

## Supplementary Material

### A. Auxillary Lemmas

**Lemma A.1** (Lemma 14.1 (Shalev-Shwartz & Ben-David, 2014)). Any update of the form

$$\mathbf{P}^{(t+1)} = \Pi_C(\mathbf{P}^{(t)} - \eta \mathbf{g}_t), \quad (13)$$

for an arbitrary sequence of matrices  $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_T$ , projection  $\Pi_C$  onto an arbitrary convex set  $\mathcal{C}$ , and initialization  $\mathbf{P}^{(1)} = 0$  satisfies

$$\sum_{t=1}^T \langle \mathbf{P}^{(t)} - \mathbf{P}, \mathbf{g}_t \rangle \leq \frac{\|\mathbf{P}\|_F^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_F^2, \quad (14)$$

for any  $\mathbf{P} \in \mathcal{C}$ .

**Lemma A.2** (Lemma 4 (Warmuth & Kuzmin, 2006)). For any PSD  $\mathbf{A}$  and any symmetric  $\mathbf{B}, \mathbf{C}$  for which  $\mathbf{B} \preceq \mathbf{C}$ , it holds that  $\text{Tr}(\mathbf{A}\mathbf{B}) \leq \text{Tr}(\mathbf{A}\mathbf{C})$ .

### B. Proofs of Section 3

*Proof of Theorem 2.1.* Using Lemma A.1, noting  $\|\mathbf{P}\|_F^2 \leq k$  for all  $\mathbf{P} \in \mathcal{C}$ , and  $\mathbf{g}_t = -\mathbf{x}_t \mathbf{x}_t^\top$  so that  $\|\mathbf{g}_t\|_F^2 = \|\mathbf{x}_t\|_2^4 \leq 1$ ,

$$\sum_{t=1}^T \mathbf{x}_t^\top \mathbf{P} \mathbf{x}_t - \sum_{t=1}^T \mathbf{x}_t^\top \mathbf{P}^{(t)} \mathbf{x}_t \leq \frac{k}{2\eta} + \frac{\eta T}{2}.$$

Choosing  $\eta = \sqrt{\frac{k}{T}}$  completes the proof.  $\square$

*Proof of Theorem 3.1.* Let  $\widehat{\mathbf{g}}_t = \mathbf{g}_t + \mathbf{E}_t$  be the noisy gradient. The analysis begins with bounding the distance between the  $t$ -th iterate and any candidate  $\mathbf{P} \in \mathcal{C}$ ,  $D_t = \|\mathbf{P}_t - \mathbf{P}\|_F$ .

$$\begin{aligned} D_{t+1}^2 &= \left\| \mathbf{P}^{(t+1)} - \mathbf{P} \right\|_F^2 = \left\| \Pi_C(\mathbf{P}^{(t)} - \eta \widehat{\mathbf{g}}_t) - \mathbf{P} \right\|_F^2 \\ &\leq \left\| \mathbf{P}^{(t)} - \eta \widehat{\mathbf{g}}_t - \mathbf{P} \right\|_F^2 \\ &= \left\| \mathbf{P}^{(t)} - \mathbf{P} \right\|_F^2 + \eta^2 \|\widehat{\mathbf{g}}_t\|_F^2 \\ &\quad - 2\eta \langle \mathbf{P}^{(t)} - \mathbf{P}, \mathbf{g}_t + \mathbf{E}_t \rangle \\ &\leq D_t^2 + \eta^2 G^2 - 2\eta \langle \mathbf{P}^{(t)} - \mathbf{P}, \mathbf{g}_t \rangle \\ &\quad - 2\eta \langle \mathbf{P}^{(t)} - \mathbf{P}, \mathbf{E}_t \rangle \\ &\leq D_t^2 + \eta^2 G^2 - 2\eta \langle \mathbf{P}^{(t)} - \mathbf{P}, \mathbf{g}_t \rangle \\ &\quad + 2\eta \left\| \mathbf{P}^{(t)} - \mathbf{P} \right\|_* \|\mathbf{E}_t\|_2 \\ &\leq D_t^2 + \eta^2 G^2 - 2\eta \langle \mathbf{P}^{(t)} - \mathbf{P}, \mathbf{g}_t \rangle + 4k\eta \|\mathbf{E}_t\|_2 \end{aligned}$$

Where the first inequality is due to projection onto a convex set being non-expanding and the third one is by Holder's inequality. The last inequality follows from triangle inequality and constraints:  $\left\| \mathbf{P}^{(t)} - \mathbf{P} \right\|_* \leq \left\| \mathbf{P}^{(t)} \right\|_* + \|\mathbf{P}\|_* \leq 2k$ .

Noting  $\mathbf{g}_t = -\mathbf{x}_t \mathbf{x}_t^\top$ , rearranging and dividing both sides by  $2\eta$  we get

$$\langle \mathbf{P} - \mathbf{P}^{(t)}, \mathbf{x}_t \mathbf{x}_t^\top \rangle \leq \frac{D_t^2 - D_{t+1}^2}{2\eta} + \frac{\eta}{2} G^2 + 2k \|\mathbf{E}_t\|_2 \quad (15)$$

We then sum over  $T$  iterates and use the telescopic property of the first term of the right hand side of 15, so that  $\sum_{t=1}^T D_t^2 - D_{t+1}^2 = D_1^2 - D_{T+1}^2 \leq D_1^2$ . The initial distance  $D_1$  is bounded as follows:

$$\begin{aligned} D_1^2 &= \left\| \mathbf{P}^{(1)} - \mathbf{P} \right\|_F^2 = \left\| \mathbf{P}^{(1)} \right\|_F^2 + \|\mathbf{P}\|_F^2 - 2\langle \mathbf{P}^{(1)}, \mathbf{P} \rangle \\ &\leq k + k + 2 \left\| \mathbf{P}^{(1)} \right\|_* \|\mathbf{P}\|_2 \leq 4k \end{aligned}$$

Plugging back the above bound in 15 and noting  $\sum_{t=1}^T \|\mathbf{E}_t\|_2 \leq E$ , we get

$$\sum_{t=1}^T \mathbf{x}_t^\top \mathbf{P} \mathbf{x}_t - \sum_{t=1}^T \mathbf{x}_t^\top \mathbf{P}^{(t)} \mathbf{x}_t \leq \frac{2k}{\eta} + \frac{\eta G^2 T}{2} + 2kE$$

By choosing optimal learning rate  $\eta = \frac{2\sqrt{k}}{G\sqrt{T}}$  we get the desired results.  $\square$

*Proof of Theorem 3.2.* Let  $\widehat{\mathbf{g}}_k = (\mathbf{x}_k + \mathbf{y}_k)(\mathbf{x}_k + \mathbf{y}_k)^\top$  denote the corrupted gradient and let  $\mathbf{g}_k = \mathbf{x}_k \mathbf{x}_k^\top$  denote the unbiased estimate of  $\mathbf{C}$  based on the  $k$ -th sample. Let  $\mathbf{v}$  denote the top eigenvector of  $\mathbf{C} := \mathbb{E}_{\mathcal{D}}[\mathbf{x}\mathbf{x}^\top]$  with associated eigenvalue  $\lambda$  and let  $\mathbf{u} \sim \mathcal{N}(0, \mathbf{I})$  be the initialization for Oja's algorithm. We are going to assume that  $\sum_{k=1}^T \|\mathbf{y}_k\| + \|\mathbf{y}_k\|^2 \leq \sqrt{T}$ . We note that our assumption implies that  $\sum_{k=1}^T \|\widehat{\mathbf{g}}_k - \mathbf{g}_k\| \leq O(\sqrt{T})$  and  $\sum_{k=1}^T \|\widehat{\mathbf{g}}_k\|^n \leq O(T^{(n/2)})$ . Let  $\Phi_k^{\mathbf{M}}$  and  $\Psi_k$  be defined as in the proof of Theorem 5.3. We are first going to upper bound  $\text{Tr}(\Phi_T^{\mathbf{u}\mathbf{u}^\top})$ :

$$\begin{aligned} \text{Tr}(\Phi_T^{\mathbf{u}\mathbf{u}^\top}) &= \text{Tr}(\Phi_{T-1}^{\mathbf{u}\mathbf{u}^\top}) + 2\eta \text{Tr}(\mathbf{x}_T \mathbf{x}_T^\top \Phi_{T-1}^{\mathbf{u}\mathbf{u}^\top}) \\ &\quad + 2\eta \text{Tr}((\widehat{\mathbf{g}}_T - \mathbf{g}_T) \Phi_{T-1}^{\mathbf{u}\mathbf{u}^\top}) + \eta^2 \text{Tr}(\widehat{\mathbf{g}}_T^2 \Phi_{T-1}^{\mathbf{u}\mathbf{u}^\top}) \\ &\leq \text{Tr}(\Phi_{T-1}^{\mathbf{u}\mathbf{u}^\top}) \left( 1 + 2\eta \mathbf{w}_T^\top \mathbf{g}_T \mathbf{w}_T + 2\eta \|\widehat{\mathbf{g}}_T - \mathbf{g}_T\| + \eta^2 \|\widehat{\mathbf{g}}_T\|^2 \right) \\ &\leq \text{Tr}(\Phi_{T-1}^{\mathbf{u}\mathbf{u}^\top}) \exp \left( 2\eta \mathbf{w}_T^\top \mathbf{g}_T \mathbf{w}_T + 2\eta \|\widehat{\mathbf{g}}_T - \mathbf{g}_T\| + \eta^2 \|\widehat{\mathbf{g}}_T\|^2 \right) \\ &\leq \dots \\ &\leq \|\mathbf{u}\|_2^2 \exp \left( \sum_{k=1}^T 2\eta \mathbf{w}_k^\top \mathbf{g}_k \mathbf{w}_k + 2\eta \|\widehat{\mathbf{g}}_k - \mathbf{g}_k\| + \eta^2 \|\widehat{\mathbf{g}}_k\|^2 \right), \end{aligned} \quad (16)$$

where we have used lemma 4 in (Warmuth & Kuzmin, 2006).

