

663 The Appendix is organized as follows:

- 664 1. Section A provides utility results useful in subsequent proofs.
- 665 2. Section B provides the proof of the lower bound described in Section 3.1
- 666 3. Section C proves helper lemmas for the results in Section 4
- 667 4. Section D proves Theorems 1, 2 and 3.
- 668 5. Section E proves the boosting result (Lemma 3) and end to end analysis of Algorithm 1
- 669 followed by the boosting algorithm 2.
- 670 6. Section F provides details of our experimental setup
- 671 7. Section G provides more related work.

672 A Utility Results

673 **Lemma A.1.** For any $\mathbf{y} \in \mathbb{R}^d$, let $\mathbf{Q}(\mathbf{y}, \mathcal{Q})$ denote the quantization of \mathbf{y} using Eq (5), where each
 674 coordinate, y_i , is quantized independently as $\mathbf{Q}(y_i, \mathcal{Q})$. Then, for any $x \in \mathbb{R}$ and $\mathbf{v} \in \mathbb{R}^d$,

- 675 1. $\mathbb{E}[\mathbf{Q}(x, \mathcal{Q})|x] = x$.
- 676 2. $|\mathbf{Q}(x, \mathcal{Q}) - x| \leq u - l$.
- 677 3. $\text{Var}[\mathbf{Q}(x, \mathcal{Q})|x] \leq \frac{(u-l)^2}{4}$.
- 678 4. $\text{Cov}[\mathbf{Q}(\mathbf{v}, \mathcal{Q})|\mathbf{v}]$ is a diagonal matrix.

679 *Proof.* Throughout the proof, we condition on the fixed x and treat all randomness as coming from
 680 the independent choices made by the quantizer.

681 (i) *Unbiasedness.* We have

$$\mathbb{E}[\mathbf{Q}(x, \delta) | x] = p_i(x)u + (1 - p_i(x))l = x.$$

682 (ii) *Boundedness.* By definition, after rounding, we always round any $x \in [u, l]$ to either u or l .
 683 Therefore, $|\mathbf{Q}(x, \mathcal{Q}) - x| \leq u - l$.

684 (iii) *Variance bound.* Using the variance of a Bernoulli random variable, we have,

$$\text{Var}[\mathbf{Q}(x, \mathcal{Q}) | x] = p_i(x)(1 - p_i(x))(u - l)^2 \leq \frac{(u-l)^2}{4},$$

685 since $t(1 - t) \leq \frac{1}{4}$ for all $t \in [0, 1]$.

686 (iv) *Covariance bound.* Since the elements of \mathbf{v} are quantized independently, this follows.

687 □

688 **Lemma A.2** (Choice of learning rate). Let $\eta := \frac{\alpha \log(n)}{b(\lambda_1 - \lambda_2)}$. Then, under Assumption 1 for $\theta \in (0, 1)$,
 689 η satisfies

$$b(\eta^2 \mathcal{M}^2 + \kappa^2) \leq \frac{0.008}{\log(d/\theta)}, \text{ and } \eta \in (0, 1)$$

690 for $\alpha > 1$, $b \geq 250\alpha^2 \log^2(n) \log(\frac{d}{\theta}) / (\lambda_1 - \lambda_2)^2$, and $\kappa^2 b \leq 0.004 / \log(\frac{d}{\theta})$.

691 *Proof.* For Lemma A.8 we require,

$$4b(\eta^2 \mathcal{M}^2 + \kappa^2)(1 + 2 \log(d)) \leq 1 \tag{A.13}$$

692 For Theorem A.4 we require,

$$4e^2 b(\eta^2 \mathcal{M}^2 + \kappa^2) \log\left(\frac{d}{\theta}\right) \leq \frac{1}{4} \tag{A.14}$$

where $\theta \in (0, 1)$ represents the failure probability. It is not hard to see that (A.14) implies (A.13). Therefore it suffices to ensure

$$b(\eta^2 \mathcal{M}^2 + \kappa^2) \log \left(\frac{d}{\theta} \right) \leq 0.008$$

Setting each term smaller than 0.004, it suffices to have

$$b \geq \frac{250\alpha^2 \log^2(n) \log \left(\frac{d}{\theta} \right)}{(\lambda_1 - \lambda_2)^2}, \quad \kappa^2 b \leq \frac{0.004}{\log \left(\frac{d}{\theta} \right)}$$

which completes the proof for the first condition.

The second condition on η follows by setting $\eta \leq 1$ and solving for b . This yields

$$b \geq \max \left\{ 250\alpha^2 \log^2(n) \log \left(\frac{d}{\theta} \right) / (\lambda_1 - \lambda_2)^2, \alpha \log(n) / (\lambda_1 - \lambda_2) \right\}$$

Since $\alpha > 1$, the first term is larger than the second one, which completes the proof. \square

Lemma A.3. Let \mathbf{w} and $\boldsymbol{\xi}$ be vectors in \mathbb{R}^d such that $\|\mathbf{w}\| = 1$ and $\mathbf{w} + \boldsymbol{\xi} \neq 0$. Then,

$$\sin^2(\mathbf{w}, \mathbf{w} + \boldsymbol{\xi}) \leq \left(\frac{\|\boldsymbol{\xi}\|}{\|\mathbf{w} + \boldsymbol{\xi}\|} \right)^2.$$

Proof.

$$\begin{aligned} \sin^2(\mathbf{w}, \mathbf{w} + \boldsymbol{\xi}) &= 1 - \left(\frac{\mathbf{w}^\top (\mathbf{w} + \boldsymbol{\xi})}{\|\mathbf{w} + \boldsymbol{\xi}\|} \right)^2 = \frac{(\mathbf{w} + \boldsymbol{\xi})^\top (\mathbf{w} + \boldsymbol{\xi}) - (1 + \mathbf{w}^\top \boldsymbol{\xi})^2}{\|\mathbf{w} + \boldsymbol{\xi}\|^2} \\ &= \frac{\boldsymbol{\xi}^\top \boldsymbol{\xi} - (\mathbf{w}^\top \boldsymbol{\xi})^2}{\|\mathbf{w} + \boldsymbol{\xi}\|^2} \leq \left(\frac{\|\boldsymbol{\xi}\|}{\|\mathbf{w} + \boldsymbol{\xi}\|} \right)^2. \end{aligned}$$

\square

Lemma A.4. Let \mathbf{x} and \mathbf{y} be unit vectors in \mathbb{R}^d . Then,

$$\frac{1}{2} \min(\|\mathbf{x} - \mathbf{y}\|^2, \|\mathbf{x} + \mathbf{y}\|^2) \leq \sin^2(\mathbf{x}, \mathbf{y}) \leq \min(\|\mathbf{x} - \mathbf{y}\|^2, \|\mathbf{x} + \mathbf{y}\|^2).$$

Proof. We express $\sin^2(\mathbf{x}, \mathbf{y})$ in terms of $\|\mathbf{x} - \mathbf{y}\|$ and $\|\mathbf{x} + \mathbf{y}\|$. Since $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\mathbf{x}^\top \mathbf{y} = 2 - 2\cos(\mathbf{x}, \mathbf{y})$ and $\|\mathbf{x} + \mathbf{y}\|^2 = 2 + 2\cos(\mathbf{x}, \mathbf{y})$,

$$\|\mathbf{x} - \mathbf{y}\|^2 + \|\mathbf{x} + \mathbf{y}\|^2 = 4 \text{ and } \sin^2(\mathbf{x}, \mathbf{y}) = 1 - \cos^2(\mathbf{x}, \mathbf{y}) = \frac{1}{4} \|\mathbf{x} - \mathbf{y}\|^2 \|\mathbf{x} + \mathbf{y}\|^2.$$

The upper bound on $\sin^2(\mathbf{x}, \mathbf{y})$ follows immediately from the above equations. For the lower bound, note that at least one of $\|\mathbf{x} - \mathbf{y}\|^2$ and $\|\mathbf{x} + \mathbf{y}\|^2$ is at least 2. If $\|\mathbf{x} + \mathbf{y}\|^2 \geq 2$, then $\sin^2(\mathbf{x}, \mathbf{y}) \geq \|\mathbf{x} - \mathbf{y}\|^2 / 2$. Otherwise, $\sin^2(\mathbf{x}, \mathbf{y}) \geq \|\mathbf{x} + \mathbf{y}\|^2 / 2$. \square

Lemma A.5. Let \mathbf{x}, \mathbf{y} , and \mathbf{z} be non-zero vectors in \mathbb{R}^d . Then,

$$\sin^2(\mathbf{x}, \mathbf{z}) \leq 2\sin^2(\mathbf{x}, \mathbf{y}) + 2\sin^2(\mathbf{y}, \mathbf{z}).$$

B Lower Bounds

Proof of Lemma 1

Proof. Let $\mathbf{v}_1 \in \mathbb{R}^d$ be the unit vector with $\mathbf{v}_1(i) = \delta/3$ for $i \in [d-1]$ and $\mathbf{v}_1(d) = \sqrt{1 - \frac{(d-1)\delta^2}{9}}$.

Consider any a vector $\mathbf{w} \in \mathcal{V}_L$, and let $\tilde{\mathbf{w}} = \mathbf{w}/\|\mathbf{w}\|$. Since $\mathbf{w} \in \mathcal{V}_L$, $\mathbf{w}(i) = 0$ or $|\mathbf{w}(i)| \geq \delta/2$. In particular, $|\mathbf{v}_1(i) - \mathbf{w}(i)| \geq \delta/6$ and $|\mathbf{v}_1(i) + \mathbf{w}(i)| \geq \delta/6$ for all $i \in [d-1]$. It follows that

$$\|\mathbf{v}_1 - \mathbf{w}\|^2 \geq \sum_{i=1}^{d-1} (\mathbf{v}_1(i) - \mathbf{w}(i))^2 \geq (d-1) \left(\frac{\delta}{6} \right)^2 = \frac{\delta^2(d-1)}{36}$$

and $\|\mathbf{v}_1 + \mathbf{w}\|^2 \geq \frac{\delta^2(d-1)}{36}$ similarly. The Lemma follows from A.4 \square

\square

716 **Proof of Lemma 2**

717 *Proof.* It suffices to construct two unit vectors \mathbf{v}_1 and \mathbf{v}_2 such that $\inf_{\mathbf{w} \in \mathcal{V}_{NL}} \sin^2(\mathbf{w}, \mathbf{v}_1) = \Omega(\zeta^2)$
 718 and $\inf_{\mathbf{w} \in \mathcal{V}_{NL}} \sin^2(\mathbf{w}, \mathbf{v}_2) = \Omega(\delta_0^2 d)$.

