

# Proximal Mean-field for Neural Network Quantization

Thalaiyasingam Ajanthan<sup>1</sup>, Puneet K. Dokania<sup>1</sup>, Richard Hartley<sup>2</sup>, and Philip H. S. Torr<sup>1</sup>

<sup>1</sup>University of Oxford      <sup>2</sup>Australian National University

## Abstract

*Compressing large neural networks by quantizing the parameters, while maintaining the performance is often highly desirable due to the reduced memory and time complexity. In this work, we formulate neural network quantization as a discrete labelling problem and design an efficient approximate algorithm based on the popular mean-field method. To this end, we devise a projected stochastic gradient descent algorithm and show that it is, in fact, equivalent to a proximal version of the mean-field method. Thus, we provide an MRF optimization perspective to neural network quantization, which would enable research on modelling higher-order interactions among the network parameters to design better quantization schemes. Our experiments on standard image classification datasets with convolutional and residual architectures evidence that our algorithm obtains fully-quantized networks with accuracies very close to the floating-point reference networks.*

## 1. Introduction

Despite the success of deep neural networks, they are highly overparametrized, resulting in excessive computational and memory requirements. Compressing such large networks by quantizing the parameters, while maintaining the performance is highly desirable for real-time applications, or for resource-limited devices.

In Neural Network (NN) quantization, the objective is to learn a network while restricting the parameters to take values from a small discrete set (usually binary) representing quantization levels. This can be formulated as a *discrete labelling problem* where each learnable parameter takes a label from the discrete set and the learning objective is to find the label configuration that minimizes the empirical loss. This is an extremely challenging discrete optimization problem because the number of label configurations grows exponentially with the number of parameters in the network and the loss function is highly non-convex.

Over the past 20 years, similar large-scale discrete labelling problems have been extensively studied under the

context of Markov Random Field (MRF) optimization, and many efficient approximate algorithms have been developed [2, 5, 9, 28, 38, 39]. In this work, we bridge the gap between NN quantization and MRF optimization, and design an efficient approximate algorithm based on the popular mean-field method [39] for NN quantization.

Specifically, we first formulate NN quantization as a discrete labelling problem. Then, we relax the discrete solution space to a convex polytope and introduce an algorithm to iteratively optimize the first-order Taylor approximation of the loss function over the polytope. This approach is in fact a (stochastic) gradient descent type method with an additional projection step at each iteration. For a particular choice of differentiable projection, we show that our method is equivalent to a proximal version of the well-known mean-field method. Furthermore, we prove that under certain conditions, our algorithm specializes to the popular binary connect quantization method [8].

The MRF optimization perspective to NN quantization opens up many interesting research directions. Mainly, our approach represents the simplest case where the NN parameters are assumed to be independent of each other. However, one can potentially model second-order or even high-order interactions among parameters and use efficient inference algorithms developed and well-studied in the MRF optimization literature. Therefore, we believe, many such algorithms can be transposed into this framework to design better network quantization schemes. Furthermore, in contrast to the existing NN quantization methods [19, 32], we eliminate the need for additional heuristics by quantizing all the learnable parameters in the network (including biases) and our formulation can be seamlessly extended to quantization levels beyond binary.

We evaluate the merits of our algorithm on MNIST, CIFAR-10/100, and TinyImageNet classification datasets with convolutional and residual architectures. Our experiments evidence that the quantized networks obtained by our algorithm yield accuracies very close to the floating point counterparts while consistently outperforming directly comparable baselines when all learnable parameters are quantized.

## 2. Neural Network Quantization

Neural Network (NN) quantization is defined as learning neural networks with the parameters restricted to a small discrete set representing quantization levels. This primarily relies on the hypothesis that since NNs are usually over-parametrized, perhaps it is possible to obtain a quantized network with performance comparable to the floating-point network. To this end, given a dataset  $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ , NN quantization problem can be written as:

$$\min_{\mathbf{w}} L(\mathbf{w}; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; (\mathbf{x}_i, \mathbf{y}_i)), \quad (1)$$

$$\mathbf{w} \in \mathcal{Q}^m.$$

Here,  $\ell(\cdot)$  is the standard loss function (e.g., cross-entropy loss),  $\mathbf{w}$  is the  $m$  dimensional parameter vector, and  $\mathcal{Q}$  with  $|\mathcal{Q}| = d$  is a predefined discrete set representing quantization levels (e.g.,  $\mathcal{Q} = \{-1, 1\}$  or  $\mathcal{Q} = \{-1, 0, 1\}$ ). In Eq. (1), all the learnable parameters including biases are quantized. Thus, the objective is to obtain a *fully-quantized network*, as opposed to previous methods where some parts of the network are not quantized (e.g., biases and last layer parameters) [8, 32]. This also eliminates the need for additional heuristics.

### 2.1. NN Quantization as Discrete Labelling

We observe that Eq. (1) is a *discrete labelling problem* where each learnable parameter  $w_j$  takes a label  $\lambda$  from the discrete label set  $\mathcal{Q}$ . Here, the learning objective is to find a label configuration  $\mathbf{w}^*$  that minimizes the loss function  $L(\mathbf{w})$ . Notice that, in general, optimizing Eq. (1) directly is intractable due to two inherent problems [23, 29]: 1) the solution space is discrete with exponentially many feasible points ( $d^m$  with  $m$  in the order of millions); and 2) the loss function (known as the *energy* in MRF literature) is highly non-convex and does not satisfy any regularity condition (e.g., submodularity). This hinders the use of any off-the-shelf discrete optimization algorithm.

To this end, we introduce an approximate algorithm by relaxing both the above mentioned problems. In particular, we first relax the discrete solution space to a convex polytope (a standard practice in discrete optimization) and then iteratively optimize the first-order Taylor approximation of the loss function over the polytope. Our approach, as will be shown subsequently, falls under the regime of (stochastic) gradient descent type methods and is applicable to any type of loss function. In what follows we talk about the above mentioned relaxations and related optimization in detail.

### 2.2. Continuous Relaxation of the Solution Space

The parameter vector  $\mathbf{w} \in \mathcal{Q}^m$  can be equivalently represented using indicator variables as follows. Let  $u_{j:\lambda} \in$

$\{0, 1\}$  be the indicator variable, where  $u_{j:\lambda} = 1$  if and only if  $w_j = \lambda \in \mathcal{Q}$ . Then, for any  $j \in \{1 \dots m\}$ , we can write

$$w_j = \sum_{\lambda \in \mathcal{Q}} \lambda u_{j:\lambda}, \quad (2)$$

$$\text{s.t. } \sum_{\lambda \in \mathcal{Q}} u_{j:\lambda} = 1, \quad u_{j:\lambda} \in \{0, 1\} \quad \forall \lambda \in \mathcal{Q}.$$

Note, any  $w_j$  represented using Eq. (2) belongs to  $\mathcal{Q}$ . For convenience, by denoting the vector of quantization levels as  $\mathbf{q}$ , a parameter vector  $\mathbf{w} \in \mathcal{Q}^m$  can be written in a matrix vector product as:

$$\mathbf{w} = \mathbf{u} \mathbf{q}, \quad (3)$$

$$\text{s.t. } \mathbf{u} \in \mathcal{V} = \left\{ \mathbf{u} \mid \sum_{\lambda} u_{j:\lambda} = 1, \quad \forall j, \right. \\ \left. u_{j:\lambda} \in \{0, 1\}, \quad \forall j, \lambda \right\}.$$

Here,  $\mathbf{u}$  is a  $m \times d$  matrix<sup>1</sup> (i.e., each row  $\mathbf{u}_j = \{u_{j:\lambda} \mid \lambda \in \mathcal{Q}\}$ ), and  $\mathbf{q}$  is a column vector of dimension  $d$ . Note that there is a one-to-one correspondence between the sets  $\mathcal{V}$  and  $\mathcal{Q}^m$ . Substituting Eq. (3) in the NN quantization objective (1) results in the variable change from  $\mathbf{w}$  to  $\mathbf{u}$  as:

$$\min_{\mathbf{w} \in \mathcal{Q}^m} L(\mathbf{w}; \mathcal{D}) = \min_{\mathbf{u} \in \mathcal{V}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{u} \mathbf{q}; (\mathbf{x}_i, \mathbf{y}_i)). \quad (4)$$

Even though the above variable change converts the problem from  $m$  to  $md$  dimensional space, the one-to-one correspondence between  $\mathbf{w}$  and  $\mathbf{u}$  ensures that the cardinality of the sets  $\mathcal{Q}^m$  and  $\mathcal{V}$  are exactly the same. The binary constraint  $u_{j:\lambda} \in \{0, 1\}$  together with the non-convex loss function  $L(\cdot)$  makes the problem NP-hard [29]. By relaxing the binary constraints to  $u_{j:\lambda} \in [0, 1]$ , we obtain the convex hull of the set  $\mathcal{V}$  as defined below:

$$\mathcal{S} = \text{conv}(\mathcal{V}) = \left\{ \mathbf{u} \mid \sum_{\lambda} u_{j:\lambda} = 1, \quad \forall j, \right. \\ \left. u_{j:\lambda} \geq 0, \quad \forall j, \lambda \right\}. \quad (5)$$

Furthermore, the set  $\mathcal{S}$  decomposes over each  $j$ , and it is in fact the Cartesian product of the probability simplexes of dimension  $d$ . Thus,

$$\mathcal{S} = \prod_{j=1}^m \Delta_j, \quad \text{where } \Delta_j = \left\{ \mathbf{z} \mid \sum_{\lambda} z_{\lambda} = 1, \right. \\ \left. z_{\lambda} \geq 0, \quad \forall \lambda \right\}. \quad (6)$$

Therefore, for a feasible point  $\mathbf{u} \in \mathcal{S}$ , the vector  $\mathbf{u}_j$  for each  $j$  ( $j$ -th row of matrix  $\mathbf{u}$ ) belongs to the probability simplex of dimension  $d$ . Hence, we can interpret the value  $u_{j:\lambda}$  as the probability of assigning the discrete label  $\lambda$  to the weight  $w_j$ . Now, the relaxed optimization can be written as:

$$\min_{\mathbf{u} \in \mathcal{S}} \tilde{L}(\mathbf{u}; \mathcal{D}) := L(\mathbf{u} \mathbf{q}; \mathcal{D}), \quad (7)$$

<sup>1</sup>To simplify the notation, we denote  $\mathbf{u}$  as a matrix but it can be thought of as an  $md$  dimensional vector obtained by flattening the matrix  $\mathbf{u}$ .