Next we are going to lower bound  $\mathbb{E}[\text{Tr}(\Psi_1)]$ :

$$\begin{aligned}
 \mathbb{E}[\text{Tr}(\Psi_1)] &\geq \mathbb{E}[\text{Tr}(\mathbf{v}\mathbf{v}^\top (\mathbf{I} + 2\eta\widehat{\mathbf{g}}_T) \Phi_{T-1}^1)] \\
 &= \mathbb{E}[\text{Tr}(\mathbf{v}\mathbf{v}^\top (\mathbf{I} + 2\eta\mathbf{C}) \Phi_{T-1}^1)] \\
 &+ \mathbb{E}[\text{Tr}(\mathbf{v}\mathbf{v}^\top (\mathbf{I} + 2\eta\mathbf{y}_T\mathbf{y}_T^\top) \Phi_{T-1}^1)] \\
 &\geq (1 + 2\eta\lambda)\mathbb{E}[\text{Tr}(\mathbf{v}\mathbf{v}^\top \Phi_{T-1}^1)] \\
 &\geq \exp(2\eta\lambda - 2\eta^2\lambda^2) \mathbb{E}[\text{Tr}(\mathbf{v}\mathbf{v}^\top \Phi_{T-1}^1)] \geq \dots \\
 &\geq \exp(T(2\eta\lambda - 2\eta^2\lambda^2)),
 \end{aligned} \tag{17}$$

where we used the fact that  $\mathbb{E}[\mathbf{x}_k^\top] = 0$  and that  $1 + 2x \geq \exp(2x - 2x^2)$  for  $x \in [0, 1]$ . Finally using Lieb-Thirring's inequality we have the following upper bound on  $\mathbb{E}[\text{Tr}(\Psi_1^2)]$ :

$$\begin{aligned}
 \mathbb{E}[\text{Tr}(\Psi_1^2)] &\leq \mathbb{E}[\text{Tr}((\mathbf{I} + \eta\widehat{\mathbf{g}}_1)^4 \Psi_2^2)] \\
 &= \mathbb{E}[\left(\mathbf{I} + 4\eta\mathbf{x}_1\mathbf{x}_1^\top + 4\eta\mathbf{y}_1\mathbf{y}_1^\top + 6\eta^2\widehat{\mathbf{g}}_1^2 + 4\eta^3\widehat{\mathbf{g}}_1^3\right. \\
 &\quad \left.+ \eta^4\widehat{\mathbf{g}}_1^4\right) \text{Tr}(\Psi_2^2)] \\
 &\leq \exp\left(4T\eta\lambda + 4\alpha_1\eta\sqrt{T} + 6\alpha_2\eta^2T + 4\alpha_3\eta^3T^{(3/2)}\right. \\
 &\quad \left.+ \alpha_4\eta^4T^2\right),
 \end{aligned} \tag{18}$$

where for the last inequality we have pushed the expectation inside the trace and used the fact that  $\mathbb{E}[\|\widehat{\mathbf{g}}_1^n\|]$  is bounded above by some linear combinations of  $\|\mathbf{y}_1\|^m$ , where  $m \in [2n]$ . Next we have used the fact that  $1 + x \leq \exp(x)$  and induction together with  $\sum_{k=1}^T \|\mathbf{y}_k\| \leq O(\sqrt{T})$  and  $\sum_{k=1}^T \|\mathbf{y}_k\|^{2n} \leq O(T^{n/2})$ . Let  $\alpha = 4\alpha_1 + 6\alpha_2\eta\sqrt{T} + 4\alpha_3\eta^2T + \alpha_4\eta^3T^{3/2} + 4\eta\lambda^2\sqrt{T}$ . Using Chebyshev's inequality with equations 17 and 18 we have:

$$\begin{aligned}
 \mathbb{P}[\text{Tr}(\Psi_1) \leq \exp(T(2\eta\lambda - 2\eta^2\lambda^2))] &(1) \\
 - \delta^{-1/2} \sqrt{\exp(\sqrt{T}\eta\alpha) - 1} &\leq \delta.
 \end{aligned}$$

Let  $\eta = \frac{\beta}{\sqrt{T}}$  for some constant  $\beta$  independent of  $T$ . We claim that for small enough  $\beta$  it holds that with probability at least  $1 - \delta$ ,  $\text{Tr}(\Psi_1) \geq \frac{2}{3}\exp(T(2\eta\lambda - 2\eta^2\lambda^2))$ . For the choice of  $\eta$ , this is equivalent to  $\exp(\beta\alpha) \leq 1 + \frac{\delta}{9}$ , where  $\alpha$  now becomes  $\alpha = 4\alpha_1 + 6\alpha_2\beta + 4\lambda^2\beta + 4\alpha_3\beta^2 + \alpha_4\beta^3$ . Taking  $\beta \rightarrow 0$  we can see that  $\exp(\alpha\beta) \rightarrow 1 < 1 + \delta/9$  and since  $\exp(\alpha\beta)$  is a continuous function of  $\beta$  the desired value of  $\beta > 0$  exists. From the derivation in (Allen-Zhu & Li, 2017) Appendix I, it follows that  $\|\mathbf{u}\|_2^2 \leq O(d + \log(1/\delta))$  and  $\text{Tr}(\Phi_T^{\text{uu}\top}) \geq \Omega(\delta^{-2})\text{Tr}(\Psi_1)$  with

probability  $1 - \delta$ . Thus with probability  $1 - 2\delta$  it holds that

$$\begin{aligned}
 O(d + \log(1/\delta)) \exp\left(\sum_{k=1}^T 2\eta\mathbf{w}_k^\top \widehat{\mathbf{g}}_k \mathbf{w}_k + 2\eta \|\widehat{\mathbf{g}}_k - \mathbf{g}_k\| \right. \\
 \left. + \eta^2 \|\widehat{\mathbf{g}}_k\|^2\right) &\geq \frac{2}{3\delta^2} \exp(T(2\eta\lambda - 2\eta^2\lambda^2)) \implies \\
 \sum_{k=1}^T \langle \widehat{\mathbf{g}}_k, \mathbf{P}^{(*)} - \mathbf{P}^k \rangle &\leq \frac{\log(O(d + \log(1/\delta)))}{\eta} \\
 + O(\sqrt{T}) + \eta O(T) &- \frac{\log(\frac{2}{3\delta^2})}{\eta}.
 \end{aligned}$$

Substituting  $\eta = \frac{\beta}{\sqrt{T}}$  finishes the proof.  $\square$

### C. Proofs of Section 4

**Lemma C.1.** The naive estimator  $\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top$  is not unbiased, i.e.

$$\mathbb{E}_{\mathcal{R}}[\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top | \mathbf{x}] = q^2\mathbf{x}\mathbf{x}^\top + (q - q^2) \text{diag } \mathbf{x}\mathbf{x}^\top,$$

where the expectation is taken with respect to Bernoulli random model.

*Proof of Lemma C.1.* For off-diagonal elements ( $i \neq j$ ),  $\mathbb{E}_{\mathcal{R}}[\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top | \mathbf{x}]_{i,j} = \mathbf{x}_i\mathbf{x}_j\mathbb{P}(i, j \in \{i_1, \dots, i_r\}) = \mathbf{x}_i\mathbf{x}_j\mathbb{P}(Z_i = 1)\mathbb{P}(Z_j = 1) = q^2\mathbf{x}_i\mathbf{x}_j$ . Now consider the diagonal elements.  $\mathbb{E}_{\mathcal{R}}[\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top | \mathbf{x}]_{i,i} = \mathbf{x}_i^2\mathbb{P}(Z_i = 1) = q\mathbf{x}_i^2$ . Putting together,  $\mathbb{E}_{\mathcal{R}}[\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top | \mathbf{x}] = q^2\mathbf{x}\mathbf{x}^\top + (q - q^2) \text{diag } \mathbf{x}\mathbf{x}^\top$ .  $\square$

*Proof of Lemma 4.1.* First note that

$$\mathbb{E}_{\mathcal{S}}[\mathbf{z}\mathbf{z}^\top | \tilde{\mathbf{x}}] = \frac{r - rq}{q^2} \mathbb{E}_{\mathcal{S}}[\tilde{x}_{i_s}^2 \mathbf{e}_{i_s} \mathbf{e}_{i_s}^\top | \tilde{\mathbf{x}}] = \frac{1 - q}{q^2} \text{diag } \tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top.$$

The proof simply follows from the following equalities.

$$\begin{aligned}
 \mathbb{E}_{\mathcal{S}, \mathcal{R}}[\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top - \mathbf{z}\mathbf{z}^\top | \mathbf{x}] & \\
 &= \mathbb{E}_{\mathcal{R}}[\mathbb{E}_{\mathcal{S}}[\frac{1}{q^2}\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top | \tilde{\mathbf{x}}] | \mathbf{x}] - \mathbb{E}_{\mathcal{R}}[\mathbb{E}_{\mathcal{S}}[\mathbf{z}\mathbf{z}^\top | \tilde{\mathbf{x}}] | \mathbf{x}] \\
 &= \frac{1}{q^2} \mathbb{E}_{\mathcal{R}}[\widehat{\mathbf{x}}\widehat{\mathbf{x}}^\top | \mathbf{x}] - \frac{1 - q}{q^2} \mathbb{E}_{\mathcal{R}}[\text{diag } \tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top | \mathbf{x}] \\
 &= \mathbf{x}\mathbf{x}^\top + \frac{q - q^2}{q^2} \text{diag } \mathbf{x}\mathbf{x}^\top - \frac{1 - q}{q} \text{diag } \mathbf{x}\mathbf{x}^\top \\
 &= \mathbf{x}\mathbf{x}^\top.
 \end{aligned} \tag{19}$$

