# Quantization through Piecewise-Affine Regularization: Optimization and Statistical Guarantees

Jianhao Ma
University of Pennsylvania
jianhaom@wharton.upenn.edu

Lin Xiao
Meta FAIR
linx@meta.com

August 18, 2025

## Abstract

Optimization problems over discrete or quantized variables are very challenging in general due to the combinatorial nature of their search space. Piecewise-affine regularization (PAR) provides a flexible modeling and computational framework for quantization based on continuous optimization. In this work, we focus on the setting of supervised learning and investigate the theoretical foundations of PAR from optimization and statistical perspectives. First, we show that in the overparameterized regime, where the number of parameters exceeds the number of samples, every critical point of the PAR-regularized loss function exhibits a high degree of quantization. Second, we derive closed-form proximal mappings for various (convex, quasi-convex, and non-convex) PARs and show how to solve PAR-regularized problems using the proximal gradient method, its accelerated variant, and the Alternating Direction Method of Multipliers. Third, we study statistical guarantees of PAR-regularized linear regression problems; specifically, we can approximate classical formulations of $\ell_1$-, squared $\ell_2$-, and nonconvex regularizations using PAR and obtain similar statistical guarantees with quantized solutions.

## 1 Introduction

In many machine learning and decision-making problems, we need to optimize an objective function where some variables are constrained to be discrete:

$$\min_{\boldsymbol{x} \in \mathcal{Q}^{d_1}, \boldsymbol{y} \in \mathbb{R}^{d_2}} f(\boldsymbol{x}, \boldsymbol{y}). \tag{1}$$

Here $\boldsymbol{y}$ is the continuous variable, and the elements of $\boldsymbol{x}$ are restricted to a discrete set $\mathcal{Q}$. For example, $\mathcal{Q}$ can be the set of binary values $\{0, 1\}$, a subset of integers, or a finite set of discrete real numbers. The prevalence and importance of such problems are highlighted by the following examples.

- **Quantization in machine learning model compression.** Modern deep learning models offer remarkable capabilities in vision and language processing, but they often come with substantial computational and memory requirements. Quantization, which maps model parameters from high-precision to low-precision formats, has emerged as an effective approach for model compression. It can significantly reduce memory footprint, computational cost, and inference latency [HMD16, SCYE17].

- **Communications and signal processing.** Quantization plays a fundamental role in digital communications, where continuous-amplitude signals must be converted into discrete values for transmission,

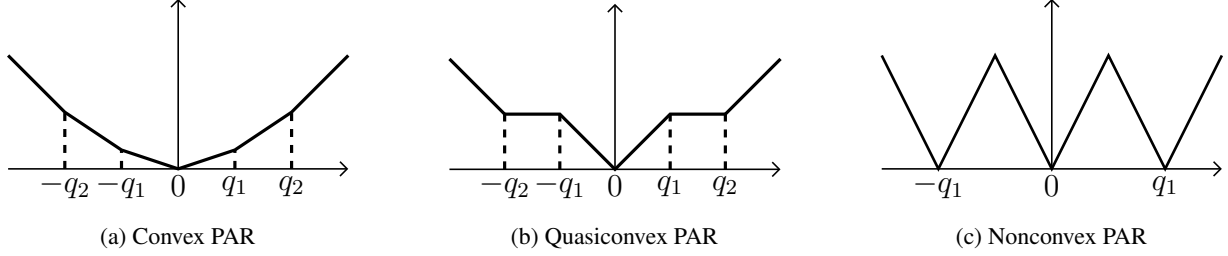|  |  |  |
|:---:|:---:|:---:|
| (a) Convex PAR | (b) Quasiconvex PAR | (c) Nonconvex PAR |

Figure 1: Three different types of PAR $\Psi(\cdot)$ for inducing quantization towards $\mathcal{Q} = \{\pm q_i\}$.

storage, and processing [PS08, Opp99, GN02]. This process underlies analog-to-digital conversion, enabling real-world analog signals to be represented with finite bit-depth. The quality of the quantization directly affects signal fidelity, bandwidth efficiency, and error rates in communication systems.

- **Mixed integer programming in operations research.** Many optimization problems in operations research require variables to take discrete values, such as facility location, production scheduling, or resource allocation [SS19]. These problems are formulated as mixed integer programs where some variables are constrained to integer values while others remain continuous, leading to challenging computational problems that combine combinatorial and continuous optimization [WN99, Wol20].

Solving Problem (1) exactly is extremely challenging due to its combinatorial nature. For instance, even the simple convex quadratic binary program $\min_{\boldsymbol{x} \in \{0,1\}^d} \boldsymbol{x}^\top \boldsymbol{A} \boldsymbol{x}$ with $\boldsymbol{A} \succeq 0$ is NP-hard [Har82].

In this paper, we focus on finding *approximate solutions* using continuous optimization methods, by adding regularizations/penalties that can induce discrete solutions to the objective function. Specifically, we consider the following unconstrained optimization problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} F_\lambda(\boldsymbol{x}) := f(\boldsymbol{x}) + \lambda \Psi(\boldsymbol{x}), \tag{2}$$

where $\Psi(\cdot)$ is a regularization that encourages the variables to be discrete (within the set $\mathcal{Q}$), and $\lambda$ controls the strength of regularization.[1] Among many possible choices, the family of Piecewise-Affine Regularizers (PARs) are particularly suited for inducing quantization due to their nature of nonsmoothness (Figure 1 illustrates three representative types of PAR). In particular, they tend to trap the optimization variables at the set of nondifferentiable breakpoints. This mirrors how the $\ell_1$-regularizer promotes sparsity through nondifferentiability at zero. For PARs, we align their set of nondifferentiable points with the target quantization values, making them effective for inducing desired quantization. We refer to the Problem (2) with such a regularizer as **Piecewise-Affine Regularized Optimization (PARO)**.

PAR has a long history in statistics. The classic example is the Lasso [Tib96, CDS98], which uses $\ell_1$-regularization to induce sparsity in linear regression and other statistical learning tasks. Subsequent work introduced other convex PARs to induce different structures; a prominent case is the graph-based total-variation penalty, which suppresses jumps and favors piecewise-constant signals [Con17, KPR16]. Several recent studies examined the geometry of solutions produced by convex PARs [EDA24, ST22, TSGS21].

Parallel advances in machine learning highlight an intrinsic link between PARs and quantization. In particular, [CBD15] introduced the straight-through estimator (STE) which has become a workhorse for quantization-aware training, and [YZL+18] draws its connection to the framework of regularized optimization.

---

[1] Here we omit the continuous part $\boldsymbol{y}$ and focus on the case where all the variables are subject to regularization. It can be easily extended to the general setting by setting $\Psi(\boldsymbol{y}) = 0$.

[BWL19] reinterpret STE as regularized dual averaging [Xia10] with a non-convex PAR (illustrated in Figure 1c), and follow-up works broaden this analysis to wider families of PARs, devising optimization algorithms with rigorous convergence guarantees [DYS+21, LYLPN24]. Most recently, [JML+25] propose to use convex PARs (Figure 1a) for quantization, and introduce an aggregated proximal stochastic gradient method with last-iterate convergence guarantee and strong empirical results on deep learning tasks.

While PARs have achieved strong empirical success, their theoretical properties remain underexplored. In particular, the mechanisms by which PARs promote discrete structure, their optimization guarantee, and statistical behavior are not well understood. In this work, we bridge this gap by analyzing PARs through the lenses of quantization, optimization, and statistics, establishing new theoretical foundations that explain and support their empirical effectiveness. Our main contributions are:

- **Quantization guarantees.** We provide theoretical backing for PAR's ability to induce high quantization rate in supervised learning models. We show that all critical points of PARO have a high proportion of quantized entries. The quantization rate is directly linked to the ratio of parameter dimension to sample size, indicating that overparameterized models naturally achieve higher quantization rate.

- **Optimization methods.** We derive closed-form expressions for the proximal mappings of various PARs, spanning convex, quasiconvex, and nonconvex formulations. We show that the proximal gradient method and its accelerated variants can efficiently converge to critical points of PARO for both convex and nonconvex cases. Additionally, for structured problems like linear regression, we demonstrate that the Alternating Direction Method of Multipliers (ADMM) [BPC+11] can be effectively applied, often achieving faster convergence.

- **Statistical properties.** We demonstrate that PARs can closely approximate a wide range of conventional regularizers, including squared $\ell_2$-, $\ell_1$-, and general nonconvex regularizers. For linear regression, we prove that specific PARs can effectively mimic classic regularizers, achieving optimal statistical guarantees with quantized solutions. This highlights PARO's ability to reduce model size without sacrificing performance.

- **Numerical experiments.** We conduct extensive simulations across linear and logistic regression tasks. These experiments empirically corroborate our theoretical findings, validating the quantization, optimization, and statistical guarantees provided by the PARO framework.

The rest of this paper is organized as follows. Section 2 introduces different types of PAR along with their first-order optimality conditions. We also provide quantization guarantees of PAR for generalized linear models and supervised learning. In Section 3, we derive the proximal operators of various PARs and illustrate how to solve PARO using several standard optimization algorithms. In Section 4, we establish statistical guarantees for various PARs in the linear regression setting. Finally, in Section 5, we present numerical experiments that support our theoretical findings and demonstrate the effectiveness of our approach.

**Notations.** We use bold lowercase letters (e.g., $\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{z}$) to denote vectors, and bold uppercase letters (e.g., $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$) to denote matrices. For a vector $\boldsymbol{x} \in \mathbb{R}^d$, we define its $\ell_2$-norm as $\|\boldsymbol{x}\| = (\sum_{i=1}^d x_i^2)^{1/2}$ and its $\ell_\infty$-norm as $\|\boldsymbol{x}\|_\infty = \max_i |x_i|$. The sign function is denoted as $\mathrm{sign}(x)$, which returns 1 if $x > 0$, $-1$ if $x < 0$, and 0 otherwise. For a set $I$, we denote its cardinality by $|I|$. Given a matrix $\boldsymbol{A} \in \mathbb{R}^{m \times d}$, we use $\mathbf{A}_i$ to denote its $i$-th row. The notation $\boldsymbol{A}_I$ refers to the submatrix of $\boldsymbol{A}$ containing only the rows indexed by $I$. For a function $f : \mathbb{R}^d \to \mathbb{R}$, we denote its gradient at $\boldsymbol{x}$ by $\nabla f(\boldsymbol{x})$ if $f$ is differentiable at $\boldsymbol{x}$ and its Clarke subdifferential by $\partial f(\boldsymbol{x})$ otherwise. For a scalar $x$, we use $\lfloor x \rceil$ to denote the closest integer to $x$.

# 2 PAR and quantization guarantees

In this paper, we consider coordinate-wise piecewise-affine regularization (PAR), i.e.,

$$\Psi(\boldsymbol{x}) = \sum_{i=1}^{d} \Psi(x_i),$$

where we use $\Psi$ for both vector and scalar inputs (slight abuse of notation). For ease of presentation and broad applicability in practice, we focus on PARs that are symmetric with respect to the origin (as illustrated in Figure 1), with the definition

$$\Psi(x) = a_k(|x| - q_k) + b_k \quad \text{if } q_k \leq |x| \leq q_{k+1}, \tag{3}$$

where $\mathcal{Q} = \{0, \pm q_1, \ldots, \pm q_m\}$, with $0 = q_0 < q_1 < \cdots < q_m$, denotes the set of targeted quantization values. The slopes $\mathcal{A} = \{a_0, a_1, \cdots, a_m\}$ satisfy $a_k \neq a_{k+1}$ for all $0 \leq k \leq m - 1$, and the intercepts $\{b_0, b_1, \cdots, b_m\}$ are determined recursively by setting $b_0 = 0$ and

$$b_k = b_{k-1} + a_{k-1}(q_k - q_{k-1}), \qquad k = 1, \ldots, m. \tag{4}$$

The remainder of this section presents general optimality conditions for PARO, followed by quantization guarantees established under a supervised learning framework.

## 2.1 Optimality conditions

In this section, we examine the first-order optimality conditions for PARO.[2] Specifically, we assume that $f$ is differentiable and $\boldsymbol{x}^\star$ is a Clarke critical point [BDLS07, Definition 2] of PARO (equation 2), i.e.,

$$\boldsymbol{0} \in \nabla f(\boldsymbol{x}^\star) + \lambda \partial \Psi(\boldsymbol{x}^\star), \tag{5}$$

where $\partial \Psi(\boldsymbol{x}^\star)$ is the Clarke subdifferential of $\Psi$ at $\boldsymbol{x}^\star$. This can be rewritten as $\nabla f(\boldsymbol{x}^\star) \in -\lambda \, \partial \Psi(\boldsymbol{x}^\star)$, which yields the following conditions for each coordinate $i = 1, \cdots, d$ and quantization level $k = 1, \cdots, m$:

$$
\begin{aligned}
x_i^\star &= -q_k, &\Longleftarrow\quad& \nabla_i f(\boldsymbol{x}^\star) \in \lambda\,(a_{k-1}, a_k) \\
x_i^\star &\in (-q_k, -q_{k-1}) &\Longrightarrow\quad& \nabla_i f(\boldsymbol{x}^\star) = \lambda\,a_{k-1} \\
x_i^\star &= 0 &\Longleftarrow\quad& -\nabla_i f(\boldsymbol{x}^\star) \in \lambda\,(-a_0, a_0) \\
x_i^\star &\in (q_{k-1}, q_k) &\Longrightarrow\quad& \nabla_i f(\boldsymbol{x}^\star) = -\lambda\,a_{k-1} \\
x_i^\star &= q_k, &\Longleftarrow\quad& \nabla_i f(\boldsymbol{x}^\star) \in \lambda\,(-a_k, -a_{k-1}).
\end{aligned}
$$

Here the symbol $\Longleftarrow$ ($\Longrightarrow$) means that the left-hand side expression is a necessary (sufficient) condition for the right-hand side expression.

We immediately recognize that the sufficient condition for $x_i^\star = 0$ is the same as for the $\ell_1$-regularization $\Psi(x) = \lambda a_0 \cdot |x|$. Further examination reveals that for any parameter not clustered at a discrete value in $\mathcal{Q}$, i.e., if $x_i^\star \in (q_{k-1}, q_k)$, the corresponding partial derivative of $f$ must equal to one of the $2m$ discrete values $\{\pm \lambda a_{k-1}\}_{k=1}^m$. Conversely, almost all values of the partial derivatives of $f$, except for these $2m$ discrete values, can be balanced by assigning the model parameters at the $2m + 1$ discrete values in $\mathcal{Q}$. Intuitively, this implies that the model parameters at optimality are more likely to be clustered at these discrete values.

---

[2]The results for convex PARs were previously presented in [JML$^+$25]; here, we extend them to general PARs and include the full statements for completeness.