Note that, for any quantization set  $\mathcal{Q}$ , if  $\mathbf{u} \in \mathcal{V}$ , then the loss function  $\tilde{L}(\mathbf{u})$  has the same value as the original loss function  $L(\mathbf{w})$ . Furthermore, the relaxation of  $\mathbf{u}$  from  $\mathcal{V}$  to  $\mathcal{S}$ , translates into relaxing  $\mathbf{w}$  from  $\mathcal{Q}^m$  to the convex region  $[q_{\min}, q_{\max}]^m$ . Here,  $q_{\min}$  and  $q_{\max}$  represent the minimum and maximum quantization levels, respectively.

In fact,  $\mathbf{u}$  is an overparametrized representation of  $\mathbf{w}$ , and the mapping  $\mathbf{u} \rightarrow \mathbf{w}$  defined by Eq. (3) is many-to-one, and precisely an *onto*<sup>2</sup> mapping. However, this representation has an interesting probabilistic interpretation that learning  $\mathbf{u}$  can be interpreted as learning a discrete probability distribution over the NN parameters  $\mathbf{w}$ . This interpretation would be useful in drawing the connection between our algorithm and the mean-field method later in Sec. 3.

In addition, it can be shown that any local minimum in  $\mathcal{S}$  (the relaxed  $\mathbf{u}$ -space) is also a local minimum in  $[q_{\min}, q_{\max}]^m$  (the relaxed  $\mathbf{w}$ -space) and vice versa (Proposition 2.1). This essentially means that the variable change from  $\mathbf{w}$  to  $\mathbf{u}$  does not alter the optimization problem and a local minimum in the  $\mathbf{w}$ -space can be obtained by optimizing in the  $\mathbf{u}$ -space.

**Proposition 2.1.** Let  $f(\mathbf{w})$  be a continuous function with  $\mathbf{w} = g(\mathbf{u}) = \mathbf{u}\mathbf{q}$ . Then a point  $\mathbf{u}^k \in \mathcal{S}$  is a local minimum of  $f \circ g$ , if and only if  $\mathbf{w}^k = \mathbf{u}^k\mathbf{q}$  is a local minimum of  $f$  in the region  $[q_{\min}, q_{\max}]^m$ .

*Proof.* It can be easily proved using the properties that the function  $g : \mathcal{S} \rightarrow [q_{\min}, q_{\max}]^m$  is onto and continuous. See Appendix A.  $\square$

Finally, we would like to point out that the relaxation used while moving from  $\mathbf{w}$  to  $\mathbf{u}$  space is well studied in the MRF optimization literature and has been used to prove bounds on the quality of the solutions [7, 22]. In the case of NN quantization, in addition to the connection to mean-field (Sec. 3), we believe that this relaxation allows for more exploration, which would be useful in the stochastic setting.

### 2.3. First-order Approximation and Optimization

Here we talk about the optimization of  $\tilde{L}(\mathbf{u})$  over  $\mathcal{S}$ , discuss how our optimization scheme allows more exploration in the parameter space, and also discuss the conditions when this optimization will lead to a quantized solution in the  $\mathbf{w}$  space, which is our prime objective.

Stochastic Gradient Descent (SGD)<sup>3</sup> [34] is the de facto method of choice for optimizing neural networks. In this section, we interpret SGD as a proximal method, which will be useful later to show its difference to our final algorithm. In particular, SGD (or gradient descent) can be interpreted as

<sup>2</sup>A mapping  $f(x) : \mathcal{X} \rightarrow \mathcal{Y}$  is onto if  $\forall y \in \mathcal{Y}, \exists x \in \mathcal{X}$  such that  $f(x) = y$ .

<sup>3</sup>The difference between SGD and gradient descent is that the gradients are approximated using a stochastic oracle in the former case.

iteratively minimizing the first-order Taylor approximation of the loss function augmented by a proximal term [30]. In our case, the objective function is the same as SGD but the feasible points are now constrained to form a convex polytope. Thus, at each iteration  $k$ , the first-order objective can be written as:

$$\mathbf{u}^{k+1} = \underset{\mathbf{u} \in \mathcal{S}}{\operatorname{argmin}} \tilde{L}(\mathbf{u}^k) + \langle \mathbf{g}^k, \mathbf{u} - \mathbf{u}^k \rangle_F + \frac{1}{2\eta} \|\mathbf{u} - \mathbf{u}^k\|_F^2, \quad (8)$$

where  $\eta > 0$  is the learning rate and  $\mathbf{g}^k := \nabla_{\mathbf{u}} \tilde{L}^k$  is the stochastic (or mini-batch) gradient of  $\tilde{L}$  with respect to  $\mathbf{u}$  evaluated at  $\mathbf{u}^k$ . Here,  $\langle \cdot, \cdot \rangle_F$  is the Frobenius inner product<sup>4</sup> and  $\|\cdot\|_F$  is the Frobenius norm, respectively. In the unconstrained case, by setting the derivative with respect to  $\mathbf{u}$  to zero, one can easily verify that the above formulation leads to standard SGD updates. For constrained optimization (as in our case (8)), it is natural to use the stochastic version of Projected Gradient Descent (PGD) [35]. Specifically, at iteration  $k$ , the projected stochastic gradient update can be written as:

$$\mathbf{u}^{k+1} = P_{\mathcal{S}}(\mathbf{u}^k - \eta \mathbf{g}^k), \quad (9)$$

where  $P_{\mathcal{S}}(\cdot)$  denotes the projection to the polytope  $\mathcal{S}$ . Even though this type of problems can be optimized using projection-free algorithms [3, 11, 24], by relying on PGD, we enable the use of any off-the-shelf first-order optimization algorithms (e.g., Adam [21]). Furthermore, for a particular choice of projection, we show that the PGD update is equivalent to a proximal version of the mean-field method.

#### 2.3.1 Projection to the Polytope $\mathcal{S}$

As pointed out in Eq. (6), the polytope  $\mathcal{S}$  can be decomposed into  $m$  probability simplexes of dimension  $d$ . From Eq. (8), it is clear that the objective function is also separable for each  $j$ . Therefore, projection to  $\mathcal{S}$  can be decomposed into  $m$  independent projections to the  $d$ -dimensional probability simplexes. This decomposition amounts to an independence assumption that the probability of parameter  $w_j$  taking a label  $\lambda$  (represented by  $u_{j:\lambda}$ ) is independent for each  $j$ .

Under this assumption, a simple and differentiable way of projecting to the probability simplex is via the softmax function. Thus, for a given update  $\tilde{\mathbf{u}}^{k+1} = \mathbf{u}^k - \eta \mathbf{g}^k$  and a scalar  $\beta > 0$  (usually  $\beta = 1$ ), the softmax projection for each  $j \in \{1 \dots m\}$  can be written as:

$$\mathbf{u}_j^{k+1} = \operatorname{softmax}(\beta \tilde{\mathbf{u}}_j^{k+1}), \quad \text{where} \quad (10)$$

$$u_{j:\lambda}^{k+1} = \frac{e^{\beta(u_{j:\lambda}^k - \eta g_{j:\lambda}^k)}}{\sum_{\mu \in \mathcal{Q}} e^{\beta(u_{j:\mu}^k - \eta g_{j:\mu}^k)}} \quad \forall \lambda \in \mathcal{Q}.$$

<sup>4</sup>This is equivalent to vectorizing the matrices and applying the standard inner product.

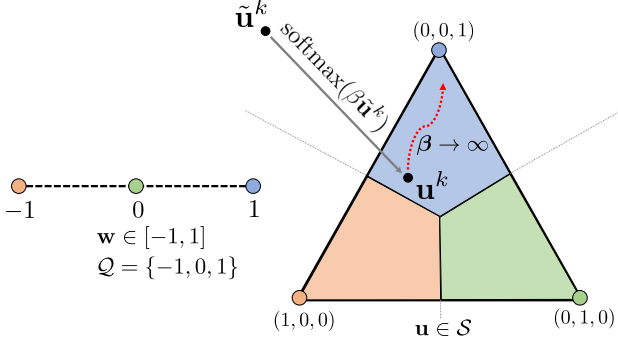


Figure 1: Illustration of  $\mathbf{w}$  and  $\mathbf{u}$ -spaces and exploration with softmax projection for an example where  $m = 1$ . Here each vertex of the simplex corresponds to a discrete quantization level in the  $\mathbf{w}$ -space and the simplex is partitioned based on its vertex association. Given an infeasible point  $\tilde{\mathbf{u}}^k$ , it is projected to the simplex via softmax and when  $\beta \rightarrow \infty$ , the projected point would move towards the associated vertex.

It can be easily verified that  $\mathbf{u}^{k+1} \in \mathcal{S}$ . For brevity, we use  $\mathbf{u} = \text{softmax}(\beta \tilde{\mathbf{u}})$  to denote the softmax projection applied for each  $j \in \{1 \dots m\}$  independently.

### 2.3.2 Exploration and Quantization

Recall that, even though we iteratively optimize over the relaxed polytope  $\mathcal{S}$ , our objective is to obtain a quantized solution in the  $\mathbf{w}$  space. This is equivalent to obtaining a solution  $\mathbf{u}$  that is a vertex of the polytope  $\mathcal{S}$ . We achieve this using an increasing schedule for  $\beta$ , which enables both exploration and quantization.

More precisely, the update Eq. (10) can be interpreted as an *exploration mechanism* over the probability simplexes. In particular, once the gradient step is computed, the resulting point is projected to the polytope via a “noisy” operator (softmax) with noise controlled by the hyperparameter  $\beta$  (lower the  $\beta$  the more noise). By noise we mean how far is the projected point from a vertex, *i.e.*, farther the projection, the noisier is the operator. It is easy to see that, when  $\beta \rightarrow \infty$ , the resulting projection tends to place all the probability mass in one dimension (for a vector  $\mathbf{u}_j$ ), and zeros everywhere else. Thus, this limiting case leads to a vertex of the polytope  $\mathcal{S}$ , which is the zero-noise case. Hence, a monotonically increasing schedule for  $\beta$  results in exploration over the polytope and finally reaches a vertex. This interpretation is illustrated in Fig. 1, and an entropy based view of the same phenomena is provided in Sec. 3. Note that, similar to ours, constraining the solution space through a hyperparameter is extensively studied in the optimization literature and one such example is the barrier method [6].

