

## A Matrix Stochastic Gradient for CCA

Throughout this section, we denote the error in the gradient at time  $t$  by  $E_t = g_t - \partial_t$ . First, we introduce the following structural results, which give a lower bound on the smallest eigenvalue of the empirical auto-covariance matrices, which holds with high probability for any iterate (A.2), and uniformly over all iterates (2.1). We will use Matrix Bernstein [20] inequality in proof of Lemma A.2

**Theorem A.1** (Matrix Bernstein [20]). *consider a finite sequence  $\{X_k\}$  of independent, random, self-adjoint matrices with dimension  $d$ . Assume that each random matrix satisfies  $\mathbb{E}[X_k] = 0$  and  $\lambda_{\max}(X_k) \leq R$  almost surely. Then, for all  $\epsilon \geq 0$ ,*

$$\mathbb{P} \left\{ \lambda_{\max} \left( \sum_k X_k \right) \geq \epsilon \right\} \leq d \exp \left( \frac{-\epsilon^2/2}{\sigma^2 + R\epsilon/3} \right)$$

where  $\sigma^2 := \|\sum_k \mathbb{E}[X_k^2]\|$ .

**Lemma A.2.** *With probability at least  $1 - \delta'$  with respect to training data drawn i.i.d. from  $\mathcal{D}$ , it holds that  $\lambda_{\min}(C_{x,\tau}) \geq \frac{r_x}{2}$  and  $\lambda_{\min}(C_{y,\tau}) \geq \frac{r_y}{2}$ , whenever*

$$\tau \geq \max \left\{ \frac{2 \log \left( \frac{d_x}{\delta'} \right) B^2}{r_x^2} + \frac{\log \left( \frac{d_x}{\delta'} \right) B}{3r_x}, \frac{2 \log \left( \frac{d_y}{\delta'} \right) B^2}{r_y^2} + \frac{\log \left( \frac{d_y}{\delta'} \right) B}{3r_y} \right\}.$$

*Proof of Lemma A.2* Set  $X_k = \frac{1}{t} (x_k x_k^\top - C_x)$ , so that  $\mathbb{E} \left[ \left\| \sum_{k=1}^t X_k \right\|_2 \right] = \mathbb{E} [\|C_{x,t} - C_x\|_2]$ . Define

$$\sigma^2 = \sum_{k=1}^{\tau} \mathbb{E} [X_k^2] \preceq \frac{1}{\tau^2} \sum_k \{ \mathbb{E} [B x_k x_k^\top] - C_x^2 \} \preceq \frac{B}{\tau^2} \mathbb{E} \left[ \sum_{k=1}^{\tau} x_k x_k^\top \right] \leq \frac{B^2}{\tau}.$$

Using Theorem A.1 with  $\epsilon = \frac{r_x}{2}$  we get that with probability at least  $1 - \delta'$ , it holds that  $\|C_x - C_{x,t}\| \leq \frac{r_x}{2}$ . By Weyl's inequality, we have that

$$|\lambda_{\min}(C_{x,\tau}) - r_x| = |\lambda_{\min}(C_{x,\tau}) - \lambda_{\min}(C_x)| \leq \|C_x - C_{x,\tau}\| \leq \frac{r_x}{2}.$$

A similar derivation for  $\lambda_{\min}(C_{y,\tau})$  completes the proof.  $\square$

*Proof of Lemma 2.1* We show the result for  $C_{x,t}$  with  $c = c_x = \frac{3r_x^2}{6B^2 + Br_x}$ . Proof for  $C_{y,t}$  is symmetric. By Lemma A.2, we have that for every  $t$ :

$$\mathbb{P} \{ \|C_x - C_{x,t}\| \leq \frac{r_x}{2} \} \geq 1 - d_x e^{-ct}.$$

Probability that  $\lambda_{\min}(C_{x,t}) \geq \frac{r_x}{2}$  uniformly for all  $t \geq \tau + 1$  is  $\prod_{t=\tau+1}^{T+\tau} (1 - d_x e^{-ct})$ . Taking the logarithm, we have

$$\begin{aligned} \log \left( \prod_{t=\tau+1}^{T+\tau} (1 - d_x e^{-ct}) \right) &= \sum_{t=\tau+1}^{T+\tau} \log (1 - d_x e^{-ct}) \\ &\geq \sum_{t=\tau+1}^{T+\tau} -2d_x e^{-ct} \quad (\log(1-z) \geq -2z \text{ for } z \in (0, 0.5)) \\ &= -2d_x \frac{(e^{-c})^{\tau+1} - (e^{-c})^{T+\tau+1}}{1 - e^{-c}} \\ &\geq -2d_x \frac{e^{-c(\tau+1)}}{1 - e^{-c}} \end{aligned}$$