## 2.2 Quantization guarantee

In this section, we study the quantization properties of solutions obtained using the PARO framework

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} F_\lambda(\boldsymbol{x}) = f(\boldsymbol{x}) + \lambda \Psi(\boldsymbol{x}). \tag{6}$$

We start with a simple generalized linear model setting, and later generalize it to general supervised learning setting. Given a specific PAR $\Psi(\cdot)$ with quantization values $\mathcal{Q} = \{0, \pm q_1, \dots, \pm q_m\}$, for an arbitrary point $\boldsymbol{x} \in \mathbb{R}^d$, we define its **quantization rate** $\mathrm{qr}(\boldsymbol{x})$ as the fraction of coordinates with quantized values, i.e.,

$$\mathrm{qr}(\boldsymbol{x}) := \frac{|\{i : x_i \in \mathcal{Q}\}|}{d}. \tag{7}$$

**Generalized linear model.** We consider objective functions of the form $f(\boldsymbol{x}) = g(\boldsymbol{A}\boldsymbol{x})$. Here $g(\cdot)$ is a smooth function, and $\boldsymbol{A} \in \mathbb{R}^{n \times d}$ is a matrix representing the input data. This formulation includes various generalized linear models, such as linear regression and logistic regression. The following result provides a quantization guarantee for the critical points of this problem when regularized by a class of PARs.

**Theorem 1** (Quantization guarantee for generalized linear models). *Consider a PAR $\Psi(\cdot)$ where the slopes $\mathcal{A}$ satisfies $0 \notin \mathcal{A}$. Suppose each row of the design matrix $\boldsymbol{A} \in \mathbb{R}^{n \times d}$ with $n \leq d$ is i.i.d. drawn from a distribution $\mathcal{D}_{\boldsymbol{a}}$ that is absolutely continuous with respect to the $p$-dimensional Lebesgue measure. Then, with probability one, any critical point $\boldsymbol{x}^\star$ of PARO (equation 2) satisfies $\mathrm{qr}(\boldsymbol{x}^\star) \geq 1 - n/d$.*

Below, we highlight the implications and significance of Theorem 1:

- **Effect of overparameterization.** Theorem 1 establishes a lower bound on the quantization rate, namely $1 - \frac{n}{d}$. Thus, in highly overparameterized models where $d \gg n$, the quantization rate approaches $1$. This suggests that larger models are inherently easier to quantize, which is consistent with empirical observations reported in [CZL$^+$25, CZG$^+$24].

- **Guarantees for critical points.** This result applies to all critical points of the regularized objective, not just global minimizers. This is particularly useful, as we show in Section 3 that standard optimization algorithms such as the proximal gradient method can efficiently find such critical points.

- **Tightness of the quantization guarantee.** The lower bound on the quantization rate depends only on the ratio of sample size to data dimension and is independent of the regularization strength. While this may appear weak, extensive simulations in Section 5 demonstrate that this bound is nearly tight, particularly when the regularization strength is moderate.

**General supervised learning.** We consider a training dataset $S_n = \{(\boldsymbol{a}_i, b_i)\}_{i=1}^n$ consisting of $n$ data points, where $\boldsymbol{a}_i \in \mathbb{R}^p$ represents the input and $b_i \in \mathbb{R}$ represents the output. We assume that each data point $(\boldsymbol{a}_i, b_i)$ is i.i.d. generated from an underlying distribution $\mathcal{D} = \mathcal{D}_{\boldsymbol{a}} \times \mathcal{D}_b$. We consider a family of machine learning models $f(\boldsymbol{a}; \boldsymbol{x})$ indexed by parameter $\boldsymbol{x} \in \mathbb{R}^d$. Then, we optimize the PAR-regularized empirical risk function

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} F_\lambda(\boldsymbol{x}) := \frac{1}{n} \sum_{i=1}^n \ell\left(f(\boldsymbol{a}_i; \boldsymbol{x}), b_i\right) + \lambda \Psi(\boldsymbol{x}). \tag{8}$$

Here $\Psi(\cdot)$ is a PAR with quantization values $\mathcal{Q}$ and the slopes $\mathcal{A}$.

We now introduce a set of conditions on the data distribution, the model, and the PAR.

**Assumption 1.** *The marginal data distribution $\mathcal{D}_{\boldsymbol{a}}$, the model $f(\boldsymbol{a}; \boldsymbol{x})$, and the PAR regularizer $\Psi(\cdot)$ satisfy the following conditions:*

- *(Continuous data distribution) The marginal data distribution $\mathcal{D}_{\boldsymbol{a}}$ is absolutely continuous with respect to the $p$-dimensional Lebesgue measure.*

- *(Real analytic model) The model $f : \mathbb{R}^p \times \mathbb{R}^d \to \mathbb{R}$ is a real analytic function.*

- *(Non-degenerate Jacobian) For any critical point $\boldsymbol{x}^\star$ and any index set $I$ with $|I| = n + 1$, and for any vectors $\boldsymbol{v} \in \mathbb{R}^d$, we all have*

$$\det \left( [\nabla_{\boldsymbol{x}} f(\boldsymbol{a}_1; \boldsymbol{x}^\star), \cdots, \nabla_{\boldsymbol{x}} f(\boldsymbol{a}_n; \boldsymbol{x}^\star), \boldsymbol{v}]_I \right) \not\equiv 0. \tag{9}$$

- *(Non-degenerate PAR) For the PAR $\Psi(\cdot)$, $0$ is not in its slope set $\mathcal{A}$.*

These conditions are relatively mild. First, we only assume absolute continuity of the data distribution and impose no further requirements such as sub-Gaussianity or bounded moments. Second, the real analyticity of the model $f$ is a standard assumption in the machine learning theory literature, and is satisfied by many common architectures, including feed-forward neural networks with analytic activation functions such as sigmoid or tanh; see, e.g., [Sch23, NH17, WO10, KS21]. Third, the non-degenerate Jacobian condition rules out trivial models (e.g., constant functions) and ensures sufficient expressiveness. A similar assumption is also used in [KS21]. Lastly, the non-degenerate PAR condition is necessary: if $0 \in \mathcal{A}$, the PAR family includes the zero function $\Psi(x) \equiv 0$, which yields no quantization effect.

**Theorem 2** (Quantization guarantee). *Suppose Assumption 1 holds. Then, with probability one, the quantization rate of any critical point $\boldsymbol{x}^\star$ of the objective in Equation (8) is at least $1 - n/d$, i.e.,*

$$\mathrm{qr}(\boldsymbol{x}^\star) \geq 1 - \frac{n}{d}. \tag{10}$$

Finally, we prove Theorems 1 and 2. We begin with the proof of Theorem 2, and then verify that the generalized linear model setting satisfies Assumption 1. This allows us to apply Theorem 2 and thereby establish Theorem 1 as a special case.

*Proof of Theorem 2.* For any critical point $\boldsymbol{x}^\star$ of the objective in Equation (8), it satisfies the first-order optimality condition

$$\boldsymbol{0} \in \boldsymbol{J}\boldsymbol{\delta} + \lambda \partial \Psi(\boldsymbol{x}^\star). \tag{11}$$

Here $\boldsymbol{J}$ is the Jacobian matrix defined by

$$\boldsymbol{J} = [\nabla_{\boldsymbol{x}} f(\boldsymbol{a}_1; \boldsymbol{x}^\star), \cdots, \nabla_{\boldsymbol{x}} f(\boldsymbol{a}_n; \boldsymbol{x}^\star)] \in \mathbb{R}^{d \times n}, \tag{12}$$

and $\boldsymbol{\delta} = [\nabla_f \ell(f(\boldsymbol{a}_1; \boldsymbol{x}^\star), b_1), \cdots, \nabla_f \ell(f(\boldsymbol{a}_n; \boldsymbol{x}^\star), b_n)]^\top$ is the residual vector.

Then, we partition $\boldsymbol{x}^\star$ into two parts, i.e., $\boldsymbol{x}^\star = [\boldsymbol{x}_I^\star; \boldsymbol{x}_{I^c}^\star]$ where $\boldsymbol{x}_I^\star \in \mathcal{Q}$ stands for the quantized part and $\boldsymbol{x}_{I^c}^\star$ stands for the non-quantized part. Similarly, we partition the Jacobian matrix $\boldsymbol{J}$ as $\boldsymbol{J} = [\boldsymbol{J}_I; \boldsymbol{J}_{I^c}]$. Suppose that $|I| \geq d - n$, then we are done. Hence, we assume that $|I| < d - n$. Now let us consider the optimality condition over the index set $[d] - I$, which is given by

$$-\frac{1}{\lambda} \boldsymbol{J}_{I^c} \boldsymbol{\delta} \in \partial \Psi \left( \boldsymbol{x}_{I^c}^\star \right). \tag{13}$$

Since $\boldsymbol{x}_{I^c}$ is the non-quantized part, we have $\partial\Psi(x_i^\star)\in\mathcal{A}$ for all $i\in I^c$, which implies that

$$-\frac{1}{\lambda}\boldsymbol{J}_{I^c}\boldsymbol{\delta}\in\mathcal{A}^{d-|I|}.\tag{14}$$

Now, we show that this event happens with probability zero. To this end, for an arbitrary $\boldsymbol{v}\in\mathcal{A}^{d-|I|}$, we have

$$\mathbb{P}\left(-\frac{1}{\lambda}\boldsymbol{J}_{I^c}\boldsymbol{\delta}=\boldsymbol{v}\right)\leq\mathbb{P}\left(\boldsymbol{v}\in\mathrm{col}\left(\boldsymbol{J}_{I^c}\right)\right)\leq\mathbb{P}\left(\det\left([\boldsymbol{J}_{I^c},\boldsymbol{v}]\right)=0\right).\tag{15}$$

To proceed, according to Assumption 1, we know that the function

$$H(\boldsymbol{a}_1,\cdots,\boldsymbol{a}_n)=\det\left([\boldsymbol{J}_{I^c},\boldsymbol{v}]\right)=\det\left([\nabla_{\boldsymbol{x}}f(\boldsymbol{a}_1;\boldsymbol{x}^\star),\cdots,\nabla_{\boldsymbol{x}}f(\boldsymbol{a}_n;\boldsymbol{x}^\star),\boldsymbol{v}]_{I^c}\right)\tag{16}$$

is also a real analytic function and $H(\boldsymbol{a}_1,\cdots,\boldsymbol{a}_n)\not\equiv 0$. As a consequence, the set $S=\{(\boldsymbol{a}_1,\cdots,\boldsymbol{a}_n):H(\boldsymbol{a}_1,\cdots,\boldsymbol{a}_n)=0\}$ has zero Lebesgue measure [Mit15]. Additionally, since $\boldsymbol{a}_i\sim\mathcal{D}_{\boldsymbol{a}}$ is absolutely continuous with respect to the $p$-dimensional Lebesgue measure, we have

$$\mathbb{P}\left(\det\left([\boldsymbol{J}_{I^c},\boldsymbol{v}]\right)=0\right)=\mathbb{P}(S)=0.\tag{17}$$

Therefore, we show that with probability 1, $-\frac{1}{\lambda}\boldsymbol{J}_{I^c}\boldsymbol{\delta}\notin\mathcal{A}^{d-|I|}$, which leads to the contradiction. $\square$

*Proof of Theorem 1.* It suffices to check the *Non-degenerate Jacobian* condition in Assumption 1. Note that

$$[\nabla_{\boldsymbol{x}}f(\boldsymbol{a}_1;\boldsymbol{x}^\star),\cdots,\nabla_{\boldsymbol{x}}f(\boldsymbol{a}_n;\boldsymbol{x}^\star),\boldsymbol{v}]=[\boldsymbol{a}_1,\cdots,\boldsymbol{a}_n,\boldsymbol{v}].\tag{18}$$

Therefore, for any index set $I\subset[d]$ satisfying $|I|=n+1$ and $\boldsymbol{v}\in\mathcal{A}^d$, we can always choose $[\boldsymbol{a}_1,\cdots,\boldsymbol{a}_n]_I$ such that $[\boldsymbol{a}_1,\cdots,\boldsymbol{a}_n,\boldsymbol{v}]_I$ is non-degenerate. This completes the proof. $\square$

# 3 Optimization algorithms

In this section, we first derive the proximal mappings for different variants of PARs and then introduce optimization algorithms that can leverage these mappings to efficiently solve the PARO problem.

## 3.1 Proximal mapping of PARs

In this section, we study the **proximal mappings** for different variants of PARs, including convex, quasiconvex, and nonconvex formulations, which are defined by

$$\mathbf{prox}_{\lambda\Psi}(\boldsymbol{x})=\arg\min_{\boldsymbol{z}}\left\{\lambda\Psi(\boldsymbol{z})+\frac{1}{2}\|\boldsymbol{x}-\boldsymbol{z}\|^2\right\}.\tag{19}$$

Since we consider coordinate-wise PARs, the proximal mapping can be decomposed across coordinates, so we only need to analyze the scalar case

$$\mathbf{prox}_{\lambda\Psi}(x)=\arg\min_{z}\left\{\lambda\Psi(z)+\frac{1}{2}(x-z)^2\right\}.\tag{20}$$

Proximal mapping is a key tool for optimizing regularized objectives and exhibits a strong connection to quantization with PAR: it often maps inputs to predefined discrete levels, effectively acting as a quantizer.
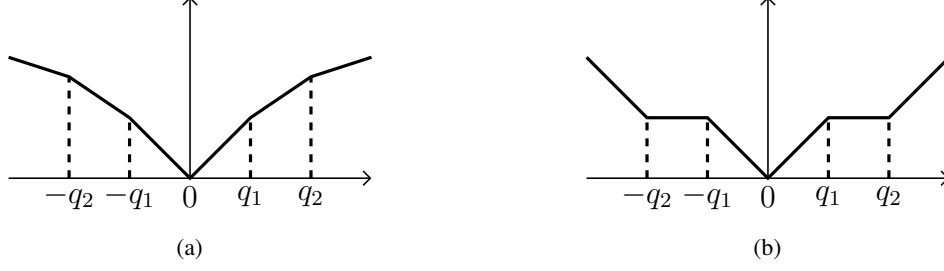
Figure 2: Two different forms of quasiconvex PARs.

For any PAR, the proximal mapping is also piecewise affine, though it may exhibit discontinuities. Within each linear segment of the regularizer, the composite objective $\lambda\Psi(z) + \frac{1}{2}(x-z)^2$ becomes a strongly convex quadratic function, whose minimizer either coincides with an endpoint of the segment or corresponds to the unconstrained minimizer, which is also an affine function of $x$, provided it lies within the interval. In the remainder of this section, we derive explicit formulas for the proximal mappings of several representative PARs and visualize them in Figure 3. Notably, the proximal mappings act as (soft) quantizers, further revealing the quantization effect induced by PAR.

**Convex PAR.** Suppose that the set of target quantization values is given as $\mathcal{Q} = \{0, \pm q_1, \ldots, \pm q_m\}$ with $0 = q_0 < q_1 < \cdots < q_m$. A convex PAR can be defined as

$$\Psi(x) = \max_{k \in \{0,\ldots,m\}} \{a_k(|x| - q_k) + b_k\}, \tag{21}$$

where the slopes $\{a_k\}_{k=0}^m$ are free parameters that satisfy $0 \le a_0 < a_1 < \cdots < a_m = +\infty$, and $\{b_k\}_{k=0}^m$ are determined by setting $b_0 = 0$, and

$$b_k = b_{k-1} + a_{k-1}(q_k - q_{k-1}), \qquad k = 1, \ldots, m. \tag{22}$$

Its proximal mapping is provided by

$$\mathbf{prox}_{\lambda\Psi}(x) = \begin{cases} \text{sign}(x)\, q_k & \text{if } |x| \in [\lambda a_{k-1} + q_k, \ \lambda a_k + q_k], \\ x - \text{sign}(x)\lambda a_k & \text{if } |x| \in [\lambda a_k + q_k, \ \lambda a_k + q_{k+1}]. \end{cases} \tag{23}$$

As shown in Figure 3 (top row), when $x$ falls within the intervals $\pm [\lambda a_{k-1} + q_k, \ \lambda a_k + q_k]$, the proximal mapping quantizes it to the corresponding quantization value $\text{sign}(x)q_k$. Moreover, as the regularization strength $\lambda$ increases, these quantization intervals become wider, making it more likely for the proximal mapping to produce a quantized solution when incorporated into an optimization algorithm.