It is interesting to note that, if we always project to the

---

### Algorithm 1 Proximal Mean-Field (PMF)

---

**Require:**  $K, b, \{\eta^k\}, \rho > 1, \mathcal{D}, \tilde{L}$

**Ensure:**  $\mathbf{w}^* \in \mathcal{Q}^m$

- 1:  $\tilde{\mathbf{u}}^0 \in \mathbb{R}^{m \times d}, \beta \leftarrow 1$  ▷ Initialization
  - 2: **for**  $k \leftarrow 0 \dots K$  **do**
  - 3:    $\mathbf{u}^k \leftarrow \text{softmax}(\beta \tilde{\mathbf{u}}^k)$  ▷ Projection (Eq. (10))
  - 4:    $\mathcal{D}^b = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^b \sim \mathcal{D}$  ▷ Sample a mini-batch
  - 5:    $\mathbf{g}_{\mathbf{u}}^k \leftarrow \nabla_{\mathbf{u}} \tilde{L}(\mathbf{u}; \mathcal{D}^b) \big|_{\mathbf{u}=\mathbf{u}^k}$  ▷ Gradient w.r.t.  $\mathbf{u}$  at  $\mathbf{u}^k$
  - 6:    $\mathbf{g}_{\tilde{\mathbf{u}}}^k \leftarrow \mathbf{g}_{\mathbf{u}}^k \frac{\partial \mathbf{u}}{\partial \tilde{\mathbf{u}}} \big|_{\tilde{\mathbf{u}}=\tilde{\mathbf{u}}^k}$  ▷ Gradient w.r.t.  $\tilde{\mathbf{u}}$  at  $\mathbf{u}^k$
  - 7:    $\tilde{\mathbf{u}}^{k+1} \leftarrow \tilde{\mathbf{u}}^k - \eta^k \mathbf{g}_{\tilde{\mathbf{u}}}^k$  ▷ Gradient descent on  $\tilde{\mathbf{u}}$
  - 8:    $\beta \leftarrow \rho \beta$  ▷ Increase  $\beta$
  - 9: **end for**
  - 10:  $\mathbf{w}^* \leftarrow \text{hardmax}(\tilde{\mathbf{u}}^K) \mathbf{q}$  ▷ Quantization (Eq. (13))
- 

set of vertices  $\mathcal{V}$  (as in existing methods [19]), the gradient descent may get stuck, as single gradient step may not be sufficient to move from one vertex to the next. However, in our case, the feasible domain is a convex polytope where the quantization error on the gradients are not catastrophic, meaning that the projected gradient steps are significant enough to move from one feasible point to the next.

Furthermore, we would like to point out that, the PGD update (10) in fact yields an approximate solution to Eq. (8). However, in the following section, we show that this update is exactly minimizing a similar first-order objective function but augmented by an entropy term. Hence, the overall PGD algorithm can be shown to be equivalent to a proximal version of the popular mean-field method.

## 3. Softmax based PGD as Proximal Mean-field

Here we discuss the connection between softmax based PGD and the well-known mean-field method [39]. More precisely, we show that the update in (10) is actually an *exact fixed point update* of a modified mean-field objective function. This connection bridges the gap between the MRF optimization and the NN quantization literature. We now begin with a brief review of the mean-field method and then proceed with our proof.

**Mean-field Method.** A self-contained overview is provided in the Appendix B, but here we review the important details. Given an energy (or loss) function  $L(\mathbf{w})$  and the corresponding probability distribution of the form  $P(\mathbf{w}) = e^{-L(\mathbf{w})}/Z$ , mean-field approximates  $P(\mathbf{w})$  using a fully-factorized distribution  $U(\mathbf{w}) = \prod_{j=1}^m U_j(w_j)$ . Here, the distribution  $U$  is obtained by minimizing the KL-divergence  $\text{KL}(U \parallel P)$ . Note that, from the probabilistic interpretation of  $\mathbf{u} \in \mathcal{S}$  (see Sec. 2.2), for each  $j \in \{1 \dots m\}$ , the probability  $U_j(w_j = \lambda) = u_{j:\lambda}$ . Therefore, the distribution  $U$  can be represented using the variables  $\mathbf{u} \in \mathcal{S}$ , and hence,



the mean-field objective can be written as:

$$\operatorname{argmin}_{\mathbf{u} \in \mathcal{S}} \text{KL}(\mathbf{u} \| P) = \operatorname{argmin}_{\mathbf{u} \in \mathcal{S}} \mathbb{E}_{\mathbf{u}}[L(\mathbf{w})] - H(\mathbf{u}), \quad (11)$$

where  $\mathbb{E}_{\mathbf{u}}[\cdot]$  denotes the expectation over  $\mathbf{u}$  and  $H(\mathbf{u}) = -\sum_{j=1}^m \sum_{\lambda \in \mathcal{Q}} u_{j:\lambda} \log u_{j:\lambda}$  is the entropy.

In fact, mean-field is extensively studied in the context of MRF literature where the energy function  $L(\mathbf{w})$  factorizes over small subsets of variables  $\mathbf{w}$ . This leads to efficient minimization of the KL-divergence as the expectation  $\mathbb{E}_{\mathbf{u}}[L(\mathbf{w})]$  can be computed efficiently. However, in a standard neural network, the function  $L(\mathbf{w})$  does not have an explicit factorization and direct minimization of the KL-divergence is not straight forward. To simplify the NN loss function one can approximate it using its first-order Taylor approximation which discards the interactions between the NN parameters altogether.

Interestingly, in Theorem 3.1, we show that our softmax based PGD iteratively applies a proximal version of mean-field to the first-order approximation of  $L(\mathbf{w})$ . At iteration  $k$ , let  $\hat{L}^k(\mathbf{w})$  be the first-order Taylor approximation of  $L(\mathbf{w})$ . Then, since there are no interactions among parameters in  $\hat{L}^k(\mathbf{w})$ , and it is linear, our proximal mean-field objective has a closed form solution, which is exactly the softmax based PGD update.

**Theorem 3.1.** *At iteration  $k$ , let  $\mathbf{u}^{k+1}$  be the point from the softmax based PGD update (10). Then,*

$$\mathbf{u}^{k+1} = \operatorname{argmin}_{\mathbf{u} \in \mathcal{S}} \eta \mathbb{E}_{\mathbf{u}}[\hat{L}^k(\mathbf{w})] - \langle \mathbf{u}^k, \mathbf{u} \rangle_F - \frac{1}{\beta} H(\mathbf{u}), \quad (12)$$

where  $\hat{L}^k(\mathbf{w})$  is the first-order Taylor approximation of  $L$  at  $\mathbf{w}^k = \mathbf{u}^k \mathbf{q}$  and  $\eta > 0$  is the learning rate.

*Proof.* First we will show that,  $\mathbb{E}_{\mathbf{u}}[\hat{L}^k(\mathbf{w})] = \langle \mathbf{g}^k, \mathbf{u} \rangle_F + c$  for some constant  $c$ . Then, the softmax update can be derived by writing the Lagrangian and setting the derivatives with respect to  $\mathbf{u}$  to zero. See Appendix C.  $\square$

The objective function Eq. (12) is the same as the mean-field objective for  $\hat{L}^k(\mathbf{w})$  (refer Eq. (11)) except for the term  $\langle \mathbf{u}^k, \mathbf{u} \rangle_F$ . This, in fact, acts as a proximal term. Note, it is the cosine similarity but subtracted from the loss to enforce proximity. Therefore, it encourages the resulting  $\mathbf{u}^{k+1}$  to be closer to the current point  $\mathbf{u}^k$  and its influence relative to the loss term is governed by the learning rate  $\eta$ . Since gradient estimates are stochastic in our case, such a proximal term is highly desired as it encourages the updates to make a smooth transition.

Furthermore, the negative entropy term acts as a convex regularizer and when  $\beta \rightarrow \infty$  its influence becomes negligible and the update results in a binary labelling  $\mathbf{u} \in \mathcal{V}$ . It is interesting to note that the softmax projection in Eq. (10) translates into an entropy term in the objective

function (12), and for small values of  $\beta$ , it allows the iterative procedure to explore the optimization landscape. Such an explorative behaviour (at least in the beginning of optimization) is desirable especially for a stochastic optimization algorithm. In addition, the entropy term in Eq. (12) captures the (in)dependency between the parameters. To encode dependency, the entropy of the fully-factorized distribution can perhaps be replaced with a more complex entropy such as a tree-structured entropy, following the idea of [33]. Furthermore, in place of  $\hat{L}^k$ , a higher-order approximation can be used. However, such explorations go beyond the scope of this paper.

Algorithm 1 summarizes our approach. Here, similar to the existing methods [19] we store the auxiliary variables  $\tilde{\mathbf{u}} \in \mathbb{R}^{m \times d}$  and perform gradient descent on them. In contrast to the existing methods, this is not a necessity but empirically it improves the performance. Finally, since  $\beta$  can never be  $\infty$ , to ensure a fully-quantized network, the final quantization is performed using the hardmax projection (defined below in Eq. (13)). This is equivalent to performing Maximum a Posteriori (MAP) estimate on the learned probability distribution  $\mathbf{u} \in \mathcal{S}$ .

### 3.1. Proximal ICM as a Special Case

In Algorithm 1, if the softmax projection is replaced by the hardmax projection, the resulting update can be shown to be the same as a proximal version of Iterative Conditional Modes (ICM) [4]. To see this, let us define the hardmax projection at iteration  $k$  as:

$$\begin{aligned} \mathbf{u}_j^k &= \text{hardmax}(\tilde{\mathbf{u}}_j^k), \\ u_{j:\lambda}^k &= \begin{cases} 1 & \text{if } \lambda = \operatorname{argmax}_{\mu \in \mathcal{Q}} \tilde{u}_{j:\mu}^k \\ 0 & \text{otherwise} \end{cases} \quad \forall \lambda \in \mathcal{Q}. \end{aligned} \quad (13)$$

where  $\tilde{\mathbf{u}}^k = \mathbf{u}^{k-1} - \eta \mathbf{g}^{k-1}$ . Here, following the proof of Theorem 3.1, it can be shown that this update yields a fixed point of the following equation:

$$\min_{\mathbf{u} \in \mathcal{S}} \eta \langle \mathbf{g}^k, \mathbf{u} \rangle_F - \langle \mathbf{u}^k, \mathbf{u} \rangle_F. \quad (14)$$

Notice, this is exactly the same as the ICM objective augmented by the proximal term. In this case,  $\mathbf{u} \in \mathcal{V} \subset \mathcal{S}$ , meaning, the feasible domain is restricted to be the vertices of the polytope  $\mathcal{S}$ . Moreover, when  $\beta \rightarrow \infty$ , the softmax function approaches hardmax, and hence, this is a special case of proximal mean-field.