$\square$

*Proof of Theorem 4.2.* We first bound  $\|\widehat{\mathbf{g}}_t\|_F^2$  as follows:

$$\begin{aligned}\|\widehat{\mathbf{g}}_t\|_F^2 &= \left\| \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top - \mathbf{z}_t \mathbf{z}_t^\top \right\|_F^2 \\ &= \left\| \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top \right\|_F^2 + \left\| \mathbf{z}_t \mathbf{z}_t^\top \right\|_F^2 - 2 \langle \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top, \mathbf{z}_t \mathbf{z}_t^\top \rangle \\ &\leq \left\| \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top \right\|_F^2 + \left\| \mathbf{z}_t \mathbf{z}_t^\top \right\|_F^2 \\ &\leq \|\widehat{\mathbf{x}}_t\|^4 + \|\mathbf{z}_t\|^4 \\ &\leq \frac{1}{q^2} \|\tilde{\mathbf{x}}\|^4 + \frac{r_t^2(1-q)^2}{q^4} \|\mathbf{x}\|_\infty^4 \\ &\leq \frac{1}{q^2} + \frac{r_t^2(1-q)^2}{q^4}\end{aligned}$$

where the first inequality follows from the fact that inner product of PSD matrices is non-negative, the third inequality is by definition of  $\widehat{\mathbf{x}}$  and  $\mathbf{z}$ , and the last inequality holds since  $\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\| \leq 1$ . Taking the expectation of both sides, noting the for Binomial distribution  $\mathbb{E}_{\mathcal{R}}[r_t^2] = \text{var}(r_t) + \mathbb{E}_{\mathcal{R}}[r_t]^2 = dq(1-q) + d^2q^2$ , we get  $\mathbb{E}_{\mathcal{R}}[\|\widehat{\mathbf{g}}\|_F^2] \leq \frac{1}{q^2} + \frac{dq(1-q)^3 + d^2q^2(1-q)^2}{q^4}$ . Now, using Lemma A.1 with  $\{\widehat{\mathbf{g}}_t\}_{t=1}^T$ , noting  $\|\mathbf{P}\|_F^2 \leq k$  for all  $\mathbf{P} \in \mathcal{C}$

$$\begin{aligned}\mathbb{E}\left[\sum_{t=1}^T \left\langle \mathbf{P}^* - \mathbf{P}^{(t)}, -\widehat{\mathbf{g}}_t \right\rangle\right] \\ \leq \frac{k}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \frac{q^2 + dq(1-q)^3 + d^2q^2(1-q)^2}{q^4} \\ \leq \frac{k}{2\eta} + \frac{\eta}{2} T \frac{q^2 + dq(1-q)^3 + d^2q^2(1-q)^2}{q^4}.\end{aligned}$$

Setting  $\eta = \frac{q^2}{\sqrt{q^2 + dq(1-q)^3 + d^2q^2(1-q)^2}} \sqrt{\frac{k}{T}}$  and taking expectation with respect to the internal randomization  $\mathcal{S}$  and the distribution of the missing data  $\mathcal{R}$  finishes the proof.  $\square$

*Proof of Theorem Oja with missing entries 4.3.* We begin by bounding  $\mathbb{E}_{\mathcal{S}, \mathcal{R}}[\|\widehat{\mathbf{g}}_k\|^2]$ ,  $\mathbb{E}_{\mathcal{S}, \mathcal{R}}[\|\widehat{\mathbf{g}}_k\|^3]$  and  $\mathbb{E}_{\mathcal{S}, \mathcal{R}}[\|\widehat{\mathbf{g}}_k\|^4]$ . Let  $a := \left\| \widehat{\mathbf{x}}_k \widehat{\mathbf{x}}_k^\top \right\|$  and  $b := \left\| \mathbf{z}_k \mathbf{z}_k^\top \right\|$ . Using Jensen's inequality we have  $\|\widehat{\mathbf{g}}_k\|^2 \leq 2a^2 + 2b^2$ ,  $\|\widehat{\mathbf{g}}_k\|^3 \leq 4a^3 + 4b^3$ ,  $\|\widehat{\mathbf{g}}_k\|^4 \leq 8a^4 + 8b^4$ . Since  $\|\mathbf{x}\|_2 \leq 1$  it holds that  $a^n \leq \frac{1}{q^n}$  and  $b^n \leq \frac{r_k^n(1-q)^n}{q^{2n}}$ . Let

$$\begin{aligned}\mu_2 &= dq(1-q) + dq \\ \mu_3 &= dq(1-3q+3dq+2q^2-3dq^2+d^2q^2) \\ \mu_4 &= dq(1-7q+7dq+12q^2-18dq^2+6d^2q^2 \\ &\quad -6q^3+11dq^3-6d^2q^3+d^3q^3).\end{aligned}$$

These are the second, third and fourth moment of the binomial random variable  $r_k$ , respectively. Also let  $\tilde{\alpha}_2^k =$

$\frac{2}{q^2} + \frac{2r_k^2(1-q)^2}{q^4}$ . Thus we have

$$\mathbb{E}_{\mathcal{S}, \mathcal{R}}[\|\widehat{\mathbf{g}}_k\|^n] \leq \alpha_n = \frac{2^{n-1}}{q^n} + \frac{2^{n-1}\mu_n(1-q)^n}{q^{2n}}.$$

Define  $\Phi_k^M$ ,  $\Psi_k$ ,  $\mathbf{w}_k$ ,  $\lambda, \nu$  and  $\mathbf{u}$  as in the proof of theorem 5.3. We have:

$$\begin{aligned}\text{Tr}\left(\Phi_T^{\mathbf{u}\mathbf{u}\top}\right) \\ &= \text{Tr}\left(\Phi_{T-1}^{\mathbf{u}\mathbf{u}\top}\right) + 2\eta \text{Tr}\left(\widehat{\mathbf{g}}_T \Phi_{T-1}^{\mathbf{u}\mathbf{u}\top}\right) + \eta^2 \text{Tr}\left(\widehat{\mathbf{g}}_T^2 \Phi_{T-1}^{\mathbf{u}\mathbf{u}\top}\right) \\ &\leq \text{Tr}\left(\Phi_{T-1}^{\mathbf{u}\mathbf{u}\top}\right) + 2\eta \text{Tr}\left(\widehat{\mathbf{g}}_T \Phi_{T-1}^{\mathbf{u}\mathbf{u}\top}\right) + \eta^2 \tilde{\alpha}_2^T \text{Tr}\left(\Phi_{T-1}^{\mathbf{u}\mathbf{u}\top}\right) \\ &\leq \text{Tr}\left(\Phi_{T-1}^{\mathbf{u}\mathbf{u}\top}\right) \exp\left(\eta^2 \tilde{\alpha}_2^T + 2\eta \mathbf{w}_T^\top \widehat{\mathbf{g}}_T \mathbf{w}_T\right) \\ \dots &\leq \|\mathbf{u}\|_2^2 \exp\left(\sum_{t=1}^T \eta^2 \tilde{\alpha}_2^t + 2\eta \sum_{t=1}^T \mathbf{w}_t^\top \widehat{\mathbf{g}}_t \mathbf{w}_t\right).\end{aligned}\tag{20}$$

Next we need a lower bound on  $\mathbb{E}[\Psi_1]$ . This is done in the same way as in Theorem 5.3:

$$\mathbb{E}[\text{Tr}(\Psi_1)] \geq \exp(T(2\eta\lambda - 2\eta^2\lambda^2)),\tag{21}$$

Finally, we need to bound  $\mathbb{E}[\Psi_1^2]$ . It holds that:

$$\begin{aligned}\mathbb{E}[\text{Tr}(\Psi_1^2)] &= \mathbb{E}[\text{Tr}((I + \eta\widehat{\mathbf{g}}_1)^2 \Psi_2 (I + \eta\widehat{\mathbf{g}}_1)^2 \Psi_2)] \\ &\leq \mathbb{E}[\text{Tr}((I + \eta\widehat{\mathbf{g}}_1)^4 \Psi_2^2)] \\ &= \mathbb{E}[\text{Tr}((I + 4\eta^3\widehat{\mathbf{g}}_1^3 + 6\eta^2\widehat{\mathbf{g}}_1^2 + 4\eta\widehat{\mathbf{g}}_1 + \eta^4\widehat{\mathbf{g}}_1^4) \Psi_2^2)] \\ &\leq \text{Tr}((I + \eta^4\alpha_4 + 4\eta^3\alpha_3 + 6\eta^2\alpha_2 + 4\eta\lambda) \Psi_2^2) \\ &\leq \text{Tr}((I + \eta^2\alpha + 4\eta\lambda) \mathbb{E}[\Psi_2^2]) \\ &\leq \exp(\alpha\eta^2 + 4\eta\lambda) \mathbb{E}[\text{Tr}(\Psi_2^2)],\end{aligned}\tag{22}$$

where  $\alpha = \alpha_4 + 4\alpha_3 + 6\alpha_2$ . Using Chebyshev's inequality it holds that

$$\begin{aligned}\mathbb{P}[\text{Tr}(\Psi_1) \leq \exp(T(2\eta\lambda - 2\eta^2\lambda^2))] \\ - \delta^{-1/2} \sqrt{\exp(\alpha\eta^2 + 4\eta^2\lambda^2)}] \leq \delta.\end{aligned}$$

As long as  $\eta \leq \frac{\log(1+\delta/9)}{\sqrt{T(\alpha+4\lambda^2)}}$ , this implies that with probability  $1 - \delta$ ,  $\text{Tr}(\Psi_1) \geq \frac{2}{3} \exp(T(2\eta\lambda - 2\eta^2\lambda^2))$ . From the derivation in (Allen-Zhu & Li, 2017) Appendix I, it follows that  $\|\mathbf{u}\|_2^2 \leq O(d + \log(1/\delta))$  and  $\text{Tr}(\Phi_T^{\mathbf{u}\mathbf{u}\top}) \geq \Omega(\delta^{-2}) \text{Tr}(\Psi_1)$  with probability  $1 - \delta$ . Thus with probability  $1 - 2\delta$  it holds that

$$\begin{aligned}O(d + \log(1/\delta)) \exp\left(\sum_{t=1}^T \eta^2 \tilde{\alpha}_2^t + 2\eta \sum_{t=1}^T \mathbf{w}_t^\top \widehat{\mathbf{g}}_t \mathbf{w}_t\right) \\ \geq \frac{2}{3\delta^2} \exp(T(2\eta\lambda - 2\eta^2\lambda^2)).\end{aligned}$$