**Quasiconvex PAR.** In this section, we study the proximal mapping for a class of quasiconvex PARs. A function $f : \mathbb{R} \to \mathbb{R}$ is said to be quasiconvex if, for any $x, y \in \mathbb{R}$ and $\lambda \in [0, 1]$, it satisfies

$$f(\lambda x + (1 - \lambda)y) \le \max\{f(x), f(y)\}.$$

Figure 2 illustrates two representative quasiconvex PARs. In this paper, we primarily focus on the variant illustrated in Figure 2b, as it empirically promotes quantization not only at zero but also at nonzero levels. In

contrast, the version in Figure 2a exhibits local concavity around all breakpoints except zero, which makes it less effective at inducing quantization to nonzero levels.

Although Figure 2b technically violates the slope condition stated in Theorem 1, since it contains flat regions with zero slope, this issue is not problematic in practice. One can slightly perturb the zero slopes (for example, to $\pm\epsilon$ for a small $\epsilon > 0$) without significantly affecting the solution or its quantization guarantees. For clarity of exposition, we continue with the unperturbed version shown in Figure 2b.

We consider the following quasiconvex PAR with uniformly spaced quantization gaps

$$\Psi(x) = \begin{cases} |x| - \frac{k}{2}q & \text{if} \quad kq \le |x| \le \frac{2k+1}{2}q, \\ \frac{k+1}{2}q & \text{if} \quad \frac{2k+1}{2}q \le |x| \le (k+1)q. \end{cases} \tag{24}$$

We focus on this equal-gap setting since general nonuniform cases yield more complex proximal mappings. We now characterize the proximal mapping for this PAR:

- When the regularization strength $\lambda \le q$, we have

$$\mathbf{prox}_{\lambda\Psi}(x) = \begin{cases} \text{sign}(x)kq & \text{if} \quad kq \le |x| \le kq + \lambda, \\ \text{sign}(x)\left(|x| - \lambda\right) & \text{if} \quad kq + \lambda \le |x| \le \frac{2k+1}{2}q + \frac{\lambda}{2}, \\ \text{sign}(x)|x| & \text{if} \quad \frac{2k+1}{2}q + \frac{\lambda}{2} \le |x| \le (k+1)q. \end{cases} \tag{25}$$

- When the regularization strength $\lambda \ge q$, we have

$$\mathbf{prox}_{\lambda\Psi}(x) = \text{sign}(x)\left\lfloor \frac{|x| - \frac{\lambda}{2}}{q} \right\rceil q. \tag{26}$$

Unlike the convex PAR case, when $\lambda$ exceeds a certain threshold (e.g., $\lambda \ge q$), the proximal operator becomes a hard quantizer, mapping inputs exactly to discrete levels in the quantization set $\mathcal{Q}$.

**Nonconvex PAR.** Given the quantization values $\mathcal{Q} = \{q_1, \cdots, q_m\}$ satisfying $q_1 < q_2 < \cdots < q_m$ (here we allow general asymmetric quantization values), we consider the following nonconvex PAR

$$\Psi(x) = \begin{cases} x - q_k & \text{if } q_k \le x \le \frac{q_k + q_{k+1}}{2}, \\ -x + q_{k+1} & \text{if } \frac{q_k + q_{k+1}}{2} \le x \le q_{k+1}. \end{cases} \tag{27}$$

Its proximal mapping $\mathbf{prox}_{\lambda\Psi}(\cdot)$ admits the following closed-form solution:

$$\mathbf{prox}_{\lambda\Psi}(x) = \begin{cases} \text{clip}\left(x - \lambda, q_k, \frac{q_k + q_{k+1}}{2}\right) & \text{if } q_k \le x \le \frac{q_k + q_{k+1}}{2}, \\ \text{clip}\left(x + \lambda, \frac{q_k + q_{k+1}}{2}, q_{k+1}\right) & \text{if } \frac{q_k + q_{k+1}}{2} \le x \le q_{k+1}. \end{cases} \tag{28}$$

Here the clip function $\text{clip}(x, a, b)$ is defined by $\text{clip}(x, a, b) = \min\{\max\{x, a\}, b\}$. In particular, if $\lambda \ge \frac{1}{2}\max_{1 \le k \le m-1}(q_{k+1} - q_k)$, the proximal mapping reduces to an exact quantization mapping

$$\mathbf{prox}_{\lambda\Psi}(x) = \underset{q \in \mathcal{Q}}{\arg\min} |x - q|. \tag{29}$$

Similar to the quasiconvex PAR, this nonconvex PAR induces hard quantization once $\lambda$ exceeds a certain threshold. However, unlike the quasiconvex case, where the input magnitude is first reduced before snapping to the nearest quantization level, the proximal mapping of the nonconvex PAR remains fixed once $\lambda$ surpasses the threshold.

Another interesting feature to notice is that (see Figure 3), for convex and nonconvex PARs, the magnigude of their output is no larger than that of the input; but for nonconvex PARs, this is no longer true.

**Convex PAR**



(a) small $\lambda$      (b) medium $\lambda$      (c) large $\lambda$

**Quasiconvex PAR**



(d) small $\lambda$      (e) medium $\lambda$      (f) large $\lambda$

**Nonconvex PAR**



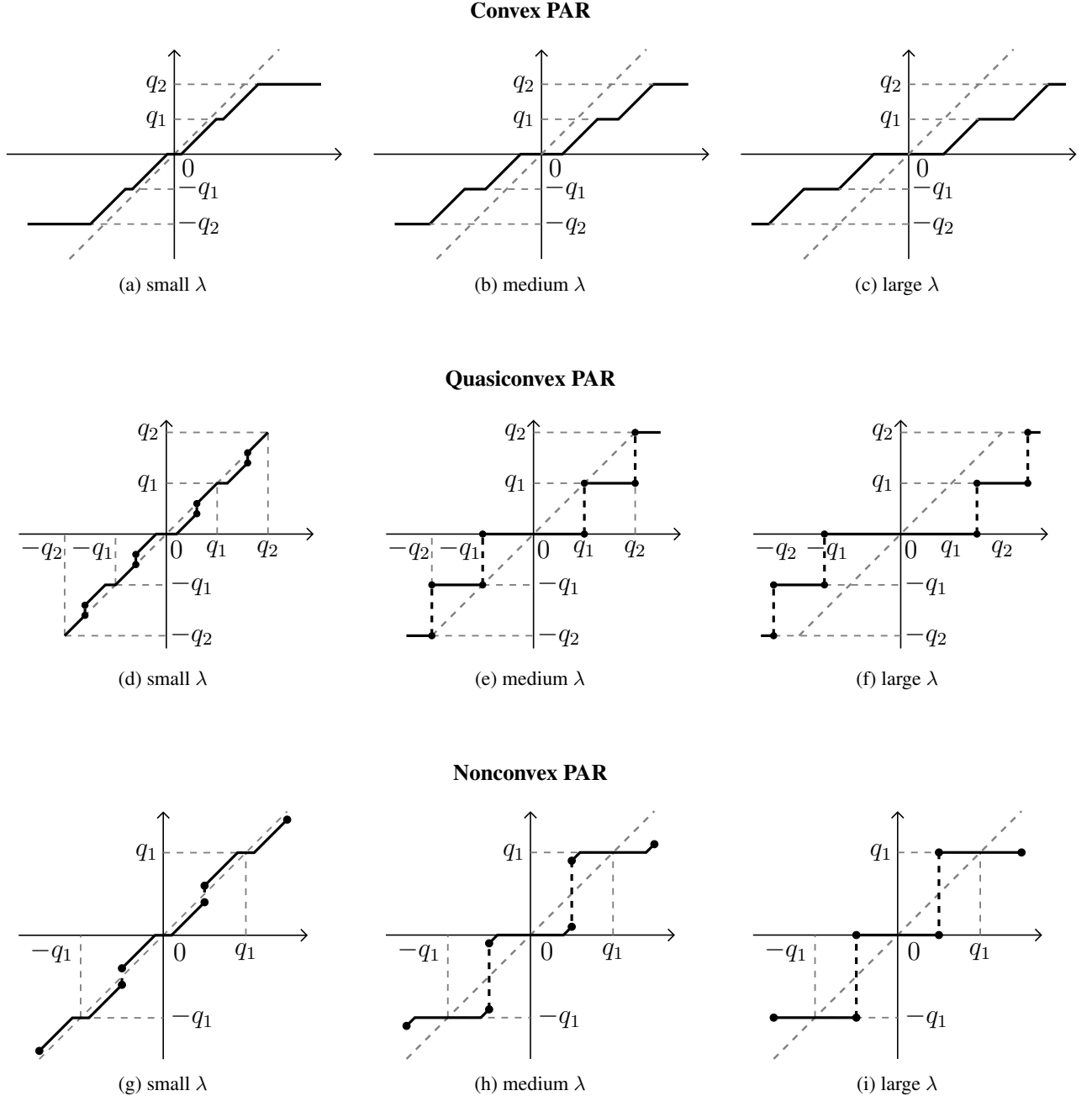(g) small $\lambda$      (h) medium $\lambda$      (i) large $\lambda$

Figure 3: Proximal mappings for different PARs. Each row corresponds to a class of PARs: convex (top), quasiconvex (middle), and nonconvex (bottom). Each column illustrates the proximal mapping under a different regularization strength: small, medium, and large $\lambda$.

## 3.2 Proximal gradient method

To minimize the PARO objective $F_\lambda(\boldsymbol{x}) = f(\boldsymbol{x}) + \lambda\Psi(\boldsymbol{x})$ where $f(\cdot)$ is a smooth loss function and $\Psi(\cdot)$ is a PAR, we first consider the classic proximal gradient method

$$\boldsymbol{x}^{t+1} = \mathbf{prox}_{\eta_t\lambda\Psi}\left(\boldsymbol{x}^t - \eta_t\nabla f(\boldsymbol{x}^t)\right). \tag{30}$$

Here $\eta_t > 0$ is the stepsize. In practice, the stepsize $\eta_t$ is often selected adaptively using a line search scheme. We refer the reader to [PB$^+$14, Bec17, WR22] for a comprehensive overview of proximal gradient methods. The following result gives the convergence guarantee of the proximal gradient algorithm (30) for PARO.

**Theorem 3.** *Suppose $f(\cdot)$ is L-smooth and the stepsize $\eta_t \equiv \eta \leq \frac{1}{2L}$. Then, the following arguments hold*

- *Convex case: If both $f(\cdot)$ and $\Psi(\cdot)$ are convex, then we have*

$$F_\lambda\left(\boldsymbol{x}^T\right) - F_\lambda^\star \leq \frac{\left\|\boldsymbol{x}^0 - \boldsymbol{x}^\star\right\|^2}{2\eta T}. \tag{31}$$

- *General case: If $F_\lambda(\boldsymbol{x}) = f(x) + \lambda\Psi(\boldsymbol{x})$ is proper and coercive, then the iterates $\{\boldsymbol{x}^t\}_{t=0}^\infty$ generated by (30) are bounded. Moreover, any accumulation point $\boldsymbol{x}^\star$ of $\{\boldsymbol{x}^t\}$ satisfies $\mathbf{0} \in \partial F_\lambda(\boldsymbol{x}^\star)$, i.e., $\boldsymbol{x}^\star$ is a critical point.*

These results demonstrate that the proximal gradient method efficiently converges to a critical point of the regularized loss function $F_\lambda(\cdot)$, which corresponds to the global minimum in the convex case. This is particularly desirable since, under mild conditions, all critical points in general supervised learning problems are highly quantized (Theorem 2). Combining these two findings, we conclude that the proximal gradient method efficiently converges to a highly quantized solution.

*Proof of Theorem 3.* The convex case follows directly from [BT09, Theorem 3.1]. For the general nonconvex case, our proof parallels the argument used in the proof of Theorem 1 in [LL15]. We include it here for completeness.

To begin, observe that

$$\begin{aligned}
&F_\lambda(\boldsymbol{x}^{t+1}) \\
&= f(\boldsymbol{x}^{t+1}) + \lambda\Psi(\boldsymbol{x}^{t+1}) \\
&\overset{(a)}{\leq} f(\boldsymbol{x}^t) + \left\langle\nabla f(\boldsymbol{x}^t), \boldsymbol{x}^{t+1} - \boldsymbol{x}^t\right\rangle + \frac{L}{2}\left\|\boldsymbol{x}^{t+1} - \boldsymbol{x}^t\right\|^2 + \lambda\Psi(\boldsymbol{x}^{t+1}) \\
&= f(\boldsymbol{x}^t) + \left\langle\nabla f(\boldsymbol{x}^t), \boldsymbol{x}^{t+1} - \boldsymbol{x}^t\right\rangle + \frac{1}{2\eta}\left\|\boldsymbol{x}^{t+1} - \boldsymbol{x}^t\right\|^2 + \lambda\Psi(\boldsymbol{x}^{t+1}) + \left(\frac{L}{2} - \frac{1}{2\eta}\right)\left\|\boldsymbol{x}^{t+1} - \boldsymbol{x}^t\right\|^2 \tag{32} \\
&\overset{(b)}{\leq} f(\boldsymbol{x}^t) + \lambda\Psi(\boldsymbol{x}^t) + \left(\frac{L}{2} - \frac{1}{2\eta}\right)\left\|\boldsymbol{x}^{t+1} - \boldsymbol{x}^t\right\|^2 \\
&\overset{(c)}{=} F_\lambda(\boldsymbol{x}^t) - \frac{1}{4\eta}\left\|\boldsymbol{x}^{t+1} - \boldsymbol{x}^t\right\|^2.
\end{aligned}$$

Here $(a)$ comes from the fact that $f(\cdot)$ is L-smooth; $(b)$ is due to the definition of $\boldsymbol{x}^{t+1}$; and in $(c)$ we use the condition that $\eta \leq \frac{1}{2L}$. As a result, $\{F_\lambda(\boldsymbol{x}^t)\}$ is nonincreasing. Since $\{F_\lambda(\cdot)\}$ is assumed to be proper and coercive, the iterates $\{\boldsymbol{x}^t\}$ is bounded. Hence, $\{\boldsymbol{x}^t\}$ has at least one accumulation point.

