

Effective Sparsification of Neural Networks with Global Sparsity Constraint

A. Appendix

In this appendix, we present additional MobileNetV1 [15] experiment on ImageNet-1K, the general experimental configurations, proof for equation 3, proof for theorem 1, analysis on the effect of temperature annealing and PyTorch code snippets of ProbMask.

A.1. MobileNetV1 on ImageNet-1K

Dataset	ImageNet		
	ProbMask	STR	GMP
Ratio			
89%	65.19	62.10	61.80
94.1%	60.10	23.61	-

Table 4. ProbMask surpasses state-of-the-art methods by 3.09% and 38.49% Top-1 Accuracy, demonstrating the effectiveness and generalizability of ProbMask on lightweight MobileNetV1 [15] architectures. Following the setting of [18], 89% and 94.1% sparsity is chosen to compare at the same pruning rate.

A.2. Experimental Configurations

Dataset	CIFAR	ImageNet
GPUs	1	4
Batch Size	256	256
Epochs	300	100
Weight Optimizer	SGD	SGD
Weight Learning Rate	0.1	0.256
Weight Momentum	0.9	0.875
Probability Optimizer	Adam	Adam
Probability Learning Rate	6e-3	6e-3
t_1	48	16
t_2	180	60
Warmup	✗	✓
Label Smoothing	✗	0.1

Table 5. The bold-face probability learning rate 6e-3 is the **only** hyperparameter obtained by grid search on CIFAR-10 experiments on a small size network Conv-4 [6] and applied directly to larger datasets and networks. This demonstrates the generality of our proposed ProbMask to different datasets, different networks and different tasks, i.e., pruning networks and finding supermasks. Other hyperparameters are applied following the same practice of previous works [29, 18, 22, 41]. The channels of ResNet32 for CIFAR experiments are doubled following the same practice of [33]. The temperature annealing scheme follows the same practice of [37]

A.3. Proof for equation 3

Proof. The PDF (probability density function) of Gumbel($\mu, 1$) is

$$f(z; \mu) = e^{-(z-\mu)-e^{-(z-\mu)}}. \quad (9)$$

The CDF (cumulative distribution function) of Gumbel($\mu, 1$) is

$$F(z; \mu) = e^{-e^{-(z-\mu)}}. \quad (10)$$

We just need to prove that

$$\forall i, P(\log(s_i) - \log(1 - s_i) + g_{1,i} - g_{2,i} \geq 0) = s_i. \quad (11)$$

$g_{1,i}$ and $g_{2,i}$ are two Gumbel(0, 1) random variables sampled for s_i . The probability is taken with respect to $g_{1,i}$ and $g_{2,i}$. s_i can be seen as a constant in the following proof.

Let $z_1 = \log(s_i) + g_{1,i}$, $z_2 = \log(1 - s_i) + g_{2,i}$. Then $z_1 \sim \text{Gumbel}(\log(s_i), 1)$, $z_2 \sim \text{Gumbel}(\log(1 - s_i), 1)$.

$$P(\log(s_i) - \log(1 - s_i) + g_{1,i} - g_{2,i} \geq 0) \quad (12)$$

$$= P(z_2 \leq z_1) \quad (13)$$

$$= \int_{-\infty}^{+\infty} \int_{-\infty}^{z_1} f(z_2; \log(1 - s_i)) f(z_1; \log(s_i)) dz_2 dz_1 \quad (14)$$

$$= \int_{-\infty}^{+\infty} F(z_1; \log(1 - s_i)) f(z_1; \log(s_i)) dz_1 \quad (15)$$

$$= \int_{-\infty}^{+\infty} e^{-e^{-(z_1 - \log(1 - s_i))}} \cdot e^{-(z_1 - \log(s_i)) - e^{-(z_1 - \log(s_i))}} dz_1 \quad (16)$$

$$= \int_{-\infty}^{+\infty} e^{-e^{-z_1(1-s_i)-z_1+\log s_i} - e^{-z_1 s_i}} dz_1 \quad (17)$$

$$= s_i \int_{-\infty}^{+\infty} e^{-e^{-z_1} - z_1} dz_1 \quad (18)$$

$$= s_i \quad (19)$$

$\int_{-\infty}^{+\infty} e^{-e^{-z_1} - z_1} dz_1$ is the integral of a Gumbel(0,1) random variable. \square

A.4. Proof for Theorem 1

Proof. The projection from \mathbf{z} to set C can be formulated in the following optimization problem:

$$\min_{\mathbf{s} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{s} - \mathbf{z}\|^2, \\ \text{s.t. } \mathbf{1}^\top \mathbf{s} \leq K \text{ and } 0 \leq \mathbf{s}_i \leq 1.$$

Then we solve the problem with Lagrangian multiplier method.