Taking logarithms on both sides implies that

$$2\eta\lambda T - 2\eta \sum_{t=1}^T \text{Tr} \left( \mathbb{E}_{\mathcal{S}, \mathcal{R}}[\mathbf{w}_t \mathbf{w}_t^\top] \mathbb{E}_{\mathcal{S}, \mathcal{R}}[\widehat{\mathbf{g}}_t] \right) \leq$$

$$\log(O(d + \log(1/\delta))) - \log\left(\frac{2}{3\delta^2}\right) + \eta^2 \sum_{t=1}^T \mathbb{E}_{\mathcal{S}, \mathcal{R}}[\tilde{\alpha}_2^t],$$

where we have used the independence between  $\widehat{\mathbf{g}}_t$  and  $\mathbf{w}_t$ . This implies

$$\sum_{t=1}^T \langle \widehat{\mathbf{g}}_t, \mathbf{P}^{(*)} - \mathbb{E}_{\mathcal{S}, \mathcal{R}}[\mathbf{P}^{(t)}] \rangle$$

$$\frac{\log(O(d + \log(1/\delta))) - \log\left(\frac{2}{3\delta^2}\right)}{2\eta} + T\eta\alpha_2$$

Setting  $\eta = \frac{\log(1+\delta/9)}{\sqrt{T}(\alpha+4\lambda^2)}$  finishes the proof.  $\square$

## D. Proofs of Section 5

**Lemma D.1.** The naive estimator  $\tilde{x}_i \tilde{x}_i^\top$  is not unbiased, i.e.

$$\mathbb{E}_{\mathcal{R}}[\tilde{x} \tilde{x}^\top | \mathbf{x}] = \frac{r(r-1)}{d(d-1)} \mathbf{x} \mathbf{x}^\top + \frac{r(d-r)}{d(d-1)} \text{diag} \mathbf{x} \mathbf{x}^\top,$$

where expectation is taken with respect to the uniform sampling model.

*Proof of Lemma D.1.* For off-diagonal elements ( $i \neq j$ ),  $\mathbb{E}_{\mathcal{R}}[\tilde{x} \tilde{x}^\top | \mathbf{x}]_{i,j} = x_i x_j \mathbb{P}(i, j \in \{i_1, \dots, i_r\})$  so we need to compute the probability of the event that two fixed element  $i, j \in [1 \dots d]$  are both in the subset  $\{i_1, \dots, i_r\}$  sampled uniformly at random from all subsets of size  $r$ . The number of sets of  $r$  elements which contain  $i$  and  $j$  is exactly  $\binom{d-2}{r-2}$  and the total number of  $r$  sets is  $\binom{d}{r}$  thus  $\mathbb{P}(i, j \in \{i_1, \dots, i_r\}) = \frac{\binom{d-2}{r-2}}{\binom{d}{r}} = \frac{r(r-1)}{d(d-1)}$  and so  $\mathbb{E}_{\mathcal{R}}[\tilde{x} \tilde{x}^\top | \mathbf{x}]_{i,j} = x_i x_j \frac{r(r-1)}{d(d-1)}$ . Now consider the diagonal elements.  $\mathbb{E}_{\mathcal{R}}[\tilde{x} \tilde{x}^\top | \mathbf{x}]_{i,i} = x_i^2 \mathbb{P}(i \in \{i_1, \dots, i_r\})$  so we need to compute the probability of the event  $i \in \{i_1, \dots, i_r\}$ . The number of sets with size  $r$  containing  $i$  is  $\binom{d-1}{r-1}$  and thus  $\mathbb{E}_{\mathcal{R}}[\tilde{x} \tilde{x}^\top | \mathbf{x}]_{i,i} = x_i^2 \frac{r}{d}$ . Putting together,

$$\mathbb{E}_{\mathcal{R}}[\tilde{x} \tilde{x}^\top | \mathbf{x}] = \frac{r(r-1)}{d(d-1)} \mathbf{x} \mathbf{x}^\top + \left( \frac{r}{d} - \frac{r(r-1)}{d(d-1)} \right) \text{diag} \mathbf{x} \mathbf{x}^\top$$

$$= \frac{r(r-1)}{d(d-1)} \mathbf{x} \mathbf{x}^\top + \frac{r(d-r)}{d(d-1)} \text{diag} \mathbf{x} \mathbf{x}^\top.$$

$\square$

*Proof of Lemma 5.1.* First note that

$$\mathbb{E}_{\mathcal{S}}[\mathbf{z} \mathbf{z}^\top | \tilde{\mathbf{x}}] = \frac{dr - r^2}{r-1} \mathbb{E}_{\mathcal{S}}[\tilde{x}_{i_s}^2 \mathbf{e}_{i_s} \mathbf{e}_{i_s}^\top | \tilde{\mathbf{x}}] = \frac{d-r}{r-1} \text{diag} \tilde{x} \tilde{x}^\top.$$

The proof simply follows from the following equalities.

$$\mathbb{E}_{\mathcal{S}, \mathcal{R}}[\widehat{\mathbf{x}} \widehat{\mathbf{x}}^\top - \mathbf{z} \mathbf{z}^\top | \mathbf{x}]$$

$$= \mathbb{E}_{\mathcal{R}}[\mathbb{E}_{\mathcal{S}}[\widehat{\mathbf{x}} \widehat{\mathbf{x}}^\top | \tilde{\mathbf{x}}] | \mathbf{x}] - \mathbb{E}_{\mathcal{R}}[\mathbb{E}_{\mathcal{S}}[\mathbf{z} \mathbf{z}^\top | \tilde{\mathbf{x}}] | \mathbf{x}]$$

$$= \frac{d(d-1)}{r(r-1)} \mathbb{E}_{\mathcal{R}}[\tilde{x} \tilde{x}^\top | \mathbf{x}] - \frac{d-r}{r-1} \mathbb{E}_{\mathcal{R}}[\text{diag} \tilde{x} \tilde{x}^\top] \quad (23)$$

$$= \mathbf{x} \mathbf{x}^\top + \frac{d-r}{r-1} \text{diag} \mathbf{x} \mathbf{x}^\top - \frac{d-r}{r-1} \text{diag} \mathbf{x} \mathbf{x}^\top$$

$$= \mathbf{x} \mathbf{x}^\top.$$

where  $\text{diag} \mathbf{A}$ ,  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is the  $d \times d$  matrix consisting of the diagonal of  $\mathbf{A}$ .  $\square$

*Proof of Theorem 5.2.* We first bound  $\|\widehat{\mathbf{g}}_t\|_F^2$  as follows:

$$\|\widehat{\mathbf{g}}_t\|_F^2 = \left\| \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top - \mathbf{z}_t \mathbf{z}_t^\top \right\|_F^2$$

$$= \left\| \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top \right\|_F^2 + \left\| \mathbf{z}_t \mathbf{z}_t^\top \right\|_F^2 - 2 \langle \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top, \mathbf{z}_t \mathbf{z}_t^\top \rangle$$

$$\leq \left\| \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top \right\|_F^2 + \left\| \mathbf{z}_t \mathbf{z}_t^\top \right\|_F^2$$

$$\leq \|\widehat{\mathbf{x}}_t\|^4 + \|\mathbf{z}_t\|^4$$

$$\leq \frac{d^2(d-1)^2}{r^2(r-1)^2} + \frac{r^2(d-r)^2}{(r-1)^2} \|\mathbf{x}\|_\infty^4$$

$$\leq \frac{d^2(d-1)^2 + r^4(d-r)^2}{r^2(r-1)^2}.$$

where the first inequality follows from  $\langle \widehat{\mathbf{x}}_t \widehat{\mathbf{x}}_t^\top, \mathbf{z}_t \mathbf{z}_t^\top \rangle \geq 0$ , since this is an inner product of PSD matrices, the third inequality follows the definition of  $\widehat{\mathbf{x}}$  and  $\mathbf{z}$ , and the fourth inequality holds since  $\|\tilde{\mathbf{x}}\| \leq \|\mathbf{x}\| \leq 1$ . Now, using Lemma A.1 with  $\{\widehat{\mathbf{g}}_t\}_{t=1}^T$ , noting  $\|\mathbf{P}\|_F^2 \leq k$  for all  $\mathbf{P} \in \mathcal{C}$ , we have that

$$\sum_{t=1}^T \left\langle \mathbf{P}^* - \mathbf{P}^{(t)}, -\widehat{\mathbf{g}}_t \right\rangle$$

$$\leq \frac{k}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \frac{d^2(d-1)^2 + r^4(d-r)^2}{r^2(r-1)^2}$$

$$= \frac{k}{2\eta} + \frac{\eta}{2} T \frac{d^2(d-1)^2 + r^4(d-r)^2}{r^2(r-1)^2}.$$

Since  $\mathbb{E}_{\widehat{\mathbf{g}}_t}[\widehat{\mathbf{g}}_t] = \mathbf{x}_t \mathbf{x}_t^\top$ , setting  $\eta = \frac{r(r-1)}{\sqrt{d^2(d-1)^2 + r^4(d-r)^2}} \sqrt{\frac{k}{2T}}$  and taking expectation with respect to the internal randomization  $\mathcal{S}$  and the distribution of the missing data  $\mathcal{R}$  finishes the proof.  $\square$