### 3.2. Binary Connect as Proximal ICM

In this section, considering binary neural networks, *i.e.*,  $\mathcal{Q} = \{-1, 1\}$ , and non-stochastic setting, we show that the Proximal Iterative Conditional Modes (PICM) algorithm is equivalent to the deterministic version of the popular Binary Connect (BC) method [8]. In these algorithms, the gradients

---

**Algorithm 2** Non-stochastic Binary Connect (BC) [8]

---

**Require:**  $K, \eta_w, \mathcal{D}, L, \mathcal{Q}$ **Ensure:**  $\mathbf{w}^* \in \mathcal{Q}^m$ 

```
1:  $\tilde{\mathbf{w}}^0 \in \mathbb{R}^m$  ▷ Initialization
2: for  $k \leftarrow 0 \dots K$  do
3:    $\mathbf{w}^k \leftarrow \text{sign}(\tilde{\mathbf{w}}^k)$  ▷ Projection
4:    $\mathbf{g}_{\mathbf{w}}^k \leftarrow \nabla_{\mathbf{w}} L(\mathbf{w}; \mathcal{D})|_{\mathbf{w}=\mathbf{w}^k}$  ▷ Gradient w.r.t.  $\mathbf{w}$ 
5:    $\mathbf{g}_{\tilde{\mathbf{w}}}^k \leftarrow \mathbf{g}_{\mathbf{w}}^k \frac{\partial \mathbf{w}}{\partial \tilde{\mathbf{w}}}|_{\tilde{\mathbf{w}}=\tilde{\mathbf{w}}^k}$  ▷ Gradient w.r.t.  $\tilde{\mathbf{w}}$ 
6:    $\tilde{\mathbf{w}}^{k+1} \leftarrow \tilde{\mathbf{w}}^k - \eta_w \mathbf{g}_{\tilde{\mathbf{w}}}^k$  ▷ Gradient descent
7: end for
8:  $\mathbf{w}^* \leftarrow \text{sign}(\tilde{\mathbf{w}}^K)$  ▷ Final quantization
```

---

are computed in two different spaces and therefore to alleviate any discrepancy due to the stochasticity, we assume that gradients are computed using the full dataset, *i.e.*, the non-stochastic setting.

Let  $\tilde{\mathbf{w}} \in \mathbb{R}^m$  and  $\mathbf{w} \in \mathcal{Q}^m$  be the infeasible and feasible points of binary connect. Similarly,  $\tilde{\mathbf{u}} \in \mathbb{R}^{m \times d}$  and  $\mathbf{u} \in \mathcal{V} \subset \mathcal{S}$  be the infeasible and feasible points of our PICM method, respectively. For convenience, we summarize binary connect (Algorithm 2) as a projected gradient descent method in  $\mathbf{w}$ -space. Recall that, PICM is the same as Algorithm 1 except, in PICM, the softmax projection is replaced by hardmax. Now, we show that the update steps in both BC and PICM are equivalent.

**Proposition 3.1.** Consider BC and PICM with  $\mathbf{q} = [-1, 1]^T$  and  $\eta_w > 0$ . For an iteration  $k > 0$ , if  $\tilde{\mathbf{w}}^k = \tilde{\mathbf{u}}^k \mathbf{q}$  then,

1. the projections in BC:  $\mathbf{w}^k = \text{sign}(\tilde{\mathbf{w}}^k)$  and PICM:  $\mathbf{u}^k = \text{hardmax}(\tilde{\mathbf{u}}^k)$  satisfy  $\mathbf{w}^k = \mathbf{u}^k \mathbf{q}$ .
2. if  $\eta_u = \eta_w/2$ , then the updated points after the gradient descent step in BC and PICM satisfy  $\tilde{\mathbf{w}}^{k+1} = \tilde{\mathbf{u}}^{k+1} \mathbf{q}$ .

*Proof.* Case (1) is simply applying  $\tilde{\mathbf{w}}^k = \tilde{\mathbf{u}}^k \mathbf{q}$  whereas case (2) can be proved by writing  $\mathbf{w}^k$  as a function of  $\tilde{\mathbf{u}}^k$  (from case (1)) and then applying chain rule. See Appendix D.  $\square$

Note that, in the implementation of BC, the auxiliary variables  $\tilde{\mathbf{w}}$  are clipped between  $[-1, 1]$  as it does not affect the sign function. In the  $\mathbf{u}$ -space, this clipping operation would translate into a projection to the polytope  $\mathcal{S}$ , meaning  $\tilde{\mathbf{w}} \in [-1, 1]$  implies  $\tilde{\mathbf{u}} \in \mathcal{S}$ , where  $\tilde{\mathbf{w}}$  and  $\tilde{\mathbf{u}}$  are related according to  $\tilde{\mathbf{w}} = \tilde{\mathbf{u}} \mathbf{q}$ . Even in this case, Proposition 3.1 holds, as the assumption  $\tilde{\mathbf{w}}^k = \tilde{\mathbf{u}}^k \mathbf{q}$  is still satisfied.

Since hardmax is a non-differentiable operation, the partial derivative  $\partial \mathbf{u} / \partial \tilde{\mathbf{u}} = \partial \text{hardmax} / \partial \tilde{\mathbf{u}}$  is not defined. However, to allow backpropagation, we write hardmax function in terms of the sign function, and used the straight-through-estimator [15] to allow gradient flow similar to binary connect. For details of the derivation please refer to Appendix D.1.

## 4. Related Work

There are many works on NN quantization focusing on different aspects such as quantizing parameters [8], activations [18], loss aware quantization [16] and quantization for specialized hardware [10], to name a few. Here we give a brief summary of latest works and for a comprehensive survey we refer the reader to [13].

In this work, we consider parameter quantization, which can be treated as a post-processing scheme [12] or incorporated into the learning process. Popular methods [8, 19] falls into the latter category and optimize the constrained problem using some form of projected stochastic gradient descent. In contrast to projection, quantization can also be enforced using a penalty term [40]. Even though, our method is based on projected gradient descent, by optimizing in the  $\mathbf{u}$ -space, we provide theoretical insights based on mean-field and bridge the gap between NN quantization and MRF optimization literature.

In contrast, the variational approach can also be used for quantization, where the idea is to learn a posterior probability on the network parameters in a Bayesian framework. In this family of methods, the quantized network can be obtained either via a quantizing prior [1] or using the MAP estimate on the learned posterior [37]. Interestingly, the learned posterior distribution can be used to estimate the model uncertainty and in turn determine the required precision for each network parameter [26]. Note that, even in our seemingly different method, we learn a probability distribution over the parameters (see Sec. 2.2) and it would be interesting to understand the connection between Bayesian methods and our algorithm.

## 5. Experiments

Since neural network binarization is popular, following the previous works [8, 32], we set the quantization levels to be binary, *i.e.*,  $\mathcal{Q} = \{-1, 1\}$ . However, our formulation is applicable to any predefined set of quantization levels given sufficient resources at training time. We would like to point out that, we consider quantizing all learnable parameters, meaning, all quantization algorithms result in 32 times less memory compared to the floating point counterparts.

We evaluate our Proximal Mean-Field (PMF) algorithm on MNIST, CIFAR-10, CIFAR-100 and TinyImageNet<sup>5</sup> classification datasets with convolutional and residual architectures and compare against the BC method [8]. Note that BC constitutes the closest and directly comparable baseline to PMF as both of them consider quantizing the parameters. Furthermore, many other methods have been developed based on BC by relaxing some of the constraints, *e.g.*, layer-wise scalars [32], and we believe, similar extensions are possible with our method as well. Our results show that

<sup>5</sup><https://tiny-imagenet.herokuapp.com/>

Dataset	Image	# class	Train / Val.	$b$	$K$
MNIST	$28 \times 28$	10	50k / 10k	100	20k
CIFAR-10	$32 \times 32$	10	45k / 5k	128	100k
CIFAR-100	$32 \times 32$	100	45k / 5k	128	100k
TinyImageNet	$64 \times 64$	200	100k / 10k	128	100k

Table 1: *Experiment setup. Here,  $b$  denotes the batch size and  $K$  denotes the total number of iterations used for all the methods.*

the binary networks obtained by PMF yield accuracies very close to the floating point counterparts while consistently outperforming the BC method.

### 5.1. Experiment Setup

The details of the datasets and their corresponding experiment setups are given in Table 1. In all the experiments, standard multi-class cross-entropy loss is minimized. MNIST is tested using LeNet-300 and LeNet-5, where the former consists of three fully-connected (FC) layers while the latter is composed of two convolutional and two FC layers. For CIFAR and TinyImageNet, VGG-16 [36] and ResNet-18 [14] architectures adapted for CIFAR dataset are used. In particular, for CIFAR experiments, similar to [25], the size of the FC layers of VGG-16 is set to 512 and no dropout layers are employed. For TinyImageNet, the stride of the first convolutional layer of ResNet-18 is set to 2 to handle the image size [17]. In all the models, batch normalization [20] (with no learnable parameters) and ReLU non-linearity are used. Except for MNIST, standard data augmentation is used (*i.e.*, random crop and horizontal flip) and weight decay is set to 0.0001 unless stated otherwise.

For all the algorithms, the hyperparameters such as the optimizer and the learning rate (also its schedule) are cross-validated using the validation set<sup>6</sup> and the chosen parameters are given in the supplementary material. The growth-rate  $\rho$  in Algorithm 1 (the multiplicative factor used to increase  $\beta$ ) is set to 1.05 except for MNIST (where it is 1.2) and  $\beta$  is increased (line 8 in Algorithm 1) every 100 iterations. Furthermore, since the original implementation of BC do not binarize all the learnable parameters, for fair comparison, we implemented BC in our experiment setting based on the publicly available code<sup>7</sup>. All the methods are trained from a random initialization and the model with the best validation performance is chosen for each method to report the performance on the test set. Our algorithm is implemented in PyTorch [31] and the code will be released upon publication.

<sup>6</sup>For TinyImageNet, since the ground truth labels for the test set were not available, validation set is used for both cross-validation and testing.

<sup>7</sup><https://github.com/itayhubara/BinaryNet.pytorch>

## 5.2. Results

The classification accuracies on the test set (both top-1 and top-5 accuracies) of both versions of our algorithm, namely, Proximal Mean-Field (PMF) and Proximal Iterative Conditional Modes (PICM), the baseline Binary Connect (BC) and the floating point Reference Network (REF) are reported in Table 2. The training curves for CIFAR-10 and CIFAR-100 with ResNet-18 are shown in Fig. 2. Note that our PMF algorithm consistently produces better results than other binarization methods and the degradation in performance to the full floating point reference network is minimal especially for small datasets. For larger datasets (*e.g.*, CIFAR-100), binarizing ResNet-18 results in much smaller degradation compared to VGG-16.