$$\begin{aligned} L(\mathbf{s}, v) &= \frac{1}{2} \|\mathbf{s} - \mathbf{z}\|^2 + v(\mathbf{1}^\top \mathbf{s} - K) \\ &= \frac{1}{2} \|\mathbf{s} - (\mathbf{z} - v\mathbf{1})\|^2 + v(\mathbf{1}^\top \mathbf{z} - K) - \frac{n}{2} v^2. \end{aligned} \quad (20)$$

with $v \geq 0$ and $0 \leq s_i \leq 1$. Minimize the problem with respect to \mathbf{s} , we have

$$\tilde{\mathbf{s}} = \mathbf{1}_{\mathbf{z}-v\mathbf{1} \geq 1} + (\mathbf{z} - v\mathbf{1})_{\mathbf{1} > \mathbf{z}-v\mathbf{1} > 0} \quad (22)$$

Then we have

$$\begin{aligned} g(v) &= L(\tilde{\mathbf{s}}, v) \\ &= \frac{1}{2} \|[z - v\mathbf{1}]_- + [z - (v+1)\mathbf{1}]_+\|^2 \\ &\quad + v(\mathbf{1}^\top \mathbf{z} - s) - \frac{n}{2} v^2 \\ &= \frac{1}{2} \|[z - v\mathbf{1}]_-\|^2 + \frac{1}{2} \|[z - (v+1)\mathbf{1}]_+\|^2 \\ &\quad + v(\mathbf{1}^\top \mathbf{z} - s) - \frac{n}{2} v^2, v \geq 0. \\ g'(v) &= \mathbf{1}^\top [v\mathbf{1} - \mathbf{z}]_+ + \mathbf{1}^\top [(v+1)\mathbf{1} - \mathbf{z}]_- \\ &\quad + (\mathbf{1}^\top \mathbf{z} - s) - nv \\ &= \mathbf{1}^\top \min(1, \max(0, \mathbf{z} - v\mathbf{1})) - K, v \geq 0. \end{aligned}$$

It is easy to verify that $g'(v)$ is a monotone decreasing function with respect to v and we can use a bisection method solve the equation $g'(v) = 0$ with solution v_1^* . Then we get that $g(v)$ increases in the range of $(-\infty, v_1^*]$ and decreases in the range of $[v_1^*, +\infty)$. The maximum of $g(v)$ is achieved at 0 if $v_1^* \leq 0$ and v_1^* if $v_1^* > 0$. Then we set $v_2^* = \max(0, v_1^*)$. Finally we have

$$\begin{aligned} \mathbf{s}^* &= \mathbf{1}_{\mathbf{z}-v_2^*\mathbf{1} \geq 1} + (\mathbf{z} - v_2^*\mathbf{1})_{\mathbf{1} > \mathbf{z}-v_2^*\mathbf{1} > 0} \\ &= \min(1, \max(0, \mathbf{z} - v_2^*\mathbf{1})). \end{aligned} \quad (23)$$

□

A.5. Temperature Annealing

Thanks to the ℓ_1 norm and cube $[0, 1]^n$ in our constraint, most probabilities will converge to 0 or 1 at the end of training, which is shown in

$$\mathbf{s} = \min(1, \max(0, \mathbf{z} - v_2^*\mathbf{1})).$$

Traditional temperature annealing starts with a relative high value, i.e., 1 to have a smooth relaxation and gradually decrease to a small value to make relaxation close to the original objective function. In this section we analyze how the temperature annealing contributes to the training process, especially helping probabilities converge to 0 or 1.

Firstly consider the gradient:

$$\nabla_{s_i} \mathcal{L} \left(\mathbf{w}, \sigma \left(\frac{\log(\frac{s}{1-s}) + \mathbf{g}_1 - \mathbf{g}_0}{\tau} \right) \right) \quad (25)$$

$$= \nabla_{\sigma_i} \mathcal{L} \left(\mathbf{w}, \sigma \left(\frac{\log(\frac{s}{1-s}) + \mathbf{g}_1 - \mathbf{g}_0}{\tau} \right) \right) S, \quad (26)$$

where $S = \nabla_{s_i} \sigma \left(\frac{\log(\frac{s}{1-s}) + \mathbf{g}_1 - \mathbf{g}_0}{\tau} \right)$. We can see that the larger the magnitude $|S|$, z_i (step 8 in Algorithm 1) will vary more greatly.