where we require  $d_x e^{-ct} \leq \frac{1}{2}$ , which holds for  $\tau \geq \frac{1}{c} \log(2d_x)$ . We want  $\exp\left(-2d_x \frac{e^{-c(\tau+1)}}{1-e^{-c}}\right) \geq 1 - \delta$ , which gives the following

$$\begin{aligned} \exp\left(-2d_x \frac{e^{-c(\tau+1)}}{1-e^{-c}}\right) &\geq 1 - \delta \\ \iff -e^{-c(\tau+1)} &\geq \frac{1-e^{-c}}{2d_x} \log(1-\delta) \\ \iff -c(\tau+1) &\leq \log\left(\frac{1-e^{-c}}{2d_x} \log\left(\frac{1}{1-\delta}\right)\right) \\ \iff \tau &\geq \frac{1}{c} \log\left(\frac{2d_x}{\log\left(\frac{1}{1-\delta}\right)}\right) - 1 \end{aligned}$$

so that the algorithm succeeds whenever  $\tau \geq \max\{\frac{1}{c} \log\left(\frac{1-e^{-c}}{2d_x} \log\left(\frac{1}{1-\delta}\right)\right) - 1, \frac{1}{c} \log(2d_x)\}$ .  $\square$

**Remark A.3.** Throughout the Appendix,  $\delta$  and  $\tau$  are as defined in statement of Lemma 2.1

Next, we introduce a result on perturbations of matrix square roots which is used in proof of Lemma 2.2

**Lemma A.4** (Perturbation Bounds for Matrix Square Roots [17]). *Let  $A_j \in \mathbb{R}^{n \times n}$  with  $A_j \succeq \mu_j^2 I$  in the positive semi-definite order where  $j = 1, 2$ . Then  $A_j$  has a square root satisfying  $A_j^{\frac{1}{2}} \succeq \mu_j I$  and  $\|A_1^{\frac{1}{2}} - A_2^{\frac{1}{2}}\|_2 \leq \frac{1}{\mu_1 + \mu_2} \|A_1 - A_2\|_2$ .*

Next, we present the following bound on convergence of the empirical covariance matrix to the population covariance matrix.

**Lemma A.5.** Under the same assumptions as Lemma 2.2

$$\mathbb{E}_{\mathcal{D}} [\|C_{x,t} - C_x\|_2] \leq \sqrt{\frac{2B^2}{t} \log(d_x)} + \frac{B}{3t} \log(d_x).$$

*Proof.* We bound the quantity by applying the Matrix Bernstein Inequality ([21], Theorem 6.6.1).

Set  $X_k = \frac{1}{t} (x_k x_k^\top - C_x)$ , so that  $\mathbb{E} \left[ \left\| \sum_{k=1}^t X_k \right\|_2 \right] = \mathbb{E} [\|C_{x,t} - C_x\|_2]$ . To apply the inequality we need to verify that  $\mathbb{E}[X_k] = 0$  and  $\|X_k\|_2 \leq R$  and bound  $\sigma^2 := \|\sum_k \mathbb{E}[X_k^2]\|_2$ . It follows from the definition that  $\mathbb{E}[X_k] = 0$ . To bound  $\|X_k\|_2$ , note that

$$\|X_k\|_2 = \frac{1}{t} \|x_k x_k^\top - C_x\|_2 \leq \frac{1}{t} (\|x_k x_k^\top\|_2 + \|\mathbb{E}[x_k x_k^\top]\|_2) \leq \frac{1}{t} (B + \mathbb{E}[\|x_k x_k^\top\|_2]) \leq \frac{2B}{t}.$$

Finally, we bound  $\sigma^2$  by observing that

$$\sum_{k=1}^t \mathbb{E}[X_k^2] \preceq \frac{1}{t^2} \sum_k \{\mathbb{E}[B x_k x_k^\top] - C_x^2\} \preceq \frac{B}{t^2} \mathbb{E} \left[ \sum_{k=1}^t x_k x_k^\top \right],$$

which implies  $\sigma^2 \leq \frac{B^2}{t}$ . By Matrix Bernstein's Inequality we have

$$\mathbb{E} \left[ \left\| \sum_{k=1}^t X_k \right\|_2 \right] = \mathbb{E} [\|C_{x,t} - C_x\|_2] \leq \sqrt{\frac{2B^2}{t} \log(d_x)} + \frac{B}{3t} \log(d_x)$$

which completes the proof.  $\square$

*Proof of lemma 2.2* Let  $A = W_x x_t$ ,  $B = W_y y_t$ ,  $\hat{A} = W_{x,t} x_t$ ,  $\hat{B} = W_{y,t} y_t$ . From the lower bound assumption on the spectrum of the population auto-covariance matrices and Lemma 2.1 we have

