

---

# PROXQUANT: Quantized Neural Networks via Proximal Operators

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 To make deep neural networks feasible in resource-constrained environments (such  
2 as mobile devices), it is beneficial to quantize models by using low-precision  
3 weights. One common technique for quantizing neural networks is the straight-  
4 through gradient method, which enables back-propagation through the quantization  
5 mapping. Despite its empirical success, little is understood about why the straight-  
6 through gradient method works.

7 Building upon a novel observation that the straight-through gradient method is in  
8 fact *identical* to the well-known Nesterov’s dual-averaging algorithm on a quanti-  
9 zation constrained optimization problem, we propose a more principled alternative  
10 approach, called PROXQUANT, that formulates quantized network training as a  
11 regularized learning problem instead and optimizes it via the prox-gradient method.  
12 PROXQUANT does back-propagation on the underlying full-precision vector and  
13 applies an efficient prox-operator in between stochastic gradient steps to encourage  
14 quantizedness. For quantizing ResNets and LSTMs, PROXQUANT outperforms  
15 state-of-the-art results on binary quantization and is on par with state-of-the-art on  
16 multi-bit quantization. For binary quantization, our analysis shows both theoret-  
17 ically and experimentally that PROXQUANT is more stable than the straight-through  
18 gradient method (i.e. BinaryConnect), challenging the indispensability of the  
19 straight-through gradient method and providing a powerful alternative.

## 20 1 Introduction

21 In this paper, we formulate the problem of model quantization as a regularized learning problem and  
22 propose to solve it with a proximal gradient method. Our contributions are summarized as follows.

- 23 • We present a unified framework for defining regularization functionals that encourage  
24 binary, ternary, and multi-bit quantized parameters, through penalizing the distance to  
25 quantized sets. For binary quantization, the resulting regularizer is a  $W$ -shaped non-smooth  
26 regularizer, which shrinks parameters towards either  $-1$  or  $1$  in the same way that the  $L_1$   
27 norm regularization shrinks parameters towards  $0$ . We demonstrate that the prox-operators  
28 for regularizers that come out of our framework often admit linear-time solutions (or linear  
29 time approximation heuristics) which result in numerically *exact* quantized parameters.
- 30 • We propose training quantized networks using PROXQUANT (Algorithm 1) — a stochastic  
31 proximal gradient method with a homotopy scheme. Compared with the straight-through gra-  
32 dient method, PROXQUANT has access to additional gradient information at non-quantized  
33 points, and its homotopy scheme prevents potential overshoot early in the training. Algo-  
34 rithmically, PROXQUANT involves just adding a simple proximal step with respect to a  
35 quantization-inducing regularizer after each stochastic gradient step, thus can be efficiently

36 implemented under any major deep learning frameworks without incurring significant system overhead and be used as a modular component to add to the training pipeline of any  
 37 deep networks to result in a quantized network.  
 38

- 39 • We demonstrate the effectiveness and flexibility of PROXQUANT through systematic experiments on (1) image classification with ResNets (Section 3.1); (2) language modeling with  
 40 LSTMs (Section 3.2). The PROXQUANT method outperforms the state-of-the-art results on binary quantization and is comparable with the state-of-the-art on ternary and multi-bit  
 41 quantization.  
 42
- 43 • For binary nets, we show that BinaryConnect suffers from more optimization instability than PROXQUANT through (1) a theoretical characterization of convergence for BinaryConnect  
 44 (Section 4.1) and (2) a sign change experiment on CIFAR-10 (Section G). Experimentally, PROXQUANT finds better binary nets that is also closer to the initialization in the sign  
 45 change metric.  
 46

47 We present the main ingredients of our contribution in this extended abstract. See the Appendices B  
 48 for the prior work, C for the notation, A and D for the motivation and preliminary discussions about the straight-through gradient method and prox operators.  
 49

## 52 2 Quantized net training via regularized learning

53 We propose the PROXQUANT algorithm, which adds a quantization-inducing regularizer onto the  
 54 loss and optimizes via the (non-lazy) prox-gradient method with a finite  $\lambda$ . The prototypical version  
 55 of PROXQUANT is described in Algorithm 1.

---

**Algorithm 1** PROXQUANT: Prox-gradient method for quantized net training

---

**Require:** Regularizer  $R$  that induces desired quantizedness, initialization  $\theta_0$ , learning rates  $\{\eta_t\}_{t \geq 0}$ ,  
 regularization strengths  $\{\lambda_t\}_{t \geq 0}$

**while** not converged **do**

Perform the prox-gradient step

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \left\{ L(\theta_t) + \langle \theta - \theta_t, \tilde{\nabla} L(\theta_t) \rangle + \frac{1}{2\eta_t} \|\theta - \theta_t\|_2^2 + \lambda_t R(\theta) \right\} \quad (1)$$

$$= \text{prox}_{\eta_t \lambda_t R} \left( \theta_t - \eta_t \tilde{\nabla} L(\theta_t) \right). \quad (2)$$

The inner SGD step in eq. (2) can be replaced by any preferred stochastic optimization method such as Momentum SGD or Adam [Kingma and Ba, 2014].

**end while**

---

56 Compared to usual full-precision training, PROXQUANT only adds a prox step after each stochastic  
 57 gradient step, hence can be implemented straightforwardly upon existing full-precision training. As  
 58 the prox step does not need to know how the gradient step is performed, our method adapts to other  
 59 stochastic optimizers as well such as Adam. Further, each iteration is a prox-gradient step over the  
 60 objective  $L(\theta) + \lambda_t R(\theta)$  with learning rates  $\eta_t$ , and by choosing  $(\eta_t, \lambda_t)$  we obtain a joint control  
 61 over the speed of training and falling onto the quantized set.