Next, we show that any accumulation point $\boldsymbol{x}^\star$ is a critical point of $F_\lambda(\cdot)$. By telescoping equation 32, we obtain

$$\frac{1}{T}\sum_{t=0}^{T-1}\left\|\boldsymbol{x}^{t+1}-\boldsymbol{x}^t\right\|^2 \leq \frac{\left(F_\lambda(\boldsymbol{x}^0)-F_\lambda^\star\right)}{4\eta T}, \tag{33}$$

which implies $\left\|\boldsymbol{x}^{t+1}-\boldsymbol{x}^t\right\| \to 0$ as $t \to \infty$. From the proximal update rule,

$$\boldsymbol{x}^{t+1} \in \arg\min_{\boldsymbol{z}}\left\{\eta_t\lambda\Psi(\boldsymbol{z}) + \frac{1}{2}\left\|\boldsymbol{z}-\left(\boldsymbol{x}^t-\eta_t\nabla f(\boldsymbol{x}^t)\right)\right\|^2\right\}. \tag{34}$$

The optimality condition gives

$$\boldsymbol{0} \in \left(\eta_t\lambda\partial\Psi(\boldsymbol{x}^{t+1}) + \boldsymbol{x}^{t+1} - \left(\boldsymbol{x}^t-\eta_t\nabla f(\boldsymbol{x}^t)\right)\right). \tag{35}$$

Let's define

$$\boldsymbol{v}^t := \boldsymbol{x}^t - \boldsymbol{x}^{t+1} + \eta_t\left(\nabla f(\boldsymbol{x}^t)-\nabla f(\boldsymbol{x}^{t+1})\right),$$

then (35) implies

$$\boldsymbol{v}^t \in \partial F_\lambda(\boldsymbol{x}^{t+1}). \tag{36}$$

Using the Lipschitz continuity of $\nabla f$, we have

$$\begin{aligned}
\left\|\boldsymbol{v}^t\right\| &= \left\|\boldsymbol{x}^t - \boldsymbol{x}^{t+1} + \eta_t\left(\nabla f(\boldsymbol{x}^t)-\nabla f(\boldsymbol{x}^{t+1})\right)\right\| \\
&\leq (1+\eta_t L)\left\|\boldsymbol{x}^t-\boldsymbol{x}^{t+1}\right\| \\
&= \frac{3}{2}\left\|\boldsymbol{x}^t-\boldsymbol{x}^{t+1}\right\|,
\end{aligned}$$

which converges to $0$ as $t \to \infty$. Therefore, $\left\|\boldsymbol{v}^t\right\| \to 0$ as $t \to \infty$.

For any accumulation point $\boldsymbol{x}^\star$, there exists a subsequence $\{\boldsymbol{x}^{t_k}\}$ such that $\boldsymbol{x}^{t_k} \to \boldsymbol{x}^\star$. Since $\boldsymbol{v}^{t_k} \in \partial F_\lambda(\boldsymbol{x}^{t_k})$ and $\left\|\boldsymbol{v}^{t_k}\right\| \to 0$, we conclude that $\boldsymbol{0} \in \partial F_\lambda(\boldsymbol{x}^\star)$ by the closedness of the limiting subdifferential. This completes the proof.

$\square$

**Accelerated proximal gradient methods.** To further accelerate convergence, one may apply the accelerated proximal gradient method, which incorporates a momentum term

$$\begin{aligned}
\boldsymbol{y}^{t+1} &= \boldsymbol{x}^t + \beta^t\left(\boldsymbol{x}^t-\boldsymbol{x}^{t-1}\right), \\
\boldsymbol{x}^{t+1} &= \mathbf{prox}_{\eta_t\lambda\Psi}\left(\boldsymbol{y}^{t+1}-\eta_t\nabla f\left(\boldsymbol{y}^{t+1}\right)\right).
\end{aligned} \tag{37}$$

Here $\beta_t \in (0,1)$ is the momentum coefficient and $\eta_t > 0$ is the stepsize. See [PB+14, Section 4.3] for more details about the choices of the stepsize and momentum coefficient. In the convex setting, classical results show that the accelerated method achieves a convergence rate of $O(1/T^2)$ for the regularized objective when appropriate hyperparameters are used [Van10, BT10]. This rate extends to nonconvex objectives with convex penalties, as shown in [Nes13]. For general nonconvex problems, more sophisticated variants have been developed; see [LL15].

## 3.3 Alternating direction method of multipliers (ADMM)

When the loss function $f(\cdot)$ exhibits certain structures, such as in linear or logistic regression, the alternating direction method of multipliers (ADMM) offers an effective alternative to the proximal gradient method, often yielding faster convergence in practice [BPC$^+$11].

To apply ADMM, we reformulate the original problem $\min_{\boldsymbol{x}} F_\lambda(\boldsymbol{x}) = f(\boldsymbol{x}) + \lambda \Psi(\boldsymbol{x})$ as the constrained problem

$$\min_{\boldsymbol{x}, \boldsymbol{z}} f(\boldsymbol{x}) + \lambda \Psi(\boldsymbol{z}) \quad \text{such that} \quad \boldsymbol{x} - \boldsymbol{z} = 0. \tag{38}$$

The corresponding augmented Lagrangian with penalty parameter $\rho$ is given by

$$L_\rho(\boldsymbol{x}, \boldsymbol{z}, \boldsymbol{y}) = f(\boldsymbol{x}) + \lambda \Psi(\boldsymbol{z}) + \rho \langle \boldsymbol{y}, \boldsymbol{x} - \boldsymbol{z} \rangle + \frac{\rho}{2} \|\boldsymbol{x} - \boldsymbol{z}\|^2, \tag{39}$$

where $\boldsymbol{y}$ is the dual variable. ADMM performs the following updates at each iteration

$$\boldsymbol{x}^{t+1} = \arg\min_{\boldsymbol{x}} L_\rho(\boldsymbol{x}, \boldsymbol{z}^t, \boldsymbol{y}^t), \quad \text{($x$-minimization)}$$

$$\boldsymbol{z}^{t+1} = \arg\min_{\boldsymbol{z}} L_\rho(\boldsymbol{x}^{t+1}, \boldsymbol{z}, \boldsymbol{y}^t) = \mathbf{prox}_{(\lambda/\rho) \cdot \Psi} \left( \boldsymbol{x}^{t+1} + \boldsymbol{y}^t \right), \quad \text{($z$-minimization)}$$

$$\boldsymbol{y}^{t+1} = \boldsymbol{y}^t + \rho \left( \boldsymbol{x}^{t+1} - \boldsymbol{z}^{t+1} \right). \quad \text{(dual update)}$$

When $f$ admits a convenient structure, the $x$-minimization can often be computed efficiently. For example, if $f(\boldsymbol{x}) = \frac{1}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|^2$, it admits a closed-form solution $\boldsymbol{x}^{t+1} = \left( \boldsymbol{A}^\top \boldsymbol{A} + \rho \boldsymbol{I} \right)^{-1} \left( \boldsymbol{A}^\top \boldsymbol{b} + \rho \left( \boldsymbol{z}^t - \boldsymbol{y}^t \right) \right)$. ADMM is known to achieve convergence rates comparable to the proximal gradient method in both convex and certain nonconvex settings [YH16]. Furthermore, when additional structural conditions are satisfied, including linear and logistic regression with convex PARs as special cases, ADMM enjoys a linear convergence rate [HL17]. In Section 5, we empirically compare the performance of ADMM against other benchmark methods on linear regression tasks.

# 4 Statistical guarantees of PAR for linear regression

In practice, we seek quantized solutions that not only achieve low training loss but also exhibit strong statistical guarantees. This section investigates the statistical properties of PAR-regularized solutions in the context of linear regression. Our main result shows that specific PAR formulations closely approximate widely used regularizers, including $\ell_2$-, $\ell_1$-, and more general nonconvex regularizers; see Figure 4 for an illustration. As a result, the corresponding PAR-regularized estimators inherit similar statistical guarantees as their classical counterparts.

Specifically, we consider the following PAR-regularized linear regression problem:

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} F_{\mathrm{PAR}}(\boldsymbol{x}) = \frac{1}{2n} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|^2 + \lambda \Psi(\boldsymbol{x}). \tag{40}$$

Here, $\boldsymbol{A} \in \mathbb{R}^{n \times d}$ is the design matrix and $\boldsymbol{b} \in \mathbb{R}^n$ is the response vector. We will specify the data model and the PAR formulation in each subsection. Throughout this section, we denote the PAR-regularized solution by $\boldsymbol{x}_{\mathrm{PAR}}^\star = \arg\min_{\boldsymbol{x} \in \mathbb{R}^d} F_{\mathrm{PAR}}(\boldsymbol{x})$.

(a) PAR as $x^2$      (b) PAR as $|x|$      (c) PAR as $\sqrt{|x|}$
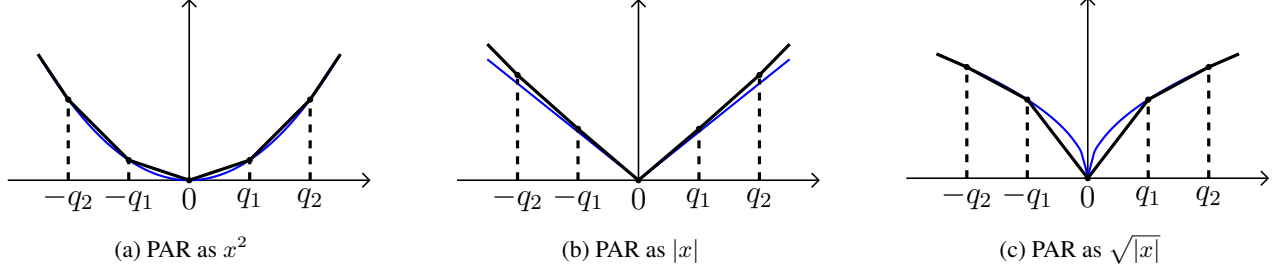
Figure 4: PARs as approximations to classical regularizers. The PARs in black color approximate the standard regularizers shown in blue.

## 4.1 PAR as ridge regression

In this section, we demonstrate that a special class of PARs can effectively approximate ridge regression

$$F_{\text{ridge}}(\boldsymbol{x}) = \frac{1}{2n}\|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|^2 + \frac{\lambda}{2}\|\boldsymbol{x}\|^2, \tag{41}$$

which admits a closed-form solution

$$\boldsymbol{x}_{\text{ridge}}^{\star} = \arg\min F_{\text{ridge}}(\boldsymbol{x}) = (\boldsymbol{A}^\top \boldsymbol{A} + n\lambda \boldsymbol{I})^\top \boldsymbol{A}^\top \boldsymbol{b}.$$

**PAR formulation.** We consider a special class of PARs, illustrated in Figure 4a. For this class, the quantization set is $\mathcal{Q} = \{0, \pm q, \pm 2q, \cdots\}$, and the slope set is $\mathcal{A} = \{q, 2q, \cdots\}$ where $q$ is the distance between adjacent quantization levels.[3]

**Data model.** We consider the fixed design regime, where the design matrix $\boldsymbol{A}$ is fixed. The response vector $\boldsymbol{b}$ is generated by $\boldsymbol{b} = \boldsymbol{A}\boldsymbol{x}^{\star} + \boldsymbol{\epsilon}$, with $\boldsymbol{\epsilon} = [\epsilon_1, \ldots, \epsilon_n]^\top$ where each $\epsilon_i$ is an independent random variable with zero mean and variance $\sigma^2$.[4]

We denote the sample covariance matrix by $\widehat{\boldsymbol{\Sigma}} := \frac{1}{n}\boldsymbol{A}^\top \boldsymbol{A}$. For any estimator $\boldsymbol{x}$, we define the in-sample risk and excess risk as

$$\mathcal{R}(\boldsymbol{x}) := \mathbb{E}_{\tilde{\boldsymbol{b}} \sim \mathcal{D}}\left[\frac{1}{n}\left\|\boldsymbol{A}\boldsymbol{x} - \tilde{\boldsymbol{b}}\right\|^2\right], \quad \mathcal{E}(\boldsymbol{x}) = \mathcal{R}(\boldsymbol{x}) - \mathcal{R}^{\star} = \|\boldsymbol{x} - \boldsymbol{x}^{\star}\|_{\widehat{\boldsymbol{\Sigma}}}^2. \tag{42}$$

Let $\mathcal{R}^{\star} = \min_{\boldsymbol{x} \in \mathbb{R}^d} \mathcal{R}(\boldsymbol{x})$ represents the optimal risk. Our next theorem characterizes the distances between the PAR and the ridge solutions.

**Theorem 4.** *The distances between $\boldsymbol{x}_{\text{PAR}}^{\star}$ and $\boldsymbol{x}_{\text{ridge}}^{\star}$ are characterized as follows*

$$\left\|\boldsymbol{x}_{\text{PAR}}^{\star} - \boldsymbol{x}_{\text{ridge}}^{\star}\right\| \leq \sqrt{\frac{d}{2}}q, \quad \text{and} \quad \left\|\boldsymbol{x}_{\text{PAR}}^{\star} - \boldsymbol{x}_{\text{ridge}}^{\star}\right\|_{\widehat{\boldsymbol{\Sigma}}} \leq \sqrt{\frac{d\lambda}{2}}q. \tag{43}$$

The above two bounds both scale with the quantization level $q$. However, unlike the distance evaluated in the $\ell_2$-norm, the one in the Mahalanobis norm $\|\cdot\|_{\widehat{\boldsymbol{\Sigma}}}$ depends on the regularization strength $\lambda$. Specifically, as the regularization strength $\lambda$ decreases, the two estimators become closer and closer in the Mahalanobis norm

---

[3]We can also design PARs with nonuniform quantization intervals to approximate the $\ell_2$-regularizer.

[4]It may be possible to extend these results to the random design regime by incorporating results from [HKZ14].

$\|\cdot\|_{\widehat{\boldsymbol{\Sigma}}}$. This discrepancy arises because the least-square loss $f(\boldsymbol{x}) = \frac{1}{2n}\|\boldsymbol{Ax} - \boldsymbol{b}\|^2$ is not strongly convex with respect to the $\ell_2$-norm when $n \ll d$, whereas it remains strongly convex in the Mahalanobis norm. This result is particularly appealing, as our next theorem establishes that the PAR-regularized solution enjoys statistical guarantees comparable to the ridge estimator, provided that the quantization gap $q$ is sufficiently small.

*Proof.* To start with, we first show that the two loss functions $F_{\mathrm{PAR}}(\boldsymbol{x})$ and $F_{\mathrm{ridge}}(\boldsymbol{x})$ are uniformly close to each other in the sense that

$$\sup_{\boldsymbol{x}\in\mathbb{R}^d}\left\{F_{\mathrm{PAR}}(\boldsymbol{x}) - F_{\mathrm{ridge}}(\boldsymbol{x})\right\} \leq \frac{d\lambda q^2}{8}. \tag{44}$$

To this end, we first consider the 1-dimensional function $\Psi(x) - \frac{1}{2}x^2$ for $x \in \mathbb{R}$. Suppose that $x \in [kq, (k+1)q)$ for some $k \in \mathbb{Z}$. Without loss of generality, we assume $k \geq 0$. Then, we have

$$0 \leq \Psi(x) - \frac{1}{2}x^2 = \left(k + \frac{1}{2}\right)q(x - kq) + \frac{k^2q^2}{2} - \frac{1}{2}x^2. \tag{45}$$

It attains its maximum at the point $x^\star = \left(k + \frac{1}{2}\right)q$ with the optimal value $\frac{1}{8}q^2$. Then, the argument follows by taking summation over all the coordinates.