The superior performance of PMF against BC and PICM empirically validates the hypothesis that performing “noisy” projection via softmax and annealing the noise is indeed beneficial in the stochastic setting. Furthermore, even though PICM and BC are theoretically equivalent in the non-stochastic setting, PICM yields slightly better accuracies in all our experiments. We conjecture that this is due to the fact that in PICM, the training is performed on a larger network (*i.e.*, in the  $\mathbf{u}$ -space). Similar behaviour is empirically observed in [27].

To further consolidate our implementation of BC, we quote the accuracies reported in the original papers here. In [8], the top-1 accuracy on CIFAR-10 with a modified VGG type network is 90.10%. In the same setting, even with additional layer-wise scalars, (the Binary Weight Network (BWN) method [32]), the corresponding accuracy is 90.12%. Note that, in both the above cases, the last layer parameters and biases in all layers were not binarized.

### 5.3. Proximal Mean-field Analysis

To analyse the effect of storing the auxiliary variables  $\tilde{\mathbf{u}}$  in Algorithm 1, we evaluate PMF with and without storing  $\tilde{\mathbf{u}}$ . The results are reported in Table 3. Storing the auxiliary variables and updating them is in fact improves the overall performance. However, even without storing  $\tilde{\mathbf{u}}$ , PMF obtains reasonable performance, indicating the usefulness of our continuous relaxation. Note that, if the auxiliary variables are not stored in BC, it is impossible to train the network as the quantization error in the gradients are catastrophic and single gradient step is not sufficient to move from one discrete point to the next.

## 6. Discussion

In this work, we have formulated NN quantization as a discrete labelling problem and introduced a projected stochastic gradient descent algorithm to optimize it. By showing our approach as a proximal mean-field method, we have also provided an MRF optimization perspective to the

Dataset	Architecture	REF (Float) Top-1/5 (%)	BC Top-1/5 (%)	PICM Top-1/5 (%)	PMF Top-1/5 (%)	REF - PMF Top-1 (%)
MNIST	LeNet-300	98.55/99.93	98.05/99.93	98.18/99.91	<b>98.24</b> /99.97	+0.31
	LeNet-5	99.39/99.98	99.30/99.98	99.31/99.99	<b>99.44</b> /100.0	−0.05
CIFAR-10	VGG-16	93.01/99.38	86.40/98.43	88.96/99.17	<b>90.51</b> /99.56	+2.50
	ResNet-18	94.64/99.78	91.60/99.74	92.02/99.71	<b>92.55</b> /99.80	+2.09
CIFAR-100	VGG-16	70.33/88.58	43.70/73.43	45.65/74.70	<b>61.52</b> /85.83	+8.81
	ResNet-18	73.85/92.49	69.93/90.75	70.85/91.46	<b>71.85</b> /91.88	+2.00
TinyImageNet	ResNet-18	56.41/79.75	49.33/74.13	49.66/74.54	<b>50.78</b> /75.01	+5.63

Table 2: Classification accuracies on the test set for different methods. Note that our PMF algorithm consistently produces better results than other binarization methods and the degradation in performance to the full floating point network (last column) is minimal especially for small datasets. For larger datasets (e.g., CIFAR-100), binarizing ResNet-18 results in much smaller degradation compared to VGG-16. Even though, PICM and BC are theoretically equivalent in the non-stochastic setting, PICM yields slightly better accuracies. Note, all PMF, PICM, and BC require **32** times less memory compared to single-precision floating points networks at test time.

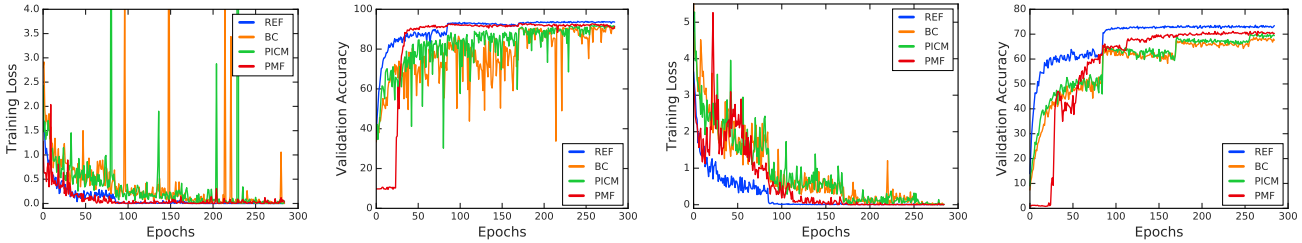


Figure 2: Training curves for CIFAR-10 (first two) and CIFAR-100 (last two) with ResNet-18. For quantization methods, the validation accuracy is always measured with the quantized networks. Specifically, in PMF, the hardmax projection is applied before the evaluation. The validation accuracies of PMF are the worst at the beginning (until 25 epochs). We believe, this is the exploration phase of PMF and hardmax projection introduces errors in this phase. While BC and PICM are extremely noisy, PMF training curves are fairly smooth and closely resembles the high-precision reference network (especially on CIFAR-10). Note that, in CIFAR-10, PMF surpasses REF for epochs between 30 – 80 and closely follows it afterwards.

Dataset	Architecture	PMF wo $\tilde{\mathbf{u}}$ Top-1/5 (%)	PMF Top-1/5 (%)
MNIST	LeNet-300	96.74/99.92	<b>98.24</b> /99.97
	LeNet-5	98.78/99.95	<b>99.44</b> /100.0
CIFAR-10	VGG-16	80.18/98.24	<b>90.51</b> /99.56
	ResNet-18	87.36/99.50	<b>92.55</b> /99.80

Table 3: Comparison of PMF with and without storing the auxiliary variables  $\tilde{\mathbf{u}}$ . Storing the auxiliary variables and updating them is in fact improves the overall performance. However, even without storing  $\tilde{\mathbf{u}}$ , PMF obtains reasonable performance, indicating the usefulness of our relaxation.

NN quantization problem. This connection to MRF opens up many interesting research directions primarily on consider-

ing dependency between the neural network parameters to derive better network quantization schemes. Furthermore, our PMF approach learns a probability distribution over the network parameters, which is similar in spirit to Bayesian deep learning methods. Therefore, we believe, it is interesting to explore the connection between Bayesian methods and our algorithm, which can potentially drive research in both the fields.

## 7. Acknowledgements

This work was supported by the ERC grant ERC-2012-AdG 321162-HELIOS, EPSRC grant Seebibyte EP/M013774/1 and EPSRC/MURI grant EP/N019474/1. We would also like to acknowledge the Royal Academy of Engineering and FiveAI.



# Appendices

Here, we provide the proofs of propositions and theorems stated in the main paper and a self-contained overview of the mean-field method. Later in Sec. E, we give the experimental details to cater reproducibility and also provide additional training curves to understand the convergence behaviour of our algorithm. In short, the behaviour observed in the main paper holds.

## A. Relationship between $\mathbf{w}$ -space and $\mathbf{u}$ -space

**Proposition A.1.** Let  $f(\mathbf{w})$  be a continuous function with  $\mathbf{w} = g(\mathbf{u}) = \mathbf{u}\mathbf{q}$ . Then a point  $\mathbf{u}^k \in \mathcal{S}$  is a local minimum of  $f \circ g$ , if and only if  $\mathbf{w}^k = \mathbf{u}^k\mathbf{q}$  is a local minimum of  $f$  in the region  $[q_{\min}, q_{\max}]^m$ .

*Proof.* We will prove this by contradiction. Let  $\bar{\mathbf{w}}$  be a local minimum around the neighbourhood of  $\mathbf{w}^k$ . Since, the function  $g : \mathcal{S} \rightarrow [q_{\min}, q_{\max}]^m$  is onto and continuous, there exists a  $\bar{\mathbf{u}}$  such that  $\bar{\mathbf{w}} = \bar{\mathbf{u}}\mathbf{q}$  in the neighbourhood of  $\mathbf{u}^k$ , and it satisfies  $f \circ g(\bar{\mathbf{u}}) < f \circ g(\mathbf{u}^k)$ . This is a contradiction, hence, if  $\mathbf{u}^k$  is a local minimum of  $f \circ g$ , then  $\mathbf{w}^k$  is a local minimum of  $f$  in the region  $[q_{\min}, q_{\max}]^m$ . Similarly, from  $\mathbf{w}$ -space to the  $\mathbf{u}$ -space can be proved.  $\square$

## B. Mean-field Method

For completeness we briefly review the underlying theory of the mean-field method. For in-depth details, we refer the interested reader to the Chapter 5 of [39]. Furthermore, for background on Markov Random Field (MRF), we refer the reader to the Chapter 2 of [2]. In this section, for better understanding, we use the notations from the main paper and highlight the similarities wherever possible.

**Markov Random Field.** Let  $\mathcal{W} = \{W_1, \dots, W_m\}$  be a set of random variables, where each random variable  $W_j$  takes a label  $w_j \in \mathcal{Q}$ . For a given labelling  $\mathbf{w} \in \mathcal{Q}^m$ , the energy associated with an MRF can be written as:

$$L(\mathbf{w}) = \sum_{C \in \mathcal{C}} L_C(\mathbf{w}), \quad (15)$$

where  $\mathcal{C}$  is the set of subsets (cliques) of  $\mathcal{W}$  and  $L_C(\mathbf{w})$  is a positive function (factor or clique potential) that depends only on the values  $w_j$  for  $j \in C$ . Now, the joint probability distribution over the random variables can be written as:

$$P(\mathbf{w}) = \frac{1}{Z} e^{-L(\mathbf{w})}, \quad (16)$$

where the normalization constant  $Z$  is usually referred to as the partition function. From Hammersley-Clifford theorem, for the factorization given in Eq. (15), the joint probability distribution  $P(\mathbf{w})$  can be shown to factorize over

each clique  $C \in \mathcal{C}$ , which is essentially the Markov property. However, this Markov property is not necessary to write Eq. (16) and in turn for our formulation, but since mean-field is usually described in the context of MRFs we provide it here for completeness. The objective of mean-field is to obtain the most probable configuration, which is equivalent to minimizing the energy  $L(\mathbf{w})$ .