*Proof of Theorem 5.3.* First we bound  $\|\widehat{\mathbf{g}}_k\|_2$ . Since  $\widehat{\mathbf{g}}_k = \widehat{\mathbf{x}}_k \widehat{\mathbf{x}}_k^\top - \mathbf{z}_k \mathbf{z}_k^\top$  is the difference of two PSD matrices it holds that

$$\|\widehat{\mathbf{g}}_k\| = \max \left( \left\| \widehat{\mathbf{x}}_k \widehat{\mathbf{x}}_k^\top \right\|, \left\| \mathbf{z}_k \mathbf{z}_k^\top \right\| \right) \leq \alpha. \quad (24)$$

Let  $\mathbf{v}$  be the top eigenvector of  $\mathbf{C}$  with associated eigenvalue  $\lambda$ . Also let  $\mathbf{g}_k = \mathbf{x}_k \mathbf{x}_k^\top$  and  $\mathbf{w}_k$  be the  $k$ -th iterate of Oja's algorithm, such that  $\mathbf{P}^{(k)} = \mathbf{w}_k \mathbf{w}_k^\top$ . Define

$$\begin{aligned}\Phi_k^{\mathbf{M}} &= (\mathbf{I} + \eta \widehat{\mathbf{g}}_k) \cdots (\mathbf{I} + \eta \widehat{\mathbf{g}}_1) \mathbf{M} (\mathbf{I} + \eta \widehat{\mathbf{g}}_1) \cdots (\mathbf{I} + \eta \widehat{\mathbf{g}}_k) \\ \Psi_k &= (\mathbf{I} + \eta \widehat{\mathbf{g}}_k) \cdots (\mathbf{I} + \eta \widehat{\mathbf{g}}_T) \mathbf{v} \mathbf{v}^\top (\mathbf{I} + \eta \widehat{\mathbf{g}}_T) \cdots (\mathbf{I} + \eta \widehat{\mathbf{g}}_k).\end{aligned}$$

Following the calculations in (Allen-Zhu & Li, 2017) Appendix I, we have:

$$\begin{aligned}\text{Tr}(\Phi_T^{\mathbf{uu}^\top}) &= \text{Tr}((\mathbf{I} + \eta \widehat{\mathbf{g}}_T) \Phi_{T-1}^{\mathbf{uu}^\top} (\mathbf{I} + \eta \widehat{\mathbf{g}}_T)) \\ &= \text{Tr}(\Phi_{T-1}^{\mathbf{uu}^\top}) + 2\eta \text{Tr}(\widehat{\mathbf{g}}_T \Phi_{T-1}^{\mathbf{uu}^\top}) + \eta^2 \text{Tr}(\widehat{\mathbf{g}}_T \Phi_{T-1}^{\mathbf{uu}^\top} \widehat{\mathbf{g}}_T) \\ &= \text{Tr}(\Phi_{T-1}^{\mathbf{uu}^\top}) + 2\eta \text{Tr}(\widehat{\mathbf{g}}_T \Phi_{T-1}^{\mathbf{uu}^\top}) + \eta^2 \text{Tr}(\widehat{\mathbf{g}}_T^2 \Phi_{T-1}^{\mathbf{uu}^\top}) \\ &\leq \text{Tr}(\Phi_{T-1}^{\mathbf{uu}^\top}) + 2\eta \text{Tr}(\widehat{\mathbf{g}}_T \Phi_{T-1}^{\mathbf{uu}^\top}) + \eta^2 \alpha^2 \text{Tr}(\Phi_{T-1}^{\mathbf{uu}^\top}) \\ &= \text{Tr}(\Phi_{T-1}^{\mathbf{uu}^\top}) (1 + \eta^2 \alpha^2 + 2\eta \mathbf{w}_T^\top \widehat{\mathbf{g}}_T \mathbf{w}_T) \\ &\leq \text{Tr}(\Phi_{T-1}^{\mathbf{uu}^\top}) \exp(\eta^2 \alpha^2 + 2\eta \mathbf{w}_T^\top \widehat{\mathbf{g}}_T \mathbf{w}_T) \\ \cdots &\leq \|\mathbf{u}\|_2^2 \exp\left(T\eta^2 \alpha^2 + 2\eta \sum_{k=1}^T \mathbf{w}_k^\top \widehat{\mathbf{g}}_k \mathbf{w}_k\right),\end{aligned}\tag{25}$$

where for the first inequality we have used the fact that  $\Phi_{T-1}^{\mathbf{uu}^\top}$  is PSD and  $\|\widehat{\mathbf{g}}_k\|_2 \leq \alpha$  and for the second inequality we have used  $1 + x \leq \exp(x)$ . We have also used the fact that  $\Phi_{T-1}^{\mathbf{uu}^\top} = \text{Tr}(\Phi_{T-1}^{\mathbf{uu}^\top}) \mathbf{w}_T \mathbf{w}_T^\top$ . Next we have

$$\begin{aligned}\mathbb{E}[\Psi_1] &= \mathbb{E}[\text{Tr}(\mathbf{v} \mathbf{v}^\top (\mathbf{I} + \eta \widehat{\mathbf{g}}_T) \Phi_{T-1}^{\mathbf{I}} (\mathbf{I} + \eta \widehat{\mathbf{g}}_T))] \\ &= \mathbb{E}[\text{Tr}(\mathbf{v} \mathbf{v}^\top (\mathbf{I} + \eta \widehat{\mathbf{g}}_T) \Phi_{T-1}^{\mathbf{I}}) + \eta^2 \mathbf{v}^\top \widehat{\mathbf{g}}_T \Phi_{T-1}^{\mathbf{I}} \widehat{\mathbf{g}}_T \mathbf{v}] \\ &\geq \mathbb{E}[\text{Tr}(\mathbf{v} \mathbf{v}^\top (\mathbf{I} + \eta \mathbf{C}) \Phi_{T-1}^{\mathbf{I}})] \\ &= (1 + 2\eta\lambda) \mathbb{E}[\mathbf{v}^\top \Phi_{T-1}^{\mathbf{I}} \mathbf{v}] \\ &\geq \exp(2\eta\lambda - 2\eta^2 \lambda^2) \mathbb{E}[\mathbf{v}^\top \Phi_{T-1}^{\mathbf{I}} \mathbf{v}] \geq \cdots \\ &\geq \exp((2\eta\lambda - 2\eta^2 \lambda^2)T),\end{aligned}\tag{26}$$

where the expectation is both with respect to  $\mathcal{R}$  and  $\mathcal{D}$  and we have used  $1 + 2x \geq \exp(2x - 2x^2)$ , for  $0 \leq x \leq 1$ . Finally we have:

$$\begin{aligned}\mathbb{E}[\text{Tr}(\Psi_1^2)] &= \mathbb{E}[\text{Tr}((\mathbf{I} + \eta \widehat{\mathbf{g}}_1)^2 \Psi_2 (\mathbf{I} + \eta \widehat{\mathbf{g}}_1)^2 \Psi_2)] \\ &\leq \mathbb{E}[\text{Tr}((\mathbf{I} + \eta \widehat{\mathbf{g}}_1)^4 \Psi_2^2)] \\ &= \mathbb{E}[\text{Tr}((\mathbf{I} + 4\eta^3 \widehat{\mathbf{g}}_1^3 + 6\eta^2 \widehat{\mathbf{g}}_1^2 + 4\eta \widehat{\mathbf{g}}_1 + \eta^4 \widehat{\mathbf{g}}_1^4) \Psi_2^2)] \\ &\leq \mathbb{E}[\text{Tr}((\mathbf{I} + 11\eta^2 \alpha^4 \mathbf{I} + 4\eta \mathbf{g}_1) \Psi_2^2)] \\ &\leq \exp(11\eta^2 \alpha^4 + 4\eta\lambda) \mathbb{E}[\text{Tr}(\Psi_2^2)] \leq \cdots \\ &\leq \exp(T(11\eta^2 \alpha^4 + 4\eta\lambda)).\end{aligned}\tag{27}$$

Combining equations 26 and 27 with Chebyshev's inequality

we get:

$$\begin{aligned}\mathbb{P}[\text{Tr}(\Psi_1) \leq \exp(T(2\eta\lambda - 2\alpha^2 \eta^2)) - \\ \exp(T(2\eta\lambda - 2\alpha^2 \eta^2)) \sqrt{\frac{(\exp(T(11\eta^2 \alpha^4 + 4\eta^2 \lambda^2)) - 1)}{\delta}}] \\ \leq \delta.\end{aligned}$$

Thus with probability  $1 - \delta$  it holds

$$\begin{aligned}\text{Tr}(\Psi_1) &\geq \exp(T(2\eta\lambda - 2\alpha^2 \eta^2)) (1 - \\ &\quad \delta^{-1/2} \sqrt{\exp(T(11\eta^2 \alpha^4 + 4\eta^2 \lambda^2)) - 1}) \\ &\geq \frac{2}{3} \exp(T(2\eta\lambda - 2\alpha^2 \eta^2))\end{aligned}$$

as long as  $\eta \leq \frac{\log(1+\delta/9)}{(11\alpha^2+4\lambda^2)\sqrt{T}}$ . Following the derivations (Allen-Zhu & Li, 2017), we have that with probability  $1 - \delta$  it holds that  $\text{Tr}(\Phi_T^{\mathbf{uu}^\top}) \geq \Omega(\delta^{-2}) \text{Tr}(\Psi_1)$  and that  $\|\mathbf{u}\|_2^2 \leq O(d + \log(1/\delta^2))$ . Union bound, together with inequality 25 implies that with probability  $1 - 2\delta$  we have:

$$\begin{aligned}O(d + \log(1/\delta^2)) \exp\left(T\eta^2 \alpha^2 + 2\eta \sum_{k=1}^T \mathbf{w}_k^\top \widehat{\mathbf{g}}_k \mathbf{w}_k\right) \\ \geq \frac{2}{3\delta^2} \exp(T(2\eta\lambda - 2\alpha^2 \eta^2)) \implies \\ \log(O(d + \log(1/\delta^2))) + 2\eta \sum_{k=1}^T \mathbf{w}_k^\top \widehat{\mathbf{g}}_k \mathbf{w}_k \\ \geq \log(2/3\delta^2) + 2T\eta\lambda - 2T\alpha^2 \eta^2 \implies \\ \sum_{k=1}^T \lambda - \mathbf{w}_k^\top \widehat{\mathbf{g}}_k \mathbf{w}_k \leq \\ \frac{\log(O(d + \log(1/\delta^2))) - \log(2/3\delta^2)/2}{\eta} + 2T\alpha^2 \eta \implies \\ \sum_{k=1}^T \langle \mathbf{g}_k, \mathbf{P}^* - \mathbb{E}_{\mathcal{S}, \mathcal{R}}[\mathbf{P}_k] \rangle \leq \\ \frac{\sqrt{T}(11\alpha^2 + 4\lambda^2)}{\log(1 + \delta/9)} (\log(O(d + \log(1/\delta^2))) - \log(2/3\delta^2)/2) \\ + 2 \frac{\sqrt{T}}{\log(1 + \delta/9)},\end{aligned}$$

where in the last implication we have substituted  $\eta = \frac{\log(1+\delta/9)}{(11\alpha^2+4\lambda^2)\sqrt{T}}$  and used the fact that  $\alpha^2 < 11\alpha^2 + 4\lambda^2$ .  $\square$