Provided this result, we can establish the following relationship between the losses $F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{PAR}})$ and $F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{ridge}})$:

$$F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{PAR}}) \leq F_{\mathrm{PAR}}(\boldsymbol{x}^\star_{\mathrm{PAR}}) + \frac{d\lambda q^2}{8} \leq F_{\mathrm{PAR}}(\boldsymbol{x}^\star_{\mathrm{ridge}}) + \frac{d\lambda q^2}{8} \leq F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{ridge}}) + \frac{d\lambda q^2}{4}. \tag{46}$$

Now, we first consider the distance in the $\ell_2$-norm. Note that the regularizer $\frac{\lambda}{2}x^2$ is $\lambda$-strongly convex with respect to the $\ell_2$-norm and the loss function $\frac{1}{2n}\|\boldsymbol{Ax} - \boldsymbol{b}\|^2$ is convex. Hence, we have

$$\begin{aligned}
F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{PAR}}) &\geq F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{ridge}}) + \left\langle \nabla F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{ridge}}), \boldsymbol{x}^\star_{\mathrm{PAR}} - \boldsymbol{x}^\star_{\mathrm{ridge}} \right\rangle + \frac{\lambda}{2}\left\|\boldsymbol{x}^\star_{\mathrm{PAR}} - \boldsymbol{x}^\star_{\mathrm{ridge}}\right\|^2 \\
&= F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{ridge}}) + \frac{\lambda}{2}\left\|\boldsymbol{x}^\star_{\mathrm{PAR}} - \boldsymbol{x}^\star_{\mathrm{ridge}}\right\|^2.
\end{aligned} \tag{47}$$

Here we use the optimality condition that $\nabla F_{\mathrm{ridge}}(\boldsymbol{x}^\star_{\mathrm{ridge}}) = \boldsymbol{0}$. Substituting this inequality into Equation (46), we derive the desired result

$$\left\|\boldsymbol{x}^\star_{\mathrm{PAR}} - \boldsymbol{x}^\star_{\mathrm{ridge}}\right\| \leq \sqrt{\frac{d}{2}}q. \tag{48}$$

To control the Mahalanobis norm, we note that the least-square loss $f(\boldsymbol{x}) = \frac{1}{2n}\|\boldsymbol{Ax} - \boldsymbol{b}\|^2$ is 1-strongly convex with respect to the norm $\|\cdot\|_{\widehat{\boldsymbol{\Sigma}}}$. To show this, it suffices to prove that for any $t \in [0, 1]$ and any $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d$, the following inequality holds

$$f(t\boldsymbol{x} + (1-t)\boldsymbol{y}) \leq tf(\boldsymbol{x}) + (1-t)f(\boldsymbol{y}) - \frac{t(1-t)}{2}\|\boldsymbol{x} - \boldsymbol{y}\|^2_{\widehat{\boldsymbol{\Sigma}}}. \tag{49}$$

This follows by the following equality

$$\begin{aligned}
f(t\boldsymbol{x} + (1-t)\boldsymbol{y}) &= \frac{1}{2n}\|t(\boldsymbol{Ax} - \boldsymbol{b}) + (1-t)(\boldsymbol{Ay} - \boldsymbol{b})\|^2 \\
&= tf(\boldsymbol{x}) + (1-t)f(\boldsymbol{y}) - \frac{t(1-t)}{2n}\|\boldsymbol{A}(\boldsymbol{x} - \boldsymbol{y})\|^2 \\
&= tf(\boldsymbol{x}) + (1-t)f(\boldsymbol{y}) - \frac{t(1-t)}{2}\|\boldsymbol{x} - \boldsymbol{y}\|^2_{\widehat{\boldsymbol{\Sigma}}}.
\end{aligned} \tag{50}$$

15

Therefore, leveraging the property of the strong convexity, we have

$$
\begin{aligned}
F_{\mathrm{ridge}}(\boldsymbol{x}^{\star}_{\mathrm{PAR}}) &\geq F_{\mathrm{ridge}}(\boldsymbol{x}^{\star}_{\mathrm{ridge}}) + \left\langle \nabla F_{\mathrm{ridge}}(\boldsymbol{x}^{\star}_{\mathrm{ridge}}), \boldsymbol{x}^{\star}_{\mathrm{PAR}} - \boldsymbol{x}^{\star}_{\mathrm{ridge}} \right\rangle + \frac{1}{2} \left\| \boldsymbol{x}^{\star}_{\mathrm{PAR}} - \boldsymbol{x}^{\star}_{\mathrm{ridge}} \right\|^{2}_{\widehat{\boldsymbol{\Sigma}}} \\
&= F_{\mathrm{ridge}}(\boldsymbol{x}^{\star}_{\mathrm{ridge}}) + \frac{1}{2} \left\| \boldsymbol{x}^{\star}_{\mathrm{PAR}} - \boldsymbol{x}^{\star}_{\mathrm{ridge}} \right\|^{2}_{\widehat{\boldsymbol{\Sigma}}}.
\end{aligned}
\tag{51}
$$

Substituting this inequality into Equation (46), we derive that

$$
\left\| \boldsymbol{x}^{\star}_{\mathrm{PAR}} - \boldsymbol{x}^{\star}_{\mathrm{ridge}} \right\|^{2}_{\widehat{\boldsymbol{\Sigma}}} \leq \frac{d \lambda q^{2}}{2}.
\tag{52}
$$

This completes the proof. $\qquad \square$

Our next theorem characterizes the excess risk of the PAR-regularized solution.

**Theorem 5** (Excess risk). *Consider the data model and PAR formulation described above. If the regularization strength in (40) is chosen as* $\lambda = \Theta\left( \frac{\sigma}{\|\boldsymbol{x}^{\star}\| + \sqrt{d}q} \sqrt{\frac{\operatorname{tr}(\widehat{\boldsymbol{\Sigma}})}{n}} \right)$, *the excess risk of the PAR-regularized solution* $\boldsymbol{x}^{\star}_{\mathrm{PAR}}$ *satisfies*

$$
\mathcal{E}(\boldsymbol{x}^{\star}_{\mathrm{PAR}}) \lesssim \sigma \left( \|\boldsymbol{x}^{\star}\| + \sqrt{d}q \right) \sqrt{\frac{\operatorname{tr}(\widehat{\boldsymbol{\Sigma}})}{n}}.
\tag{53}
$$

**Statistical guarantee.** If the quantization level is chosen as $q \leq \frac{\|\boldsymbol{x}^{\star}\|}{\sqrt{d}}$, the excess risk of the PAR-regularized solution is $O\left( \sigma \|\boldsymbol{x}^{\star}\| \sqrt{\operatorname{tr}(\widehat{\boldsymbol{\Sigma}})/n} \right)$, which is the same as that of the ridge regression estimator.

**Quantization guarantee.** Theorem 1 implies that at least $d - n$ coordinates of $\boldsymbol{x}^{\star}_{\mathrm{PAR}}$ are quantized. Moreover, if the covariate spectrum is concentrated in only a few directions, which means that the effective rank $\operatorname{tr}(\widehat{\boldsymbol{\Sigma}})/\|\widehat{\boldsymbol{\Sigma}}\| \ll d$, then the required sample complexity can be significantly smaller than the dimensionality, i.e., $n \ll d$. In this scenario, most coordinates of $\boldsymbol{x}^{\star}_{\mathrm{PAR}}$ are quantized.

**Comparison with simple estimators.** One might argue that the quantization result derived in Theorem 1 is too weak. For instance, we can easily construct an estimator with the same quantization guarantee: we simply set $d - n$ coordinates to be zero and pick the remaining $n$ coordinates by solving a linear equation $\boldsymbol{A}_{:n}\boldsymbol{x}_{:n} = \boldsymbol{b}$. While this approach achieves zero training loss, it fails to generalize well, highlighting the advantage of the PAR estimator, which maintains both quantization efficiency and strong statistical guarantees.

**Storage advantage.** The PAR estimator offers significant storage benefits compared to the ridge estimator. To quantify this, we first consider the ridge estimator, whose parameters are typically dense and stored in single-precision floating-point format (FP32), requiring 32 bits to store each parameter. In contrast, the PAR estimator reduces storage via quantization. To match the statistical accuracy of ridge regression, it suffices to select the quantization gap $q \leq \frac{\|\boldsymbol{x}^{\star}\|}{\sqrt{d}}$. A crude bound on the maximum entry of the PAR estimator gives

$$
\left\| \boldsymbol{x}^{\star}_{\mathrm{PAR}} \right\|_{\infty} \leq \left\| \boldsymbol{x}^{\star}_{\mathrm{PAR}} - \boldsymbol{x}^{\star}_{\mathrm{ridge}} \right\| + \left\| \boldsymbol{x}^{\star}_{\mathrm{ridge}} \right\| \leq \sqrt{\frac{d}{2}} q + \|\boldsymbol{x}^{\star}\| \leq 2 \|\boldsymbol{x}^{\star}\|.
\tag{54}
$$

Therefore, storing the magnitude of each parameter requires at most $\lceil \ln(2\sqrt{d}) \rceil$ bits, plus one additional bit for the sign. For dimensions up to $1 \times 10^8$, this amounts to roughly 16 bits per parameter, resulting in a $2\times$ reduction in storage compared to the ridge estimator.

Before proving Theorem 5, we first introduce the following classic textbook result on Ridge regression.

**Proposition 1** (Proposition 3.7 in [Bac24]). *For ridge estimator $\boldsymbol{x}_{\mathrm{ridge}}^\star$, its excess risk is characterized as*

$$\mathcal{E}(\boldsymbol{x}_{\mathrm{ridge}}^\star) \leq \frac{\lambda}{2} \|\boldsymbol{x}^\star\|^2 + \frac{\sigma^2 \mathrm{tr}(\widehat{\boldsymbol{\Sigma}})}{2\lambda n}. \tag{55}$$

We now proceed to prove Theorem 5.

*Proof of Theorem 5.* First, we have the following decomposition

$$\mathcal{E}(\boldsymbol{x}_{\mathrm{PAR}}^\star) = \|\boldsymbol{x}_{\mathrm{PAR}}^\star - \boldsymbol{x}^\star\|_{\widehat{\boldsymbol{\Sigma}}}^2 \leq 2 \|\boldsymbol{x}_{\mathrm{PAR}}^\star - \boldsymbol{x}_{\mathrm{ridge}}^\star\|_{\widehat{\boldsymbol{\Sigma}}}^2 + 2 \|\boldsymbol{x}_{\mathrm{ridge}}^\star - \boldsymbol{x}^\star\|_{\widehat{\boldsymbol{\Sigma}}}^2 . \tag{56}$$

By Theorem 4, the first term can be bounded as

$$\|\boldsymbol{x}_{\mathrm{PAR}}^\star - \boldsymbol{x}_{\mathrm{ridge}}^\star\|_{\widehat{\boldsymbol{\Sigma}}} \leq \sqrt{\frac{d\lambda}{2}} q. \tag{57}$$

For the second term, applying Proposition 1 yields

$$\mathcal{E}(\boldsymbol{x}_{\mathrm{ridge}}^\star) = \|\boldsymbol{x}_{\mathrm{ridge}}^\star - \boldsymbol{x}^\star\|_{\widehat{\boldsymbol{\Sigma}}}^2 \leq \frac{\lambda}{2} \|\boldsymbol{x}^\star\|^2 + \frac{\sigma^2 \mathrm{tr}(\widehat{\boldsymbol{\Sigma}})}{2\lambda n}. \tag{58}$$

Combining the above two recipes, we obtain that

$$\mathcal{E}(\boldsymbol{x}_{\mathrm{PAR}}^\star) \leq \lambda \left( dq^2 + \|\boldsymbol{x}^\star\|^2 \right) + \frac{\sigma^2 \mathrm{tr}(\widehat{\boldsymbol{\Sigma}})}{\lambda n}. \tag{59}$$

Therefore, upon setting $\lambda = \Theta\left( \frac{\sigma}{\|\boldsymbol{x}^\star\| + \sqrt{d}q} \sqrt{\frac{\mathrm{tr}(\widehat{\boldsymbol{\Sigma}})}{n}} \right)$, we derive that

$$\mathcal{E}(\boldsymbol{x}_{\mathrm{PAR}}^\star) \lesssim \sigma \left( \|\boldsymbol{x}^\star\| + \sqrt{d}q \right) \sqrt{\frac{\mathrm{tr}(\widehat{\boldsymbol{\Sigma}})}{n}}. \tag{60}$$

This completes the proof. $\square$

## 4.2 PAR as Lasso and nonconvex regularizers

In this section, we demonstrate that a special class of PAR regularizers can effectively approximate the Lasso objective

$$F_{\mathrm{Lasso}}(\boldsymbol{x}) = \frac{1}{2n} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|^2 + \frac{\lambda}{2} \|\boldsymbol{x}\|_1 . \tag{61}$$

as well as the nonconvex regularizers, including several commonly used in sparse linear regression, such as the bridge regularizer ($\Psi(\boldsymbol{x}) = \|\boldsymbol{x}\|_p^p$ for $0 < p < 1$), smoothly clipped absolute deviations (SCAD) penalty [FL01], and minimax concave penalty (MCP) [Zha10]. Nonconvex regularizers are introduced to better approximate the $\ell_0$-norm and to address the bias issue of the Lasso, which penalizes large coefficients more heavily. See [ZZ12] for a comprehensive survey of nonconvex regularizers.

While the $\ell_1$-regularizer $\Psi(\boldsymbol{x}) = \|\boldsymbol{x}\|_1$ is itself a special case of a PAR, it primarily encourages sparsity by promoting solutions concentrated at zero. In contrast, we propose a richer class of PARs that not only approximate the $\ell_1$-penalty but also introduce additional quantization levels beyond zero. This structure enables the regularizer to promote parameter values near multiple predefined levels, thereby facilitating quantization while retaining statistical properties comparable to those of the standard Lasso solution.

**Data model.** We consider the following sparse linear regression setting:

**Assumption 2.** *The true solution $\boldsymbol{x}^\star$ is $s$-sparse, i.e., its support $S = \operatorname{supp}(\boldsymbol{x}^\star)$ satisfies $|S| = s$.*

**Assumption 3** (Restricted eigenvalue). *We assume the design matrix $\boldsymbol{A}$ satisfies the restricted eigenvalue (RE) condition over $S = \operatorname{supp}(\boldsymbol{x}^\star)$ with parameters $(\alpha, \gamma)$, that is*

$$\frac{1}{n}\left\|\boldsymbol{A}\boldsymbol{v}\right\|^2 \geq \gamma \left\|\boldsymbol{v}\right\|^2, \quad \forall \boldsymbol{v} \in \mathcal{C}_\alpha(S) := \{\boldsymbol{v} \in \mathbb{R}^d : \left\|\boldsymbol{v}_{S^c}\right\|_1 \leq \alpha \left\|\boldsymbol{v}_S\right\|_1\}. \tag{62}$$

This condition holds for a broad class of random design matrices, particularly those with sub-Gaussian or isotropic rows, provided that the sample size is sufficiently large relative to the sparsity level $s$. In particular, for a Gaussian design matrix, where $\boldsymbol{A} \in \mathbb{R}^{n \times d}$ has i.i.d. $\mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$ rows, the restricted eigenvalue condition with parameters $(\alpha, \gamma)$ holds with probability at least $1 - \exp\{-\Omega(n)\}$ as long as the sample size $n \gtrsim \frac{\|\boldsymbol{\Sigma}\|^2(1+\alpha)^2}{\gamma}s\log(d)$ [RWY10, Corollary 1].

