};
```

类似的，在include/caffe/layers/inner\_product\_layer.hpp也添加该函数声明。

在src/caffe/layers/conv\_layer.cpp 添加该函数的声明，用于初始化mask和对权值进行聚类。

```
9 template <typename Dtype>  
10 void ConvolutionLayer<Dtype>::ComputeBlobMask(float ratio)  
11 {  
12     // LOG(INFO) << "conv blob mask" << endl;  
13     int count = this->blobs_[0]->count();  
14     this->masks_.resize(count);  
15     //this->dmasks_ = new Dtype[count];  
16     //this->indices_.resize(count);  
17     this->indices_.resize(count);  
18     this->centroids_.resize(CONV_QUNUM);  
19     //calculate min max value of weight  
20     const Dtype* weight = this->blobs_[0]->cpu_data();  
21     Dtype min_weight = weight[0], max_weight = weight[0];  
22     vector<Dtype> sort_weight(count);  
23     for (int i = 0; i < count; ++i)  
24     {  
25         //this->masks_[i] = 1; //initialize  
26         sort_weight[i] = fabs(weight[i]);  
27     }  
28     sort(sort_weight.begin(), sort_weight.end());  
29     max_weight = sort_weight[count - 1];  
30     //cout << sort_weight[0] << " " << sort_weight[count - 1] << endl;  
31     int index = int(count*ratio); //int(count*(1- max_weight));  
32     Dtype thr;  
33     Dtype* muweight = this->blobs_[0]->mutable_cpu_data();  
34     float rat = 0;  
35     if (index > 0) {  
36         //thr = ratio;  
37         thr = sort_weight[index-1];  
38         LOG(INFO) << "CONV THR: " << thr << " " << ratio << endl;  
39         for (int i = 0; i < count; ++i)  
40         {  
41             this->masks_[i] = ((weight[i] >= thr || weight[i] < -thr) ? 1 : 0);  
42             //this->masks_[i] = (weight[i] > thr ? 1 : 0); // (weight[i] == 0 ? 0 : 1); // (weight[i] > thr ? 1 : 0);  
43             muweight[i] *= this->masks_[i];  
44             //this->dmasks_[i] = this->masks_[i];  
45             rat += (1-this->masks_[i]);  
46         }  
47     }  
48     else {  
49         for (int i = 0; i < count; ++i)  
50         {  
51             this->masks_[i] = (weight[i] == 0 ? 0 : 1); //keep unchanged  
52             rat += (1-this->masks_[i]);  
53         }  
54     }  
55     LOG(INFO) << "sparsity: " << rat/count << endl;  
56     min_weight = sort_weight[index];  
57     int ncentroid = CONV_QUNUM;  
58     kmeans_cluster(this->indices_, this->centroids_, muweight, count, this->masks_, /* max_weight, min_weight,*/ ncentroid, 1000);  
59 }  
60 }
```

同时，修改前向和后向函数。

在前向函数中，需要将权值用其聚类中心表示，红框部分为添加部分：

```
82 template <typename Dtype>  
83 void ConvolutionLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,  
84 const vector<Blob<Dtype>*>& top) {  
85     Dtype* muweight = this->blobs_[0]->mutable_cpu_data();  
86     int count = this->blobs_[0]->count();  
87     for (int i = 0; i < count; ++i)  
88     {  
89         if (this->masks_[i])  
90             muweight[i] = this->centroids_[this->indices_[i]];  
91     }  
92     const Dtype* weight = this->blobs_[0]->cpu_data();  
93     for (int i = 0; i < bottom.size(); ++i) {  
94         const Dtype* bottom_data = bottom[i]->cpu_data();  
95         Dtype* top_data = top[i]->mutable_cpu_data();  
96         for (int n = 0; n < this->num_; ++n) {  
97             this->forward_cpu_gemm(bottom_data + n * this->bottom_dim_, weight,  
98 top_data + n * this->top_dim_,  
99 if (this->bias_term_) {  
100 const Dtype* bias = this->blobs_[1]->cpu_data();  
101 this->forward_cpu_bias(top_data + n * this->top_dim_, bias);  
102 }  
103 }  
104 }  
105 }  
106 }
```

在后向函数中，需要添加两部分，一是对mask为0，即屏蔽掉的权值不再进行更新，即将其weight\_diff设为0，另一个则是统计每一类内的梯度差值均值，并将其反传回去，红框内为添加部分。