719 Let \mathbf{v}_1 be the vector in \mathbb{R}^d with coordinates

$$\mathbf{v}_1(1) = \frac{1}{\sqrt{1 + (1 + \zeta/2)^2}}, \quad \mathbf{v}_1(2) = \frac{1 + \zeta/2}{\sqrt{1 + (1 + \zeta/2)^2}}, \quad \mathbf{v}_1(i) = 0 \quad \forall i \geq 3.$$

720 For the sake of contradiction, suppose there exists $\mathbf{w}_1 \in \mathcal{V}_{NL}$ such that $\sin^2(\mathbf{w}_1, \mathbf{v}_1) \leq \zeta^2/100$.

721 Let $\tilde{\mathbf{w}}_1 := \mathbf{w}_1 / \|\mathbf{w}_1\|$. By Lemma A.4

$$\min(\|\mathbf{v}_1 - \tilde{\mathbf{w}}_1\|_2^2, \|\mathbf{v}_1 + \tilde{\mathbf{w}}_1\|_2^2) \leq 2 \sin^2(\mathbf{v}_1, \mathbf{w}_1) \leq \frac{\zeta^2}{50}.$$

722 Flipping the sign of \mathbf{w}_1 if necessary, we may assume $\|\mathbf{v}_1 - \tilde{\mathbf{w}}_1\|_2^2 \leq \zeta^2/50$. So,

$$|\mathbf{v}_1(i) - \tilde{\mathbf{w}}_1(i)| \leq \zeta/7 \quad \forall i \in [d]. \quad (\text{A.15})$$

723 The bound $\zeta \leq 0.1$ ensures $\mathbf{v}_1(1) \geq 20/29$ and $\mathbf{v}_1(2) - \mathbf{v}_1(1) \geq \zeta/3$, which also implies $\tilde{\mathbf{w}}_1(2) -$
 724 $\tilde{\mathbf{w}}_1(1) \geq \zeta/21 > 0$. It follows that

$$\begin{aligned} \frac{\mathbf{w}_1(2) + \delta_0/\zeta}{\mathbf{w}_1(1) + \delta_0/\zeta} &= \frac{\tilde{\mathbf{w}}_1(2) + \delta_0/\zeta \cdot 1/\|\mathbf{w}_1\|}{\tilde{\mathbf{w}}_1(1) + \delta_0/\zeta \cdot 1/\|\mathbf{w}_1\|} \leq \frac{\mathbf{v}_1(2) + \zeta/7 + \delta_0/2\zeta}{\mathbf{v}_1(1) - \zeta/7 + \delta_0/2\zeta} \\ &= 1 + \frac{\zeta}{2} + \frac{\delta_0/2\zeta + \zeta/7 - (1 + \zeta/2)(\delta_0/2\zeta - \zeta/7)}{\mathbf{v}_1(1) + \delta_0/2\zeta - \zeta/7} \\ &= 1 + \frac{\zeta}{2} + \frac{2\zeta/7 + \zeta^2/14 - \delta_0/4}{\mathbf{v}_1(1) - \zeta/7 + \delta_0/2\zeta} \\ &\leq 1 + \frac{\zeta}{2} + \frac{2\zeta/7}{2/3} < 1 + \zeta, \end{aligned}$$

725 and

$$\begin{aligned} \frac{\mathbf{w}_1(2) + \delta_0/\zeta}{\mathbf{w}_1(1) + \delta_0/\zeta} &= \frac{\tilde{\mathbf{w}}_1(2) + \delta_0/\zeta \cdot 1/\|\mathbf{w}_1\|}{\tilde{\mathbf{w}}_1(1) + \delta_0/\zeta \cdot 1/\|\mathbf{w}_1\|} \geq \frac{\mathbf{v}_1(2) - \zeta/7 + 2\delta_0/\zeta}{\mathbf{v}_1(1) + \zeta/7 + 2\delta_0/\zeta} \\ &= 1 + \frac{\zeta}{2} + \frac{2\delta_0/\zeta - \zeta/7 - (1 + \zeta/2)(2\delta_0/\zeta + \zeta/7)}{\mathbf{v}_1(1) + \zeta/7 + 2\delta_0/\zeta} \\ &= 1 + \frac{\zeta}{2} - \frac{2\zeta/7 + \zeta^2/14 + \delta_0}{\mathbf{v}_1(1) + \zeta/7 + 2\delta_0/\zeta} \\ &> 1 + \frac{\zeta}{2} - \frac{\zeta \mathbf{v}_1(1)/2 + \zeta^2/14 + \delta_0}{\mathbf{v}_1(1) + \zeta/7 + 2\delta_0/\zeta} = 1. \end{aligned}$$

726 Under the logarithmic quantization scheme, it can be inductively shown that $q_k + \delta_0/\zeta = (\delta_0/\zeta) \cdot (1 +$
 727 $\zeta)^k$ for all non-negative integers k such that $q_k \in \mathcal{Q}_{NL}$. In particular, $\frac{\mathbf{w}_1(2) + \delta_0/\zeta}{\mathbf{w}_1(1) + \delta_0/\zeta}$ is a non-negative
 728 integer power of $1 + \zeta$, contradicting

$$1 < \frac{\mathbf{w}_1(2) + \delta_0/\zeta}{\mathbf{w}_1(1) + \delta_0/\zeta} < 1 + \zeta.$$

729 Therefore, $\inf_{\mathbf{w}_1 \in \mathcal{V}_{NL}} \sin^2(\mathbf{w}_1, \mathbf{v}_1) \geq \zeta^2/100$.

730 For the other bound, let $\mathbf{v}_2(i) = \delta_0/3$ for $i \in d-1$ and $\mathbf{v}_2(d) = \sqrt{1 - (d-1)\delta_0^2/9}$. Any $\mathbf{w}_2 \in \mathcal{V}_{NL}$
 731 satisfies $\mathbf{w}_2(i) = 0$ or $|\mathbf{w}_2(i)| \geq \delta_0$ for all $i \in [d]$. Since $\|\mathbf{w}_2\| \in [1/2, 2]$, the normalized vector
 732 $\tilde{\mathbf{w}}_2 = \mathbf{w}_2/\|\mathbf{w}_2\|$ satisfies $|\tilde{\mathbf{w}}_2(i)| = 0$ or $|\tilde{\mathbf{w}}_2(i)| \geq \delta_0/2$ for all $i \in [d]$.

733 In particular $|\mathbf{v}_2(i) - \tilde{\mathbf{w}}_2(i)| \geq \delta_0/6$ and $|\mathbf{v}_2(i) + \tilde{\mathbf{w}}_2(i)| \geq \delta_0/6$ for all $i \in [d]$. By Lemma A.4

$$\sin^2(\mathbf{w}_2, \mathbf{v}_2) \geq \frac{1}{2} \min(\|\mathbf{w}_2 - \mathbf{v}_2\|^2, \|\mathbf{w}_2 + \mathbf{v}_2\|^2) \geq \frac{\delta_0^2(d-1)}{72}.$$

734

□

C Proof of Results in Section 4

For ease of exposition, all results in this section are stated with a generic number of datapoints n . We apply these results with different choices of n (e.g. number of batches b) for proving the main theorems (Theorem 1, 2, 3). Consider Oja's Algorithm applied to the matrices $\mathbf{A}_i \in \mathbb{R}_{d \times d}$, such that $\mathbf{A}_i = \eta \mathbf{D}_i + \boldsymbol{\Xi}_i$ where \mathbf{D}_i are independent with $\mathbb{E}[\mathbf{D}_i] = \boldsymbol{\Sigma}$. Let \mathcal{S}_i be the set of all random vectors $\boldsymbol{\xi}$ resulting from the quantizations in the first i iterations of the algorithm, and let \mathcal{F}_{i-} denote the σ -field generated by $\mathbf{D}_1, \dots, \mathbf{D}_i$ and \mathcal{S}_{i-1} , and denote $\mathbb{E}_i[\cdot] := \mathbb{E}[\cdot | \mathcal{F}_{i-}]$. We assume the noise term $\boldsymbol{\Xi}_i$ is conditionally unbiased, i.e., $\mathbb{E}_i[\boldsymbol{\Xi}_i] = \mathbf{0}_{d \times d}$.

$$\mathcal{F}_{i-} := \sigma(\{\mathbf{D}_1, \dots, \mathbf{D}_i, \mathcal{S}_{i-1}\}), \quad \mathcal{F}_i := \sigma(\{\mathbf{D}_1, \dots, \mathbf{D}_i, \mathcal{S}_i\}).$$

Recall the update rule

$$\mathbf{u}_i = (\mathbf{I} + \mathbf{A}_i) \mathbf{w}_{i-1}; \quad \mathbf{w}_i = \frac{\prod_{t=i}^1 (\mathbf{I} + \mathbf{A}_t) \mathbf{u}_0}{\|\prod_{t=i}^1 (\mathbf{I} + \mathbf{A}_t) \mathbf{u}_0\|}. \quad (\text{A.16})$$

We bound the numerator and denominator in (A.16) separately.

For the numerator, we will show that $\|\prod_{t=n}^1 (\mathbf{I} + \mathbf{A}_t) - (\mathbf{I} + \eta \boldsymbol{\Sigma})^n\|$ is small. Let $\mathbf{Y}_i = \mathbf{I} + \mathbf{A}_i$ for $i \in [n]$, and let $\{\mathbf{Z}_i\}_{0 \leq i \leq n}$ be defined as

$$\mathbf{Z}_i := \mathbf{Y}_i \mathbf{Z}_{i-1}, \quad \mathbf{Z}_0 := \mathbf{I}. \quad (\text{A.17})$$

Note that \mathbf{Z}_{i-1} is measurable w.r.t \mathcal{F}_{i-} .