$\frac{1}{\sqrt{r_x}}\mathbf{I} \succeq \mathbf{W}_x, \sqrt{\frac{2}{r_x}}\mathbf{I} \succeq \mathbf{W}_{x,t}, \frac{1}{\sqrt{r_y}}\mathbf{I} \succeq \mathbf{W}_y, \sqrt{\frac{2}{r_y}}\mathbf{I} \succeq \mathbf{W}_{y,t}$ . Therefore,

$$\begin{aligned} \mathbb{E}[\|\mathbf{E}_t\|_2 | \mathcal{A}_t] &= \mathbb{E}[\|\mathbf{g}_t - \partial_t\|_2 | \mathcal{A}_t] = \mathbb{E}[\|\mathbf{W}_x \mathbf{x}_t \mathbf{y}_t^\top \mathbf{W}_y - \mathbf{W}_{x,t} \mathbf{x}_t \mathbf{y}_t^\top \mathbf{W}_{y,t}\|_2 | \mathcal{A}_t] \\ &= \mathbb{E}\left[\left\|\mathbf{A}\mathbf{B}^\top - \widehat{\mathbf{A}}\widehat{\mathbf{B}}^\top\right\|_2 | \mathcal{A}_t\right] \\ &= \mathbb{E}\left[\left\|\mathbf{A}\mathbf{B}^\top - \mathbf{A}\widehat{\mathbf{B}}^\top + \mathbf{A}\widehat{\mathbf{B}}^\top - \widehat{\mathbf{A}}\widehat{\mathbf{B}}^\top\right\|_2 | \mathcal{A}_t\right] \\ &\leq \mathbb{E}\left[\|\mathbf{A}\|_2 \|\mathbf{B} - \widehat{\mathbf{B}}\|_2 + \|\widehat{\mathbf{B}}\|_2 \|\mathbf{A} - \widehat{\mathbf{A}}\|_2 | \mathcal{A}_t\right]. \end{aligned} \quad (14)$$

where the inequality is due to the triangle inequality and sub-multiplicativity of the operator norm. We first bound  $\|\mathbf{A}\|_2$  and  $\|\widehat{\mathbf{B}}\|_2$

$$\|\mathbf{A}\|_2 \leq \|\mathbf{W}_x\|_2 \|\mathbf{x}_t\| \leq \sqrt{\frac{B}{r_x}}, \quad \|\widehat{\mathbf{B}}\|_2 \leq \|\mathbf{W}_{y,t}\|_2 \|\mathbf{y}_t\| \leq \sqrt{\frac{2B}{r_y}}.$$

This implies that (14) is bounded by

$$\sqrt{\frac{B}{r_x}} \mathbb{E}[\|\mathbf{B} - \widehat{\mathbf{B}}\|_2 | \mathcal{A}_t] + \sqrt{\frac{2B}{r_y}} \mathbb{E}[\|\mathbf{A} - \widehat{\mathbf{A}}\|_2 | \mathcal{A}_t]. \quad (15)$$

We now bound  $\mathbb{E}[\|\mathbf{A} - \widehat{\mathbf{A}}\|_2 | \mathcal{A}_t]$