62 Details of choosing the regularizer  $R$ , deriving the prox-operator  $\text{prox}_{\lambda R}$ , and choosing the regular-  
 63 ization strength are deferred to Appendix E.

## 64 3 Experiments

65 We evaluate the performance of PROXQUANT on two tasks: image classification with ResNets, and  
 66 language modeling with LSTMs. On both tasks, we show that the default straight-through gradient  
 67 method is not the only choice, and our PROXQUANT can achieve the same and often better results.

### 68 3.1 Image classification on CIFAR-10

69 **Problem setup** We perform image classification on the CIFAR-10 dataset, which contains 50000  
 70 training images and 10000 test images of size 32x32. We apply a commonly used data augmentation  
 71 strategy (pad by 4 pixels on each side, randomly crop to 32x32, do a horizontal flip with probability  
 72 0.5, and normalize). Our models are ResNets [He et al., 2016] of depth 20, 32, and 44 with weights  
 73 quantized to binary or ternary.

74 **Method** We use PROXQUANT with suitable regularizers in the binary case and the ternary case,  
 75 which we respectively denote as PQ-B and PQ-T. The training is initialized at pre-trained full-  
 76 precision nets (warm-start). For the regularization strength we use the homotopy method  $\lambda_t = \lambda \cdot t$   
 77 with  $\lambda = 10^{-4}$ . We initialize at pre-trained full-precision networks and use the Adam optimizer  
 78 with constant learning rate 0.01. To accelerate training in the final stage, we do a hard quantization  
 79  $\theta \mapsto q(\theta)$  at epoch 400 and keeps training till the 600-th epoch to stabilize the BatchNorm layers.

80 We compare with BinaryConnect (BC) for binary nets and Trained Ternary Quantization (TTQ) [Zhu  
 81 et al., 2016] for ternary nets. For BinaryConnect, we haven’t found reported results with ResNets  
 82 on CIFAR-10, and we train with the recommended Adam optimizer with learning rate decay [Cour-  
 83 bariaux et al., 2015] (initial learning rate 0.01, multiply by 0.1 at epoch 81 and 122, hard-quantize at  
 84 epoch 400), which we find leads to the best result for BinaryConnect.

85 **Result** The top-1 classification errors are reported in Table 1. For binary nets, our PROXQUANT-  
 86 Binary consistently yields better results than BinaryConnect. For ternary nets, our results are  
 87 comparable with the reported results of TTQ, and the best performance of our method over 4 runs  
 88 (from the same initialization) is slightly better than TTQ.

Table 1: Top-1 classification error of quantized ResNets on CIFAR-10. Performance is reported in mean(std) over 4 runs, where for PQ-T we report in addition the best of 4 (Bo4).

Model (Bits)	Full-Precision (32)	BC (1)	PQ-B (ours) (1)	TTQ (2)	PQ-T (ours) (2)	PQ-T (Bo4) (2)
ResNet-20	8.06	9.49 (0.22)	<b>9.15</b> (0.21)	8.87	<b>8.40</b> (0.13)	8.22
ResNet-32	7.25	8.66 (0.36)	<b>8.40</b> (0.23)	7.63	7.65 (0.15)	7.53
ResNet-44	6.96	8.26 (0.24)	<b>7.79</b> (0.06)	7.02	7.05 (0.08)	6.98

### 89 3.2 Language modeling with LSTMs

90 See Appendix F for details.

## 91 4 Stability analysis of binary quantization

### 92 4.1 Convergence characterization for BinaryConnect

93 We now show that BinaryConnect has a very stringent convergence condition. Consider the Bina-  
 94 ryConnect method with batch gradients:

$$s_t = \text{sign}(\theta_t), \quad \theta_{t+1} = \theta_t - \eta_t \nabla L(s_t). \quad (3)$$

95 **Definition 4.1** (Fixed point and convergence). *We say that  $s \in \{\pm 1\}^d$  is a **fixed point** of the*  
 96 *BinaryConnect algorithm, if  $s_0 = s$  in eq. (3) implies that  $s_t = s$  for all  $t = 1, 2, \dots$ . We say that the*  
 97 *BinaryConnect algorithm **converges** if there exists  $t < \infty$  such that  $s_t$  is a fixed point.*

98 **Theorem 4.1.** *Assume that the learning rates satisfy  $\sum_{t=0}^{\infty} \eta_t = \infty$ , then  $s \in \{\pm 1\}^d$  is a fixed*  
 99 *point for BinaryConnect eq. (3) if and only if  $\text{sign}(\nabla L(s)[i]) = -s[i]$  for all  $i \in [d]$  such that*  
 100  *$\nabla L(\theta)[i] \neq 0$ . Such a point may not exist, in which case BinaryConnect does not converge for any*  
 101 *initialization  $\theta_0 \in \mathbb{R}^d$ .*

102 In a sign change experiment on CIFAR-10 (see Appendix G), we are going to see that BinaryConnect  
 103 indeed fails to converge to a fixed sign pattern, corroborating Theorem 4.1.