We are now ready to state our first result. Note that

$$\mathbf{Z}_n = \prod_{i=n}^1 (\mathbf{I} + \mathbf{A}_i).$$

where $\mathbf{A}_i = \eta \mathbf{D}_i + \boldsymbol{\Xi}_i$ and \mathbf{D}_i are independent $d \times d$ random matrices with mean $\boldsymbol{\Sigma}$.

C.1 Proof of Proposition 1

Proposition A.1. [Proposition 1 in main paper] Under Assumption 1 for $\eta \in (0, 1)$, we have:

$$\begin{aligned} \|\mathbf{Z}_n\|_{p,q}^2 &\leq (\alpha + C_p \gamma)^n \|\mathbf{Z}_0\|_{p,q}^2, \text{ and} \\ \|\mathbf{Z}_n - (\mathbf{I} + \eta \boldsymbol{\Sigma})^n\|_{p,q}^2 &\leq \alpha^n (\exp(C_p n \gamma) - 1) \|\mathbf{Z}_0\|_{p,q}^2 \end{aligned}$$

for $\mathbf{Z}_0 = \mathbf{I}$, $\alpha := (1 + \eta \lambda_1)^2$, $\gamma := 2(\eta^2 \mathcal{M}^2 + \kappa^2)$ and $C_p := p - 1$.

Proof. Note that

$$\mathbb{E}[\mathbf{Y}_i | \mathcal{F}_{i-1}] = \mathbf{I} + \eta \boldsymbol{\Sigma} + \mathbb{E}[\boldsymbol{\Xi}_i | \mathcal{F}_{i-1}] = \mathbf{I} + \eta \boldsymbol{\Sigma}$$

Note that $m_i = 1 + \eta \lambda_1$ and

$$\|\mathbf{Y}_i - \mathbb{E}[\mathbf{Y}_i | \mathcal{F}_{i-1}]\| = \|\eta(\mathbf{D}_i - \boldsymbol{\Sigma}) + \boldsymbol{\Xi}_i\| \leq \eta \mathcal{M} + \kappa$$

The last line uses Eq 9. Thus $\sigma_i = \frac{\eta \mathcal{M} + \kappa}{1 + \eta \lambda_1}$. Note that $\nu \leq 2(\eta^2 \mathcal{M}^2 + \kappa^2)$. The same argument as in Theorem 7.4 in [HNW20] gives the bound. \square

Lemma A.6. Under Assumption 1 and η set according to Lemma A.2 with $b = n$,

$$\mathbb{P}(\|\mathbf{Z}_n - (\mathbf{I} + \eta \boldsymbol{\Sigma})^n\| \geq t(1 + \eta \lambda_1)^n) \leq \max(d, e) \exp\left(-\frac{t^2}{2e^2 n \gamma}\right) \quad \forall t \leq e.$$

where $\gamma := 2(\eta^2 \mathcal{M}^2 + \kappa^2)$ and $e = 2.718 \dots$ is the Napier's constant.

Proof. By Proposition A.1,

$$\begin{aligned} \mathbb{P}(\|\mathbf{Z}_n - (\mathbf{I} + \eta \boldsymbol{\Sigma})^n\| \geq t(1 + \eta \lambda_1)^n) &\leq \frac{\mathbb{E}[\|\mathbf{Z}_n - (\mathbf{I} + \eta \boldsymbol{\Sigma})^n\|^p]}{t^p (1 + \eta \lambda_1)^p} \leq \frac{\|\mathbf{Z}_n - (\mathbf{I} + \eta \boldsymbol{\Sigma})^n\|_{p,p}^p}{t^p (1 + \eta \lambda_1)^p} \\ &\leq \frac{\alpha^{\frac{p}{2}} (\exp(C_p n \gamma) - 1)^{p/2} d}{t^p (1 + \eta \lambda_1)^p} \leq d (t^{-2} (\exp(C_p n \gamma) - 1))^{p/2}. \end{aligned}$$

760 If $\frac{t^2}{e^2 n \gamma} < 2$, then $e \cdot \exp\left(-\frac{t^2}{2e^2 n \gamma}\right) \geq 1$. Otherwise, let $p := \frac{t^2}{e^2 n \gamma} \geq 2$. Since $t \leq e$, $C_p n \gamma \leq$
761 $p n \gamma \leq \frac{t^2}{e^2} \leq 1$. Therefore, $\exp(C_p n \gamma) - 1 \leq e C_p n \gamma \leq \frac{t^2}{e}$. Therefore,

$$\mathbb{P}(\|\mathbf{Z}_n - (\mathbf{I} + \eta \boldsymbol{\Sigma})^n\| \geq t(1 + \eta \lambda_1)^n) \leq d \cdot (t^{-2} \cdot (\exp(C_p n \gamma) - 1))^{p/2} \leq d \exp\left(-\frac{t^2}{2e^2 n \gamma}\right).$$

762

□

763 **Lemma A.7.** Under Assumption [I](#) and η set according to Lemma [A.2](#) with $b = n$,

$$\mathbb{E}[\|\mathbf{Z}_n\|^2] \leq \exp\left(2\sqrt{2n\gamma \max\{2n\gamma, \log(d)\}}\right) (1 + \eta \lambda_1)^{2n},$$

764 where $\gamma = 2(\eta^2 \mathcal{M}^2 + \kappa^2)$. Moreover, if $2n\gamma(1 + 2\log(d)) \leq 1$, then

$$\mathbb{E}[\|\mathbf{Z}_n - \mathbb{E}[\mathbf{Z}_n]\|^2] \leq 2e^2 n \gamma (1 + 2\log(d)) (1 + \eta \lambda_1)^{2n}.$$

765 *Proof.* Using Proposition [A.1](#) $\alpha := (1 + \eta \lambda_1)^2$, and $\gamma := 2(\eta^2 \mathcal{M}^2 + \kappa^2)$,

$$\mathbb{E}[\|\mathbf{Z}_n\|^2] \leq \|\mathbf{Z}_n\|_{p,2}^2 \leq (\alpha + C_p \gamma)^n \|\mathbf{Z}_0\|_{p,2}^2 \leq (1 + \eta \lambda_1)^{2n} \exp(C_p n \gamma) \|\mathbf{Z}_0\|_{p,2}^2.$$

766 Set $p := \max\left(2, \sqrt{\frac{2\log d}{n\gamma}}\right)$. Then, $\|\mathbf{Z}_0\|_{p,2} = d^{\frac{1}{p}} \leq \exp\left(\frac{pn\gamma}{2}\right)$ and

$$\mathbb{E}[\|\mathbf{Z}_n\|^2] \leq (1 + \eta \lambda_1)^{2n} \exp(2pn\gamma) = \exp\left(2\sqrt{2n\gamma \max\{2n\gamma, \log(d)\}}\right) (1 + \eta \lambda_1)^{2n}.$$

767 For the second result, set $p := 2(1 + \log(d))$. Then, $C_p n \gamma \leq 1$ and $\|\mathbf{Z}_0\|_p \leq \sqrt{e}$. By Proposi-
768 tion [A.1](#),

$$\begin{aligned} \mathbb{E}[\|\mathbf{Z}_n - \mathbb{E}[\mathbf{Z}_n]\|^2] &\leq \|\mathbf{Z}_n - \mathbb{E}[\mathbf{Z}_n]\|_{p,2}^2 \leq (\exp(C_p n \gamma) - 1) (1 + \eta \lambda_1)^n \|\mathbf{Z}_0\|_p^2 \\ &\leq e^2 C_p n \gamma (1 + \eta \lambda_1)^n \\ &\leq 2e^2 n \gamma (1 + \log(d)) (1 + \eta \lambda_1)^n. \end{aligned}$$

769

□

770 C.2 Proof of Lemma [4](#)

771 **Lemma A.8** (Lemma [4](#) in main paper). Let Assumption [I](#) hold and η be set according to Lemma [A.2](#)
772 with $b = n$. Define $\gamma := 2(\eta^2 \mathcal{M}^2 + \kappa^2)$. If $2n\gamma(1 + 2\log(d)) \leq 1$, then

$$\mathbb{E}\left[\text{Tr}\left(\mathbf{V}_\perp^\top \mathbf{Z}_n \mathbf{Z}_n^\top \mathbf{V}_\perp\right)\right] \leq \exp(2\eta n \lambda_1 + \eta^2 n (\mathcal{V}_0 + \lambda_1^2)) \left[\frac{d}{\exp(2\eta n (\lambda_1 - \lambda_2))} + \frac{5(\eta^2 \mathcal{V}_0 + \kappa_1)}{\eta(\lambda_1 - \lambda_2)} \right].$$

773 *Proof.* Let $\beta_i := \mathbb{E}\left[\text{Tr}\left(\mathbf{V}_\perp^\top \mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{V}_\perp\right)\right]$ for all $0 \leq i \leq n$. Then, for any $i \in [n]$,

$$\begin{aligned} \beta_i &= \mathbb{E}\left[\text{Tr}\left(\mathbf{V}_\perp^\top (\mathbf{I} + \mathbf{A}_i) \mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top (\mathbf{I} + \mathbf{A}_i^\top) \mathbf{V}_\perp\right)\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\text{Tr}\left(\mathbf{V}_\perp^\top (\mathbf{I} + \mathbf{A}_i) \mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top (\mathbf{I} + \mathbf{A}_i^\top) \mathbf{V}_\perp\right) \middle| \mathcal{F}_{i-}\right]\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\text{Tr}\left(\mathbf{V}_\perp^\top (\mathbf{I} + \eta \mathbf{Y}_i) \mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top (\mathbf{I} + \eta \mathbf{Y}_i) \mathbf{V}_\perp\right) \middle| \mathcal{F}_{i-}\right]\right] \\ &\quad + \mathbb{E}\left[\mathbb{E}\left[\text{Tr}\left(\mathbf{V}_\perp^\top \boldsymbol{\Xi}_i \mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top \boldsymbol{\Xi}_i^\top \mathbf{V}_\perp\right) \middle| \mathcal{F}_{i-}\right]\right]. \end{aligned}$$