$$\begin{aligned} \mathbb{E}[\|\mathbf{A} - \widehat{\mathbf{A}}\|_2 | \mathcal{A}_t] &\leq \mathbb{E}[\|\mathbf{W}_x - \mathbf{W}_{x,t}\|_2 \|\mathbf{x}_t\| | \mathcal{A}_t] \\ &\leq \sqrt{B} \mathbb{E}\left[\left\|\mathbf{C}_x^{-\frac{1}{2}} - \mathbf{C}_{x,t}^{-\frac{1}{2}}\right\|_2 | \mathcal{A}_t\right] \\ &= \sqrt{B} \mathbb{E}\left[\left\|\mathbf{C}_x^{-\frac{1}{2}} \left(\mathbf{C}_{x,t}^{\frac{1}{2}} - \mathbf{C}_x^{\frac{1}{2}}\right) \mathbf{C}_{x,t}^{-\frac{1}{2}}\right\|_2 | \mathcal{A}_t\right] \\ &\leq \sqrt{B} \mathbb{E}\left[\left\|\mathbf{C}_x^{-\frac{1}{2}}\right\|_2 \left\|\mathbf{C}_{x,t}^{-\frac{1}{2}}\right\|_2 \left\|\mathbf{C}_{x,t}^{\frac{1}{2}} - \mathbf{C}_x^{\frac{1}{2}}\right\|_2 | \mathcal{A}_t\right] \\ &\leq \frac{\sqrt{2B}}{r_x} \mathbb{E}\left[\left\|\mathbf{C}_{x,t}^{\frac{1}{2}} - \mathbf{C}_x^{\frac{1}{2}}\right\|_2 | \mathcal{A}_t\right] \\ &\leq \frac{\sqrt{2B}}{(1 + \sqrt{2}/2)r_x^{3/2}} \mathbb{E}[\|\mathbf{C}_{x,t} - \mathbf{C}_x\|_2 | \mathcal{A}_t] \leq \frac{\sqrt{B}}{r_x^{3/2}} \mathbb{E}[\|\mathbf{C}_{x,t} - \mathbf{C}_x\|_2 | \mathcal{A}_t], \end{aligned} \quad (16)$$

where the last inequality follows from Lemma A.4. By Lemma A.5  $\mathbb{E}[\|\mathbf{C}_{x,t} - \mathbf{C}_x\|_2] \leq \sqrt{\frac{2B^2}{t} \log(d_x)} + \frac{2B}{3t} \log(d_x)$  and thus by equation (16),

$$\begin{aligned} \mathbb{E}[\|\mathbf{A} - \widehat{\mathbf{A}}\|_2 | \mathcal{A}_t] &\leq \frac{\sqrt{B}}{r_x^{3/2}} \left\{ \sqrt{\frac{2B^2}{t} \log(d_x)} + \frac{B}{3t} \log(d_x) \right\} \\ \mathbb{E}[\|\mathbf{B} - \widehat{\mathbf{B}}\|_2 | \mathcal{A}_t] &\leq \frac{\sqrt{B}}{r_y^{3/2}} \left\{ \sqrt{\frac{2B^2}{t} \log(d_y)} + \frac{B}{3t} \log(d_y) \right\} \end{aligned} \quad (17)$$

Finally (15) together with (17) implies that

$$\mathbb{E}[\|\mathbf{E}_t\|_2 | \mathcal{A}_t] \leq \frac{2B^2}{\sqrt{r_x r_y}} \left\{ \frac{1}{r_y} \left\{ \sqrt{\frac{2 \log(d_y)}{t}} + \frac{3}{t} \log(d_y) \right\} + \frac{1}{r_x} \left\{ \sqrt{\frac{2 \log(d_x)}{t}} + \frac{3}{t} \log(d_x) \right\} \right\}.$$

Let  $r := \min\{r_x, r_y\}$  and  $d := \max\{d_x, d_y\}$ . As long as  $t > \frac{9 \log(d)}{2}$ , we have that  $\mathbb{E}[\|\mathbf{E}_t\|_2 | \mathcal{A}_t] \leq \frac{\kappa}{\sqrt{t}}$ , where  $\kappa := \frac{8B^2 \sqrt{2 \log(d)}}{r^2}$ .  $\square$

**Lemma A.6.** Assume that the event  $\mathcal{A}_T$  occurs. Let  $\kappa$  be a constant such that for all iterates  $\mathbb{E}_{\mathcal{Q}}[\|\mathbf{g}_t - \partial_t\|_2 | \mathcal{A}_t] \leq \frac{\kappa}{\sqrt{t}}$ . Then, we have that  $\sum_{t=1}^T \mathbb{E}[\|\mathbf{E}_t\|_2 | \mathcal{A}_t] \leq 2\kappa\sqrt{T}$ .

*Proof.* We note that  $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq \int_{t=1}^T \frac{1}{\sqrt{t}} dt + 1$ . Substituting  $z = \sqrt{t}$  and noting  $dt = 2z dz$  we get

$$\int_{t=1}^T \frac{1}{\sqrt{t}} dt + 1 = \int_{z=1}^{\sqrt{T}} \frac{1}{z} 2z dz + 1 = 2\sqrt{T} - 1 \leq 2\sqrt{T}.$$