```

107 template <typename Dtype>
108 void ConvolutionLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
109 const vector<bool>& propagate_down, const vector<Blob<Dtype>*>& bottom) {
110     const Dtype* weight = this->blobs_[0]->cpu_data();
111     Dtype* weight_diff = this->blobs_[0]->mutable_cpu_diff();
112     int count = this->blobs_[0]->count();
113     for (int i = 0; i < top.size(); ++i) {
114         const Dtype* top_diff = top[i]->cpu_diff();
115         const Dtype* bottom_data = bottom[i]->cpu_data();
116         Dtype* bottom_diff = bottom[i]->mutable_cpu_diff();
117         // Bias gradient, if necessary.
118         if (this->bias_term_ && this->param_propagate_down_[1]) {
119             Dtype* bias_diff = this->blobs_[1]->mutable_cpu_diff();
120             for (int n = 0; n < this->num_; ++n) {
121                 this->backward_cpu_bias(bias_diff, top_diff + n * this->top_dim_);
122             }
123         }
124         if (this->param_propagate_down_[0] || propagate_down[i]) {
125             for (int n = 0; n < this->num_; ++n) {
126                 // gradient w.r.t. weight. Note that we will accumulate diffs.
127                 if (this->param_propagate_down_[0]) {
128                     this->weight_cpu_gemm(bottom_data + n * this->bottom_dim_,
129                         top_diff + n * this->top_dim_, weight_diff);
130                 }
131                 Dtype* weight_diff = this->blobs_[0]->mutable_cpu_diff();
132                 for (int j = 0; j < count; ++j) {
133                     weight_diff[j] *= this->masks_[j];
134                     vector<Dtype> tmpDiff(CONV_QUNUM);
135                     vector<int> freq(CONV_QUNUM);
136                     for (int j = 0; j < count; ++j) {
137                         {
138                             if (this->masks_[j])
139                             {
140                                 tmpDiff[this->indices_[j]] += weight_diff[j];
141                                 freq[this->indices_[j]]++;
142                             }
143                         }
144                     }
145                     for (int j = 0; j < count; ++j) {
146                         if (this->masks_[j])
147                             weight_diff[j] = tmpDiff[this->indices_[j]] / freq[this->indices_[j]];
148                     }
149                 }
150                 // gradient w.r.t. bottom data, if necessary.
151                 if (propagate_down[i]) {
152                     this->backward_cpu_gemm(top_diff + n * this->top_dim_, weight,
153                         bottom_diff + n * this->bottom_dim_);
154                 }
155             }
156         }
157     }
158 }

```

kmeans的实现如下，当然也可以用Opencv自带的，速度会更快些。



```

1  template<typename Dtype>
2  void kmeans_cluster(vector<int> &cLabel, vector<Dtype> &cCentro, Dtype *cWeights, int nWeig
3  {
4      //find min max
5      Dtype maxWeight=numeric_limits<Dtype>::min(), minWeight=numeric_limits<Dtype>::max();
6      for(int k = 0; k < nWeights; ++k)
7      {
8          if(mask[k])
9          {
10             if(cWeights[k] > maxWeight)
11                 maxWeight = cWeights[k];
12             if(cWeights[k] < minWeight)
13                 minWeight = cWeights[k];
14         }
15     }
16     // generate initial centroids linearly
17     for (int k = 0; k < nCluster; k++)
18         cCentro[k] = minWeight + (maxWeight - minWeight)*k / (nCluster - 1);
19
20     //initialize all label to -1
21     for (int k = 0; k < nWeights; ++k)
22         cLabel[k] = -1;
23
24     const Dtype float_max = numeric_limits<Dtype>::max();
25     // initialize
26     Dtype *cDistance = new Dtype[nWeights];
27     int *cClusterSize = new int[nCluster];
28
29     Dtype *pCentroPos = new Dtype[nCluster];
30     int *pClusterSize = new int[nCluster];
31     memset(pClusterSize, 0, sizeof(int)*nCluster);
32     memset(pCentroPos, 0, sizeof(Dtype)*nCluster);
33     Dtype *ptrC = new Dtype[nCluster];
34     int *ptrS = new int[nCluster];
35
36     int iter = 0;
37     //Dtype tk1 = 0.f, tk2 = 0.f, tk3 = 0.f;
38     double mCurDistance = 0.0;
39     double mPreDistance = numeric_limits<double>::max();
40
41     // clustering
42     while (iter < max_iter)
43     {
44         // check convergence
45         if (fabs(mPreDistance - mCurDistance) / mPreDistance < 0.01) break;
46         mPreDistance = mCurDistance;
47         mCurDistance = 0.0;
48
49         // select nearest cluster
50
51         for (int n = 0; n < nWeights; n++)
52         {
53             if (!mask[n])
54                 continue;
55             Dtype distance;

```