774 The last line used $\mathbb{E}[\boldsymbol{\Xi}_i | \mathcal{F}_{i-}] = \mathbf{0}$ and that \mathbf{Z}_{i-1} is measurable with respect to \mathcal{F}_{i-} . In other words,

$$\beta_i = \mathbb{E}\left[\text{Tr}\left(\mathbf{V}_\perp^\top (\mathbf{I} + \eta \mathbf{Y}_i) \mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top (\mathbf{I} + \eta \mathbf{Y}_i) \mathbf{V}_\perp\right)\right] + \mathbb{E}\left[\text{Tr}\left(\mathbf{Z}_{i-1}^\top \mathbb{E}\left[\boldsymbol{\Xi}_i^\top \mathbf{V}_\perp \mathbf{V}_\perp^\top \boldsymbol{\Xi}_i \middle| \mathcal{F}_{i-}\right] \mathbf{Z}_{i-1}\right)\right].$$

775 For the first term, following the analysis of Lemma 10 from [JK⁺16],

$$\begin{aligned} \mathbb{E} \left[\text{Tr} \left(\mathbf{V}_\perp^\top (\mathbf{I} + \eta \mathbf{Y}_i) \mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top (\mathbf{I} + \eta \mathbf{Y}_i) \mathbf{V}_\perp \right) \right] &\leq (1 + 2\eta\lambda_2 + \eta^2 (\mathcal{V}_0 + \lambda_1^2)) \beta_{i-1} + \eta^2 \mathcal{V}_0 \mathbb{E} [\|\mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top\|_2] \\ &\leq (1 + 2\eta\lambda_2 + \eta^2 (\mathcal{V}_0 + \lambda_1^2)) \beta_{i-1} + \eta^2 \mathcal{V}_0 \mathbb{E} [\|\mathbf{Z}_{i-1}\|_2^2]. \end{aligned} \quad (\text{A.18})$$

776 The second term can be bounded as follows:

$$\begin{aligned} \mathbb{E} \left[\text{Tr} \left(\mathbf{Z}_{i-1}^\top \mathbb{E} \left[\boldsymbol{\Xi}_i^\top \mathbf{V}_\perp \mathbf{V}_\perp^\top \boldsymbol{\Xi}_i | \mathcal{F}_{i-} \right] \mathbf{Z}_{i-1} \right) \right] &= \mathbb{E} \left[\text{Tr} \left(\mathbb{E} \left[\boldsymbol{\Xi}_i^\top \mathbf{V}_\perp \mathbf{V}_\perp^\top \boldsymbol{\Xi}_i | \mathcal{F}_{i-} \right] \mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top \right) \right] \\ &\leq \mathbb{E} \left[\mathbb{E} \left[\text{Tr} \left(\boldsymbol{\Xi}_i^\top \mathbf{V}_\perp \mathbf{V}_\perp^\top \boldsymbol{\Xi}_i \right) | \mathcal{F}_{i-} \right] \|\mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top\|_2 \right] \\ &\leq \kappa_1 \mathbb{E} [\|\mathbf{Z}_{i-1} \mathbf{Z}_{i-1}^\top\|_2]. \end{aligned} \quad (\text{A.19})$$

777 Combining (A.18) and (A.19), we obtain the recurrence

$$\beta_i \leq (1 + 2\eta\lambda_2 + \eta^2 (\mathcal{V}_0 + \lambda_1^2)) \beta_{i-1} + (\eta^2 \mathcal{V}_0 + \kappa_1) \mathbb{E} [\|\mathbf{Z}_{i-1}\|_2^2].$$

778 By Lemma A.7, we have for $\gamma := 2(\eta^2 \mathcal{M}^2 + \kappa^2)$,

$$\begin{aligned} \beta_i &\leq (1 + 2\eta\lambda_2 + \eta^2 (\mathcal{V}_0 + \lambda_1^2)) \beta_{i-1} + (\eta^2 \mathcal{V}_0 + \kappa_1) \exp \left(2\sqrt{2n\gamma \log d} \right) (1 + \eta\lambda_1)^{2(i-1)} \\ &\leq \exp (2\eta\lambda_2 + \eta^2 (\mathcal{V}_0 + \lambda_1^2)) \beta_{i-1} + s \exp (2\eta\lambda_1 (i-1) + \eta^2 (\mathcal{V}_0 + \lambda_1^2) (i-1)), \end{aligned}$$

779 where $s = (\eta^2 \mathcal{V}_0 + \kappa_1) \exp (2\sqrt{2n\gamma \log d})$. Unrolling the recursion,

$$\begin{aligned} \beta_n &\leq \exp (2\eta n \lambda_1 + \eta^2 n (\mathcal{V}_0 + \lambda_1^2)) \left[\exp (-2\eta n (\lambda_1 - \lambda_2)) \beta_0 + s \cdot \sum_{t=0}^{n-1} \left(\frac{\exp (2\eta \lambda_2 + \eta^2 (\mathcal{V}_0 + \lambda_1^2))}{\exp (2\eta \lambda_1 + \eta^2 (\mathcal{V}_0 + \lambda_1^2))} \right)^{2(n-1-t)} \right] \\ &\leq \exp (2\eta n \lambda_1 + \eta^2 n (\mathcal{V}_0 + \lambda_1^2)) \left[\exp (-2\eta n (\lambda_1 - \lambda_2)) \beta_0 + \frac{s}{1 - \exp (-2\eta (\lambda_1 - \lambda_2))} \right] \\ &\leq \exp (2\eta n \lambda_1 + \eta^2 n (\mathcal{V}_0 + \lambda_1^2)) \left[\exp (-2\eta n (\lambda_1 - \lambda_2)) \beta_0 + \frac{2.35s}{2\eta (\lambda_1 - \lambda_2)} \right]. \end{aligned}$$

780 where the last line follows since $\frac{1}{1-e^{-x}} \leq \frac{2.35}{x}$ for $x \leq 2$. The proof follows since $\beta_0 \leq d$ and
781 $\exp (2\sqrt{2n\gamma \log d}) \leq \exp(\sqrt{2}) < 4.12$. \square

782 D Proofs of Theorems 1, 2, and 3

783 D.1 Proof of Theorem 1

784 We are now ready to present the proof of Theorem 1, which follows from the following Theorem A.4
785 with $\theta = 0.9$ and $n = b$.

786 **Theorem A.4.** Fix $\theta \in (0, 1)$. Then, for \mathbf{w}_b being the output of Algorithm 1 under assumption 1, η
787 set as $\frac{\alpha \log n}{b(\lambda_1 - \lambda_2)}$, α set as in Lemma A.2, $\kappa_1 \leq 1/2$, and

$$\sqrt{2e^2 b \gamma \log(d/\theta)} \leq \frac{1}{2},$$

788 for $\gamma := 2(\eta^2 \mathcal{M}^2 + \kappa^2)$. Then, with probability at least $1 - 3\theta$,

$$\sin^2(\mathbf{w}, \mathbf{v}_1) \leq \frac{24 \log(1/\theta)}{\theta^3} \left[\frac{d}{\exp(2\alpha \log(n))} + \frac{5(\eta^2 \mathcal{V}_0 + \kappa_1)}{\eta(\lambda_1 - \lambda_2)} \right] + 8\kappa_1.$$

789 *Proof.* Note that by Algorithm 1 and the definition of \mathbf{Z} in (A.17),

$$\mathbf{u}_b = \frac{\mathbf{Z}_b \mathbf{u}_0}{\|\mathbf{Z}_b \mathbf{u}_0\|}.$$

790 Since $\mathbf{v}_1 \mathbf{v}_1^\top + \mathbf{V}_\perp \mathbf{V}_\perp^\top = \mathbf{I}_d$,

$$\sin^2(\mathbf{u}_b, \mathbf{v}_1) = 1 - (\mathbf{u}_b^\top \mathbf{v}_1)^2 = \left\| \frac{\mathbf{V}_\perp \mathbf{V}_\perp^\top \mathbf{Z}_b \mathbf{u}_0}{\|\mathbf{Z}_b \mathbf{u}_0\|} \right\|^2.$$

791 By Lemma 6 from [JJK⁺16], with probability at least $1 - \theta$,

$$\sin^2(\mathbf{u}_b, \mathbf{v}_1) \leq \frac{2.5 \log(1/\theta)}{\theta^2} \frac{\text{Tr}(\mathbf{V}_\perp^\top \mathbf{Z}_b \mathbf{Z}_b^\top \mathbf{V}_\perp)}{\mathbf{v}_1^\top \mathbf{Z}_b \mathbf{Z}_b^\top \mathbf{v}_1}.$$

792 We apply Lemma A.7 with $q = 2$ and $p = 2(1 + \log(d))$. Since $b\gamma(1 + 2\log(d)) \leq 1$,

$$\mathbb{E}[\|\mathbf{Z}_b - (\mathbf{I} + \eta\mathbf{\Sigma})^b\|] \leq \|\mathbf{Z}_b - (\mathbf{I} + \eta\mathbf{\Sigma})^b\|_{p,2} \leq \sqrt{e^2 b\gamma(1 + 2\log(d))} (1 + \eta\lambda_1)^b, \quad (\text{A.20})$$

793 For the numerator, we use Lemma A.8 and Markov's inequality to get

$$\text{Tr}(\mathbf{V}_\perp^\top \mathbf{Z}_b \mathbf{Z}_b^\top \mathbf{V}_\perp) \leq \frac{1}{\theta} \exp(2\eta b\lambda_1 + \eta^2 b(\mathcal{V}_0 + \lambda_1^2)) \left[\frac{d}{\exp(2\eta b(\lambda_1 - \lambda_2))} + \frac{5(\eta^2 \mathcal{V}_0 + \kappa_1)}{\eta(\lambda_1 - \lambda_2)} \right]. \quad (\text{A.21})$$

794 with probability at least $1 - \theta$. The denominator can be bounded as

$$\|\mathbf{Z}_b^\top \mathbf{v}_1\| \geq \|(\mathbf{I} + \eta\mathbf{\Sigma})^b \mathbf{v}_1\| - \left\| \left(\mathbf{Z}_b - (\mathbf{I} + \eta\mathbf{\Sigma})^b \right)^\top \mathbf{v}_1 \right\| \geq (1 + \eta\lambda_1)^b - \|\mathbf{Z}_b - (\mathbf{I} + \eta\mathbf{\Sigma})^b\|.$$