104 **References**

- 105 A. G. Anderson and C. P. Berg. The high-dimensional geometry of binary neural networks. *arXiv*  
106 *preprint arXiv:1705.07199*, 2017.
- 107 M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- 108 M. A. Carreira-Perpinán. Model compression as constrained optimization, with application to neural  
109 nets. part i: General framework. *arXiv preprint arXiv:1707.01209*, 2017.
- 110 M. A. Carreira-Perpinán and Y. Idelbayev. Model compression as constrained optimization, with  
111 application to neural nets. part ii: Quantization. *arXiv preprint arXiv:1707.04319*, 2017.
- 112 M. Courbariaux, Y. Bengio, and J.-P. David. BinaryConnect: Training deep neural networks with  
113 binary weights during propagations. In *Advances in neural information processing systems*, pages  
114 3123–3131, 2015.
- 115 Y. Ding, J. Liu, and Y. Shi. On the universal approximability of quantized relu neural networks. *arXiv*  
116 *preprint arXiv:1802.03646*, 2018.
- 117 S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with  
118 pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- 119 S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally. EIE: Efficient inference  
120 engine on compressed deep neural network. In *Computer Architecture (ISCA), 2016 ACM/IEEE*  
121 *43rd Annual International Symposium on*, pages 243–254. IEEE, 2016.
- 122 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings*  
123 *of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 124 S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780,  
125 1997.
- 126 L. Hou and J. T. Kwok. Loss-aware weight quantization of deep networks. In *International Conference*  
127 *on Learning Representations*, 2018. URL <https://openreview.net/forum?id=BkrSv01A->.
- 128 I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Quantized neural networks:  
129 Training neural networks with low precision weights and activations. *Journal of Machine Learning*  
130 *Research*, 18:187–1, 2017.
- 131 D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
132 2014.
- 133 F. Li and B. Liu. Ternary weight networks. *arXiv preprint arXiv:1605.04711*, 2016.
- 134 H. Li, S. De, Z. Xu, C. Studer, H. Samet, and T. Goldstein. Training quantized nets: A deeper  
135 understanding. In *Advances in Neural Information Processing Systems*, pages 5811–5821, 2017.
- 136 M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini. Building a large annotated corpus of english:  
137 The penn treebank. *Computational linguistics*, 19(2):313–330, 1993.
- 138 N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends® in Optimization*, 1(3):  
139 127–239, 2014.
- 140 M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. Xnor-net: Imagenet classification using binary  
141 convolutional neural networks. In *European Conference on Computer Vision*, pages 525–542.  
142 Springer, 2016.
- 143 J. Sun and X. Sun. Adversarial probabilistic regularization. Unpublished draft, 2018.
- 144 R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*  
145 *Society. Series B (Methodological)*, pages 267–288, 1996.
- 146 L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal*  
147 *of Machine Learning Research*, 11(Oct):2543–2596, 2010.

- 148 C. Xu, J. Yao, Z. Lin, W. Ou, Y. Cao, Z. Wang, and H. Zha. Alternating multi-bit quantization for  
 149 recurrent neural networks. In *International Conference on Learning Representations*, 2018. URL  
 150 <https://openreview.net/forum?id=S19dR9x0b>.
- 151 S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou. Dorefa-net: Training low bitwidth convolutional  
 152 neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160*, 2016.
- 153 C. Zhu, S. Han, H. Mao, and W. J. Dally. Trained ternary quantization. *arXiv preprint*  
 154 *arXiv:1612.01064*, 2016.

## 155 A Drawback of the straight-through gradient method

156 Typically, training a quantized network involves (1) the design of a *quantizer*  $q$  that maps a  
 157 full-precision parameter to a  $k$ -bit quantized parameter, and (2) the *straight-through gradient*  
 158 *method* [Courbariaux et al., 2015] that enables back-propagation from the quantized parameter  
 159 back onto the original full-precision parameter, which is critical to the success of quantized network  
 160 training. With quantizer  $q$ , an iterate of the straight-through gradient method (see Figure 1a) proceeds  
 161 as  $\theta_{t+1} = \theta_t - \eta_t \widetilde{\nabla} L(\theta)|_{\theta=q(\theta_t)}$ , and  $q(\widehat{\theta})$  (for the converged  $\widehat{\theta}$ ) is taken as the output model. For  
 162 training binary networks, choosing  $q(\cdot) = \text{sign}(\cdot)$  gives the BinaryConnect method [Courbariaux  
 163 et al., 2015].

164 Though appealingly simple and empirically effective, it is information-theoretically rather mysterious  
 165 why the straight-through gradient method works well, at least in the binary case: while the goal is  
 166 to find a parameter  $\theta \in \{\pm 1\}^d$  with low loss, the algorithm only has access to stochastic gradients  
 167 at  $\{\pm 1\}^d$ . As this is a discrete set, *a priori*, gradients in this set do not necessarily contain any  
 168 information about the function values. Indeed, a simple one-dimensional example (Figure 1b) shows  
 169 that BinaryConnect fails to find the minimizer of fairly simple convex Lipschitz functions in  $\{\pm 1\}$ ,  
 170 due to a lack of gradient information in between.

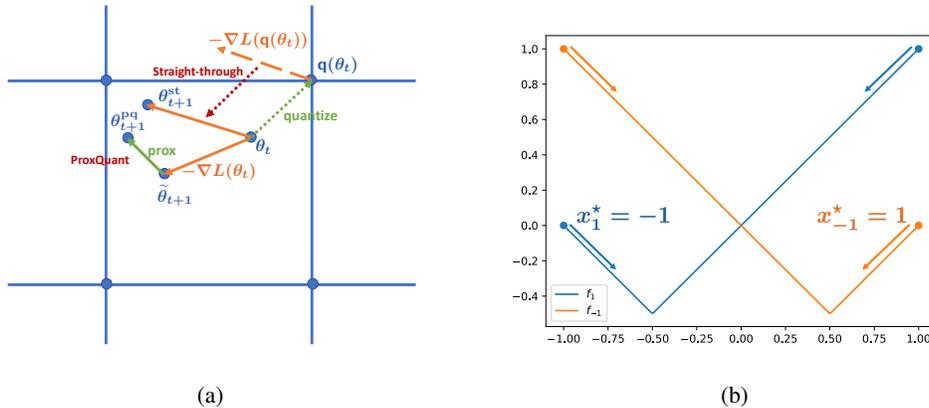


Figure 1: (a) Comparison of the straight-through gradient method and our PROXQUANT method. The straight-through method computes the gradient at the quantized vector and performs the update at the original real vector; PROXQUANT performs a gradient update at the current real vector followed by a prox step which encourages quantizedness. (b) A two-function toy failure case for BinaryConnect. The two functions are  $f_1(x) = |x + 0.5| - 0.5$  (blue) and  $f_{-1}(x) = |x - 0.5| - 0.5$  (orange). The derivatives of  $f_1$  and  $f_{-1}$  coincide at  $\{-1, 1\}$ , so any algorithm that only uses this information will have identical behaviors on these two functions. However, the minimizers in  $\{\pm 1\}$  are  $x_1^* = -1$  and  $x_{-1}^* = 1$ , so the algorithm must fail on one of them.