## E. Proofs of Section 6

*Proof of Theorem 6.1.* Denoting  $(\mathbf{P}^{(t)} - \mathbf{I})\mathbf{x}_t$  by  $\mathbf{y}_t$ , we can rewrite  $\mathbf{g}_t$  as  $\mathbf{g}_t = -\frac{\mathbf{x}_t\mathbf{y}_t^\top + \mathbf{y}_t\mathbf{x}_t^\top}{2\|\mathbf{y}_t\|_2}$ . Thus we have

$$\begin{aligned}\|g_t\|_F^2 &= \frac{1}{4\|\mathbf{y}_t\|_2^2} \text{Tr} \left( (\mathbf{x}_t\mathbf{y}_t^\top + \mathbf{y}_t\mathbf{x}_t^\top)^\top (\mathbf{x}_t\mathbf{y}_t^\top + \mathbf{y}_t\mathbf{x}_t^\top) \right) \\ &\leq \frac{4\|\mathbf{x}_t\|_2^2\|\mathbf{y}_t\|_2^2}{4\|\mathbf{y}_t\|_2^2} = \|\mathbf{x}_t\|_2^2 \leq 1.\end{aligned}$$

By convexity of norms,  $\|\mathbf{x} - \mathbf{P}^{(t)}\mathbf{x}\|_2 - \|\mathbf{x} - \mathbf{P}^*\mathbf{x}\|_2 \leq \langle \mathbf{P}^{(t)} - \mathbf{P}^*, \mathbf{g}_t \rangle$ . Using Lemma A.1, noting  $\|\mathbf{P}\|_F^2 \leq k$  for all  $\mathbf{P} \in C$ , and  $\|\mathbf{g}_t\|_F^2 \leq 1$ ,

$$\sum_{t=1}^T \|\mathbf{x}_t - \mathbf{P}^{(t)}\mathbf{x}_t\|_2 - \sum_{t=1}^T \|\mathbf{x}_t - \mathbf{P}^*\mathbf{x}_t\|_2 \leq \frac{k}{2\eta} + \frac{\eta T}{2}.$$

Choosing  $\eta = \sqrt{\frac{k}{T}}$  completes the proof.  $\square$

## F. Proofs of Section 8

*Proof of Theorem 8.1.* Our proof follows (Zinkevich, 2003). We start the analysis by bounding the distance of our iterates at time  $t$  from any dynamic competitor  $\mathbf{P}_*^{(t-1)}$  at time  $t-1$

$$\begin{aligned}\|\mathbf{P}^{(t+1)} - \mathbf{P}_*^{(t)}\|_F^2 &= \|\Pi_C \left( \mathbf{P}^{(t)} - \eta\mathbf{g}_t \right) - \mathbf{P}_*^{(t)}\|_F^2 \\ &\leq \|\mathbf{P}^{(t)} - \eta\mathbf{g}_t - \mathbf{P}_*^{(t)}\|_F^2 \\ &= \|\mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)} + \mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)} - \eta\mathbf{g}_t\|_F^2 \\ &= \|\mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)}\|_F^2 + \|\mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)}\|_F^2 \\ &\quad + 2\langle \mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)}, \mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)} \rangle + \eta^2 \|\mathbf{g}_t\|_F^2 \\ &\quad - 2\eta \langle \mathbf{g}_t, \mathbf{P}^{(t)} - \mathbf{P}_*^{(t)} \rangle \\ &\leq \|\mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)}\|_F^2 + \|\mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)}\|_F^2 \\ &\quad + 2\|\mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)}\|_F \cdot \|\mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)}\|_F + \eta^2 \\ &\quad + 2\eta \left( \mathbf{x}_t^\top \mathbf{P}^{(t)} \mathbf{x}_t - \mathbf{x}_t^\top \mathbf{P}_*^{(t)} \mathbf{x}_t \right) \\ &\leq \|\mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)}\|_F^2 + 2\sqrt{k} \|\mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)}\|_F \\ &\quad + 4\sqrt{k} \|\mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)}\|_F + \eta^2 \\ &\quad + 2\eta \left( \mathbf{x}_t^\top \mathbf{P}^{(t)} \mathbf{x}_t - \mathbf{x}_t^\top \mathbf{P}_*^{(t)} \mathbf{x}_t \right)\end{aligned}$$

Rearranging and dividing both sides by  $2\eta$  we get

$$\begin{aligned}\mathbf{x}_t^\top \mathbf{P}_*^{(t)} \mathbf{x}_t - \mathbf{x}_t^\top \mathbf{P}^{(t)} \mathbf{x}_t &\leq \frac{3\sqrt{k}}{\eta} \|\mathbf{P}_*^{(t-1)} - \mathbf{P}_*^{(t)}\|_F + \frac{\eta}{2} \\ &\quad + \frac{\|\mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)}\|_F^2 - \|\mathbf{P}^{(t+1)} - \mathbf{P}_*^{(t)}\|_F^2}{2\eta}\end{aligned}$$

Summing over  $t$ , observe the telescopic sum on the right hand side is bounded by

$$\begin{aligned}\sum_{t=1}^T \left( \|\mathbf{P}^{(t)} - \mathbf{P}_*^{(t-1)}\|_F^2 - \|\mathbf{P}^{(t+1)} - \mathbf{P}_*^{(t)}\|_F^2 \right) \\ = \|\mathbf{P}^{(1)} - \mathbf{P}_*^{(0)}\|_F^2 - \|\mathbf{P}^{(T+1)} - \mathbf{P}_*^{(T)}\|_F^2 \\ \leq \|\mathbf{P}_*^{(0)}\|_F^2 \leq k\end{aligned}$$

Plugging the definition of the total shift  $S$  into the sum

$$\begin{aligned}\sum_{t=1}^T \mathbf{x}_t^\top \mathbf{P}_*^{(t)} \mathbf{x}_t - \sum_{t=1}^T \mathbf{x}_t^\top \mathbf{P}^{(t)} \mathbf{x}_t &\leq \frac{3\sqrt{k}S}{\eta} + \frac{\eta T}{2} + \frac{k}{2\eta} \\ &= \frac{6\sqrt{k}S + k}{2\eta} + \frac{\eta T}{2}\end{aligned}$$

Choosing  $\eta = \sqrt{\frac{6\sqrt{k}S+k}{T}}$  completes the proof.  $\square$

## F.1. Experimental Results

We provide additional experiments in Figures 3, 4, 5 and 6 for PCA with partial observations and for PCA with missing entries. As stated in the main text, all our observations from Figure 1 and Figure 2 hold for the results presented here.

Further, in Figure 7, we present experimental results for the proposed *Absolute Subspace Deviation Model* algorithm of Section 6 (referred to as *Robust GD* in the plots). The algorithms against which we compare are *Robust Online PCA* proposed in (Goes et al., 2014) (referred to as *Robust MEG* in the plots) which is an instance of online mirror descent with choice of potential function being entropy and simple batch PCA on the whole dataset. As ground truth we take the  $k$ -dimensional subspace returned by the proposed method in (Lerman et al., 2012). In the following discussion the ground truth is represented by the projection matrix  $\mathbf{P}^* \in \mathbb{R}^{d \times d}$  and the estimated projection matrix returned by an algorithm is denoted by  $\mathbf{P}^{(t)}$ . Since our experiments are ran on a synthetic dataset of fixed size, we denote this dataset by  $\mathbf{X} \in \mathbb{R}^{d \times n}$ . The criteria against which we evaluate are average angle between subspaces given by  $\text{Tr} \left( (\mathbf{P}^{(t)})^\top \mathbf{P}^* \right) / k$ , reconstruction error on the whole dataset given by  $\left\| (\mathbf{U}^{(t)})^\top \mathbf{X} - (\mathbf{U}^*)^\top \mathbf{X} \right\|_F / n$  where a rank- $k$  projection matrix  $\mathbf{P}$  is decomposed as  $\mathbf{P} = \mathbf{U}\mathbf{U}^\top$  for orthogonal  $\mathbf{U} \in \mathbb{R}^{d \times k}$ . We also evaluate on the total regret incurred. The data set is generated as in (Goes et al., 2014). First inliers are sampled from a  $k$ -dimensional Multivariate Normal distribution with 0 mean and covariance matrix  $\frac{1}{k}\mathbf{I}_k$ . Next the inliers are embedded in a  $d$ -dimensional space via the linear transformation  $\mathbf{U} \in \mathbb{R}^{d \times k}$  with entries  $U_{ii} = 1$  and  $U_{i \neq j} = 0$ . Finally outliers are