795 Using Lemma A.6, with probability atleast $1 - \theta$,

$$\begin{aligned} \|\mathbf{Z}_b \mathbf{v}_1\| &\geq (1 + \eta\lambda_1)^b - \sqrt{2e^2 b\gamma \log(d/\theta)} (1 + \eta\lambda_1)^b \\ &= (1 + \eta\lambda_1)^b \left(1 - \sqrt{2e^2 b\gamma \log(d/\theta)} \right) \\ &\geq \exp(\eta\lambda_1 b - \eta^2 \lambda_1^2 b) \left(1 - \sqrt{2e^2 b\gamma \log(d/\theta)} \right). \end{aligned} \quad (\text{A.22})$$

796 where the last line follows since $(1 + x) \geq \exp(x - x^2)$ for all $x \geq 0$. From equations (A.21), (A.22), and the assumption $\sqrt{2e^2 b\gamma \log(d/\theta)} \leq 1/2$, it follows that with probability
797 $1 - 3\theta$,
798

$$\sin^2(\mathbf{u}_b, \mathbf{v}_1) \leq \frac{12 \log(1/\theta)}{\theta^3} \left[\frac{d}{\exp(2\alpha \log(n))} + \frac{5(\eta^2 \mathcal{V}_0 + \kappa_1)}{\eta(\lambda_1 - \lambda_2)} \right]. \quad (\text{A.23})$$

799 Since $\mathbf{w} \leftarrow \mathbf{Q}(\mathbf{u}_b, \mathcal{Q})$, by Lemma A.9 and using $\|\xi\| \leq \kappa_1 \leq 0.5$,

$$\sin^2(\mathbf{w}, \mathbf{u}_b) \leq \frac{\|\xi\|^2}{\|\mathbf{u}_b + \xi\|^2} \leq \frac{\|\xi\|^2}{(\|\mathbf{u}_b\| - \|\xi\|)^2} \leq \frac{\kappa_1}{0.5^2} \leq 4\kappa_1. \quad (\text{A.24})$$

800 The result follows by using equations (A.23), (A.24), and Lemma A.5. \square

801 D.2 Proofs of Theorems 2 and 3

802 Next, we apply Theorem A.4 to analyze the quantized version of Oja's algorithm as described in
803 Algorithm 1. The idea is to show that the error from the rounding operation can be incorporated into
804 the noise in the iterates of Oja's algorithm, which, given an appropriately chosen σ -field, will be
805 mean zero. For this subsection, we will use:

$$\mathbf{D}_i = \sum_{j \in B_i} \mathbf{X}_j \mathbf{X}_j^\top / (n/b),$$

806 where $\mathbf{A}_i = \eta(\mathbf{D}_i + \xi_{a,i} \mathbf{u}_{i-1}^\top) + \xi_{2,i} \mathbf{u}_{i-1}^\top + (\mathbf{I} + \eta \mathbf{D}_i) \xi_{1,i} \mathbf{u}_{i-1}^\top$.

807 We first state and prove some intermediate results needed for the proof of Theorem 2.

Theorem A.5. Let $d, n, b \in \mathbb{N}$, and let $\{\mathbf{X}_i\}_{i \in [n]}$ be a set of n IID vectors in \mathbb{R}^d satisfying assumption [1](#). Let $\eta := \frac{\alpha \log n}{b(\lambda_1 - \lambda_2)}$ be the learning rate set as in Lemma [A.2](#). Suppose the quantization grid $\mathcal{Q} = \mathcal{Q}_L$, and $\sqrt{4e^2 b(4\eta^2 + 9\delta^2 d) \log(d/\theta)} \leq \frac{1}{2}$. Then, with probability at least 0.9, the output \mathbf{w} of Algorithm [7](#) satisfies

$$\sin^2(\mathbf{w}, \mathbf{v}_1) \leq \frac{24 \log(1/\theta)}{\theta^3} \left[\frac{d}{n^{2\alpha}} + \frac{5\alpha \mathcal{V} \log n}{n(\lambda_1 - \lambda_2)^2} + \frac{30b\delta^2 d}{\alpha \log n} \right] + 48\delta^2 d.$$

Proof. In order to apply Theorem [1](#) we come up with valid choices of \mathcal{V}_0 , κ , and κ_1 .

Since each \mathbf{D}_i is symmetric and $\{\mathbf{X}_i\}_{i \in [n]}$ are independent,

$$\|\mathbb{E}[(\mathbf{D}_i - \Sigma)(\mathbf{D}_i - \Sigma)^T]\| = \left\| \frac{1}{n/b} \mathbb{E}[(\mathbf{X}_1 \mathbf{X}_1^T - \Sigma)^2] \right\| \leq \frac{b\mathcal{V}}{n} =: \mathcal{V}_0. \quad (\text{A.25})$$

Next,

$$\Xi_i = \eta \xi_{a,i} \mathbf{u}_{i-1}^T + \xi_{2,i} \mathbf{u}_{i-1}^T + (\mathbf{I} + \eta \mathbf{D}_i) \xi_{1,i} \mathbf{u}_{i-1}^T.$$

Also observe that

$$\mathbb{E}[\xi_{1,i} | \mathcal{F}_{i-}] = 0, \quad \mathbb{E}[\xi_{a,i} | \mathcal{F}_{i-}] = 0, \quad \mathbb{E}[\xi_{2,i} | \xi_{a,i}, \xi_{1,i}, \mathcal{F}_{i-}] = 0, \quad (\text{A.26})$$

By equation [A.26](#),

$$\begin{aligned} \mathbb{E}[\Xi_i^T \Xi_i | \mathcal{F}_{i-}] &= \mathbb{E}[\eta^2 \mathbf{u}_{i-1} \xi_{a,i}^T \xi_{a,i} \mathbf{u}_{i-1}^T + \mathbf{u}_{i-1} \xi_{2,i}^T \xi_{2,i} \mathbf{u}_{i-1}^T + \mathbf{u}_{i-1} \xi_{1,i}^T (\mathbf{I} + \eta \mathbf{D}_i) (\mathbf{I} + \eta \mathbf{D}_i)^T \xi_{2,i} \mathbf{u}_{i-1}^T | \mathcal{F}_{i-}] \\ &\implies \|\mathbb{E}[\Xi_i^T \Xi_i | \mathcal{F}_{i-}]\|_F \leq \eta^2 \delta^2 d + \delta^2 d + (1 + \eta)^2 \delta^2 d \leq 6\delta^2 d =: \kappa_1. \end{aligned}$$

As for κ , we have

$$\|\Xi_i\| \leq 2(1 + \eta)\delta\sqrt{d} \leq 3\delta\sqrt{d} =: \kappa$$

We are now ready to obtain the sin-squared error. Note that $\mathcal{M} \leq 2$, since $\|\mathbf{X}_i\| \leq 1$ almost surely, for all $i \in [n]$. By Theorem [A.4](#), with probability at least $1 - 3\theta$,

$$\sin^2(\mathbf{w}, \mathbf{v}_1) \leq \frac{24 \log(1/\theta)}{\theta^3} \left[\frac{d}{\exp(2\alpha \log(n))} + \frac{5(\eta^2 \mathcal{V}_0 + \kappa_1)}{\eta(\lambda_1 - \lambda_2)} \right] + 8\kappa_1.$$

as long as $\sqrt{2e^2 b \gamma \log(d/\theta)} \leq \frac{1}{2}$. Our parameter choices are $\mathcal{V}_0 = \frac{b\mathcal{V}}{n}$, $\kappa = 3\delta\sqrt{d}$, and $\kappa_1 = 6\delta^2 d$.

$$\sin^2(\mathbf{w}, \mathbf{v}_1) \leq \frac{24 \log(1/\theta)}{\theta^3} \left[\frac{d}{n^{2\alpha}} + \frac{5\alpha \mathcal{V} \log n}{n(\lambda_1 - \lambda_2)^2} + \frac{30b\delta^2 d}{\alpha \log n} \right] + 48\delta^2 d.$$

□

Lemma A.9. Let $\mathbf{u} = \mathbf{Q}(\mathbf{w}, \mathcal{Q}_{NL})$, where $\mathbf{u} \in \mathbb{R}^d$ and \mathcal{Q}_{NL} is defined in equation [4](#). Then,

$$\|\mathbf{w} - \mathbf{Q}(\mathbf{w}, \mathcal{Q}_{NL})\| \leq \delta_0 \sqrt{d} + \|\mathbf{w}\| \zeta$$

Proof. Let $\xi = \mathbf{Q}(\mathbf{w}, \mathcal{Q}_{NL}) - \mathbf{w}$. Say $\mathbf{w}_i > 0$. Let k be the unique integer such that $\mathbf{w}_i \in [q_k, q_{k+1}]$. Equivalently for negative \mathbf{w}_i , say the bin is $[-q_{k+1}, -q_k]$. We have:

$$|\xi_i| \leq q_{k+1} - q_k \leq \delta_0 + \zeta q_k \leq |\mathbf{w}_i| \zeta + \delta_0$$

Thus we have:

$$\|\xi\| \leq \delta_0 \sqrt{d} + \|\mathbf{w}\| \zeta.$$

□

Theorem A.6. Fix $\theta \in (0, 1)$. Let the initial vector $\mathbf{u}_0 \sim \mathcal{N}(0, \mathbf{I})$. Let the number of batches b and quantization scale δ be such that $\sqrt{4e^2 b(4\eta^2 + 32\delta_0^2 d + 98\zeta^2) \log(d/\theta)} \leq 1/2$. Then, under assumption [1](#) with η set as $\frac{\alpha \log n}{b(\lambda_1 - \lambda_2)}$, where α is set as in Lemma [A.2](#), $\delta_0 \sqrt{d} \leq 0.25$, and $\zeta \leq 0.25$, with probability at least $1 - 3\theta$, the output \mathbf{w}_b of Algorithm [7](#) gives:

$$\sin^2(\mathbf{w}, \mathbf{v}_1) \leq \frac{24 \log(1/\theta)}{\theta^3} \left[\frac{d}{n^{2\alpha}} + \frac{5\alpha \mathcal{V} \log n}{n(\lambda_1 - \lambda_2)^2} + \frac{5b(4\delta_0 \sqrt{d} + 7\zeta)^2}{\alpha \log n} \right] + 8(4\delta_0 \sqrt{d} + 7\zeta)^2.$$

830 *Proof.* In order to apply Theorem [1](#) we need to bound \mathcal{V} , κ and κ_1 . We start with the first. For us,
 831 \mathbf{D}_i is defined in Eq [9](#). Let \mathcal{R}_i denote the random variables in the quantization up to and including the
 832 i^{th} update.