```

56         Dtype mindistance = float_max;
57         int clostCluster = -1;
58         for (int k = 0; k < nCluster; k++)
59         {
60             distance = fabs(cWeights[n] - cCentro[k]);
61             if (distance < mindistance)
62             {
63                 mindistance = distance;
64                 clostCluster = k;
65             }
66         }
67         cDistance[n] = mindistance;
68         cLabel[n] = clostCluster;
69     }
70
71
72     // calc new distance/inertia
73
74     for (int n = 0; n < nWeights; n++)
75     {
76         if (mask[n])
77             mCurDistance = mCurDistance + cDistance[n];
78     }
79
80
81     // generate new centroids
82     // accumulation(private)
83
84     for (int k = 0; k < nCluster; k++)
85     {
86         ptrC[k] = 0.f;
87         ptrS[k] = 0;
88     }
89
90     for (int n = 0; n < nWeights; n++)
91     {
92         if (mask[n])
93         {
94             ptrC[cLabel[n]] += cWeights[n];
95             ptrS[cLabel[n]] += 1;
96         }
97     }
98
99     for (int k = 0; k < nCluster; k++)
100    {
101        pCentroPos[ k] = ptrC[k];
102        pClusterSize[k] = ptrS[k];
103    }
104
105    //reduction(global)
106    for (int k = 0; k < nCluster; k++)
107    {
108
109        cCentro[k] = pCentroPos[k];
110        cClusterSize[k] = pClusterSize[k];

```

```

111
112         cCentro[k] /= cClusterSize[k];
113     }
114
115     iter++;
116     // cout << "Iteration: " << iter << " Distance: " << mCurDistance << endl;
117 }
118 //gather centroids
119 //#pragma omp parallel for
120 //for(int n=0; n<nNode; n++)
121 //    cNodes[n] = cCentro[cLabel[n]];
122
123 delete[] cDistance;
124 delete[] cClusterSize;
125 delete[] pClusterSize;
126 delete[] pCentroPos;
127 delete[] ptrC;
128 delete[] ptrS;
129 }

```

全连接层的修改和卷积层的一致不再赘述。同样的，可以把对应的.cu文件中的gpu前向和后向函数实现也修改了，方便gpu训练。

最后，在src/caffe/net.cpp的CopyTrainedLayersFrom(const NetParameter& param)函数中调用我们定义的函数，即在读入已经训练好的模型权值时，对每一层做需要的权值mask初始化和权值聚类。