□

**Lemma A.7.** *With the same assumptions as in Lemma 2.2 we have  $\|\partial_t\|_F \leq \frac{2B}{\sqrt{r_x r_y}}$ .*

*Proof.*  $\|\partial_t\|_F = \|\mathbf{W}_{x,t} \mathbf{x}_t \mathbf{y}_t^\top \mathbf{W}_{y,t}^\top\|_F = \|\mathbf{W}_{x,t} \mathbf{x}_t\|_2 \|\mathbf{W}_{y,t} \mathbf{y}_t\|_2 \leq B \|\mathbf{W}_{x,t}\|_2 \|\mathbf{W}_{y,t}\|_2 \leq \frac{2B}{\sqrt{r_x r_y}}$ . □

*Proof of Theorem 2.3* Analysis is done by conditioning on the events that  $\lambda_{\min}(\mathbf{C}_{x,t}) \geq \frac{r_x}{2}$  and  $\lambda_{\min}(\mathbf{C}_{y,t}) \geq \frac{r_y}{2}$ . By Lemma 2.1 we know that this event occurs with probability at least  $1 - \delta$ , where  $\tau \geq \max\{\frac{1}{c_x} \log\left(\frac{1-e^{-c}}{2d_x} \log(1-\delta)\right) - 1, \frac{1}{c_x} \log(2d_x), \frac{1}{c_y} \log\left(\frac{1-e^{-c}}{2d_y} \log(1-\delta)\right) - 1, \frac{1}{c_y} \log(2d_y)\}$ . The expectations are taken by conditioning on the above events and for ease of notation we set  $r_x = \frac{r_x}{2}, r_y = \frac{r_y}{2}$ . We start the analysis by measuring the distance between the  $t$ -th iterate and the optimum,  $D_t = \|\mathbf{M}_t - \mathbf{M}_*\|_F$ .

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

where the first inequality follows since projection onto a convex set in a Hilbert space is contractive, the second inequality follows since  $G = 2B/\sqrt{r_x r_y}$  is an upper bound on  $\|\partial_t\|_F$  as given in Lemma A.7, the third inequality follows using Holder's inequality, and the last inequality follows since  $\|\mathbf{M}_t - \mathbf{M}_*\|_* \leq \|\mathbf{M}_t\|_* + \|\mathbf{M}_*\|_* \leq 2k$ . Rearranging, dividing both sides by  $2\eta$ , and taking expectation on both sides, we get

$$\mathbb{E}[\langle \mathbf{M}_* - \mathbf{M}_t, \mathbf{g}_t \rangle | \mathcal{A}_t] \leq \frac{D_t^2 - D_{t+1}^2}{2\eta} + \frac{\eta}{2} G^2 + 2k\mathbb{E}[\|\mathbf{E}_t\|_2 | \mathcal{A}_t]$$

where  $\mathbf{M}_t$  and  $\mathbf{g}_t$  are conditionally independent given  $\mathcal{A}_t$ . We average over  $T$  iterates, and note that  $\sum_{t=1}^T D_t^2 - D_{t+1}^2 = D_1^2 - D_{T+1}^2 \leq D_1^2$ , where the initial distance is bounded as follows:

$$D_1^2 = \|\mathbf{M}_1 - \mathbf{M}_*\|_F^2 = \|\mathbf{M}_1\|_F^2 + \|\mathbf{M}_*\|_F^2 - 2\langle \mathbf{M}_1, \mathbf{M}_* \rangle \leq k + k + 2\|\mathbf{M}_1\|_* \|\mathbf{M}_*\|_2 \leq 4k.$$

We get:

$$\mathbb{E}[\langle \mathbf{M}_* - \tilde{\mathbf{M}}, \mathbf{C}_x^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_y^{-\frac{1}{2}} \rangle | \mathcal{A}_T] \leq \frac{2k}{\eta T} + \frac{\eta G^2}{2} + \frac{2k\kappa\sqrt{T}}{T},$$

where we used Lemma A.6 to bound  $\sum_{t=1}^T \mathbb{E}[\|\mathbf{E}_t\|_2 | \mathcal{A}_t] \leq 2\kappa\sqrt{T}$ . Finally write

$$\begin{aligned} \mathbb{E}[\langle \mathbf{M}_* - \tilde{\mathbf{M}}, \mathbf{C}_x^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_y^{-\frac{1}{2}} \rangle] &= \mathbb{E}[\langle \mathbf{M}_* - \tilde{\mathbf{M}}, \mathbf{C}_x^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_y^{-\frac{1}{2}} \rangle | \mathcal{A}_T] (1 - \delta) + \mathbb{E}[\langle \mathbf{M}_* - \tilde{\mathbf{M}}, \mathbf{C}_x^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_y^{-\frac{1}{2}} \rangle | \bar{\mathcal{A}}_T] \delta \\ &\leq \mathbb{E}[\langle \mathbf{M}_* - \tilde{\mathbf{M}}, \mathbf{C}_x^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_y^{-\frac{1}{2}} \rangle | \mathcal{A}_T] + \delta \mathbb{E}[\langle \mathbf{M}_* - \tilde{\mathbf{M}}, \mathbf{C}_x^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_y^{-\frac{1}{2}} \rangle | \bar{\mathcal{A}}_T] \\ &\leq \frac{2k}{\eta T} + \frac{\eta G^2}{2} + \frac{2k\kappa\sqrt{T}}{T} + \delta \frac{Bk}{r}, \end{aligned}$$