## 171 B Prior work

172 **Methodologies** Han et al. [2015] propose Deep Compression, which compresses a DNN via  
 173 sparsification, nearest-neighbor clustering, and Huffman coding. This architecture is then made into  
 174 a specially designed hardware for efficient inference [Han et al., 2016]. In a parallel line of work,  
 175 Courbariaux et al. [2015] propose BinaryConnect that enables the training of binary neural networks,

176 and Li and Liu [2016], Zhu et al. [2016] extend this method into ternary quantization. Training and  
 177 inference on quantized nets can be made more efficient by also quantizing the activation [Hubara  
 178 et al., 2017, Rastegari et al., 2016, Zhou et al., 2016], and such networks have achieved impressive  
 179 performance on large-scale tasks such as ImageNet classification [Rastegari et al., 2016, Zhu et al.,  
 180 2016]. In the NLP land, quantized language models have been successfully trained using alternating  
 181 multi-bit quantization [Xu et al., 2018].

182 **Theories** Li et al. [2017] prove the convergence rate of stochastic rounding and BinaryConnect  
 183 on convex problems and demonstrate the advantage of BinaryConnect over stochastic rounding on  
 184 non-convex problems. Anderson and Berg [2017] demonstrate the effectiveness of binary networks  
 185 through the observation that the angles between high-dimensional vectors are approximately preserved  
 186 when binarized, and thus high-quality feature extraction with binary weights is possible. Ding et al.  
 187 [2018] show a universal approximation theorem for quantized ReLU networks.

188 **Principled methods** Hou and Kwok [2018] propose a proximal Newton algorithm for model quan-  
 189 tization, which makes use of the additional Hessian information but is hence slightly more expensive  
 190 in each iteration. Sun and Sun [2018] perform model quantization through a Wasserstein regulariza-  
 191 tion term and minimize via the adversarial representation, similar as in Wasserstein GANs [Arjovsky  
 192 et al., 2017]. Their method has the potential of generalizing to other generic requirements on the  
 193 parameter, but might be hard to tune due to the instability of the inner maximization problem.

194 While preparing this manuscript, we discovered the independent work of Carreira-Perpinán [2017],  
 195 Carreira-Perpinán and Idelbayev [2017]. They formulate quantized network training as a constrained  
 196 optimization problem and propose to solve them via augmented Lagrangian methods. From an  
 197 optimization perspective, our views are largely complementary: they treat the quantization as a  
 198 constraint, whereas we encourage quantization through a regularizer. Due to time constraints, we  
 199 did not do experimental comparison (they only reported results on VGG whereas we focus on  
 200 ResNets) – as they solve a full augmented Lagrangian minimization in between each compression  
 201 step, successful training of their LC algorithm will at least require a careful tuning of this inner  
 202 optimization procedure.

## 203 C Notation

204 Throughout the paper, we let  $\theta \in \mathbb{R}^d$  denote the parameters of a neural network,  $L(\theta)$  denote the loss  
 205 function over the entire dataset (the empirical risk), and  $\tilde{\nabla}L$  denote the stochastic gradient of  $L$  (e.g.  
 206 over a minibatch). For the regularization method, we denote the set of quantized parameters by  $\mathcal{Q}$ ,  
 207 the regularizer by  $R(\theta)$  and the regularization strength by  $\lambda$ . The learning rates are denoted by  $\eta_t$  for  
 208  $t \geq 0$ . We let  $\text{Proj}_S : \mathbb{R}^d \rightarrow \mathbb{R}^d$  denote the standard Euclidean projection onto a set  $S \subset \mathbb{R}^d$  and  
 209  $\text{prox}_f$  denote the proximal operator w.r.t. function  $f$  (details in Appendix D.3). The  $p$ -Wasserstein  
 210 distance is denoted as  $W_p$ . We will restrict attention to Wasserstein distances on  $\mathbb{R}$  and  $p \in \{1, 2\}$ .

## 211 D Preliminaries

212 The optimization difficulty of training quantized models is that they involve a discrete parameter space  
 213 and hence efficient local-search methods are often prohibitive. For example, the problem of training  
 214 a binary neural network is to minimize  $L(\theta)$  for  $\theta \in \{\pm 1\}^d$ . Projected SGD on this set will not  
 215 move unless with an unreasonably large stepsize [Li et al., 2017], whereas greedy nearest-neighbor  
 216 search requires  $d$  forward passes which is intractable for neural networks where  $d$  is on the order  
 217 of millions. Alternatively, quantized training can also be cast as minimizing  $L(q(\theta))$  for  $\theta \in \mathbb{R}^d$   
 218 and an appropriate *quantizer*  $q$  that maps a real vector to a nearby quantized vector, but  $\theta \mapsto q(\theta)$  is  
 219 often non-differentiable and piecewise constant (such as the binary case  $q(\cdot) = \text{sign}(\cdot)$ ), and thus  
 220 back-propagation through  $q$  does not work.