Take $x = \frac{1}{\tau} \in [1, +\infty)$, $r = \log(s_i) - \log(1 - s_i) + g_{1,i} - g_{0,i}$. We have

$$S = \frac{\sigma(rx)(1 - \sigma(rx))x}{s_i(1 - s_i)} \quad (27)$$

Since S is an even function w.r.t r , we just consider the case $r > 0$. Then we take the gradient w.r.t to x .

$$\nabla_x (\sigma(rx)(1 - \sigma(rx))x) \quad (28)$$

$$= \sigma(rx)(1 - \sigma(rx))(rx - 2rx\sigma(rx) + 1) \quad (29)$$

Take $y = rx, y > 0$ since $x > 0$ and $r > 0$. The solution to $y - 2y\sigma(y) + 1 = 0$ is around 1.55. S is a monotonically increasing function for $x \in (0, \frac{1.55}{r}]$ and monotonically decreasing function for $x \in [\frac{1.55}{r}, +\infty)$.

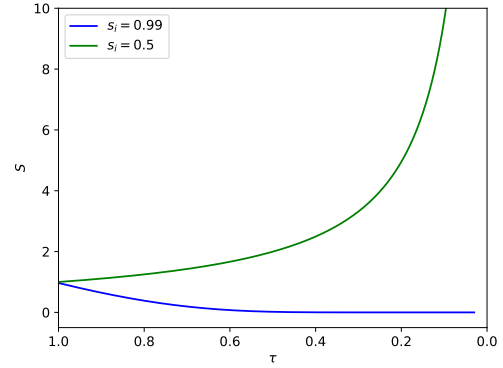


Figure 8. The value of S changes with τ approaching zero.

Now we analyze S using two special cases. Take $g = g_{1,i} - g_{2,i}, g \sim \text{Logistic}(0, 1)$. Take $s_i = 0.99$ and $g = 0.04$ for example. $r = \log(99) + 0.1 \approx 4.63$. S is monotonically decreasing function for $x \in [1, +\infty)$. Take $s_i = 0.5$ and $g = 0.04$ for example. S would increase as τ decreases to 0.03, since we set $\tau = 0.97(1 - t/T) + 0.03$. We plot the corresponding graph in Figure 8.

From the above two examples, we know that for probabilities around 0.5, S becomes larger in the training process, potentially making $|z_i|$ large and finally make probability come close to 0 or 1 after projection. For probabilities close to 0 or 1, S becomes smaller in the training process, making them stay close to 0 or 1 at the end of training.

```

1 class ProbMaskConv(nn.Conv2d):
2     def __init__(self, *args, **kwargs):
3         super().__init__(*args, **kwargs)
4         self.scores = nn.Parameter(torch.Tensor(self.weight.size())) # Probability
5         self.subnet = None # Mask
6         self.scores.data = (torch.ones_like(self.scores) * parser_args.score_init_constant)
7
8     def forward(self, x): # Sample a mask and forward propagation
9         if not parser_args.discrete: # Training
10            eps = 1e-20
11            temp = parser_args.T
12            uniform0 = torch.rand_like(self.scores)
13            uniform1 = torch.rand_like(self.scores)
14            noise = -torch.log(torch.log(uniform0 + eps) / torch.log(uniform1 + eps) + eps)
15            self.subnet = torch.sigmoid((torch.log(self.scores + eps) - torch.log(1.0 - self.scores +
16            eps) + noise) * temp)
17        else: # Testing
18            self.subnet = (torch.rand_like(self.scores) < self.scores).float()
19            w = self.weight * self.subnet
20            x = F.conv2d(x, w, self.bias, self.stride, self.padding, self.dilation, self.groups)
21            return x

```

Listing 1. PyTorch Code Snippets for ProbMaskConv.

```

1 def constrainScoreByWhole(model):
2     total = 0
3     for n, m in model.named_modules():
4         if hasattr(m, "scores"):
5             total += m.scores.nelement()
6     v = solveV(model, total) # Calculate  $v_2^*$  in Theorem 1
7     for n, m in model.named_modules():
8         if hasattr(m, "scores"):
9             m.scores.sub_(v).clamp_(0, 1) # Do the projection
10
11 def solveV(model, total): # Solve solution to Equation 7 with bisection search
12     k = total * parser_args.prune_rate
13     a, b = 0, 0
14     for n, m in model.named_modules():
15         if hasattr(m, "scores"):
16             b = max(b, m.scores.max())
17     def f(v):
18         s = 0
19         for n, m in model.named_modules():
20             if hasattr(m, "scores"):
21                 s += (m.scores - v).clamp(0, 1).sum()
22     return s - k
23 if f(0) < 0:
24     return 0
25 itr = 0
26 while (1):
27     itr += 1
28     v = (a + b) / 2
29     obj = f(v)
30     if abs(obj) < 1e-3 or itr > 20:
31         break
32     if obj < 0:
33         b = v
34     else:
35         a = v
36 v = max(0, v)
37 return v

```

Listing 2. PyTorch Code Snippets for Projection in Theorem 1