where the last inequality holds because  $\langle \mathbf{M}_*, \mathbf{C}_x^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_y^{-\frac{1}{2}} \rangle < \frac{Bk}{r}$ . To finish the proof we can set  $\delta \leq \frac{1}{\sqrt{T}}$  and choose optimal learning rate  $\eta = \frac{2\sqrt{k}}{G\sqrt{T}}$ . □

While Theorem 2.3 gives a bound on the objective of Problem 2, we can always bound the original CCA objective as given in Problem 1. Note that after rounding, we get a rank- $k$  factorization for  $\tilde{M} := UV^\top$ , such that  $U^\top U = I_k$  and  $V^\top V = I_k$ . As a result, for  $\hat{U} := C_x^{-\frac{1}{2}} U$  and  $\hat{V} := C_y^{-\frac{1}{2}} V$  it holds that  $\hat{U}^\top C_x \hat{U} = \hat{V}^\top C_y \hat{V} = I_k$ . Furthermore, it holds that:

$$\text{Tr}(U_*^\top C_{xy} V_* - \hat{U}^\top C_{xy} \hat{V}) \leq \frac{2\sqrt{k}G + 2k\kappa}{\sqrt{T}}$$

Let's denote  $\tilde{U} := C_{x,t}^{-\frac{1}{2}} U$  and  $\tilde{V} := C_{y,t}^{-\frac{1}{2}} V$ . We first give the following structural lemma which is used in Theorem 2.4 for giving generalization error bounds for  $\tilde{U}$  and  $\tilde{V}$  with respect to the original CCA problem as in 1.

**Lemma A.8.** *Assume the event  $\mathcal{A}_T$  occurs, then:*

$$\begin{aligned} \mathbb{E}[\|\tilde{U} - \hat{U}\|_2 | \mathcal{A}_T] &= \mathbb{E}\left[\left\|C_x^{-\frac{1}{2}} - C_{x,T}^{-\frac{1}{2}}\right\|_2 | \mathcal{A}_T\right] \leq \frac{1}{r_x^{3/2}} \left( \sqrt{\frac{2B^2}{T} \log(d_x)} + \frac{2B}{3T} \log(d_x) \right) \\ \mathbb{E}[\|\tilde{V} - \hat{V}\|_2 | \mathcal{A}_T] &= \mathbb{E}\left[\left\|C_y^{-\frac{1}{2}} - C_{y,T}^{-\frac{1}{2}}\right\|_2 | \mathcal{A}_T\right] \leq \frac{1}{r_y^{3/2}} \left( \sqrt{\frac{2B^2}{T} \log(d_y)} + \frac{2B}{3T} \log(d_y) \right) \end{aligned}$$

*Proof.* First observe that

$$\begin{aligned} \mathbb{E}[\|\tilde{U} - \hat{U}\|_2 | \mathcal{A}_T] &= \mathbb{E}[\|C_{x,T}^{-\frac{1}{2}} U - C_x^{-\frac{1}{2}} U\|_2 | \mathcal{A}_T] \leq \mathbb{E}[\|C_{x,T}^{-\frac{1}{2}} - C_x^{-\frac{1}{2}}\|_2 \|U\|_2 | \mathcal{A}_T] = \mathbb{E}[\|C_{x,T}^{-\frac{1}{2}} - C_x^{-\frac{1}{2}}\|_2 | \mathcal{A}_T] \\ \mathbb{E}[\|\tilde{V} - \hat{V}\|_2 | \mathcal{A}_T] &= \mathbb{E}[\|C_{y,T}^{-\frac{1}{2}} V - C_y^{-\frac{1}{2}} V\|_2 | \mathcal{A}_T] \leq \mathbb{E}[\|C_{y,T}^{-\frac{1}{2}} - C_y^{-\frac{1}{2}}\|_2 \|V\|_2 | \mathcal{A}_T] = \mathbb{E}[\|C_{y,T}^{-\frac{1}{2}} - C_y^{-\frac{1}{2}}\|_2 | \mathcal{A}_T] \end{aligned}$$