**Mean-field Inference.** The basic idea behind mean-field is to approximate the intractable probability distribution  $P(\mathbf{w})$  with a tractable one. Specifically, mean-field obtains a fully-factorized distribution (*i.e.*, each random variable  $W_j$  is independent) closest to the true distribution  $P(\mathbf{w})$  in terms of KL-divergence. Let  $U(\mathbf{w}) = \prod_{j=1}^m U_j(w_j)$  denote a fully-factorized distribution. Recall, the variables  $\mathbf{u}$  introduced in Sec. 2.2 represent the probability of each weight  $W_j$  taking a label  $\lambda$ . Therefore, the distribution  $U$  can be represented using the variables  $\mathbf{u} \in \mathcal{S}$ , where  $\mathcal{S}$  is defined as:

$$\mathcal{S} = \left\{ \mathbf{u} \mid \sum_{\lambda} u_{j:\lambda} = 1, \quad \forall j, \quad u_{j:\lambda} \geq 0, \quad \forall j, \lambda \right\}. \quad (17)$$

Now, the KL-divergence between  $U$  and  $P$  can be written as:

$$\begin{aligned} \text{KL}(U||P) &= \sum_{\mathbf{w} \in \mathcal{Q}^m} U(\mathbf{w}) \log \frac{U(\mathbf{w})}{P(\mathbf{w})}, \quad (18) \\ &= \sum_{\mathbf{w} \in \mathcal{Q}^m} U(\mathbf{w}) \log U(\mathbf{w}) - \sum_{\mathbf{w} \in \mathcal{Q}^m} U(\mathbf{w}) \log P(\mathbf{w}), \\ &= -H(U) - \sum_{\mathbf{w} \in \mathcal{Q}^m} U(\mathbf{w}) \log \frac{e^{-L(\mathbf{w})}}{Z}, \quad \text{Eq. (16)}, \\ &= -H(U) + \sum_{\mathbf{w} \in \mathcal{Q}^m} U(\mathbf{w}) L(\mathbf{w}) + \log Z. \end{aligned}$$

Here,  $H(U)$  denotes the entropy of the fully-factorized distribution, which is exactly the one used in Theorem C.1. Specifically,

$$H(U) = H(\mathbf{u}) = - \sum_{j=1}^m \sum_{\lambda \in \mathcal{Q}} u_{j:\lambda} \log u_{j:\lambda}. \quad (19)$$

Furthermore, in Eq. (18), since  $Z$  is a constant, it can be removed from the minimization. Hence the final mean-field objective can be written as:

$$\begin{aligned} \min_U F(U) &:= \sum_{\mathbf{w} \in \mathcal{Q}^m} U(\mathbf{w}) L(\mathbf{w}) - H(U), \quad (20) \\ &= \mathbb{E}_U[L(\mathbf{w})] - H(U), \end{aligned}$$

where  $\mathbb{E}_U[L(\mathbf{w})]$  denotes the expected value of the loss  $L(\mathbf{w})$  over the distribution  $U(\mathbf{w})$ . Note that, the expected value of the loss can be written as a function of the variables  $\mathbf{u}$ . In particular,

$$E(\mathbf{u}) := \mathbb{E}_U[L(\mathbf{w})] = \sum_{\mathbf{w} \in \mathcal{Q}^m} U(\mathbf{w}) L(\mathbf{w}), \quad (21)$$

$$= \sum_{\mathbf{w} \in \mathcal{Q}^m} \prod_{j=1}^m u_{j:w_j} L(\mathbf{w}).$$

Now, the mean-field objective can be written as an optimization over  $\mathbf{u}$ :

$$\min_{\mathbf{u} \in \mathcal{S}} F(\mathbf{u}) := E(\mathbf{u}) - H(\mathbf{u}). \quad (22)$$

Computing this expectation  $E(\mathbf{u})$  in general is intractable as the sum is over an exponential number of elements ( $|\mathcal{Q}|^m$  elements, where  $m$  is usually in the order millions for an image or a neural network). However, for an MRF, the energy function  $L(\mathbf{w})$  can be factorized easily as in Eq. (15) (e.g., unary and pairwise terms) and  $E(\mathbf{u})$  can be computed fairly easily as the distribution  $U$  is also fully-factorized.

In mean-field, the above objective (22) is minimized iteratively using a fixed point update. This update is derived by writing the Lagrangian and setting the derivatives with respect to  $\mathbf{u}$  to zero. The derivation is very similar to the proof of Theorem C.1, and at iteration  $k$ , the mean-field update for each  $j \in \{1 \dots m\}$  can be written as:

$$u_{j:\lambda}^{k+1} = \frac{e^{-\frac{\partial E^k}{\partial u_{j:\lambda}}}}{\sum_{\mu} e^{-\frac{\partial E^k}{\partial u_{j:\mu}}}} \quad \forall \lambda \in \mathcal{Q}. \quad (23)$$

Here,  $\frac{\partial E^k}{\partial u_{j:\lambda}}$  denotes the gradient of  $E(\mathbf{u})$  with respect to  $u_{j:\lambda}$  evaluated at  $\mathbf{u}_{j:\lambda}^k$ . This update is repeated until convergence. Once the distribution  $U$  is obtained, finding the most probable configuration is straight forward, since  $U$  is a product of independent distributions over each random variable  $W_j$ . Note that, as most probable configuration is exactly the minimum label configuration, the mean-field method iteratively minimizes the actual energy function  $L(\mathbf{w})$ .

### C. Softmax based PGD as Proximal Mean-field

Recall from Sec. 2.3.1, the softmax based PGD update can be written as:

$$\mathbf{u}_j^{k+1} = \text{softmax}(\beta \tilde{\mathbf{u}}_j^{k+1}), \quad \text{where} \quad (24)$$

$$u_{j:\lambda}^{k+1} = \frac{e^{\beta(u_{j:\lambda}^k - \eta g_{j:\lambda}^k)}}{\sum_{\mu \in \mathcal{Q}} e^{\beta(u_{j:\mu}^k - \eta g_{j:\mu}^k)}} \quad \forall \lambda \in \mathcal{Q}.$$

Here,  $\tilde{\mathbf{u}}^{k+1} = \mathbf{u}^k - \eta \mathbf{g}^k$  with  $\eta > 0$ , and  $\beta > 0$ .

**Theorem C.1.** At iteration  $k$ , let  $\mathbf{u}^{k+1}$  be the point from the softmax based PGD update (24). Then,

$$\mathbf{u}^{k+1} = \underset{\mathbf{u} \in \mathcal{S}}{\text{argmin}} \eta \mathbb{E}_{\mathbf{u}} [\hat{L}^k(\mathbf{w})] - \langle \mathbf{u}^k, \mathbf{u} \rangle_F - \frac{1}{\beta} H(\mathbf{u}), \quad (25)$$

where  $\hat{L}^k(\mathbf{w})$  is the first-order Taylor approximation of  $L$  at  $\mathbf{w}^k = \mathbf{u}^k \mathbf{q}$  and  $\eta > 0$  is the learning rate.

*Proof.* We will first prove that  $\mathbb{E}_{\mathbf{u}} [\hat{L}^k(\mathbf{w})] = \langle \mathbf{g}_{\mathbf{u}}^k, \mathbf{u} \rangle_F + c$  for some constant  $c$ . From the definition of  $\hat{L}^k(\mathbf{w})$ ,

$$\begin{aligned} \hat{L}^k(\mathbf{w}) &= L(\mathbf{w}^k) + \langle \mathbf{g}_{\mathbf{w}}^k, \mathbf{w} - \mathbf{w}^k \rangle, \\ &= \langle \mathbf{g}_{\mathbf{w}}^k, \mathbf{w} \rangle + c, \end{aligned} \quad (26)$$

where  $c$  is a constant that does not depend on  $\mathbf{w}$ . Now, the expectation over  $\mathbf{u}$  can be written as:

$$\begin{aligned} \mathbb{E}_{\mathbf{u}} [\hat{L}^k(\mathbf{w})] &= \mathbb{E}_{\mathbf{u}} [\langle \mathbf{g}_{\mathbf{w}}^k, \mathbf{w} \rangle] + c, \\ &= \langle \mathbf{g}_{\mathbf{w}}^k, \mathbb{E}_{\mathbf{u}}[\mathbf{w}] \rangle + c, \\ &= \langle \mathbf{g}_{\mathbf{w}}^k, \mathbf{u} \mathbf{q} \rangle + c, \quad \text{Definition of } \mathbf{u}. \end{aligned} \quad (27)$$

We will now show that  $\langle \mathbf{g}_{\mathbf{w}}^k, \mathbf{u} \mathbf{q} \rangle = \langle \mathbf{g}_{\mathbf{u}}^k, \mathbf{u} \rangle_F$ . To see this, let us consider an element  $j \in \{1 \dots m\}$ ,

$$\begin{aligned} g_{w_j}^k \langle \mathbf{u}_j, \mathbf{q} \rangle &= g_{w_j}^k \langle \mathbf{q}, \mathbf{u}_j \rangle, \\ &= g_{w_j}^k \mathbf{q}^T \mathbf{u}_j, \\ &= g_{\mathbf{u}_j}^k \mathbf{u}_j, \quad \mathbf{g}_{\mathbf{u}}^k = \mathbf{g}_{\mathbf{w}}^k \mathbf{q}^T. \end{aligned} \quad (28)$$

From the above equivalence, Eq. (25) can now be written as:

$$\mathbf{u}^{k+1} = \underset{\mathbf{u} \in \mathcal{S}}{\text{argmin}} \eta \langle \mathbf{g}_{\mathbf{u}}^k, \mathbf{u} \rangle_F - \langle \mathbf{u}^k, \mathbf{u} \rangle_F - \frac{1}{\beta} H(\mathbf{u}). \quad (29)$$

Now, ignoring the condition  $u_{j:\lambda} \geq 0$  for now, the Lagrangian of Eq. (29) with dual variables  $z_j$  with  $j \in \{1 \dots m\}$  can be written as:<sup>8</sup>

$$\begin{aligned} F(\mathbf{u}, \mathbf{z}) &= \beta \eta \langle \mathbf{g}^k, \mathbf{u} \rangle_F - \beta \langle \mathbf{u}^k, \mathbf{u} \rangle_F - H(\mathbf{u}) + \\ &\quad \sum_j z_j \left( 1 - \sum_{\lambda} u_{j:\lambda} \right). \end{aligned} \quad (30)$$

Note that the objective function is multiplied by  $\beta > 0$ . Now, differentiating  $F(\mathbf{u}, \mathbf{z})$  with respect to  $\mathbf{u}$  and setting the derivatives to zero:

$$\begin{aligned} \frac{\partial F}{\partial u_{j:\lambda}} &= \beta \eta g_{j:\lambda}^k - \beta u_{j:\lambda}^k + 1 + \log u_{j:\lambda} - z_j = 0, \\ \log u_{j:\lambda} &= z_j - 1 + \beta u_{j:\lambda}^k - \beta \eta g_{j:\lambda}^k, \\ u_{j:\lambda} &= e^{z_j - 1} e^{\beta(u_{j:\lambda}^k - \eta g_{j:\lambda}^k)}. \end{aligned} \quad (31)$$

<sup>8</sup>For notational clarity, we denote the gradient of  $\tilde{L}$  with respect to  $\mathbf{u}$  evaluated at  $\mathbf{u}^k$  as  $\mathbf{g}^k$ , i.e.,  $\mathbf{g}^k := \mathbf{g}_{\mathbf{u}}^k$ .