**PAR formulation.** We consider a PAR $\Psi(\cdot)$ with quantization values $\mathcal{Q} = \{0, \pm q_1, \pm q_2, \cdots\}$ where $0 < q_1 < q_2 < \cdots$ and slopes $\mathcal{A} = \{\pm a_1, \pm a_2, \cdots\}$. Let $a_{\max} = \max\{a \in \mathcal{A}\}$ denote the maximum slope magnitude. We assume that $\Psi$ satisfies a linear growth condition, i.e., there exists a universal constant $k > 0$ such that $\Psi(x) \geq \nu|x|$ for all $x \in \mathbb{R}$.

The following two examples illustrate this class of PARs.

**Example 1** (Convex PAR). *Consider a convex PAR $\Psi(\cdot)$ with arbitrary quantization values $\mathcal{Q}$ and slopes $\mathcal{A} = \{a_0, a_1, \cdots, a_m\}$ with $0 < a_0 < a_1 < \cdots < a_m < \infty$. In this case, $a_{\max} = a_m$ and $\nu = a_0$ since $\Psi(x) \geq a_0 x$ for all $x \in \mathbb{R}$.*

**Example 2** (Quasiconvex PAR). *Consider the quasiconvex PAR in Figure 2b, which is characterized by*

$$\Psi(x) = \begin{cases} |x| - \frac{k}{2}q & \text{if} \quad kq \leq |x| \leq \frac{2k+1}{2}q, \\ \frac{k+1}{2}q & \text{if} \quad \frac{2k+1}{2}q \leq |x| \leq (k+1)q, \end{cases} \tag{63}$$

*for integer $k \geq 0$ and fixed step size $q > 0$. It is straightforward to verify that this function is quasiconvex, with $a_{\max} = 1$ and $\nu = \frac{1}{2}$.*

We now characterize the statistical guarantees for the PAR-regularized solution.

**Theorem 6.** *Consider the data model and PAR formulation described above. Suppose that the regularization strength satisfies $\lambda \geq \frac{\|\boldsymbol{A}^\top \boldsymbol{\epsilon}\|_\infty}{2\nu n}$, and the design matrix $\boldsymbol{A}$ satisfies restricted eigenvalue condition with parameters $\left(\frac{3a_{\max}}{\nu}, \gamma\right)$. Then, the estimation error of the PAR-regularized solution is bounded as*

$$\left\|\boldsymbol{x}^\star_{\mathrm{PAR}} - \boldsymbol{x}^\star\right\| \leq \frac{3\lambda a_{\max}\sqrt{s}}{\gamma}. \tag{64}$$

Similar guarantees can be derived for the prediction error $\|\boldsymbol{A}(\boldsymbol{x}^\star_{\mathrm{PAR}} - \boldsymbol{x}^\star)\|$ and the estimation error in $\ell_\infty$-norm $\|\boldsymbol{x}^\star_{\mathrm{PAR}} - \boldsymbol{x}^\star\|_\infty$. However, we omit these results here and leave them for future work. This result closely resembles the classic error bound for Lasso regression [HTW15, Theorem 11.1]. In particular, if $a_{\max} \asymp k$, then the guarantees in Theorem 6 match those of Lasso up to constant factors. To illustrate this result, we consider the classical linear Gaussian model.

**Corollary 1.** *Suppose $a_{\max} = \Theta(\nu)$. Assume the design matrix $\boldsymbol{A}$ has i.i.d. standard Gaussian entries and the noise vector $\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{0}, \sigma^2 \boldsymbol{I})$ is also i.i.d. Gaussian. If the sample size satisfies $n \gtrsim s \log(d)$, then with probability at least $1 - \exp{-\Omega(\log(d))}$, the estimation error is bounded by*

$$\|\boldsymbol{x}_{\mathrm{PAR}}^\star - \boldsymbol{x}^\star\| \lesssim \sigma\sqrt{\frac{s \log(d)}{n}}. \tag{65}$$

*Proof of Theorem 6.* Note that $\boldsymbol{x}_{\mathrm{PAR}}^\star$ is the minimizer of $F_{\mathrm{PAR}}(\boldsymbol{x})$. Hence, we have $F_{\mathrm{PAR}}(\boldsymbol{x}_{\mathrm{PAR}}^\star) \leq F_{\mathrm{PAR}}(\boldsymbol{x}^\star)$, that is,

$$\frac{1}{2n}\|\boldsymbol{A}\boldsymbol{x}_{\mathrm{PAR}}^\star - \boldsymbol{b}\|^2 + \lambda\Psi(\boldsymbol{x}_{\mathrm{PAR}}^\star) \leq \frac{1}{2n}\|\boldsymbol{\epsilon}\|^2 + \lambda\Psi(\boldsymbol{x}^\star). \tag{66}$$

Rearranging this inequality and denoting $\boldsymbol{v} = \boldsymbol{x}_{\mathrm{PAR}}^\star - \boldsymbol{x}^\star$ yields

$$
\begin{aligned}
\frac{1}{2n}\|\boldsymbol{A}\boldsymbol{v}\|^2 &\leq \frac{1}{n}\left\langle \boldsymbol{A}^\top\boldsymbol{\epsilon}, \boldsymbol{v}\right\rangle + \lambda\left(\Psi(\boldsymbol{x}^\star) - \Psi(\boldsymbol{x}_{\mathrm{PAR}}^\star)\right)\\
&\overset{(a)}{\leq} \frac{\|\boldsymbol{A}^\top\boldsymbol{\epsilon}\|_\infty}{n}\|\boldsymbol{v}\|_1 + \lambda\left(\Psi(\boldsymbol{x}_S^\star) - \Psi(\boldsymbol{x}_{\mathrm{PAR},S}^\star)\right) - \lambda\Psi\left(\boldsymbol{x}_{\mathrm{PAR},S^c}^\star\right)\\
&\overset{(b)}{\leq} \frac{\|\boldsymbol{A}^\top\boldsymbol{\epsilon}\|_\infty}{n}\left(\|\boldsymbol{v}_S\|_1 + \|\boldsymbol{v}_{S^c}\|_1\right) + \lambda a_{\max}\|\boldsymbol{v}_S\|_1 - \lambda\nu\|\boldsymbol{v}_{S^c}\|_1\\
&\overset{(c)}{\leq} \frac{3}{2}\lambda a_{\max}\|\boldsymbol{v}_S\|_1 - \frac{1}{2}\lambda\nu\|\boldsymbol{v}_{S^c}\|_1.
\end{aligned}
\tag{67}
$$

Here in $(a)$, we use Hölder's inequality and the fact that $\Psi(x_i^\star) = 0$ for all $i \in S^c$. In $(b)$, we use the triangle inequality and the facts that $\Psi(\cdot)$ is $a_{\max}$-Lipschitz and $\Psi(x) \geq \nu|x|$ for all $x$. In $(c)$, we use the condition that $\lambda \geq \frac{\|\boldsymbol{A}^\top\boldsymbol{\epsilon}\|_\infty}{2\nu n}$. To proceed, note that $\|\boldsymbol{A}\boldsymbol{v}\|^2 \geq 0$, which implies that $\|\boldsymbol{v}_{S^c}\|_1 \leq \frac{3a_{\max}}{\nu}\|\boldsymbol{v}_S\|_1$. Therefore, we can apply the restricted eigenvalue condition, which yields

$$\gamma\|\boldsymbol{v}\|^2 \overset{\text{RE condition}}{\leq} \frac{1}{n}\|\boldsymbol{A}\boldsymbol{v}\|^2 \overset{\text{equation } 67}{\leq} 3\lambda a_{\max}\|\boldsymbol{v}_S\|_1 - \lambda\nu\|\boldsymbol{v}_{S^c}\|_1 \leq 3\lambda a_{\max}\sqrt{s}\|\boldsymbol{v}\|. \tag{68}$$

This accomplishes the proof. $\qquad\square$

# 5 Numerical experiments

In this section, we present numerical experiments to validate our theoretical results on quantization, optimization, and statistical performance. In Section 5.1, we demonstrate that the lower bound on the quantization rate is nearly tight in the linear regression setting. Section 5.2 evaluates the performance of various optimization algorithms across different PAR formulations. Finally, Section 5.3 examines the statistical performance of PAR-regularized models in both linear and logistic regression tasks. Our code is available at https://github.com/jianhaoma/paro.

## 5.1 Validation of the quantization guarantee

We first verify the tightness of the quantization rate lower bound given in Theorem 1, which states that the quantization rate is at least $1 - n/d$ where $n$ is the sample size and $d$ is the data dimension, and is independent of the regularization strength $\lambda$. To test this, we conduct extensive simulations on a linear regression task.
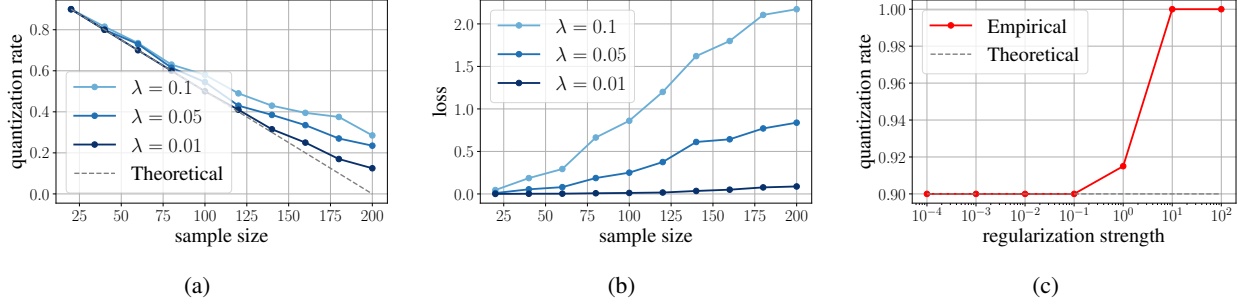
19

Figure 5: Empirical validation of the quantization guarantee on linear regression with $d = 200$. Panel (a) tests the effect of sample size; panel (b) shows the impact of $\lambda$ on training loss; and panel (c) examines the robustness of the quantization rate to $\lambda$ for the case $n = 20$ and $d = 200$.

The data dimension is set to $d = 200$, and the input matrix $A$ is generated with i.i.d. Gaussian entries. The ground truth $x^\star$ is randomly sampled, and the output is computed as $y = Ax^\star$ without additional noise.

In our experiments, we use a convex PAR with quantization values $\mathcal{Q} = \mathbb{Z}$ and slopes $\mathcal{A} = \{\ldots, -2, -1, 1, 2, \ldots\}$. Figure 5a reports the observed quantization rate across varying sample sizes and regularization strengths. The results closely match the theoretical lower bound and show that the quantization rate remains largely unaffected by $\lambda$, particularly when the sample size is small. This observation is further supported by Figure 5c, where $\lambda$ is varied from $10^{-4}$ to $100$: the quantization rate consistently exceeds $0.9$, aligning with the theoretical bound of $1 - n/d = 1 - 20/200 = 0.9$.

Additionally, Figure 5b shows that increasing $\lambda$ leads to higher training loss. Therefore, to balance quantization and model performance, we recommend using a relatively small regularization strength in practice.

## 5.2 Comparison of different optimization algorithms and PAR variants

**Convergence Behavior Across Optimization Algorithms.** We evaluate the optimization performance of three algorithms: proximal gradient (`PG`), accelerated proximal gradient (`acc_PG`), and ADMM (`ADMM`). The task is a linear regression problem regularized by three types of PARs: convex, quasiconvex, and nonconvex. The implementations of these three algorithms follow Sections 3.2 and 3.3. All methods incorporate a backtracking line search to select the step size.

We use synthetic data with feature dimension $d = 200$ and sample size $n = 20$. The true parameter $x^\star$ is generated randomly, and the design matrix $A$ is drawn from a standard Gaussian distribution. The response vector is generated as $b = Ax^\star + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.01I)$.

Figure 6 summarizes the convergence behavior across different regularizers. For the convex PAR, all algorithms exhibit linear convergence, with `acc_PG` achieving the fastest rate, followed by `ADMM` and then `PG`. For the quasiconvex PAR, `ADMM` substantially outperforms both `PG` and `acc_PG`. In the case of the nonconvex PAR, all three algorithms show comparable convergence patterns. Interestingly, the optimal objective value is typically achieved early, prior to full convergence to a critical point, which might be attributed to the nonconvex nature of the regularizer.

**Effect of PAR structure on quantization performance.** In this simulation, we examine how the structure of the PAR, whether convex, quasiconvex, or nonconvex, influences the quality of the final solution. We generate synthetic data with $d = 1000$ features and $n = 100$ samples. To ensure a fair comparison, all methods use the
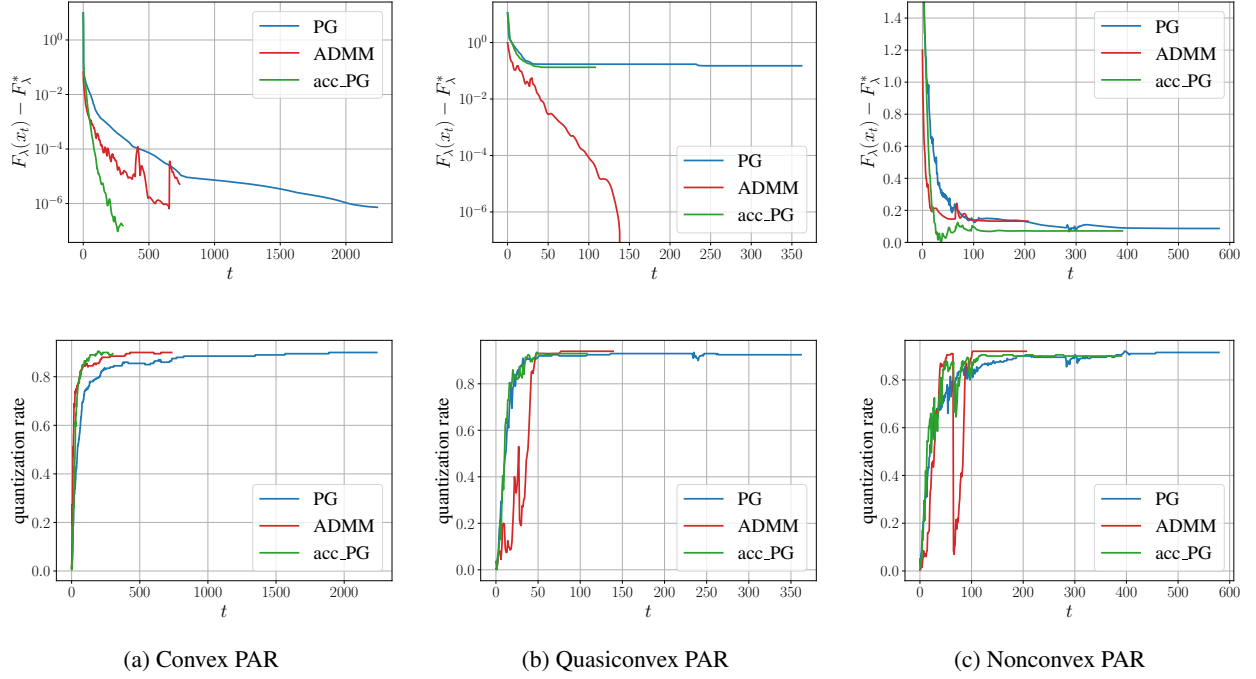
Figure 6: Comparison of optimization algorithms `PG`, `acc_PG`, and `ADMM` for linear regression with convex (left), quasiconvex (middle), and nonconvex (right) PARs. The problem dimension is 200 and the sample size is 20. All algorithms determine the step size using line search.
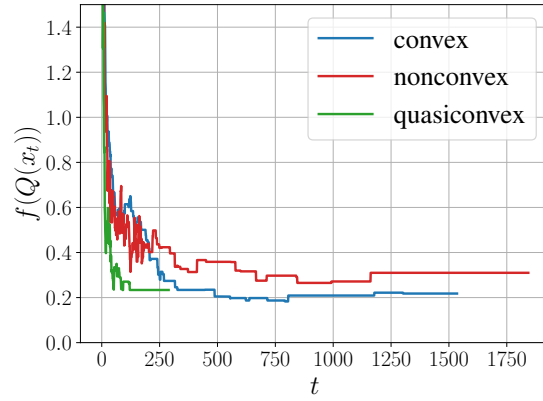


Figure 7: Comparison of three different PARs on a linear regression task. Here dimension is 1000 and the sample size is 100. To ensure a fair comparison, we use the same set of quantization values $\mathcal{Q}$ and evaluate the objective value of the fully quantized solutions. Specifically, at each iteration, the current solution $\boldsymbol{x}^t$ is projected onto $\mathcal{Q}$ to obtain $Q(\boldsymbol{x}^t)$, and we report the unregularized objective value $f(Q(\boldsymbol{x}^t))$. For each PAR, we report the best performance achieved over the regularization strengths $\{0.01, 0.015, 0.02, 0.05\}$.