### 221 D.1 The straight-through gradient method

222 The pioneering work of BinaryConnect [Courbariaux et al., 2015] proposes to solve this problem via  
 223 the *straight-through gradient method*, that is, propagate the gradient with respect to  $q(\theta)$  unaltered

224 to  $\theta$ , i.e. to let  $\frac{\partial L}{\partial \theta} := \frac{\partial L}{\partial \mathfrak{q}(\theta)}$ . One iterate of the straight-through gradient method (with the SGD  
 225 optimizer) is

$$\theta_{t+1} = \theta_t - \eta_t \tilde{\nabla} L(\theta)|_{\theta=\mathfrak{q}(\theta_t)}.$$

226 This enables the real vector  $\theta$  to move in the entire Euclidean space, and taking  $\mathfrak{q}(\theta)$  at the end of  
 227 training gives a valid quantized model. Such a customized back-propagation rule yields good empiri-  
 228 cal performance in training quantized nets and has thus become a standard practice [Courbariaux  
 229 et al., 2015, Zhu et al., 2016, Xu et al., 2018]. However, as we have discussed, it is information  
 230 theoretically unclear how the straight-through method works, and it does fail on very simple convex  
 231 Lipschitz functions (Figure 1b).

## 232 D.2 Straight-through gradient as lazy projection

233 Our first observation is that the straight-through gradient method is equivalent to Nesterov’s *dual-*  
 234 *averaging* method, or a lazy projected SGD [Xiao, 2010]. In the binary case, we wish to minimize  
 235  $L(\theta)$  over  $\mathcal{Q} = \{\pm 1\}^d$ , and the lazy projected SGD proceeds as

$$\begin{cases} \tilde{\theta}_t = \text{Proj}_{\mathcal{Q}}(\theta_t) = \text{sign}(\theta_t) = \mathfrak{q}(\theta_t), \\ \theta_{t+1} = \theta_t - \eta_t \tilde{\nabla} L(\tilde{\theta}_t). \end{cases} \quad (4)$$

236 Written compactly, this is  $\theta_{t+1} = \theta_t - \eta_t \tilde{\nabla} L(\theta)|_{\theta=\mathfrak{q}(\theta_t)}$ , which is exactly the straight-through gradient  
 237 method: take the gradient at the quantized vector and perform the update on the original real vector.

## 238 D.3 Projection as a limiting proximal operator

239 We take a broader point of view that a projection is also a limiting proximal operator with a suitable  
 240 regularizer, to allow more generality and to motivate our proposed algorithm. Given any set  $\mathcal{Q}$ , one  
 241 could identify a regularizer  $R : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$  such that the following hold:

$$R(\theta) = 0, \quad \forall \theta \in \mathcal{Q} \quad \text{and} \quad R(\theta) > 0, \quad \forall \theta \notin \mathcal{Q}. \quad (5)$$

242 In the case  $\mathcal{Q} = \{\pm 1\}^d$  for example, one could take

$$R(\theta) = R_{\text{bin}}(\theta) = \sum_{j=1}^d \min\{|\theta_j - 1|, |\theta_j + 1|\}. \quad (6)$$

243 The proximal operator (or prox operator) [Parikh and Boyd, 2014] with respect to  $R$  and strength  
 244  $\lambda > 0$  is

$$\text{prox}_{\lambda R}(\theta) := \arg \min_{\tilde{\theta} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\tilde{\theta} - \theta\|_2^2 + \lambda R(\tilde{\theta}) \right\}.$$

245 In the limiting case  $\lambda = \infty$ , the argmin has to satisfy  $R(\theta) = 0$ , i.e.  $\theta \in \mathcal{Q}$ , and the prox operator is  
 246 to minimize  $\|\theta - \theta_0\|_2^2$  over  $\theta \in \mathcal{Q}$ , which is the Euclidean projection onto  $\mathcal{Q}$ . Hence, projection is  
 247 also a prox operator with  $\lambda = \infty$ , and the straight-through gradient estimate is equivalent to a lazy  
 248 proximal gradient descent with  $\lambda = \infty$ .

249 While the prox operator with  $\lambda = \infty$  corresponds to “hard” projection onto the discrete set  $\mathcal{Q}$ , when  
 250  $\lambda < \infty$  it becomes a “soft” projection that moves towards  $\mathcal{Q}$ . Compared with the hard projection,  
 251 a finite  $\lambda$  is less aggressive and has the potential advantage of avoiding overshoot early in training.  
 252 Further, as the prox operator does not strictly enforce quantizedness, it is in principle able to query  
 253 the gradients at every point in the space, and therefore has access to more information than the  
 254 straight-through gradient method.

## 255 E Details on the PROXQUANT algorithm

### 256 E.1 Regularization for model quantization

257 We define a flexible class of quantization-inducing regularizers through “distance to the quantized  
 258 set”, derive efficient algorithms of their corresponding prox operator, and propose a homotopy method

259 for choosing the regularization strengths. Our regularization perspective subsumes most existing  
 260 algorithms for model-quantization (e.g., [Courbariaux et al., 2015, Han et al., 2015, Xu et al., 2018])  
 261 as limits of certain regularizers with strength  $\lambda \rightarrow \infty$ . Our proposed method can be viewed as a  
 262 principled generalization of these methods to  $\lambda < \infty$ .

263 Let  $\mathcal{Q} \subset \mathbb{R}^d$  be a set of quantized parameter vectors. An ideal regularizer for quantization would be  
 264 to vanish on  $\mathcal{Q}$  and reflect some type of distance to  $\mathcal{Q}$  when  $\theta \notin \mathcal{Q}$ . To achieve this, we propose  $L_1$   
 265 and  $L_2$  regularizers of the form

$$R(\theta) = \inf_{\theta_0 \in \mathcal{Q}} \|\theta - \theta_0\|_1 \quad \text{or} \quad R(\theta) = \inf_{\theta_0 \in \mathcal{Q}} \|\theta - \theta_0\|_2^2. \quad (7)$$

266 This is a highly flexible framework for designing regularizers, as one could specify any  $\mathcal{Q}$  and choose  
 267 between  $L_1$  and  $L_2$ . Specifically,  $\mathcal{Q}$  encodes certain desired quantization structure. By appropriately  
 268 choosing  $\mathcal{Q}$ , we can specify which part of the parameter vector to quantize<sup>1</sup>, the number of bits to  
 269 quantize to, whether we allow adaptively-chosen quantization levels and so on.