Since  $\sum_{\mu} u_{j:\mu} = 1$ ,

$$\sum_{\mu} u_{j:\mu} = 1 = \sum_{\mu} e^{z_j-1} e^{\beta(u_{j:\mu}^k - \eta g_{j:\mu}^k)}, \quad (32)$$

$$e^{z_j-1} = \frac{1}{\sum_{\mu} e^{\beta(u_{j:\mu}^k - \eta g_{j:\mu}^k)}}.$$

Substituting in Eq. (32),

$$u_{j:\lambda} = \frac{e^{\beta(u_{j:\lambda}^k - \eta g_{j:\lambda}^k)}}{\sum_{\mu} e^{\beta(u_{j:\mu}^k - \eta g_{j:\mu}^k)}}. \quad (33)$$

Note that,  $u_{j:\lambda} \geq 0$  for all  $j \in \{1 \dots m\}$  and  $\lambda \in \mathcal{Q}$ , and therefore,  $\mathbf{u}$  is a fixed point of Eq. (25). Furthermore, this is exactly the same formula as the softmax based PGD update (24). Hence, the proof is complete.  $\square$

## D. Binary Connect as Proximal ICM

**Proposition D.1.** Consider BC and PICM with  $\mathbf{q} = [-1, 1]^T$  and  $\eta_{\mathbf{w}} > 0$ . For an iteration  $k > 0$ , if  $\tilde{\mathbf{w}}^k = \tilde{\mathbf{u}}^k \mathbf{q}$  then,

1. the projections in BC:  $\mathbf{w}^k = \text{sign}(\tilde{\mathbf{w}}^k)$  and PICM:  $\mathbf{u}^k = \text{hardmax}(\tilde{\mathbf{u}}^k)$  satisfy  $\mathbf{w}^k = \mathbf{u}^k \mathbf{q}$ .
2. if  $\eta_{\mathbf{u}} = \eta_{\mathbf{w}}/2$ , then the updated points after the gradient descent step in BC and PICM satisfy  $\tilde{\mathbf{w}}^{k+1} = \tilde{\mathbf{u}}^{k+1} \mathbf{q}$ .

*Proof.* 1. In the binary case ( $\mathcal{Q} = \{-1, 1\}$ ), for each  $j \in \{1 \dots m\}$ , the hardmax projection can be written as:

$$u_{j:-1}^k = \begin{cases} 1 & \text{if } \tilde{u}_{j:-1}^k \geq \tilde{u}_{j:1}^k \\ 0 & \text{otherwise} \end{cases},$$

$$u_{j:1}^k = 1 - u_{j:-1}^k. \quad (34)$$

Now, multiplying both sides by  $\mathbf{q}$ , and substituting  $\tilde{w}_j^k = \tilde{u}_j^k \mathbf{q}$ ,

$$\mathbf{u}_j^k \mathbf{q} = \begin{cases} -1 & \text{if } \tilde{w}_j^k = -1 \tilde{u}_{j:-1}^k + 1 \tilde{u}_{j:1}^k \leq 0 \\ 1 & \text{otherwise} \end{cases},$$

$$w_j^k = \text{sign}(\tilde{w}_j^k). \quad (35)$$

Hence,  $\mathbf{w}^k = \text{sign}(\tilde{\mathbf{w}}^k) = \text{hardmax}(\tilde{\mathbf{u}}^k) \mathbf{q}$ .

2. Since  $\mathbf{w}^k = \mathbf{u}^k \mathbf{q}$  from case (1) above, by chain rule the gradients  $\mathbf{g}_{\mathbf{w}}^k$  and  $\mathbf{g}_{\mathbf{u}}^k$  satisfy,

$$\mathbf{g}_{\mathbf{u}}^k = \mathbf{g}_{\mathbf{w}}^k \frac{\partial \mathbf{w}}{\partial \mathbf{u}} = \mathbf{g}_{\mathbf{w}}^k \mathbf{q}^T. \quad (36)$$

Similarly, from case (1) above, for each  $j \in \{1 \dots m\}$ ,

$$w_j^k = \text{sign}(\tilde{w}_j^k) = \text{sign}(\tilde{\mathbf{u}}_j^k \mathbf{q}) = \text{hardmax}(\tilde{\mathbf{u}}_j^k) \mathbf{q},$$

$$\frac{\partial w_j}{\partial \tilde{\mathbf{u}}_j} = \frac{\partial \text{sign}}{\partial \tilde{\mathbf{u}}_j} = \frac{\partial \text{sign}}{\partial \tilde{w}_j} \frac{\partial \tilde{w}_j}{\partial \tilde{\mathbf{u}}_j} = \frac{\partial \text{hardmax}}{\partial \tilde{\mathbf{u}}_j} \mathbf{q}. \quad (37)$$

Here, the partial derivatives are evaluated at  $\tilde{\mathbf{u}} = \tilde{\mathbf{u}}^k$  but omitted for notational clarity. Moreover,  $\frac{\partial w_j}{\partial \tilde{\mathbf{u}}_j}$  is a  $d$ -dimensional column vector,  $\frac{\partial \text{sign}}{\partial \tilde{w}_j}$  is a scalar, and  $\frac{\partial \text{hardmax}}{\partial \tilde{\mathbf{u}}_j}$  is a  $d \times d$  matrix. Since,  $\frac{\partial \tilde{w}_j}{\partial \tilde{\mathbf{u}}_j} = \mathbf{q}$  (similar to Eq. (36)),

$$\frac{\partial w_j}{\partial \tilde{\mathbf{u}}_j} = \frac{\partial \text{sign}}{\partial \tilde{w}_j} \mathbf{q} = \frac{\partial \text{hardmax}}{\partial \tilde{\mathbf{u}}_j} \mathbf{q}. \quad (38)$$

Now, consider the  $\mathbf{g}_{\mathbf{u}}^k$  for each  $j \in \{1 \dots m\}$ ,

$$\mathbf{g}_{\tilde{\mathbf{u}}_j}^k = \mathbf{g}_{\mathbf{u}_j}^k \frac{\partial \mathbf{u}_j}{\partial \tilde{\mathbf{u}}_j} = \mathbf{g}_{\mathbf{u}_j}^k \frac{\partial \text{hardmax}}{\partial \tilde{\mathbf{u}}_j}, \quad (39)$$

$$\mathbf{g}_{\tilde{\mathbf{u}}_j}^k \mathbf{q} = \mathbf{g}_{\mathbf{u}_j}^k \frac{\partial \text{hardmax}}{\partial \tilde{\mathbf{u}}_j} \mathbf{q}, \quad \text{multiplying by } \mathbf{q},$$

$$= g_{w_j}^k \mathbf{q}^T \frac{\partial \text{hardmax}}{\partial \tilde{\mathbf{u}}_j} \mathbf{q}, \quad \text{Eq. (36)},$$

$$= g_{w_j}^k \mathbf{q}^T \frac{\partial \text{sign}}{\partial \tilde{w}_j} \mathbf{q}, \quad \text{Eq. (38)},$$

$$= g_{w_j}^k \frac{\partial \text{sign}}{\partial \tilde{w}_j} \mathbf{q}^T \mathbf{q},$$

$$= g_{w_j}^k \mathbf{q}^T \mathbf{q}, \quad \frac{\partial \text{sign}}{\partial \tilde{w}_j} = \frac{\partial w_j}{\partial \tilde{w}_j},$$

$$= 2 g_{w_j}^k, \quad \mathbf{q} = [-1, 1]^T.$$

Now, consider the gradient descent step for  $\tilde{\mathbf{u}}$ , with  $\eta_{\mathbf{u}} = \eta_{\mathbf{w}}/2$ ,

$$\tilde{\mathbf{u}}^{k+1} = \tilde{\mathbf{u}}^k - \eta_{\mathbf{u}} \mathbf{g}_{\tilde{\mathbf{u}}}^k, \quad (40)$$

$$\begin{aligned} \tilde{\mathbf{u}}^{k+1} \mathbf{q} &= \tilde{\mathbf{u}}^k \mathbf{q} - \eta_{\mathbf{u}} \mathbf{g}_{\tilde{\mathbf{u}}}^k \mathbf{q}, \\ &= \tilde{\mathbf{w}}^k - \eta_{\mathbf{u}} 2 \mathbf{g}_{\mathbf{w}}^k, \\ &= \tilde{\mathbf{w}}^k - \eta_{\mathbf{w}} \mathbf{g}_{\mathbf{w}}^k, \\ &= \tilde{\mathbf{w}}^{k+1}. \end{aligned}$$

Hence, the proof is complete.  $\square$

## D.1. Approximate Gradients through Hardmax Projection

In previous works [8, 32], to allow back propagation through the sign function, the straight-through-estimator [15] is used. Precisely, the partial derivative with respect to the sign function is defined as:

$$\frac{\partial \text{sign}(r)}{\partial r} := \mathbb{1}[|r| \leq 1]. \quad (41)$$

To make use of this, we intend to write the projection function hardmax in terms of the sign function. To this end, from Eq. (34), for each  $j \in \{1 \dots m\}$ ,

$$u_{j:-1}^k = \begin{cases} 1 & \text{if } \tilde{u}_{j:-1}^k - \tilde{u}_{j:1}^k \geq 0 \\ 0 & \text{otherwise} \end{cases}, \quad (42)$$

$$u_{j:1}^k = 1 - u_{j:-1}^k. \quad (43)$$

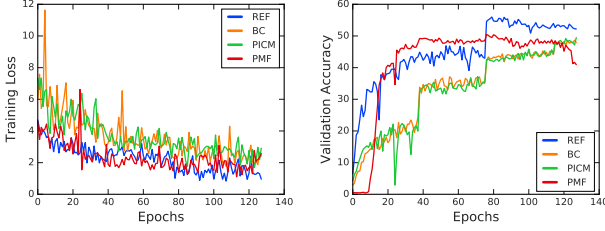


Figure 4: Training curves for TinyImageNet with ResNet-18. Compared to CIFAR-10/100 plots in Fig. 3, this plot is less noisy but the behaviour is roughly the same. PMF loss curve closely follows the loss curve of high-precision reference network and PMF even surpasses REF in validation accuracy for epochs between 20 – 80.