(a) Ridge regularizer       (b) $\ell_1$-regularizer       (c) $\ell_{0.5}$-regularizer

(d) Ridge regularizer       (e) $\ell_1$-regularizer       (f) $\ell_{0.5}$-regularizer

Figure 8: Statistical performance of Ridge (left), $\ell_1$- (middle), and $\ell_{0.5}$-regularizers (right) and their PAR approximations on linear regression tasks. The quantization rate for each PAR is also shown.

same quantization set $\mathcal{Q}$. We adopt `ADMM` as the optimization algorithm, as it consistently performs slightly better than `PG` and `acc_PG` in this setting; however, the choice of solver does not significantly impact the observed trends. At each iteration, the current iterate $\boldsymbol{x}^t$ is projected onto $\mathcal{Q}$ to obtain a fully quantized solution $Q(\boldsymbol{x}^t)$, and we evaluate the unregularized objective $f(Q(\boldsymbol{x}^t))$ to measure the solution quality. For each regularizer, we report the best performance across regularization strengths $\{0.01, 0.015, 0.02, 0.05\}$.

As shown in Figure 7, the convex PAR outperforms both quasiconvex and nonconvex counterparts, with the quasiconvex variant performing slightly better than the nonconvex one. This result highlights the potential benefits of convexity in guiding the algorithm toward high-quality solutions.

## 5.3 Statistical guarantees

We assess the statistical accuracy of three classical regularizers: $\ell_2$-, $\ell_1$-, and $\ell_{0.5}$-regularizers, against their PAR approximations. For each method, we measure the Euclidean distance $\|\hat{\boldsymbol{x}} - \boldsymbol{x}^\star\|$ between the estimated parameter $\hat{\boldsymbol{x}}$ and the true parameter $\boldsymbol{x}^\star$. All experiments are conducted with data dimension $d = 200$, where the entries of the design matrix are drawn i.i.d. from $\mathcal{N}(0, 1)$. Responses are generated from either a linear or logistic model, with additive Gaussian noise of standard deviation $\sigma = 0.1$.

For the PAR approximations, we consider three quantization gaps, $q \in \{0.1, 0.05, 0.01\}$, and construct the corresponding PARs as described in Figure 4. All methods are evaluated on the same simulated datasets to ensure a fair comparison.

Figures 8 and 9 report the parameter estimation error $\|\hat{\boldsymbol{x}} - \boldsymbol{x}^\star\|$ for each regularizer and its PAR counterpart on linear and logistic regression tasks, respectively. In both settings, PAR approximations achieve nearly identical estimation accuracy to their original counterparts across all quantization gaps. Combined with the observed quantization rates, these results confirm that PARs retain statistical accuracy while enabling quantization.
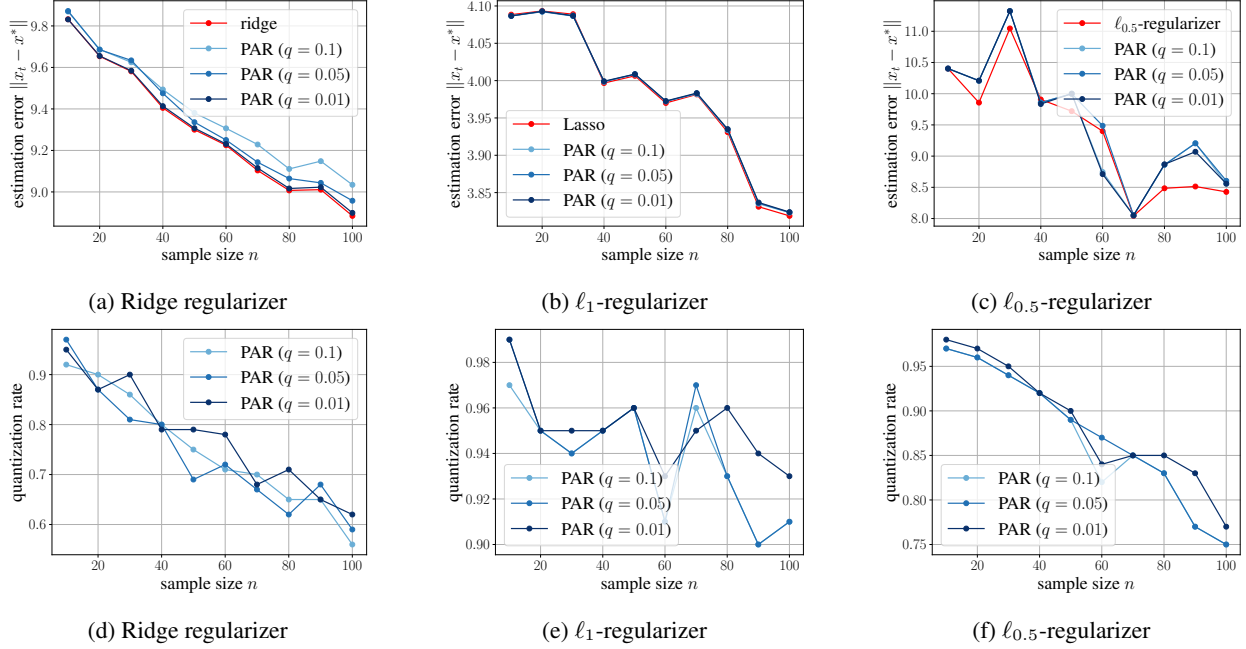
(a) Ridge regularizer     (b) $\ell_1$-regularizer     (c) $\ell_{0.5}$-regularizer

(d) Ridge regularizer     (e) $\ell_1$-regularizer     (f) $\ell_{0.5}$-regularizer

Figure 9: Statistical performance of Ridge (left), $\ell_1$- (middle), and $\ell_{0.5}$-regularizers (right) and their PAR approximations on logistic regression tasks. The quantization rate for each PAR is also shown.

# 6 Conclusion and future directions

This work introduces Piecewise-Affine Regularized Optimization (PARO), a principled and versatile framework for inducing quantization while preserving optimization and statistical guarantees. For generalized linear models, and more broadly, supervised learning models, we prove that under mild design assumptions, every critical point of the PARO objective is at least $(1 - n/d)$-quantized. Here $n$ is the sample size and $d$ is the parameter dimension, implying highly quantized solutions in the overparameterized regime where $d \gg n$. We also derive closed-form proximal mappings for three main PAR families: convex, quasiconvex, and nonconvex, and analyze the convergence of the proximal gradient method in the nonconvex setting. In the context of linear regression, we demonstrate that properly designed PARs can mimic the behavior of Ridge, Lasso, and general nonconvex penalties. They achieve comparable estimation and prediction performance while significantly reducing model storage. Extensive simulations validate our theoretical findings.

Below we point out a few interesting future directions.

- **Learnable quantization values.** Throughout this paper, we assume a fixed quantization set $\mathcal{Q}$. Some prior works have shown that jointly learning the quantization set can substantially improve model performance [EMB+19, PTT20]. Such approaches can also be interpreted through the lens of nonconvex piecewise-affine regularizers (PARs) [YZL+18]. An interesting direction is to investigate the quantization guarantees and convergence behavior under this learnable setting.

- **Stochastic gradient methods for PARO.** In this work, we focused on deterministic (full batch) proximal gradient methods and ADMM. However, in large-scale machine learning settings with high-dimensional data and massive datasets, stochastic gradient methods become necessary for scalability. Unfortunately, the standard proximal stochastic gradient method does not induce the desired

regularization effect (manifold identification) [Xia10]. This limitation has motivated the development of alternative stochastic optimization methods with strong manifold identification properties [Xia10, DYS$^+$21, JML$^+$25, QJM25]. Further investigation in this direction, by exploiting the particular structure of PARs, can be very impactful in practice.

- **Applications to combinatorial optimization.** PARs hold promise for broader applications in combinatorial optimization involving discrete variables. In integer programming, linear relaxation, where binary constraints $x \in \{0, 1\}$ are replaced with $x \in [0, 1]$, combined with (possibly stochastic) rounding, has proven both theoretically and empirically effective. This relaxation corresponds to a convex PAR defined as $\Psi(x) = 0$ for $x \in [0, 1]$ and $\Psi(x) = +\infty$ otherwise. Motivated by this, we conjecture that general PARs may induce meaningful discrete solutions for more complex combinatorial problems, including those with multi-level variables such as $x \in \{0, 1, \ldots, K\}$ for $K \geq 2$.

# Acknowledgment

# References

[Bac24]     Francis Bach. *Learning theory from first principles*. MIT press, 2024.

[BDLS07]   Jérôme Bolte, Aris Daniilidis, Adrian Lewis, and Masahiro Shiota. Clarke subgradients of stratifiable functions. *SIAM Journal on Optimization*, 18(2):556–572, 2007.

[Bec17]     Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.

[BPC$^+$11]   Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.

[BT09]      Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.

[BT10]      Amir Beck and Marc Teboulle. Gradient-based algorithms with applications to signal-recovery problems., 2010.

[BWL19]    Yu Bai, Yu-Xiang Wang, and Edo Liberty. ProxQuant: Quantized neural networks via proximal operators. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*, New Orleans, May 2019.

[CBD15]     Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. BinaryConnect: Training deep neural networks with binary weights during propagations. In *Advances in Neural Information Processing Systems*, volume 28, Montréal, Canada, December 2015.

[CDS98]   Scott Shaobing Chen, David L. Donoho, and Michael A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1998.

[Con17]   Laurent Condat. Discrete total variation: New definition and minimization. *SIAM Journal on Imaging Sciences*, 10(3):1258–1290, 2017.

[CZG⁺24]   Zeyu Cao, Cheng Zhang, Pedro Gimenes, Jianqiao Lu, Jianyi Cheng, and Yiren Zhao. Scaling laws for mixed quantization in large language models. *arXiv preprint arXiv:2410.06722*, 2024.

[CZL⁺25]   Mengzhao Chen, Chaoyi Zhang, Jing Liu, Yutao Zeng, Zeyue Xue, Zhiheng Liu, Yunshui Li, Jin Ma, Jie Huang, Xun Zhou, et al. Scaling law for quantization-aware training. *arXiv preprint arXiv:2505.14302*, 2025.

[DYS⁺21]   Tim Dockhorn, Yaoliang Yu, Eyyüb Sari, Mahdi Zolnouri, and Vahid Partovi Nia. Demystifying and generalizing binaryconnect. *Advances in Neural Information Processing Systems*, 34:13202–13216, 2021.

[EDA24]   Jasper Marijn Everink, Yiqiu Dong, and Martin Skovgaard Andersen. The geometry and well-posedness of sparse regularized linear regression. *arXiv preprint arXiv:2409.03461*, 2024.

[EMB⁺19]   Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. Learned step size quantization. In *International Conference on Learning Representations (ICLR)*, 2019.

[FL01]   Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.

[GN02]   Robert M. Gray and David L. Neuhoff. Quantization. *IEEE transactions on information theory*, 44(6):2325–2383, 2002.

[Har82]   Juris Hartmanis. Computers and intractability: a guide to the theory of np-completeness (michael r. garey and david s. johnson). *Siam Review*, 24(1):90, 1982.

[HKZ14]   Daniel Hsu, Sham M Kakade, and Tong Zhang. Random design analysis of ridge regression. *Foundations of Computational Mathematics*, 14:569–600, 2014.

[HL17]   Mingyi Hong and Zhi-Quan Luo. On the linear convergence of the alternating direction method of multipliers. *Mathematical Programming*, 162(1):165–199, 2017.

[HMD16]   Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and Huffman coding. In *Proceedings of the 4th International Conference on Learning Representations (ICLR)*, San Juan, Puerto Rico, 2016.

[HTW15]   Trevor Hastie, Robert Tibshirani, and Martin Wainwright. Statistical learning with sparsity. *Monographs on statistics and applied probability*, 143(143):8, 2015.

[JML⁺25]   Lisa Jin, Jianhao Ma, Zechun Liu, Andrey Gromov, Aaron Defazio, and Lin Xiao. Parq: Piecewise-affine regularized quantization. *arXiv preprint arXiv:2503.15748*, 2025.

[KPR16]   Vladimir Kolmogorov, Thomas Pock, and Michal Rolinek. Total variation on a tree. *SIAM Journal on Imaging Sciences*, 9(2):605–636, 2016.

[KS21]      Kenji Kawaguchi and Qingyun Sun. A recipe for global convergence guarantee in deep neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 8074–8082, 2021.

[LL15]      Huan Li and Zhouchen Lin. Accelerated proximal gradient methods for nonconvex programming. *Advances in neural information processing systems*, 28, 2015.

[LYLPN24]  Yiwei Lu, Yaoliang Yu, Xinlin Li, and Vahid Partovi Nia. Understanding neural network binarization with forward and backward proximal quantizers. *Advances in Neural Information Processing Systems*, 36, 2024.

[Mit15]     Boris Mityagin. The zero set of a real analytic function. *arXiv preprint arXiv:1512.07276*, 2015.