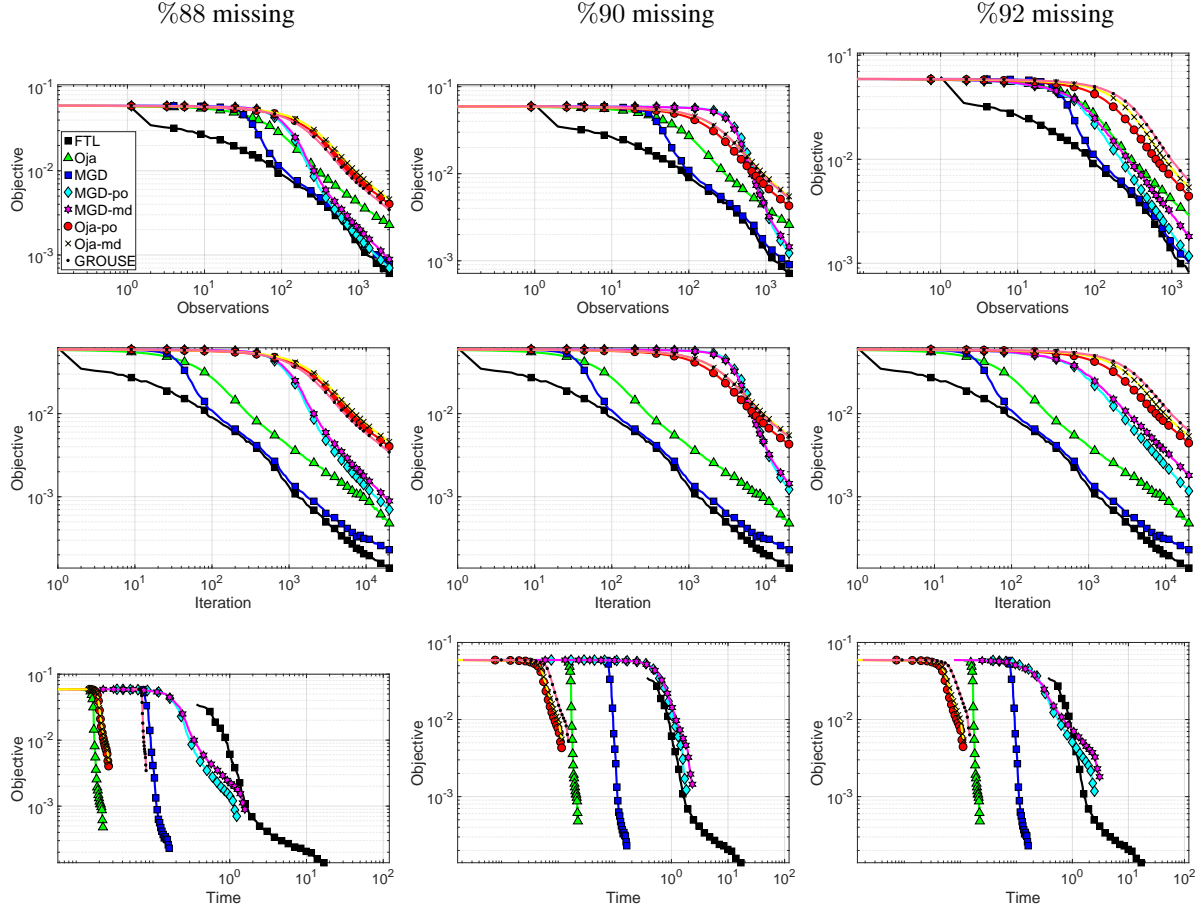


Figure 3: Comparisons of Oja, MGD, MGD-PO, MGD-MD, Oja-PO, Oja-MD, GROUSE for PCA with missing data on **XRMB** dataset, in terms of the variance captured on a test set as a function of number of observed entries for  $k=2$  (top), number of iterations (middle) and runtime (bottom).

sampled from a  $d$ -dimensional Multivariate Normal distribution with diagonal covariance  $\frac{\lambda}{d \times \text{mean}(\Sigma_{\text{out}})} \Sigma_{\text{out}}$  where  $(\Sigma_{\text{out}})_{ii} = i$ , here  $\lambda$  is a user-specified parameter which governs the SNR. In our experiments  $\lambda = 20$  the number of outliers is 40%, 60%, 80% of the total number of points,  $k = 2$  and  $d = 100$ . In our comparison we also include a “capped” version of our proposed method where the rank of each intermediate iterate is hard-capped to 15 by keeping the 15 directions associated with the top 15 singular values of our current iterate.

We now briefly discuss the more interesting points which should be observed from the provided Figure 7. First our algorithm clearly out-performs the *Robust Online PCA* and batch PCA in all the chosen criteria. Secondly, since *Robust Online PCA* is required to always keep full-rank iterates as it operates in the complementary space to the one we wish to recover it has to pay computational cost  $O(d^3)$  per iteration while our method only pays  $O(d\bar{k}^2)$  where  $\bar{k}$  is the rank of the current iterate. As can be seen from the plots

for the cases when number of outliers are 40% or 60% the intermediate rank of iterates does not grow too much so in practice our proposed method is efficient. Finally we note that the capped version of the proposed method performs as well or even better in some cases.

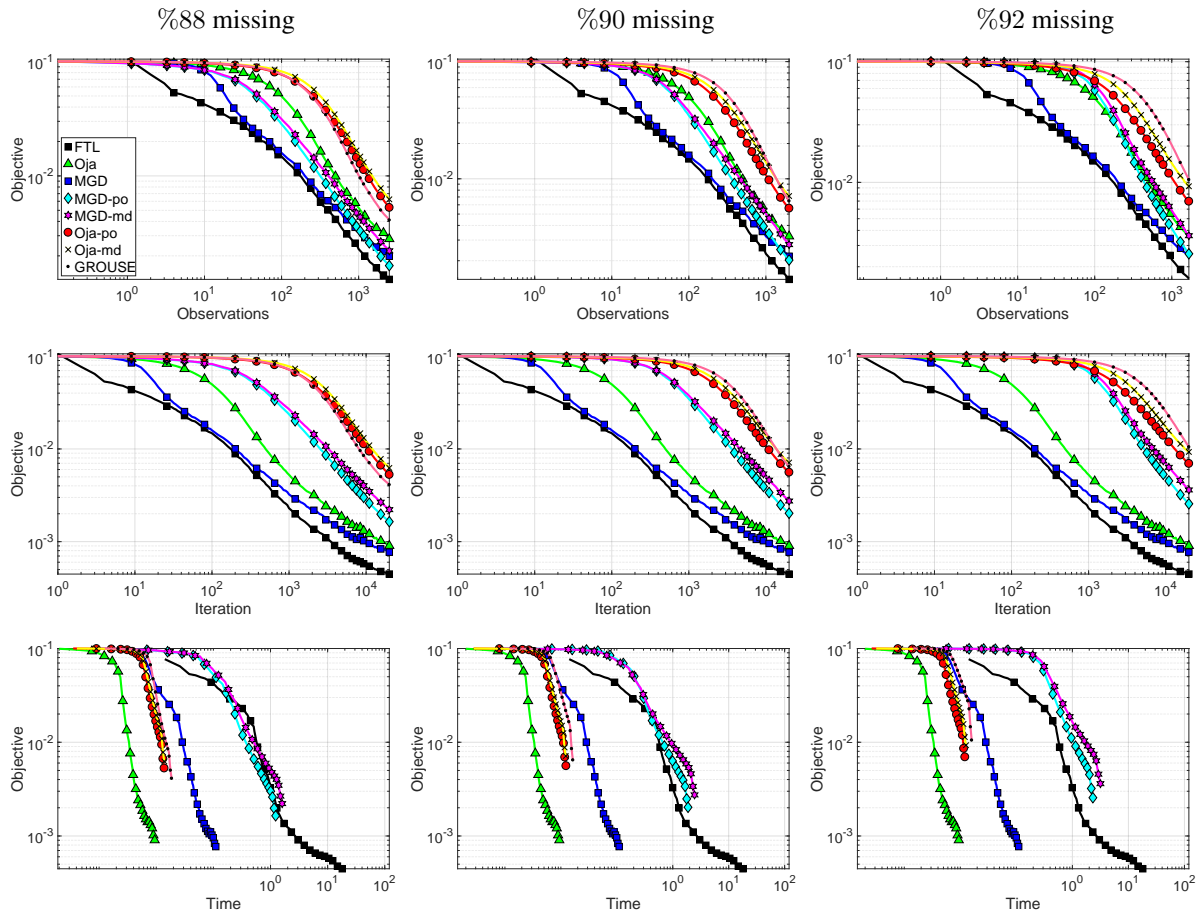


Figure 4: Comparisons of Oja, MGD, MGD-PO, MGD-MD, Oja-PO, Oja-MD, GROUSE for PCA with missing data on **XRMB** dataset, in terms of the variance captured on a test set as a function of number of observed entries for  $k=4$  (top), number of iterations (middle) and runtime (bottom).