833 Our analysis is analogous to the previous theorem. Note that the \mathcal{V}_0 parameter is as in Eq [A.25](#).

834 Now we will work out κ and κ_1 since those are the only quantities that change for the nonlinear
 835 quantization. Recall that we have,

$$\Xi_i = \eta \xi_{a,i} \mathbf{u}_{i-1}^T + \xi_{2,i} \mathbf{u}_{i-1}^T + (\mathbf{I} + \eta \mathbf{D}_i) \xi_{1,i} \mathbf{u}_{i-1}^T.$$

836 We have,

$$\begin{aligned} & \mathbb{E}[\Xi_i^T \Xi_i | \mathcal{F}_{i-}] \\ &= \eta^2 \mathbb{E}[\mathbf{u}_{i-1} \xi_{a,i}^T \xi_{a,i} \mathbf{u}_{i-1}^T | \mathcal{F}_{i-}] + \mathbb{E}[\mathbf{u}_{i-1} \xi_{2,i}^T \xi_{2,i} \mathbf{u}_{i-1}^T | \mathcal{F}_{i-}] + \mathbb{E}[\mathbf{u}_{i-1} \xi_{1,i}^T (\mathbf{I} + \eta \mathbf{D}_i) (\mathbf{I} + \eta \mathbf{D}_i)^T \xi_{2,i} \mathbf{u}_{i-1}^T | \mathcal{F}_{i-}] \end{aligned}$$

837 Now we obtain the Frobenius norm of $\xi_{a,i}$, ξ_1 , and ξ_2 under the nonlinear quantization. We start
 838 with the norm of \mathbf{w}_i , a quantized version of a unit vector \mathbf{u}_{i-1} .

839 By Lemma [A.9](#), $\|\mathbf{w}_i\| \leq 1 + \delta_0 \sqrt{d} + \zeta$. Let $\mathbf{s}_j = \mathbf{X}_j (\mathbf{X}_j^T \mathbf{w}_i)$. Then,

$$\|\mathbf{s}_j\| \leq \|\mathbf{w}_i\| \leq 1 + \delta_0 \sqrt{d} + \zeta.$$

840 Another application of Lemma [A.9](#) gives:

$$\|\xi_{a,j,i}\| = \|\mathbf{Q}(\mathbf{s}_j, \mathcal{Q}_{NL}) - \mathbf{s}_j\| \leq \delta_0 \sqrt{d} + (1 + \delta_0 \sqrt{d} + \zeta) \zeta \leq \delta_0 \sqrt{d} + 1.5\zeta$$

841 which implies $\|\xi_{a,i}\| \leq \delta_0 \sqrt{d} + 1.5\zeta$. Next, we bound $\xi_{1,i} = \mathbf{Q}(\mathbf{u}_{i-1}, \mathcal{Q}_{NL}) - \mathbf{u}_{i-1}$. By Lemma [A.9](#),

$$\|\xi_{1,i}\| \leq \delta_0 \sqrt{d} + \zeta \|\mathbf{u}_{i-1}\| = \delta_0 \sqrt{d} + \zeta.$$

842 Finally we bound $\xi_{2,i}$. Recall that:

$$\begin{aligned} \mathbf{y}_i &= \frac{\sum_{j \in \mathcal{B}_j} \mathbf{X}_j (\mathbf{X}_j^T \mathbf{w}_i)}{n/b} + \xi_{a,i} \\ \xi_{2,i} &= \mathbf{Q}(\mathbf{y}_i, \delta) - \mathbf{y}_i \end{aligned}$$

843 Since each $\|\mathbf{X}_j \mathbf{X}_j^T \mathbf{w}_i\| \leq 1 + \delta_0 \sqrt{d} + \zeta$,

$$\|\mathbf{y}_i\| \leq 1 + \delta_0 \sqrt{d} + \zeta + \|\xi_{a,i}\| \leq 1 + 2\delta_0 \sqrt{d} + 2.5\zeta \leq 3.25.$$

844 By Lemma [A.9](#)

$$\|\xi_{2,i}\| \leq \delta_0 \sqrt{d} + \zeta \|\mathbf{y}_i\| \leq \delta_0 \sqrt{d} + 3.25\zeta.$$

845 In all, it follows that

$$\|\Xi_i\| \leq \eta \|\xi_{a,i}\| + \|\xi_{2,i}\| + (1 + \eta) \|\xi_{1,i}\| \leq (\delta_0 \sqrt{d} + 1.5\zeta) + (\delta_0 \sqrt{d} + 3.25\zeta) + 2(\delta_0 \sqrt{d} + \zeta) \leq 4\delta_0 \sqrt{d} + 7\zeta =: \kappa.$$

846 We are ready to obtain the sin-squared error. Note that $\mathcal{M} \leq 2$, since $\|\mathbf{X}_i\| \leq 1$ almost surely, for all
 847 $i \in [n]$. By Theorem [A.4](#) with probability at least $1 - 3\theta$,

$$\sin^2(\mathbf{w}, \mathbf{v}_1) \leq \frac{24 \log(1/\theta)}{\theta^3} \left[\frac{d}{\exp(2\alpha \log(n))} + \frac{5(\eta^2 \mathcal{V}_0 + \kappa_1)}{\eta(\lambda_1 - \lambda_2)} \right] + 8\kappa_1.$$

848 as long as $\sqrt{2e^2 b \gamma \log(d/\theta)} \leq \frac{1}{2}$. Our parameter choices are $\mathcal{V}_0 = \frac{b\mathcal{V}}{n}$, $\kappa = 4\delta_0 \sqrt{d} + 7\zeta$, and
 849 $\kappa_1 = (4\delta_0 \sqrt{d} + 7\zeta)^2$. Therefore,

$$\sin^2(\mathbf{w}, \mathbf{v}_1) \leq \frac{24 \log(1/\theta)}{\theta^3} \left[\frac{d}{n^{2\alpha}} + \frac{5\alpha \mathcal{V} \log n}{n(\lambda_1 - \lambda_2)^2} + \frac{5b(4\delta_0 \sqrt{d} + 7\zeta)^2}{\alpha \log n} \right] + 8(4\delta_0 \sqrt{d} + 7\zeta)^2.$$

850 □

851 D.2.1 Finishing the Proofs of Theorems 2 and 3

852 *Proof of Theorem 2* For the linear quantization scheme, we apply Theorem A.5 with $\theta =$
853 $1/30$ and $b = \Theta\left(\frac{\alpha^2 \log^2 n \log d}{(\lambda_1 - \lambda_2)^2}\right)$. Moreover, since $\delta = \tilde{O}\left(\frac{\lambda_1 - \lambda_2}{\alpha \sqrt{d}}\right)$, the condition
854 $\sqrt{4e^2 b(4\eta^2 + 9\delta^2 d) \log(d/\theta)} \leq \frac{1}{2}$ holds. The Theorem follows by substituting these values into the
855 bound of Theorem A.5

856 The proof of non-linear quantization scheme follows analogously by using Theorem A.6 \square

857 *Proof of Theorem 3* We set $\theta = 1/30$. For the linear quantization scheme, we apply Theo-
858 rem A.5 with $b = n$. Moreover, since $\delta = 2^{2-\beta} = O\left(\min\left(\frac{\lambda_1 - \lambda_2}{\alpha \sqrt{d \log(n)}}, \frac{1}{\sqrt{dn}}\right)\right)$, the condition
859 $\sqrt{4e^2 b(4\eta^2 + 9\delta^2 d) \log(d/\theta)} \leq \frac{1}{2}$ holds. The Theorem follows by substituting these values into the
860 bound of Theorem A.5

861 For the non-linear scheme, the proof follows analogously by applying Theorem A.6 \square

862 E Proof of Boosting Lemma (Lemma 3)

863 In this section, we present the proof of the boosting procedure. The novelty of our analysis of the
864 boosting lemma is that we use a different quantization scale ϵ . For ease of exposition, we present all
865 proofs such that $\tilde{\rho}$ uses the extended quantization grid

$$\mathcal{Q}_L^*(\epsilon) = \{k\epsilon : k \in \mathbb{Z}\}.$$

866 which ensures $|\mathbf{Q}(\sin^2(\mathbf{u}_i, \mathbf{u}_j), \mathcal{Q}_L^*(\epsilon)) - x| \leq \epsilon$ almost surely. However, this requires a larger
867 number of bits than needed. As we note later in Remark A.1, a restricted quantization grid

$$\mathcal{Q}_L(\epsilon) = \{-2^{\beta-1}\epsilon, -(2^{\beta-1} - 1)\epsilon, \dots, -\epsilon, 0, \epsilon, \dots, (2^{\beta-1} - 1)\epsilon\}.$$

868 which only quantizes values between $[-2^{\beta-1}\epsilon, (2^{\beta-1} - 1)\epsilon]$ and projects to the ends otherwise,
869 suffices because the comparisons made in Algorithm 2 match exactly for both quantization grids. This
870 requires a modest assumption that $\beta \geq 4$ which is already assumed in Section 3.4 while optimizing
871 the parameters.