[Nes13]     Yu Nesterov. Gradient methods for minimizing composite functions. *Mathematical programming*, 140(1):125–161, 2013.

[NH17]      Quynh Nguyen and Matthias Hein. The loss surface of deep and wide neural networks. In *International conference on machine learning*, pages 2603–2612. PMLR, 2017.

[Opp99]     Alan V Oppenheim. *Discrete-time signal processing*. Pearson Education India, 1999.

[PB+14]     Neal Parikh, Stephen Boyd, et al. Proximal algorithms. *Foundations and trends® in Optimization*, 1(3):127–239, 2014.

[PS08]      John G Proakis and Masoud Salehi. *Digital communications*. McGraw-hill, 2008.

[PTT20]     Hadi Pouransari, Zhucheng Tu, and Oncel Tuzel. Least squares binary quantization of neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 698–699, 2020.

[QJM25]     Junwen Qiu, Li Jiang, and Andre Milzarek. A normal map-based proximal stochastic gradient method: Convergence and identification properties. arXiv:2305.05828, 2025.

[RWY10]     Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Restricted eigenvalue properties for correlated gaussian designs. *The Journal of Machine Learning Research*, 11:2241–2259, 2010.

[Sch23]     James Schmidt. Taylor learning. *arXiv preprint arXiv:2305.14606*, 2023.

[SCYE17]    Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, and Joel S Emer. Efficient processing of deep neural networks: A tutorial and survey. *Proceedings of the IEEE*, 105(12):2295–2329, 2017.

[SS19]      Lawrence V Snyder and Zuo-Jun Max Shen. *Fundamentals of supply chain theory*. John Wiley & Sons, 2019.

[ST22]      Ulrike Schneider and Patrick Tardivel. The geometry of uniqueness, sparsity and clustering in penalized estimation. *Journal of Machine Learning Research*, 23(331):1–36, 2022.

[Tib96]     Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.

[TSGS21]   Patrick JC Tardivel, Tomasz Skalski, Piotr Graczyk, and Ulrike Schneider. The geometry of model recovery by penalized and thresholded estimators. *HAL preprint hal-03262087*, 2021.

[Van10]   Lieven Vandenberghe. Fast proximal gradient methods. *EE236C course notes, Online, http://www. seas. ucla. edu/vandenbe C*, 236, 2010.

[WN99]   Laurence A Wolsey and George L Nemhauser. *Integer and combinatorial optimization*. John Wiley & Sons, 1999.

[WO10]   Sumio Watanabe and Manfred Opper. Asymptotic equivalence of bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of machine learning research*, 11(12), 2010.

[Wol20]   Laurence A Wolsey. *Integer programming*. John Wiley & Sons, 2020.

[WR22]   Stephen J. Wright and Benjamin Recht. *Optimization for Data Analysis*. Cambridge University Press, Cambridge, 2022.

[Xia10]   Lin Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 11(88):2543–2596, 2010.

[YH16]   Wei Hong Yang and Deren Han. Linear convergence of the alternating direction method of multipliers for a class of convex optimization problems. *SIAM journal on Numerical Analysis*, 54(2):625–640, 2016.

[YZL$^+$18]   Penghang Yin, Shuai Zhang, Jiancheng Lyu, Stanley Osher, Yingyong Qi, and Jack Xin. BinaryRelax: A relaxation approach for training deep neural networks with quantized weights. *SIAM Journal on Imaging Sciences*, 11(4):2205–2223, 2018. https://doi.org/10.1137/18M1166134.

[Zha10]   Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, pages 894–942, 2010.

[ZZ12]   Cun-Hui Zhang and Tong Zhang. A general theory of concave regularization for high-dimensional sparse estimation problems. *Statistical science*, 27(4):576–593, 2012.

# A   Proofs for Proximal Mappings

In this section, we present formal derivations of the explicit proximal mappings for the various PARs introduced in Section 3.1.

## A.1   Convex PAR

Recall that the proximal mapping of $\Psi$ is defined as

$$\mathbf{prox}_{\lambda\Psi}(x) = \arg\min_{z\in\mathbb{R}}\left\{\frac{1}{2}(z-x)^2 + \lambda\Psi(z)\right\}. \tag{69}$$

Since $\Psi(z)$ is an even function, i.e., $\Psi(z) = \Psi(-z)$, its proximal operator is an odd function. We can therefore derive the solution for $x \geq 0$, which implies the minimizer $z$ is also non-negative, and then extend

the result to all $x \in \mathbb{R}$ using the relation $\mathbf{prox}_{\lambda\Psi}(x) = \text{sign}(x)\mathbf{prox}_{\lambda\Psi}(|x|)$. For $x \geq 0$, the problem becomes

$$\mathbf{prox}_{\lambda\Psi}(x) = \arg\min_{z \geq 0} \left\{ \frac{1}{2}(z-x)^2 + \lambda\Psi(z) \right\}. \tag{70}$$

The first-order optimality condition for this convex optimization problem asserts that the minimizer $z$ must satisfy

$$0 \in z - x + \lambda\partial\Psi(z), \tag{71}$$

which can be rewritten as $x - z \in \lambda\partial\Psi(z)$. The subdifferential $\partial\Psi(z)$ for $z \geq 0$ is provided by

$$\partial\Psi(z) = \begin{cases} [-a_0, a_0] & \text{if } z = 0, \\ \{a_k\} & \text{if } z \in (q_k, q_{k+1}) \text{ for } k \in \{0, \ldots, m-1\}, \\ [a_{k-1}, a_k] & \text{if } z = q_k \text{ for } k \in \{1, \ldots, m\}. \end{cases} \tag{72}$$

We analyze the optimality condition (equation 71) for different cases.

**Case 1: Solution is zero ($z = 0$).** The optimality condition is $x - 0 \in \lambda\partial\Psi(0)$, which means $x \in \lambda[-a_0, a_0]$. Since we consider $x \geq 0$, this implies $0 \leq x \leq \lambda a_0$. Thus, for $-\lambda a_0 \leq x \leq \lambda a_0$, the minimizer is $z = 0$.

**Case 2: Solution is a non-zero quantization point ($z = q_k$ for $k \in \{1, \ldots, m\}$).** The optimality condition is $x - q_k \in \lambda\partial\Psi(q_k)$, which translates to $x - q_k \in \lambda[a_{k-1}, a_k]$. This is equivalent to

$$q_k + \lambda a_{k-1} \leq x \leq q_k + \lambda a_k.$$

Thus, if $|x|$ lies in $[q_k + \lambda a_{k-1}, q_k + \lambda a_k]$, the solution is $\mathbf{prox}_{\lambda\Psi}(x) = \text{sign}(x)q_k$.

**Case 3: Solution lies strictly between quantization points ($z \in (q_k, q_{k+1})$ for $k \in \{0, \ldots, m-1\}$).** Here, the subgradient is single-valued, $\partial\Psi(z) = \{a_k\}$. The optimality condition reduces to

$$x - z = \lambda a_k,$$

which gives the solution $z = x - \lambda a_k$. For this solution to be valid, it must lie in the assumed interval, $q_k < z < q_{k+1}$, which implies

$$q_k < x - \lambda a_k < q_{k+1}.$$

Rearranging for $x$ gives $q_k + \lambda a_k < x < q_{k+1} + \lambda a_k$. Thus, if $|x|$ lies in $(q_k + \lambda a_k, q_{k+1} + \lambda a_k)$, the solution is $\mathbf{prox}_{\lambda\Psi}(x) = x - \text{sign}(x)\lambda a_k$.

Combining the three cases above completes the proof.

## A.2 Quasiconvex PAR

Recall that the proximal mapping is given by

$$\mathbf{prox}_{\lambda\Psi}(x) = \arg\min_{z \in \mathbb{R}} \left\{ \Phi(z) := \frac{1}{2}(z-x)^2 + \lambda\Psi(z) \right\}. \tag{73}$$

Since $\Psi(x)$ is an even function, we can solve for $x \geq 0$ (which implies $z \geq 0$) and generalize using $\mathbf{prox}_{\lambda\Psi}(x) = \text{sign}(x)\mathbf{prox}_{\lambda\Psi}(|x|)$.

The first-order optimality condition for this problem is $x \in z + \lambda \partial \Psi(z)$. The Fréchet subdifferential $\partial \Psi(z)$ where $z \geq 0$ is given by[5]

$$\partial \Psi(z) = \begin{cases} [-1, 1] & \text{if } z = 0, \\ \{1\} & \text{if } z \in \left(kq, \frac{2k+1}{2}q\right), \\ \emptyset & \text{if } z = \frac{2k+1}{2}q, \\ \{0\} & \text{if } z \in \left(\frac{2k+1}{2}q, (k+1)q\right), \\ [0, 1] & \text{if } z = kq \text{ where } k \neq 0. \end{cases} \tag{74}$$

In what follows, we consider two situations: $\lambda \leq q$ and $\lambda > q$.

### A.2.1 Case I: $\lambda \leq q$

We consider the intersection between the two lines $y = x$ and $y = z + \lambda \partial \Psi(z)$. We further divide it into two cases.

**Case 1: $kq \leq x \leq kq + \lambda$.** In this case, there is only one critical point $z = kq$. Hence, we have $\mathbf{prox}_{\lambda\Psi}(x) = kq$.

**Case 2: $kq + \lambda \leq x \leq (k+1)q$.** In this case, there are two critical points: $z_1 = x - \lambda$ and $z_2 = x$. The global minimizer must be one of these candidates.

We proceed by comparing the objective function values

$$\Delta := \frac{1}{2}(x-x)^2 + \Psi(x) - \left(\frac{1}{2}(x-\lambda-x)^2 + \Psi(x-\lambda)\right) = \Psi(x) - \Psi(x-\lambda) - \frac{1}{2}\lambda^2. \tag{75}$$

- When $x \leq \frac{2k+1}{2}q$, we have

$$\Delta = \lambda^2 - \frac{1}{2}\lambda^2 = \frac{1}{2}\lambda^2 > 0. \tag{76}$$

  Hence, $\mathbf{prox}_{\lambda\Psi}(x) = x - \lambda$.

- When $\frac{2k+1}{2}q \leq x \leq \frac{2k+1}{2}q + \lambda$, we have

$$\Delta = \lambda \left(\frac{k+1}{2}q - \left(x - \lambda - \frac{k}{2}q\right)\right) - \frac{1}{2}\lambda^2. \tag{77}$$

  When $x \leq \frac{2k+1}{2}q + \frac{1}{2}\lambda$, we have $\Delta < 0$, which implies that $\mathbf{prox}_{\lambda\Psi}(x) = x - \lambda$. When $\frac{2k+1}{2}q + \frac{1}{2} \leq x \leq \frac{2k+1}{2}q + \lambda$, we have $\Delta \geq 0$, which implies $\mathbf{prox}_{\lambda\Psi}(x) = x$.

- When $\frac{2k+1}{2}q + \lambda \leq x \leq (k+1)q$, we always have $\Delta = -\frac{1}{2}\lambda^2 < 0$. Hence, $\mathbf{prox}_{\lambda\Psi}(x) = x$.

---

[5]We slightly abuse notation here by using the same symbol as for the Clarke subdifferential.

### A.2.2 Case II: $\lambda > q$

We analyze the minimizers of $\Phi(z)$ in equation 73.

First, note that when $z \in \left[\frac{(2k+1)q}{2}, (k+1)q\right]$, $\Psi(z)$ is constant and $\Phi(z)$ is a quadratic function with minimum at $z = x$. Thus, if $x \in \left[\frac{(2k^\star+1)q}{2}, (k^\star + 1)q\right]$ for some $k^\star \in \mathbb{Z}_+$, then $x$ is a candidate minimizer. In other cases, since proximal mapping is nonexpensive and $\Phi(z)$ is a decreasing function in the range $[0, x]$, the other candidates are $\{(k+1)q\}_{k<k^\star}$.

When $z \in \left[kq, \frac{(2k+1)q}{2}\right]$, $\Psi(z)$ is affine and $\Phi(z) = \frac{1}{2}(z - x + \lambda)^2 + C$ where $C$ is a constant. This quadratic achieves its minimum at $z = x - \lambda < kq$, which lies outside the interval, so the only candidates here are again grid points $\{kq\}_{k<k^\star+1}$.

Hence, all candidate minimizers belong to the set $\{(k+1)q\}_{k<k^\star+1} \cup \{x\}$. To identify the minimum, we first compare values at the grid points. For any $k$, we have

$$\Phi(kq) = \frac{1}{2}(x - kq)^2 + \frac{\lambda k q}{2} = \frac{q^2}{2}k^2 + \frac{\lambda q - 2xq}{2}k + \frac{1}{2}x^2. \tag{78}$$

This quadratic in $k$ is minimized when $k = \left\lfloor \frac{x - \frac{1}{2}\lambda}{q} \right\rfloor$. Next, we compare it with the candidate $x$. Note that

$$
\begin{aligned}
\Phi(x) - \Phi\left(\left\lfloor \frac{x - \frac{1}{2}\lambda}{q} \right\rfloor q\right) &= \frac{\lambda}{2} \cdot \left((k^\star + 1)q - \left\lfloor \frac{x - \frac{1}{2}\lambda}{q} \right\rfloor q\right) - \frac{1}{2}\left(x - \left\lfloor \frac{x - \frac{1}{2}\lambda}{q} \right\rfloor q\right)^2 \\
&\geq \frac{1}{2}\left(x - \left\lfloor \frac{x - \frac{1}{2}\lambda}{q} \right\rfloor q\right)\left(\lambda - \left(x - \left\lfloor \frac{x - \frac{1}{2}\lambda}{q} \right\rfloor q\right)\right) \\
&\geq 0.
\end{aligned}
\tag{79}
$$

where the last inequality uses $\lambda > q$. This confirms that the unique minimizer is $\left\lfloor \frac{x - \frac{1}{2}\lambda}{q} \right\rfloor q$. This completes the proof.

### A.3 Nonconvex PAR

First, we consider that $x \in \left[q_k, \frac{q_k + q_{k+1}}{2}\right]$ for some $k$. It is obvious that $\Psi_{\lambda\Psi}(x) \in \left[q_k, \frac{q_k + q_{k+1}}{2}\right]$ for any $\lambda \geq 0$. Therefore, $\Psi_{\lambda\Psi}(x)$ is indeed the minimizer of a quadratic function:

$$\mathbf{prox}_{\lambda\Psi}(x) = \underset{z \in \left[q_k, \frac{q_k + q_{k+1}}{2}\right]}{\arg\min} \left\{\frac{1}{2}(z + \lambda - x)^2\right\} = \mathrm{clip}\left(x - \lambda, q_k, \frac{q_k + q_{k+1}}{2}\right). \tag{80}$$

Similarly, if $x \in \left[\frac{q_k + q_{k+1}}{2}, q_{k+1}\right]$ for some $k$, its proximal mapping can be derived by

$$\mathbf{prox}_{\lambda\Psi}(x) = \underset{z \in \left[\frac{q_k + q_{k+1}}{2}, q_{k+1}\right]}{\arg\min} \left\{\frac{1}{2}(z - \lambda - x)^2\right\} = \mathrm{clip}\left(x + \lambda, \frac{q_k + q_{k+1}}{2}, q_{k+1}\right). \tag{81}$$

This completes the proof.