270 The choice of distance metrics will result in distinct properties in the regularized solutions. For  
 271 example, choosing the  $L_1$  version leads to non-smooth regularizers that induce exact quantizedness  
 272 in the same way that  $L_1$  norm regularization induces sparsity [Tibshirani, 1996], whereas choosing  
 273 the squared  $L_2$  version leads to smooth regularizers that induce quantizedness “softly”.

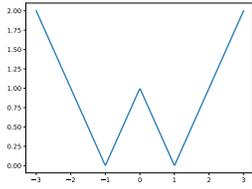
274 In the following, we present a few examples of regularizers under our framework eq. (7) which induce  
 275 binary weights, ternary weights and multi-bit quantization. We will also derive efficient algorithms  
 276 (or approximation heuristics) for solving the prox operators corresponding to these regularizers,  
 277 which generalize the projection operators used in the straight-through gradient algorithms.

278 **Binary neural nets** In a binary neural net, the entries of  $\theta$  are in  $\{\pm 1\}$ . A natural choice would be  
 279 taking  $\mathcal{Q} = \{-1, 1\}^d$ . The resulting  $L_1$  regularizer is

$$\begin{aligned} R(\theta) &= \inf_{\theta_0 \in \{\pm 1\}^d} \|\theta - \theta_0\|_1 = \sum_{j=1}^d \inf_{[\theta_0]_j \in \{\pm 1\}} |\theta_j - [\theta_0]_j| \\ &= \sum_{j=1}^d \min \{|\theta_j - 1|, |\theta_j + 1|\} = \|\theta - \text{sign}(\theta)\|_1. \end{aligned} \quad (8)$$

280 This is exactly the binary regularizer  $R_{\text{bin}}$  that we discussed earlier in eq. (6). Figure 2 plots the  
 281 W-shaped one-dimensional component of  $R_{\text{bin}}$  from which we see its effect for inducing  $\{\pm 1\}$   
 282 quantization in analog to  $L_1$  regularization for inducing exact sparsity.

283 The prox operator with respect to  $R_{\text{bin}}$ , despite being a non-convex  
 284 optimization problem, admits a simple analytical solution:



285  $\text{prox}_{\lambda R_{\text{bin}}}(\theta) = \text{SoftThreshold}(\theta, \text{sign}(\theta), \lambda)$   
 286  $= \text{sign}(\theta) + \text{sign}(\theta - \text{sign}(\theta)) \odot [|\theta - \text{sign}(\theta)| - \lambda]_+.$   
 287 (9)

288 Figure 2: W-shaped regularizer  
 289 for binary quantization.

290 We note that the choice of the  $L_1$  version is not unique: the squared  
 291  $L_2$  version works as well, whose prox operator is given by  $(\theta + \lambda \text{sign}(\theta))/(1 + \lambda)$ .

**Multi-bit quantization with adaptive levels.** Following [Xu et al., 2018], we consider  $k$ -bit quantized parameters with a structured adaptively-chosen set of quantization levels, which translates into

$$\mathcal{Q} = \left\{ \sum_{i=1}^k \alpha_i b_i : \{\alpha_1, \dots, \alpha_k\} \subset \mathbb{R}, b_i \in \{\pm 1\}^d \right\} = \left\{ \theta_0 = B\alpha : \alpha \in \mathbb{R}^k, B \in \{\pm 1\}^{d \times k} \right\}. \quad (10)$$

<sup>1</sup>Empirically, it is advantageous to keep the biases of each layers and the BatchNorm layers at full-precision, which is often a negligible fraction, say  $1/\sqrt{d}$  of the total number of parameters

292 The squared  $L_2$  regularizer for this structure is

$$R_{k\text{-bit}}(\theta) = \inf_{\alpha \in \mathbb{R}^k, B \in \{\pm 1\}^{d \times k}} \|\theta - B\alpha\|_2^2, \quad (11)$$

293 which is also the alternating minimization objective in [Xu et al., 2018].

294 We now derive the prox operator for the regularizer eq. (11). For any  $\theta$ , we have

$$\begin{aligned} \text{prox}_{\lambda R_{k\text{-bit}}}(\theta) &= \arg \min_{\tilde{\theta}} \left\{ \frac{1}{2} \|\tilde{\theta} - \theta\|_2^2 + \lambda \inf_{\alpha \in \mathbb{R}^k, B \in \{\pm 1\}^{d \times k}} \|\tilde{\theta} - B\alpha\|_2^2 \right\} \\ &= \arg \min_{\tilde{\theta}} \inf_{\alpha \in \mathbb{R}^k, B \in \{\pm 1\}^{d \times k}} \left\{ \frac{1}{2} \|\tilde{\theta} - \theta\|_2^2 + \lambda \|\tilde{\theta} - B\alpha\|_2^2 \right\}. \end{aligned} \quad (12)$$

295 This is a joint minimization problem in  $(\tilde{\theta}, B, \alpha)$ , and we adopt an alternating minimization schedule  
296 to solve it:

297 (1) Minimize over  $\tilde{\theta}$  given  $(B, \alpha)$ , which has a closed-form solution  $\tilde{\theta} = \frac{\theta + 2\lambda B\alpha}{1 + 2\lambda}$ .