872 Proof of Lemma 3

873 *Proof.* For any $i \in [r]$, define the indicator random variable $\chi_i := \mathbb{1}(\sin^2(\mathbf{u}_i, \mathbf{v}) \leq \epsilon)$. Then,
874 $\Pr(\chi_i = 1) \geq 1 - p$, where $p = 0.1$. Let $\mathcal{S} := \{i \in [r] : \chi_i = 1\}$, and define the event

$$\mathcal{E} := \{|\mathcal{S}| > 0.6r\}.$$

875 The Chernoff bound for the sum of independent Bernoulli random variables gives

$$\mathbb{P}(|\mathcal{S}| \leq (1 - \theta) \mathbb{E}[|\mathcal{S}|]) \leq \exp\left(-\frac{\theta^2 \mathbb{E}[|\mathcal{S}|]}{2}\right) \quad \forall \theta \in (0, 1).$$

876 By linearity of expectation, $\mathbb{E}[|\mathcal{S}|] \geq (1 - p)r$. Setting $\theta = 1/3$,

$$\mathbb{P}(\mathcal{E}^c) \leq \mathbb{P}(|\mathcal{S}| \leq 0.6r) \leq e^{-r/20} \leq \delta.$$

877 It suffices to show that if the event \mathcal{E} holds, then $\bar{\mathbf{u}}$ is well-defined and has small sin-squared error
878 with \mathbf{v} . Recall,

$$\bar{\mathbf{u}} := \mathbf{u}_i \text{ such that } |\{j \in [r] : \tilde{\rho}(\mathbf{u}_i, \mathbf{u}_j) \leq 5\epsilon\}| \geq 0.5r,$$

879 Conditioned on \mathcal{E} , such a $\bar{\mathbf{u}}$ always exists because \mathbf{u}_i such that $i \in \mathcal{S}$ is a valid choice. Indeed, for
880 any $i, j \in \mathcal{S}$,

$$|\tilde{\rho}(\mathbf{u}_i, \mathbf{u}_j)| \leq \sin^2(\mathbf{u}_i, \mathbf{u}_j) + \epsilon \leq 2\sin^2(\mathbf{u}_i, \mathbf{v}) + 2\sin^2(\mathbf{v}, \mathbf{u}_j) + \epsilon \leq 5\epsilon.$$

881 Since $\tilde{\rho}(\bar{\mathbf{u}}, \mathbf{u}_j) \leq 5\epsilon$ for at least $0.5r$ indices $j \in [r]$ and $|\mathcal{S}| \geq 0.6r$, there exists some index $j \in \mathcal{S}$
882 such that $\tilde{\rho}(\bar{\mathbf{u}}, \mathbf{u}_j) \leq 5\epsilon$. Therefore,

$$\sin^2(\bar{\mathbf{u}}, \mathbf{v}) \leq 2\sin^2(\bar{\mathbf{u}}, \mathbf{u}_j) + 2\sin^2(\mathbf{u}_j, \mathbf{v}) \leq 2\tilde{\rho}(\bar{\mathbf{u}}, \mathbf{u}_j) + 2\epsilon + 2\sin^2(\mathbf{u}_j, \mathbf{v}) \leq 14\epsilon.$$

883 \square

Theorem A.7. Suppose \mathcal{A} is the batched Oja’s algorithm with the setting of Theorem 2. Let ϵ be the probability 0.9 error bound guaranteed by Theorem 2, $r = \lceil 20 \log(1/\theta) \rceil$, and $m = nr$. Let $\{\mathbf{X}_i\}_{i \in [m]}$ be n IID data drawn from a distribution satisfying assumption 1, and $\mathbf{u}_j \leftarrow \mathcal{A}(\{\mathbf{X}_i\}_{(j-1)n+1 \leq i \leq jn})$ for all $j \in [r]$. Define the function $\tilde{\rho}$ as $\tilde{\rho}(\mathbf{x}, \mathbf{y}) = Q(\sin^2(\mathbf{x}, \mathbf{y}), Q_L^*(\epsilon))$, where $Q_L^*(\epsilon)$ is the linear quantization grid with spacing ϵ . Then, with probability at least $1 - \theta$, $\bar{\mathbf{u}} := \text{SuccessBoost}(\{\mathbf{u}_i\}_{i \in [r]}, \tilde{\rho}, \epsilon)$ satisfies

$$\sin^2(\bar{\mathbf{u}}, \mathbf{v}_1) \leq 14\epsilon.$$

Proof. The vectors $\mathbf{u}_1, \dots, \mathbf{u}_r$ are mutually independent. By Theorem 2, $\Pr(\sin^2(\mathbf{u}_i, \mathbf{v}_1) > \epsilon) \leq 0.1 \forall i \in [r]$. Finally, since the quantization grid while computing $\sin^2(\mathbf{x}, \mathbf{y})$ for each pair of vectors \mathbf{x} and \mathbf{y} has scale ϵ , $|\tilde{\rho}(\mathbf{u}_i, \mathbf{u}_j) - \sin^2(\mathbf{x}, \mathbf{y})| \leq \epsilon$. Therefore, all conditions of Lemma 3 apply, and the theorem follows. \square

Remark A.1. The quantization scheme for aggregation of the $\Theta(\log 1/\theta)$ vectors in the boosting algorithm is different from that of the main algorithm (which uses either a linear scheme with parameter δ or a logarithmic scheme with parameters ζ and δ_0). The quantization grid in the boosting algorithm is linear with parameter ϵ , where ϵ is the high-probability error guarantee. If the number of bits β to represent any number is fixed,

$$Q_L(\epsilon) = \{-2^{\beta-1}\epsilon, -(2^{\beta-1} - 1)\epsilon, \dots, -\epsilon, 0, \epsilon, \dots, (2^{\beta-1} - 1)\epsilon\}.$$

Since $\sin^2(\mathbf{u}_i, \mathbf{v}_1) \leq \epsilon$ with probability at least 0.9, most values of $\sin^2(\mathbf{u}_i, \mathbf{u}_j)$ will be at most 4ϵ using Lemma A.5. Any j for which $\sin^2(\mathbf{u}_i, \mathbf{u}_j)$ is outside this range of representable values has a large \sin^2 error with \mathbf{u}_i and therefore not contained in the set $\{j \in [r] : \tilde{\rho}(\mathbf{u}_i, \mathbf{u}_j) \leq 5\epsilon\}$. Rigorously, consider the extended quantization grid

$$Q_L^*(\epsilon) = \{k\epsilon : k \in \mathbb{Z}\}.$$

Then, $|Q(\sin^2(\mathbf{u}_i, \mathbf{u}_j), Q_L^*(\epsilon)) - x| \leq \epsilon$ almost surely. Moreover, if $\sin^2(\mathbf{u}_i, \mathbf{u}_j)$ is contained in the range of $Q_L(\epsilon)$, then $Q(\sin^2(\mathbf{u}_i, \mathbf{u}_j), Q_L(\epsilon)) = Q(\sin^2(\mathbf{u}_i, \mathbf{u}_j), Q_L^*(\epsilon))$. Therefore,

$$\mathbb{1}(Q(\sin^2(\mathbf{u}_i, \mathbf{u}_j), Q_L(\epsilon)) \leq 5\epsilon) = \mathbb{1}(Q(\sin^2(\mathbf{u}_i, \mathbf{u}_j), Q_L^*(\epsilon)) \leq 5\epsilon).$$

That is, by not quantizing $\tilde{\rho}(\mathbf{u}_i, \mathbf{u}_j)$ indices j for which $\sin^2(\mathbf{u}_i, \mathbf{u}_j)$ is greater than $(2^{\beta-1} - 1)\epsilon$ in magnitude but simply not counting j as a neighbor of i in the boosting procedure, the algorithm is unaffected. This requires a modest bound on β such as $\beta \geq 4$, (because $2^{\beta-1} - 1 = 7 > 5$), which is already assumed in Section 3.4 while optimizing the parameters.

F Experimental Details

Given the sample size n , dimension d , and decay exponent λ in the eigenvalues, we first draw an $n \times d$ matrix Z with independent entries uniformly distributed on $[-\sqrt{3}, \sqrt{3}]$ so that each coordinate has unit variance. We then build a kernel matrix $K \in \mathbb{R}^{d \times d}$ with entries $K_{ij} = \exp(-|i - j|^{0.01})$ and define a variance profile $\sigma_i = 5i^{-\lambda}$ for $i = 1, \dots, d$. The population covariance is formed as $\Sigma = (\sigma\sigma^\top) \circ K$, where \circ denotes the Hadamard product. Computing the eigendecomposition of Σ yields its square root $\Sigma^{1/2}$, and the observed data matrix is taken as $X = (\Sigma^{1/2}Z^\top)^\top$. We then extract the largest two eigenvalues $\lambda_1 > \lambda_2$ of Σ and the associated top eigenvector v_1 for evaluation.

G Related Work

In this section, we provide some more related work on low-precision optimization. For an excellent survey and history of quantization, see [GKD⁺22]. Dettmers *et al.* introduced QLoRA, which back-propagates through a frozen 4-bit quantized LLM into LoRA modules, enabling efficient finetuning of 65B-parameter models on a single 48 GB GPU with full 16-bit performance retention [DPHZ23]. Earlier works [XMHK23] examined the impact of stochastic roundoff errors and their bias on gradient descent convergence under low-precision arithmetic. Yu *et al.* propose Collage, a lightweight low-precision scheme for LLM training in distributed settings, combining block-wise quantization with error-feedback to stabilize large-scale pretraining [YGG⁺24]. Finally, communication-efficient distributed SGD techniques, such as 1-bit SGD with error feedback [SFD⁺14] and randomized

sketching primitives (e.g., Johnson–Lindenstrauss projections [JL84]), further underscore the broad efficacy of low-precision computation.

Low-Precision Optimization: Reducing the bit-width of model parameters and gradient updates has proven effective for alleviating communication and memory bottlenecks in large-scale learning. QSGD [AGL⁺17] uses randomized rounding to compress each coordinate to a few bits while preserving unbiasedness, incurring only an $O(\sqrt{d}/2^b)$ increase in gradient noise for b bits. TernGrad [WXY⁺17] maps gradients to $\{-1, 0, +1\}$ plus a shared scale and demonstrates negligible accuracy loss on ImageNet and CIFAR benchmarks. Suresh *et al.* [SYKM17] achieve optimal communication–accuracy trade-offs via randomized rotations and scalar quantization. More recently, “dimension-free” analyses such as Li & De Sa [LDS19] avoid scaling the required error rate with model dimension, instead depending on a suitably defined smoothness parameter.

Low-Precision PCA: Streaming PCA methods update eigenvector estimates incrementally, using $O(dp)$ memory to maintain the top p components in d dimensions. Oja’s rule [Oja82] implements a Hebbian-style one-pass update with normalization. However, to our knowledge, *no general analysis* of Oja’s algorithm under limited precision exists before our work.