```

795 layers_[target_layer_id]->ComputeBlobMask(ratio);
796

```

至此代码修改完毕，编译运行即可。

## Reference

- [1] SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size
- [2] Deep Compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding
- [3] Learning both Weights and Connections for Efficient Neural Networks
- [4] Efficient Inference Engine on Compressed Deep Neural Network

总结：

最后再提一句，几乎所有的模型压缩文章都是从Alexnet和VGG下手，一是因为他们都采用了多层较大的全连接层，而全连接层的权值甚至占到了总参数的90%以上，所以即便只对全连接层进行“开刀”，压缩效果也是显著的。另一方面，这些论文提出的结果在现在看来并不是state of art的，存在可提升的空间，而且在NIN的文章中表明，全连接层容易引起过拟合，去掉全连接层反而有助于精度提升，所以这么看来压缩模型其实是个不吃力又讨好的活，获得的好处显然是双倍的。但运用到特定的网络中，还需要不断反复试验，因地制宜，寻找适合该网络的压缩方式。

查看原文>> (<http://blog.csdn.net/may0324/article/details/52935869>)



0

#### 看过本文的人也看了：

- 深度学习知识结构图  
(<http://lib.csdn.net/base/deeplearning/structure>)
- FPGA机器学习之stanford机器学习第五堂...  
(<http://lib.csdn.net/article/deeplearning/54529>)
- 【连载】【FPGA黑金开发板】Verilog H...  
(<http://lib.csdn.net/article/deeplearning/52578>)
- xilinx fpga学习笔记5：Xst综合属性  
(<http://lib.csdn.net/article/deeplearning/52390>)
- 基于FPGA的多通道数据采集系统设计  
(<http://lib.csdn.net/article/deeplearning/54523>)
- [D-VI] my\_second\_fpga ( 1位加法器 Veri...  
(<http://lib.csdn.net/article/deeplearning/53657>)

#### 发表评论

输入评论内容

发表

#### 22个评论



([http://my.csdn.net/rill\\_zhen](http://my.csdn.net/rill_zhen))

**rill\_zhen** ([http://my.csdn.net/rill\\_zhen](http://my.csdn.net/rill_zhen))

求助一下，修改conv\_layer.cpp/conv\_layer.cu后，使用CPU/GPU进行训练时（先用CPU模式，后来尝试用GPU模式，都有错误），一旦调用forward\_cpu/forward\_gpu就会报错。尝试了各种方法后，我把muweight[i] = this->centroids\_[this->indices\_[i]]注释掉就能正常运行，backward\_cpu/backward\_gpu()也有类似情况，请问你有遇到过这种情况么？

2017-03-08 09:26:14

回复



([http://my.csdn.net/rill\\_zhen](http://my.csdn.net/rill_zhen))

**rill\_zhen** ([http://my.csdn.net/rill\\_zhen](http://my.csdn.net/rill_zhen))

回复rill\_zhen：尝试过下面评论中的方法，但没有效果。

2017-03-08 09:30:24

回复



(<http://my.csdn.net/u013478129>)



**u013478129** (<http://my.csdn.net/u013478129>)

哪位善长仁翁愿意分享一下这个修改好的代码在github呢?我一直很好奇,这个代码或者算法是不是有些很神秘的地方,或者什么潜规则,原作song han也不愿意分享这个代码?

2017-03-02 22:21:44

回复



(<http://my.csdn.net/u013478129>)

**u013478129** (<http://my.csdn.net/u013478129>)

回复u013478129： 哪位善长仁翁愿意分享一下这个修改好的代码在github呢?我一直很好奇,这个代码或者算法是不是有些很神秘的地方,或者什么潜规则,原作song han也不愿意分享这个代码?37303 9065 at qq.com

2017-03-03 16:02:33

回复



([http://my.csdn.net/qq\\_35800608](http://my.csdn.net/qq_35800608))

**qq\_35800608** ([http://my.csdn.net/qq\\_35800608](http://my.csdn.net/qq_35800608))

博主，你好，请问一下Deep Compression的具体运行步骤是什么，ReadMe里面的内容没有看懂

2017-02-21 22:00:07

回复

加载更多

---

公司简介 (<http://www.csdn.net/company/about.html>) | 招贤纳士 (<http://www.csdn.net/company/recruit.html>) | 广告服务 (<http://www.csdn.net/company/marketing.html>) | 联系方式 (<http://www.csdn.net/company/contact.html>) | 版权声明 (<http://www.csdn.net/company/statement.html>) | 法律顾问 (<http://www.csdn.net/company/layer.html>) | 问题报告 (<mailto:webmaster@csdn.net>) | 合作伙伴 (<http://www.csdn.net/friendlink.html>) | 论坛反馈 (<http://bbs.csdn.net/forums/Service>)

---

网站客服 杂志客服 (<http://wpa.qq.com/msg?d=3&uin=2251809102&site=qq&menu=yes>)

微博客服 (<http://e.weibo.com/csdnsupport/profile>) webmaster@csdn.net (<mailto:webmaster@csdn.net>) 400-600-2320 |

北京创新乐知信息技术有限公司 版权所有 | 江苏知之为计算机有限公司 | 江苏乐知网络技术有限公司

京 ICP 证 09002463 号 | Copyright © 1999-2016, CSDN.NET, All Rights Reserved



(<http://www.hd315.gov.cn/beian/view.asp?bianhao=010202001032100010>)