298 (2) Minimize over  $(B, \alpha)$  given  $\tilde{\theta}$ , which does not depend on  $\theta_0$ , and can be done via calling the  
299 alternating quantizer of [Xu et al., 2018]:  $B\alpha = \mathbf{q}_{\text{alt}}(\tilde{\theta})$ .

300 Together, the prox operator generalizes the alternating minimization procedure in [Xu et al., 2018], as  
301  $\lambda$  governs a trade-off between quantization and closeness to  $\theta$ . To see that this is a strict generalization,  
302 note that for any  $\lambda$  the solution of eq. (12) will be an interpolation between the input  $\theta$  and its Euclidean  
303 projection to  $\mathcal{Q}$ . As  $\lambda \rightarrow +\infty$ , the prox operator collapses to the projection.

304 **Ternary quantization** Ternary quantization is a variant of 2-bit quantization, in which weights are  
305 constrained to be in  $\{-\alpha, 0, \beta\}$  for real values  $\alpha, \beta > 0$ .

306 For ternary quantization, we use an approximate version of the alternating prox operator eq. (12):  
307 compute  $\tilde{\theta} = \text{prox}_{\lambda R}(\theta)$  by initializing at  $\tilde{\theta} = \theta$  and repeating

$$\hat{\theta} = \mathbf{q}(\tilde{\theta}) \quad \text{and} \quad \tilde{\theta} = \frac{\theta + 2\lambda \hat{\theta}}{1 + 2\lambda}, \quad (13)$$

308 where  $\mathbf{q}$  is the ternary quantizer defined as

$$\mathbf{q}(\theta) = \theta^+ \mathbf{1}\{\theta \geq \Delta\} + \theta^- \mathbf{1}\{\theta \leq -\Delta\}, \quad \Delta = \frac{0.7}{d} \|\theta\|_1, \quad \theta^+ = \overline{\theta|_{i:\theta_i \geq \Delta}}, \quad \theta^- = \overline{\theta|_{i:\theta_i \leq -\Delta}}. \quad (14)$$

309 This is a straightforward extension of the TWN quantizer [Li and Liu, 2016] that allows different  
310 levels for positives and negatives. We find that two rounds of alternating computation in eq. (13)  
311 achieves a good performance, which we use in our experiments.

## 312 E.2 Homotopy method for regularization strength

313 Recall that the larger  $\lambda_t$  is, the more aggressive  $\theta_{t+1}$  will move towards the quantized set. An ideal  
314 choice would be to (1) force the net to be exactly quantized upon convergence, and (2) not be too  
315 aggressive such that the quantized net at convergence is sub-optimal.

316 We let  $\lambda_t$  be a linearly increasing sequence, i.e.  $\lambda_t := \lambda \cdot t$  for some hyper-parameter  $\lambda > 0$  which  
317 we term as the *regularization rate*. With this choice, the stochastic gradient steps will start off  
318 close to full-precision training and gradually move towards exact quantizedness, hence the name  
319 “homotopy method”. The parameter  $\lambda$  can be tuned by minimizing the validation loss, and controls  
320 the aggressiveness of falling onto the quantization constraint. There is nothing special about the  
321 linear increasing scheme, but it is simple enough and works well as we shall see in the experiments.

## 322 F Experiments on LSTMs

323 **Problem setup** We perform language modeling with LSTMs Hochreiter and Schmidhuber [1997]  
324 on the Penn Treebank (PTB) dataset [Marcus et al., 1993], which contains 929K training tokens,  
325 73K validation tokens, and 82K test tokens. Our model is a standard one-hidden-layer LSTM with  
326 embedding dimension 300 and hidden dimension 300. We train quantized LSTMs with the encoder,  
327 transition matrix, and the decoder quantized to  $k$ -bits for  $k \in \{1, 2, 3\}$ . The quantization is performed  
328 in a row-wise fashion, so that each row of the matrix has its own codebook  $\{\alpha_1, \dots, \alpha_k\}$ .

329 **Method** We compare our multi-bit PROXQUANT to the state-of-the-art alternating minimization  
 330 algorithm with straight-through gradients [Xu et al., 2018]. Training is initialized at a pre-trained  
 331 full-precision LSTM. We use the SGD optimizer with initial learning rate 20.0 and decay by a factor  
 332 of 1.2 when the validation error does not improve over an epoch. We train for 80 epochs with  
 333 batch size 20, BPTT 30, dropout with probability 0.5, and clip the gradient norms to 0.25. The  
 334 regularization rate  $\lambda$  is tuned by finding the best performance on the validation set. In addition to  
 335 multi-bit quantization, we also report the results for binary LSTMs (weights in  $\{\pm 1\}$ ), comparing  
 336 BinaryConnect and our PROXQUANT-Binary.

337 **Result** We report the perplexity-per-word (PPW, lower is better) in Table 2. The performance  
 338 of PROXQUANT is comparable with the Straight-through gradient method. On Binary LSTMs,  
 339 PROXQUANT-Binary beats BinaryConnect by a large margin. These results demonstrate that PROX-  
 340 QUANT offers a powerful alternative for training recurrent networks.

Table 2: PPW of quantized LSTM on Penn Treebank.

Method / Number of Bits	1	2	3	FP (32)
BinaryConnect	419.1	-	-	88.5
PROXQUANT-Binary (ours)	<b>321.8</b>	-	-	
ALT Straight-through <sup>2</sup>	104.7	90.2	86.1	
ALT-PROXQUANT (ours)	106.2	90.0	87.2	

## 341 G Sign change experiment

342 We experimentally compare the training dynamics of PROXQUANT-Binary and BinaryConnect  
 343 through the *sign change* metric. The sign change metric between any  $\theta_1$  and  $\theta_2$  is the proportion of  
 344 their different signs, i.e. the (rescaled) Hamming distance:

$$\text{SignChange}(\theta_1, \theta_2) = \frac{\|\text{sign}(\theta_1) - \text{sign}(\theta_2)\|_1}{2d} \in [0, 1].$$

345 In  $\mathbb{R}^d$ , the space of all full-precision parameters, the sign change is a natural distance metric that  
 346 represents the closeness of the binarization of two parameters.