References

- [AGL⁺17] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pages 1707–1718, 2017.
- [AZL17] Zeyuan Allen-Zhu and Yuanzhi Li. First efficient convergence for streaming k-pca: a global, gap-free, and near-optimal rate. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 487–492. IEEE, 2017.
- [CYWZ18] Minshuo Chen, Lin Yang, Mengdi Wang, and Tuo Zhao. Dimensionality reduction for stationary time series via stochastic nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- [DPHZ23] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023. <https://arxiv.org/abs/2305.14314>
- [GKD⁺22] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In *Low-power computer vision*, pages 291–326. Chapman and Hall/CRC, 2022.
- [HCS⁺16] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. In *Advances in Neural Information Processing Systems*, 2016.
- [Heb49] Donald O. Hebb. *The Organization of Behavior: A Neuropsychological Theory*. John Wiley & Sons, New York, 1949.
- [HNWTW20] De Huang, Jonathan Niles-Weed, Joel A. Tropp, and Rachel Ward. Matrix concentration for products, 2020.
- [HW19] Amelia Henriksen and Rachel Ward. AdaOja: Adaptive Learning Rates for Streaming PCA. *arXiv e-prints*, page arXiv:1905.12115, May 2019.
- [JJK⁺16] Prateek Jain, Chi Jin, Sham Kakade, Praneeth Netrapalli, and Aaron Sidford. Streaming pca: Matching matrix bernstein and near-optimal finite sample guarantees for oja’s algorithm. In *Proceedings of The 29th Conference on Learning Theory (COLT)*, June 2016.
- [JL84] William B. Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. In *Contemporary Mathematics*, volume 26, page 189–206, 1984.
- [KLL⁺23] Jonathan Kelner, Jerry Li, Allen X Liu, Aaron Sidford, and Kevin Tian. Semi-random sparse recovery in nearly-linear time. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 2352–2398. PMLR, 2023.

- [KS24a] Syamantak Kumar and Purnamrita Sarkar. Oja’s algorithm for streaming sparse pca. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [KS24b] Syamantak Kumar and Purnamrita Sarkar. Streaming pca for markovian data. *Advances in Neural Information Processing Systems*, 36, 2024.
- [KWW⁺17] Urs Köster, Tristan J. Webb, Xin Wang, Marcel Nassar, Arjun K. Bansal, William H. Constable, Oğuz H. Elibol, Scott Gray, Stewart Hall, Luke Hornof, Amir Khosrowshahi, Carey Kloss, Ruby J. Pai, and Naveen Rao. Flexpoint: An adaptive numerical format for efficient training of deep neural networks. In *Advances in Neural Information Processing Systems*, volume 30, pages 1742–1750, 2017.
- [LD19] Zheng Li and Christopher M. De Sa. Dimension-free bounds for low-precision training. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pages 11728–11738, 2019.
- [LDS19] Zheng Li and Christopher De Sa. Dimension-free bounds for low-precision training. In *Advances in Neural Information Processing Systems*, 2019.
- [LDX⁺17] Hao Li, Soham De, Zheng Xu, Christoph Studer, Hanan Samet, and Tom Goldstein. Training quantized nets: A deeper understanding. In *Advances in Neural Information Processing Systems*, pages 5813–5823, 2017.
- [LSW21] Robert Lunde, Purnamrita Sarkar, and Rachel Ward. Bootstrapping the error of oja’s algorithm. *Advances in Neural Information Processing Systems*, 34:6240–6252, 2021.
- [MNA⁺18] Paulius Micikevicius, Sharan Narang, Gabriel Alben, Gregory Diamos, Erich Elsen, David Garcia, Dmitry Ginsburg, Michael Houston, Oleksii Kuchaiev, Sanjo Venkatesh, and Hao Wu. Mixed precision training. In *International Conference on Learning Representations*, 2018.
- [Mon22] Jean-Marie Monnez. Stochastic approximation of eigenvectors and eigenvalues of the q-symmetric expectation of a random matrix. *Communications in Statistics-Theory and Methods*, pages 1–15, 2022.
- [MP22] Nikos Mouzakis and Eric Price. Spectral guarantees for adversarial streaming pca, 2022.
- [NTSW⁺22] Miloš Nikolić, Enrique Torres Sanchez, Jiahui Wang, Ali Hadi Zadeh, Mostafa Mahmoud, Ameer Abdelhadi, Kareem Ibrahim, and Andreas Moshovos. Schrödinger’s fp: Dynamic adaptation of floating-point containers for deep learning training. *arXiv preprint arXiv:2204.13666*, 2022.
- [Oja82] Erkki Oja. Simplified neuron model as a principal component analyzer. *Journal of mathematical biology*, 15:267–273, 1982.
- [OYP25] Kaan Ozkara, Tao Yu, and Youngsuk Park. Stochastic rounding for llm training: Theory and practice. In *Proceedings of the 28th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2025. <https://arxiv.org/abs/2502.20566>.
- [Pea01] Karl Pearson. Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin philosophical magazine and journal of science*, 2(11):559–572, 1901.
- [Rie67] Bernhard Riemann. *Ueber die Darstellbarkeit einer Function durch eine trigonometrische Reihe*. Dieterich, 1867. In German.
- [SFD⁺14] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. In *Interspeech*, 2014.

- 1023 [She97] William Fleetwood Sheppard. On the calculation of the most probable values of
1024 frequency-constants for data arranged according to equidistant division of a scale.
1025 *Proceedings of the London Mathematical Society*, 1(1):353–380, 1897.
- 1026 [SLZ⁺18] Christopher De Sa, Megan Leszczynski, Jian Zhang, Alana Marzoev, Christopher R.
1027 Aberger, Kunle Olukotun, and Christopher Ré. High-accuracy low-precision training.
1028 *arXiv preprint arXiv:1803.03383*, 2018.
- 1029 [Sti86] S. M. Stigler. *The History of Statistics: The Measurement of Uncertainty before 1900*.
1030 Harvard University Press, Cambridge, 1986.
- 1031 [SYK21] Heming Sun, Lu Yu, and Jiro Katto. Learned image compression with fixed-point
1032 arithmetic. In *2021 Picture Coding Symposium (PCS)*, pages 1–5. IEEE, 2021.
- 1033 [SYKM17] Ananda Theertha Suresh, Felix X. Yu, Harsha Kumar, and H. Brendan McMa-
1034 han. Distributed mean estimation with limited communication. *arXiv preprint*
1035 *arXiv:1611.00349*, 2017.
- 1036 [SZOR15] Christopher M. De Sa, Ce Zhang, Kunle Olukotun, and Christopher Ré. Taming
1037 the wild: A unified analysis of hogwild-style algorithms. In *Advances in Neural*
1038 *Information Processing Systems*, pages 2674–2682, 2015.
- 1039 [Ver10] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices.
1040 *arXiv preprint arXiv:1011.3027*, 2010.
- 1041 [Wed72] Per-Åke Wedin. Perturbation bounds in connection with singular value decomposition.
1042 *BIT Numerical Mathematics*, 12:99–111, 1972.
- 1043 [WXY⁺17] Wei Wen, Chunpeng Xu, Felix Yan, Chunyi Wu, Yandan Wang, Yiran Chen, and Hai
1044 Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning.
1045 In *Advances in Neural Information Processing Systems*, 2017.
- 1046 [XLY⁺24] Yongqi Xu, Yujian Lee, Gao Yi, Bosheng Liu, Yucong Chen, Peng Liu, Jigang
1047 Wu, Xiaoming Chen, and Yinhe Han. Bitq: Tailoring block floating point preci-
1048 sion for improved dnn efficiency on resource-constrained devices. *arXiv preprint*
1049 *arXiv:2409.17093*, 2024.
- 1050 [XMHK23] Lu Xia, Stefano Massei, Michiel E. Hochstenbach, and Barry Koren. On the influence
1051 of stochastic roundoff errors and their bias on the convergence of the gradient descent
1052 method with low-precision floating-point computation, 2023.
- 1053 [Yat09] Randy Yates. Fixed-point arithmetic: An introduction. *Digital Signal Labs*,
1054 81(83):198, 2009.
- 1055 [YGG⁺24] Tao Yu, Gaurav Gupta, Karthick Gopalswamy, Amith R. Mamidala, Hao Zhou, Jeffrey
1056 Huynh, Youngsuk Park, Ron Diamant, Anoop Deoras, and Luke Huan. Collage: Light-
1057 weight low-precision strategy for llm training. In *Proceedings of the 41st International*
1058 *Conference on Machine Learning*, 2024.
- 1059 [YHW18] Puyudi Yang, Cho-Jui Hsieh, and Jane-Ling Wang. History pca: A new algorithm for
1060 streaming pca. *arXiv preprint arXiv:1802.05447*, 2018.
- 1061 [YIY21] Hisakatsu Yamaguchi, Makiko Ito, and Katsuhiro Yoda. Training deep neural networks
1062 in 8-bit fixed point with dynamic shared exponent management. In *Proceedings of*
1063 *the 2021 Design, Automation & Test in Europe Conference (DATE)*, 2021.
- 1064 [Zie03] Eric R Ziegel. Principal component analysis. *Technometrics*, 45(3):276–277, 2003.
- 1065 [ZLK⁺17] Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. ZipML:
1066 Training linear models with end-to-end low precision, and a little bit of deep learning.
1067 In *Proceedings of the 34th International Conference on Machine Learning*, pages
1068 4035–4043, 2017.

- 1069 [ZMK22] Sai Qian Zhang, Bradley McDanel, and T. Kung, H. Fast: Dnn training under variable
1070 precision block floating point with stochastic rounding. In *Proceedings of the 2022*
1071 *IEEE International Symposium on High-Performance Computer Architecture (HPCA)*,
1072 pages 846–860, 2022.
- 1073 [ZWG⁺23] Jiajun Zhou, Jiajun Wu, Yizhao Gao, Yuhao Ding, Chaofan Tao, Boyu Li, Fengbin
1074 Tu, Kwang-Ting Cheng, Hayden Kwok-Hay So, and Ngai Wong. Dybit: Dynamic
1075 bit-precision numbers for efficient quantized neural network inference. *arXiv preprint*
1076 *arXiv:2302.12510*, 2023.