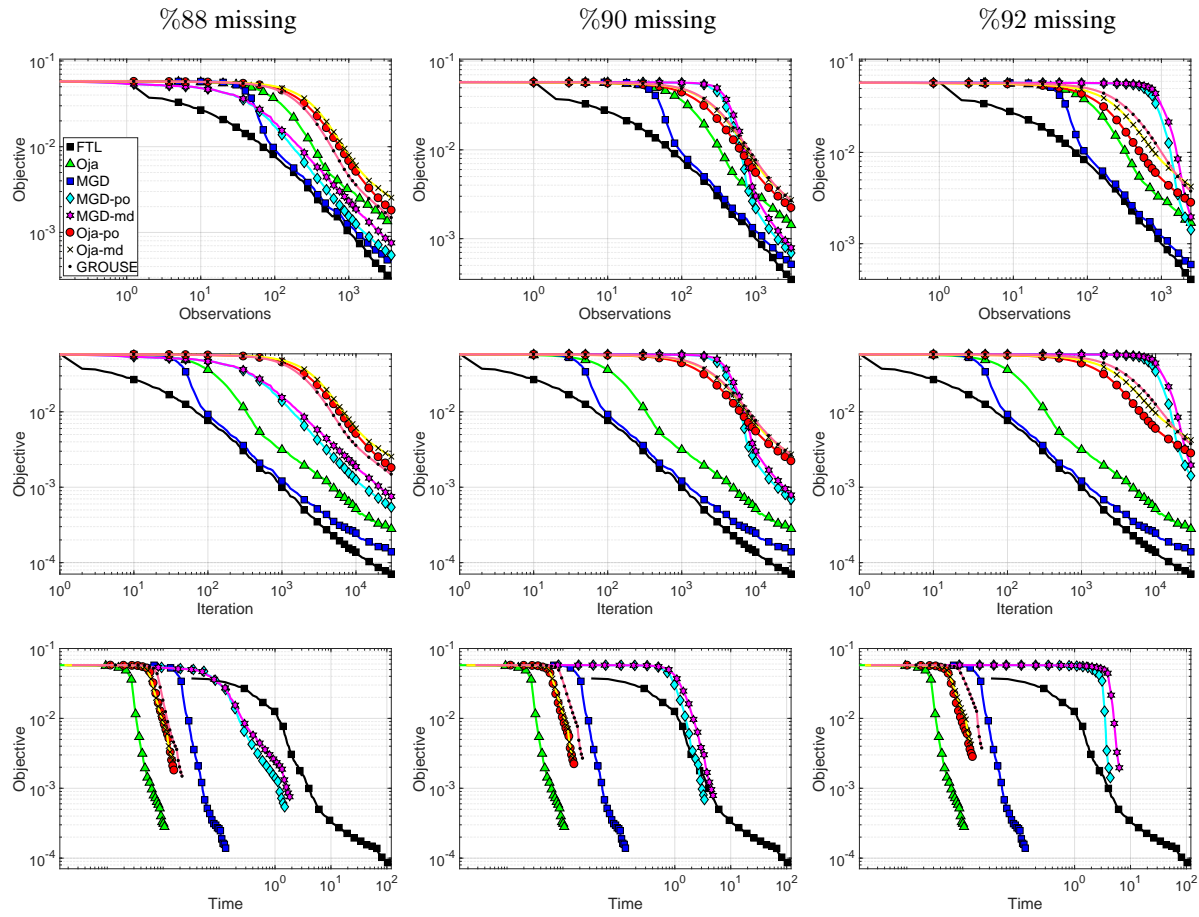


Figure 5: Comparisons of Oja, MGD, MSG-PO, MSG-MD, Oja-PO, Oja-MD and GROUSE for PCA with missing data on MNIST dataset, in terms of the variance captured on a test set as a function of number of observed entries for  $k=2$  (top), number of iterations (middle) and runtime (bottom).

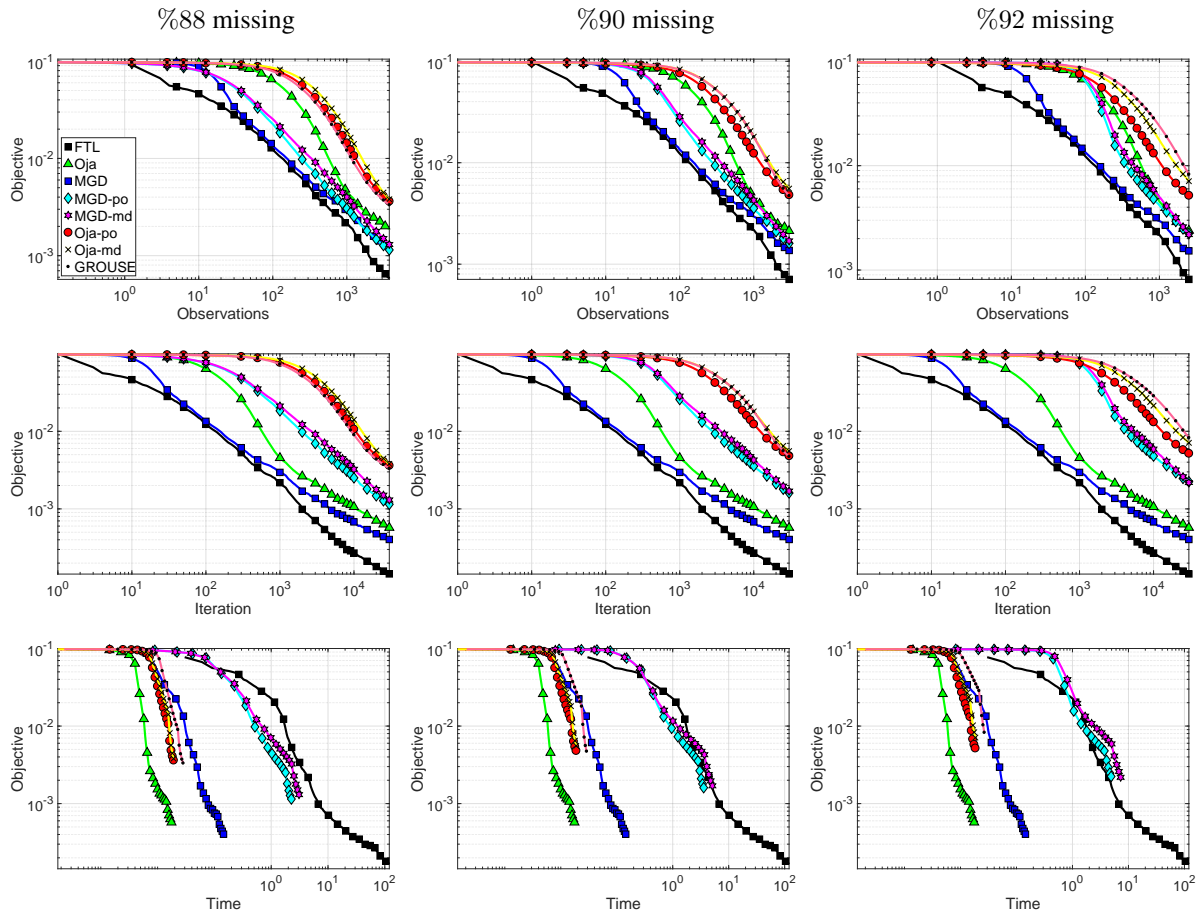


Figure 6: Comparisons of Oja, MGD, MSG-PO, MSG-MD, Oja-PO, Oja-MD and GROUSE for PCA with missing data on MNIST dataset, in terms of the variance captured on a test set as a function of number of observed entries for  $k=4$  (top), number of iterations (middle) and runtime (bottom).

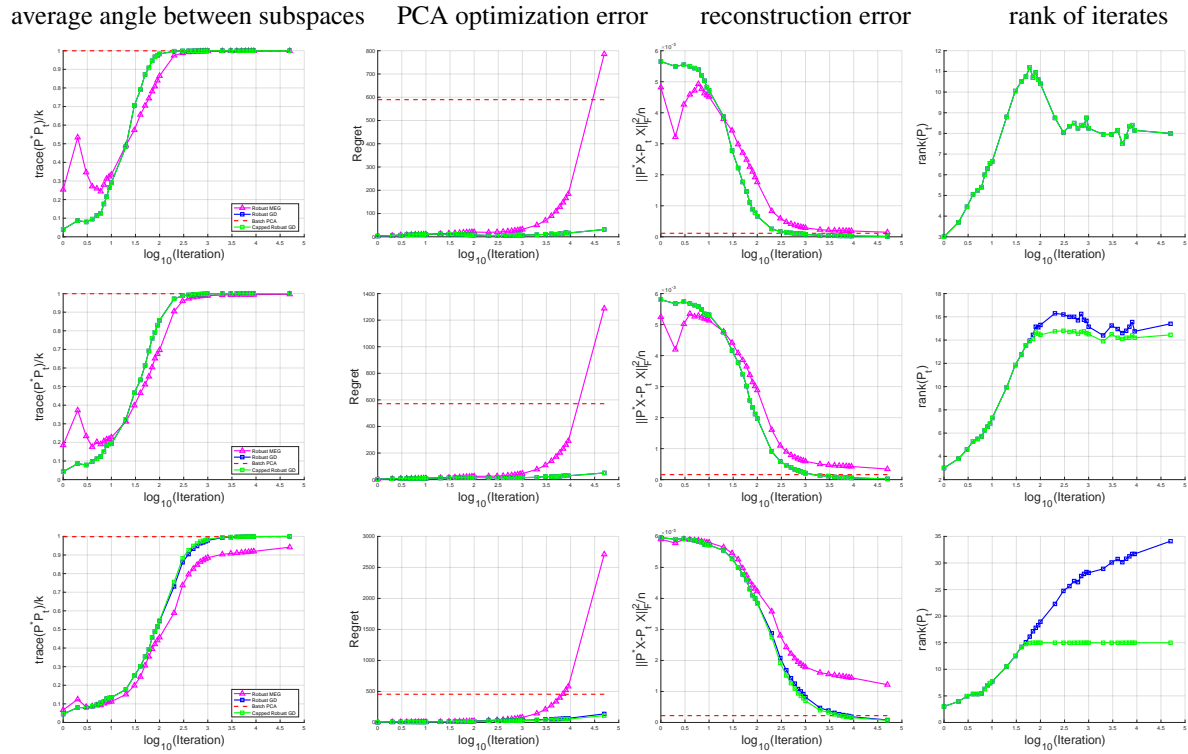


Figure 7: ( $k = 2$ ) Comparison of batch PCA, Robust Online PCA (Robust MEG), and Absolute Subspace Deviation Model (Robust GD) with 40% outliers (top row), 60% outliers (middle row) and 80% outliers (bottom row) on synthetic dataset. Experiments are in terms of the average subspace angle (left most), variance captured on a test set (left), reconstruction error (right) and the rank of iterates (right most) as a function of number of samples.