The proof simply follows from the following equations:

$$\begin{aligned} \mathbb{E}\left[\left\|C_x^{-\frac{1}{2}} - C_{x,T}^{-\frac{1}{2}}\right\|_2 | \mathcal{A}_T\right] &= \mathbb{E}\left[\left\|C_x^{-\frac{1}{2}} \left(C_{x,T}^{\frac{1}{2}} - C_x^{\frac{1}{2}}\right) C_{x,T}^{-\frac{1}{2}}\right\|_2 | \mathcal{A}_T\right] \\ &\leq \mathbb{E}\left[\left\|C_x^{-\frac{1}{2}}\right\|_2 \left\|C_{x,T}^{-\frac{1}{2}}\right\|_2 \left\|C_{x,T}^{\frac{1}{2}} - C_x^{\frac{1}{2}}\right\|_2 | \mathcal{A}_T\right] \\ &\leq \frac{\sqrt{2}}{r_x} \mathbb{E}\left[\left\|C_{x,T}^{\frac{1}{2}} - C_x^{\frac{1}{2}}\right\|_2 | \mathcal{A}_T\right] \\ &\leq \frac{\sqrt{2}}{(1 + \sqrt{2}/2)r_x^{3/2}} \mathbb{E}[\|C_{x,T} - C_x\|_2 | \mathcal{A}_T] \leq \frac{1}{r_x^{3/2}} \mathbb{E}[\|C_{x,T} - C_x\|_2 | \mathcal{A}_T] \end{aligned}$$

(by Lemma A.4)

and the fact that by Lemma A.5  $\mathbb{E}[\|C_{x,T} - C_x\|_2 | \mathcal{A}_T] \leq \sqrt{\frac{2B^2}{T} \log(d_x)} + \frac{2B}{3T} \log(d_x)$ .  $\square$

*Proof of Theorem 2.4* First note that

$$\text{Tr}(U_*^\top C_{xy} V_* - \tilde{U}^\top C_{xy} \tilde{V}) = \text{Tr}(U_*^\top C_{xy} V_* - \hat{U}^\top C_{xy} \hat{V}) + \text{Tr}(\hat{U}^\top C_{xy} \hat{V} - \tilde{U}^\top C_{xy} \tilde{V})$$

Moreover, we have that  $\text{Tr}(\hat{\mathbf{U}}^\top \mathbf{C}_{xy} \hat{\mathbf{V}} - \tilde{\mathbf{U}}^\top \mathbf{C}_{xy} \tilde{\mathbf{V}}) \leq 2k \|\hat{\mathbf{U}}^\top \mathbf{C}_{xy} \hat{\mathbf{V}} - \tilde{\mathbf{U}}^\top \mathbf{C}_{xy} \tilde{\mathbf{V}}\|_2$ . We bound the right hand side using the following equations