347 Recall in our CIFAR-10 experiments (Section 3.1), for both BinaryConnect and PROXQUANT, we  
 348 initialize at a good full-precision net  $\theta_0$  and stop at a converged binary network  $\hat{\theta} \in \{\pm 1\}^d$ . We  
 349 are interested in  $\text{SignChange}(\theta_0, \theta_t)$  along the training path, as well as  $\text{SignChange}(\theta_0, \hat{\theta})$ , i.e. the  
 350 distance of the final output model to the initialization.

351 As PROXQUANT converges to higher-performance solutions than BinaryConnect, we expect that if  
 352 we run both methods from a same warm start, the sign change of PROXQUANT should be higher than  
 353 that of BinaryConnect, as in general one needs to travel farther to find a better net.

354 However, we find that this is not the case: PROXQUANT produces binary nets with both *lower* sign  
 355 changes and *higher* performances, compared with BinaryConnect. This finding is consistent in  
 356 all layers, across different warm starts, and across different runs from each same warm start (see  
 357 Figure 3 and Table 3 in Appendix G.1). This shows that for every warm start position, there is a  
 358 good binary net nearby which can be found by PROXQUANT but not BinaryConnect, suggesting that  
 359 BinaryConnect, and in general the straight-through gradient method, suffers from higher optimization  
 360 instability than PROXQUANT. This result here is also consistent with Theorem 4.1: the signs in  
 361 BinaryConnect never stop changing until we manually freeze the signs at epoch 400.

### 362 G.1 Detailed sign change results on ResNet-20

<sup>2</sup>We thank Xu et al. [2018] for sharing the implementation of this method through a personal communication. There is a very clever trick not mentioned in their paper: after computing the alternating quantization  $q_{\text{alt}}(\theta)$ , they multiply by a constant 0.3 before taking the gradient; in other words, their quantizer is a rescaled alternating quantizer:  $\theta \mapsto 0.3q_{\text{alt}}(\theta)$ . This scaling step gives a significant gain in performance – without scaling the PPW is  $\{116.7, 94.3, 87.3\}$  for  $\{1, 2, 3\}$  bits. In contrast, our PROXQUANT does not involve a scaling step and achieves better PPW than this unscaled ALT straight-through method.

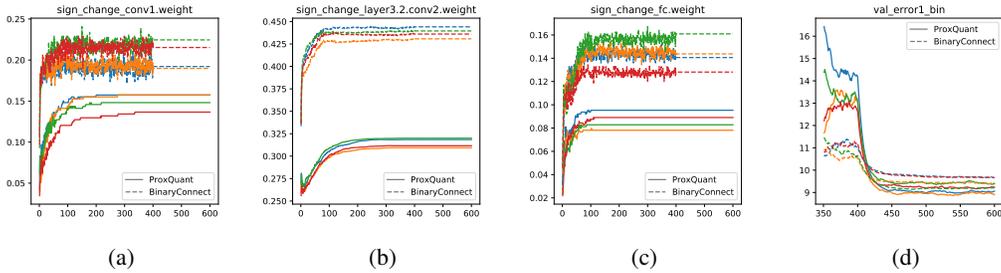


Figure 3:  $\text{SignChange}(\theta_0, \theta_t)$  against  $t$  (epoch) for BinaryConnect and PROXQUANT, over 4 runs starting from the same full-precision ResNet-20. PROXQUANT has significantly lower sign changes than BinaryConnect while converging to better models. (a) The first conv layer of size  $16 \times 3 \times 3 \times 3$ ; (b) The last conv layer of size  $64 \times 64 \times 3 \times 3$ ; (c) The fully connected layer of size  $64 \times 10$ ; (d) The validation top-1 error of the binarized nets (with moving average smoothing).

Table 3: Performances and sign changes on ResNet-20 in mean(std) over 3 full-precision initializations and 4 runs per (initialization x method). Sign changes are computed over all quantized parameters in the net.

Initialization	Method	Top-1 Error(%)	Sign change
FP-Net 1 (8.06)	BC	9.489 (0.223)	0.383 (0.006)
	PQ-B	<b>9.146</b> (0.212)	<b>0.276</b> (0.020)
FP-Net 2 (8.31)	BC	9.745 (0.422)	0.381 (0.004)
	PQ-B	<b>9.444</b> (0.067)	<b>0.288</b> (0.002)
FP-Net 3 (7.73)	BC	9.383 (0.211)	0.359 (0.001)
	PQ-B	<b>9.084</b> (0.241)	<b>0.275</b> (0.001)

Table 4: Performances and sign changes on ResNet-20 in raw data over 3 full-precision initializations and 4 runs per (initialization x method). Sign changes are computed over all quantized parameters in the net.

Initialization	Method	Top-1 Error(%)	Sign change
FP-Net 1 (8.06)	BC	9.664, 9.430, 9.198, 9.663	0.386, 0.377, 0.390, 0.381
	PQ-B	9.058, 8.901, 9.388, 9.237	0.288, 0.247, 0.284, 0.285
FP-Net 2 (8.31)	BC	9.456, 9.530, 9.623, 10.370	0.376, 0.379, 0.382, 0.386
	PQ-B	9.522, 9.474, 9.410, 9.370	0.291, 0.287, 0.289, 0.287
FP-Net 3 (7.73)	BC	9.107, 9.558, 9.538, 9.328	0.360, 0.357, 0.359, 0.360
	PQ-B	9.284, 8.866, 9.301, 8.884	0.275, 0.276, 0.276, 0.275