Hence, the projection  $\text{hardmax}(\tilde{\mathbf{u}}^k)$  for each  $j$  can be written as:

$$u_{j:-1}^k = \frac{\text{sign}(\tilde{u}_{j:-1}^k - \tilde{u}_{j:1}^k) + 1}{2}, \quad (44)$$

$$u_{j:1}^k = \frac{1 - \text{sign}(\tilde{u}_{j:-1}^k - \tilde{u}_{j:1}^k)}{2}. \quad (45)$$

Now, using Eq. (41), we can write:

$$\frac{\partial \mathbf{u}_j}{\partial \tilde{\mathbf{u}}_j} \bigg|_{\tilde{\mathbf{u}}_j = \tilde{\mathbf{u}}_j^k} = \frac{1}{2} \begin{bmatrix} \mathbb{1}[|v_j^k| \leq 1] & -\mathbb{1}[|v_j^k| \leq 1] \\ -\mathbb{1}[|v_j^k| \leq 1] & \mathbb{1}[|v_j^k| \leq 1] \end{bmatrix}, \quad (46)$$

where  $v_j^k = \tilde{u}_{j:-1}^k - \tilde{u}_{j:1}^k$ .

## E. Experimental Details

To enable reproducibility, we first give the hyperparameter settings used to obtain the results reported in the main paper in Table 4. Then, we provide additional training curves of different methods in Figs. 3 and 4 for better understanding of the convergence behaviour.

## References

- [1] J. Achterhold, J. M. Kohler, A. Schmeink, and T. Genewein. Variational network quantization. *ICLR*, 2018. 6
- [2] T. Ajanthan. *Optimization of Markov random fields in computer vision*. PhD thesis, Australian National University, 2017. 1, 9
- [3] T. Ajanthan, A. Desmaison, R. Bunel, M. Salzmann, P. H. S. Torr, and M. P. Kumar. Efficient linear programming for dense CRFs. *CVPR*, 2017. 3
- [4] J. Besag. On the statistical analysis of dirty pictures. *Journal of the Royal Statistical Society.*, 1986. 5
- [5] A. Blake, P. Kohli, and C. Rother. *Markov random fields for vision and image processing*. Mit Press, 2011. 1
- [6] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2009. 4
- [7] C. Chekuri, S. Khanna, J. Naor, and L. Zosin. A linear programming formulation and approximation algorithms for the metric labeling problem. *SIAM Journal on Discrete Mathematics*, 2004. 3
- [8] M. Courbariaux, Y. Bengio, and J.-P. David. Binaryconnect: Training deep neural networks with binary weights during propagations. *NIPS*, 2015. 1, 2, 5, 6, 7, 11
- [9] P. K. Dokania and P. K. Mudigonda. Parsimonious labeling. *ICCV*, 2015. 1
- [10] S. K. Esser, R. Appuswamy, P. A. Merolla, J. V. Arthur, and D. S. Modha. Backpropagation for energy-efficient neuromorphic computing. *NIPS*, 2015. 6
- [11] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval research logistics quarterly*, 1956. 3
- [12] Y. Gong, L. Liu, and L. Bourdev. Compressing deep convolutional networks using vector quantization. *arXiv preprint arXiv:1412.6115*, 2014. 6
- [13] Y. Guo. A survey on methods and theories of quantized neural networks. *arXiv preprint arXiv:1808.04752*, 2018. 6
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CVPR*, 2016. 7
- [15] G. Hinton. Neural networks for machine learning. *Coursera, video lectures*, 2012. 6, 11
- [16] L. Hou, Q. Yao, and J. T. Kwok. Loss-aware binarization of deep networks. *ICLR*, 2017. 6
- [17] G. Huang, Y. Li, G. Pleiss, Z. Liu, J. E. Hopcroft, and K. Q. Weinberger. Snapshot ensembles: Train 1, get m for free. *ICLR*, 2017. 7
- [18] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Binarized neural networks. *NIPS*, 2016. 6
- [19] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *JMLR*, 2017. 1, 4, 5, 6
- [20] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, 2015. 7
- [21] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 3
- [22] J. Kleinberg and E. Tardos. Approximation algorithms for classification problems with pairwise relationships: metric labeling and Markov random fields. *Journal of the ACM*, 2002. 3
- [23] V. Kolmogorov and R. Zabini. What energy functions can be minimized via graph cuts? *PAMI*, 2004. 2
- [24] S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-coordinate Frank-Wolfe optimization for structural SVMs. *ICML*, 2012. 3
- [25] N. Lee, T. Ajanthan, and P. H. S. Torr. SNIP: Single-shot network pruning based on connection sensitivity. *arXiv preprint arXiv:1810.02340*, 2018. 7
- [26] C. Louizos, K. Ullrich, and M. Welling. Bayesian compression for deep learning. *NIPS*, 2017. 6
- [27] A. Mishra, E. Nurvitadhi, J. J. Cook, and D. Marr. WRPN: Wide reduced-precision networks. *ICLR*, 2018. 7
- [28] P. K. Mudigonda. *Combinatorial and convex optimization for probabilistic models in computer vision*. PhD thesis, Oxford Brookes University, 2008. 1



Hyperparameter	MNIST with LeNet-300/5				TinyImageNet with ResNet-18			
	REF	BC	PICM	PMF	REF	BC	PICM	PMF
learning_rate	0.001	0.001	0.001	0.001	0.1	0.0001	0.0001	0.001
lr_decay	step	step	step	step	step	step	step	step
lr_interval	7k	7k	7k	7k	60k	30k	30k	30k
lr_scale	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2
momentum	-	-	-	-	0.9	-	-	-
optimizer	Adam	Adam	Adam	Adam	SGD	Adam	Adam	Adam
weight_decay	0	0	0	0	0.0001	0.0001	0.0001	0.0001

	CIFAR-10 with VGG-16				CIFAR-10 with ResNet-18			
	REF	BC	PICM	PMF	REF	BC	PICM	PMF
learning_rate	0.1	0.0001	0.0001	0.001	0.1	0.0001	0.0001	0.001
lr_decay	step	step	step	step	step	step	step	step
lr_interval	30k	30k	30k	30k	30k	30k	30k	30k
lr_scale	0.2	0.2	0.2	0.2	0.2	0.2	0.2	0.2
momentum	0.9	-	-	-	0.9	-	-	-
optimizer	SGD	Adam	Adam	Adam	SGD	Adam	Adam	Adam
weight_decay	0.0005	0.0001	0.0001	0.0001	0.0005	0.0001	0.0001	0.0001

	CIFAR-100 with VGG-16				CIFAR-100 with ResNet-18			
	REF	BC	PICM	PMF	REF	BC	PICM	PMF
learning_rate	0.1	0.01	0.01	0.0001	0.1	0.0001	0.0001	0.001
lr_decay	step	multi-step	multi-step	step	step	step	step	multi-step
lr_interval	30k	20k,30k,40k, 50k,60k,70k,80k	20k,30k,40k, 50k,60k,70k,80k	30k	30k	30k	30k	30k,40k,50k, 60k,70k,80k
lr_scale	0.2	0.5	0.5	0.2	0.1	0.2	0.2	0.5
momentum	0.9	0.9	0.9	-	0.9	-	-	0.95
optimizer	SGD	SGD	SGD	Adam	SGD	Adam	Adam	SGD
weight_decay	0.0005	0.0001	0.0001	0.0001	0.0005	0.0001	0.0001	0.0001

Table 4: Hyperparameter settings used for the experiments. Here, if  $\text{lr\_decay} == \text{step}$ , then the learning rate is multiplied by  $\text{lr\_scale}$  for every  $\text{lr\_interval}$  iterations. On the other hand, if  $\text{lr\_decay} == \text{multi-step}$ , the learning rate is multiplied by  $\text{lr\_scale}$  whenever the iteration count reaches any of the milestones specified by  $\text{lr\_interval}$ .

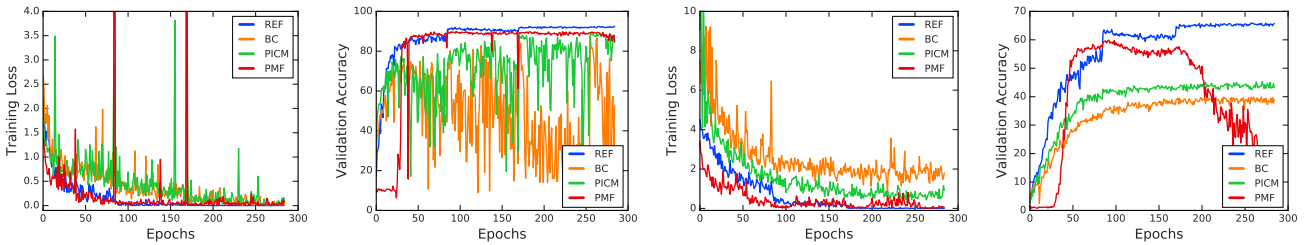


Figure 3: Training curves for CIFAR-10 (first two) and CIFAR-100 (last two) with VGG-16 (corresponding ResNet-18 plots are in Fig. 2). Similar to the main paper, while BC and PICM are extremely noisy, PMF training curves are fairly smooth and closely resembles the high-precision reference network. The validation accuracy plot for CIFAR-100 for PMF starts to decrease after 180 epochs (while training loss oscillates around a small value), this could be interpreted as overfitting to the training set.

[29] G. L. Nemhauser and L. A. Wolsey. *Integer programming and combinatorial optimization*. Springer, 1988. 2

[30] N. Parikh and S. P. Boyd. Proximal algorithms. *Foundations*

and Trends in Optimization, 2014. 3

[31] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Auto-

- matic differentiation in PyTorch. 2017. 7
- [32] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. *ECCV*, 2016. 1, 2, 6, 7, 11
  - [33] P. Ravikumar, A. Agarwal, and M. J. Wainwright. Message-passing for graph-structured linear programs: proximal projections, convergence and rounding schemes. *ICML*, 2008. 5
  - [34] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 1951. 3
  - [35] L. Rosasco, S. Villa, and B. C. Vũ. Convergence of stochastic proximal gradient algorithm. *arXiv preprint arXiv:1403.5074*, 2014. 3
  - [36] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2015. 7
  - [37] D. Soudry, I. Hubara, and R. Meir. Expectation backpropagation: Parameter-free training of multilayer neural networks with continuous or discrete weights. *NIPS*, 2018. 6
  - [38] O. Veksler. *Efficient graph-based energy minimization methods in computer vision*. PhD thesis, Cornell University New York, USA, 1999. 1
  - [39] M. J. Wainwright, M. I. Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 2008. 1, 4, 9
  - [40] P. Yin, S. Zhang, J. Lyu, S. Osher, Y. Qi, and J. Xin. Binaryrelax: A relaxation approach for training deep neural networks with quantized weights. *arXiv preprint arXiv:1801.06313*, 2018. 6