$$\begin{aligned}
\mathbb{E}[\|\hat{\mathbf{U}}^\top \mathbf{C}_{xy} \hat{\mathbf{V}} - \tilde{\mathbf{U}}^\top \mathbf{C}_{xy} \tilde{\mathbf{V}}\|_2 | \mathcal{A}_T] &= \mathbb{E}[\|\hat{\mathbf{U}}^\top \mathbf{C}_{xy} \hat{\mathbf{V}} - \tilde{\mathbf{U}}^\top \mathbf{C}_{xy} \hat{\mathbf{V}} + \tilde{\mathbf{U}}^\top \mathbf{C}_{xy} \hat{\mathbf{V}} - \tilde{\mathbf{U}}^\top \mathbf{C}_{xy} \tilde{\mathbf{V}}\|_2 | \mathcal{A}_T] \\
&\leq \mathbb{E}[\|(\hat{\mathbf{U}} - \tilde{\mathbf{U}})^\top \mathbf{C}_{xy} \hat{\mathbf{V}}\|_2 + \|\tilde{\mathbf{U}}^\top \mathbf{C}_{xy} (\hat{\mathbf{V}} - \tilde{\mathbf{V}})\|_2 | \mathcal{A}_T] \quad (\text{triangle inequality}) \\
&\leq \mathbb{E}[\|\hat{\mathbf{U}} - \tilde{\mathbf{U}}\|_2 \|\mathbf{C}_{xy}\|_2 \|\hat{\mathbf{V}}\|_2 + \|\tilde{\mathbf{U}}\|_2 \|\mathbf{C}_{xy}\|_2 \|\hat{\mathbf{V}} - \tilde{\mathbf{V}}\|_2 | \mathcal{A}_T] \quad (\text{sub-multiplicativity of the operator norm}) \\
&\leq B \mathbb{E}[\|\hat{\mathbf{U}} - \tilde{\mathbf{U}}\|_2 \|\mathbf{C}_y^{-\frac{1}{2}} \mathbf{V}\|_2 + B \|\mathbf{C}_{x,t}^{-\frac{1}{2}} \mathbf{U}\|_2 \|\hat{\mathbf{V}} - \tilde{\mathbf{V}}\|_2 | \mathcal{A}_T] \\
&\leq \mathbb{E}[\|\hat{\mathbf{U}} - \tilde{\mathbf{U}}\|_2 | \mathcal{A}_T] \frac{B}{\sqrt{r_y}} + \frac{B}{\sqrt{r_x}} \mathbb{E}[\|\hat{\mathbf{V}} - \tilde{\mathbf{V}}\|_2 | \mathcal{A}_T] \\
&\leq \frac{B}{2r_x^2} \left( \sqrt{\frac{2B^2}{T} \log(d_x)} + \frac{2B}{3T} \log(d_x) \right) \quad (\text{by Lemma A.8}) \\
&\quad + \frac{B}{2r_y^2} \left( \sqrt{\frac{2B^2}{T} \log(d_y)} + \frac{2B}{3T} \log(d_y) \right) \\
&\leq \frac{B}{r^2} \left( \sqrt{\frac{2B^2}{T} \log(d)} + \frac{2B}{3T} \log(d) \right)
\end{aligned}$$

which completes the first part of the proof. For the second part of the proof it holds:

$$\begin{aligned}
\|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \mathbf{I}\|_2 &= \|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \hat{\mathbf{U}}^\top \mathbf{C}_x \hat{\mathbf{U}}\|_2 \\
&\leq \|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \hat{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}}\|_2 + \|\hat{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \hat{\mathbf{U}}^\top \mathbf{C}_x \hat{\mathbf{U}}\|_2 \\
&\leq (\|\mathbf{C}_x \tilde{\mathbf{U}}\|_2 + \|\hat{\mathbf{U}}^\top \mathbf{C}_x\|_2) \|\tilde{\mathbf{U}} - \hat{\mathbf{U}}\|_2 \\
&\leq B \left( \|\mathbf{C}_{x,t}^{-1/2}\|_2 \|\mathbf{U}\|_2 + \|\mathbf{C}_x^{-1/2}\|_2 \|\mathbf{U}\|_2 \right) \|\tilde{\mathbf{U}} - \hat{\mathbf{U}}\|_2 \\
&\leq \frac{2B}{\sqrt{r_x}} \|\tilde{\mathbf{U}} - \hat{\mathbf{U}}\|_2.
\end{aligned}$$

Applying lemma A.8 allows us to bound  $\mathbb{E}[\|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \mathbf{I}\|_2 | \mathcal{A}_T]$ . Using Law of Total Expectation we get:

$$\begin{aligned}
\mathbb{E}[\|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \mathbf{I}\|_2] &= \mathbb{E}[\|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \mathbf{I}\|_2 | \mathcal{A}_T] (1 - \delta) + \mathbb{E}[\|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \mathbf{I}\|_2 | \bar{\mathcal{A}}_T] \delta \\
&\leq \mathbb{E}[\|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \mathbf{I}\|_2 | \mathcal{A}_T] + \delta \mathbb{E}[\|\tilde{\mathbf{U}}^\top \mathbf{C}_x \tilde{\mathbf{U}} - \mathbf{I}\|_2 | \bar{\mathcal{A}}_T] \\
&\leq \frac{B}{r_x^2} \left( \sqrt{\frac{2B^2}{T} \log(d_x)} + \frac{2B}{3T} \log(d_x) \right) + \delta(B + 1).
\end{aligned}$$

Setting  $\delta = \frac{1}{T}$  finishes the proof of the second part. The third inequality of the theorem follows similarly.  $\square$

## B Matrix Exponentiated Gradient for CCA

To make analysis easier, in this section we decide to analyze Algorithm 2 for solving a rescaled version of problem 12. In particular, we rescale the constraints in 12 so that the feasible set becomes the set of density matrices,  $\{\mathbf{M} : \text{Tr}(\mathbf{M}) = 1 \text{ and } 0 \preceq \mathbf{M} \preceq \frac{1}{k} \mathbf{I}\}$ . The results in section 3 are recovered from the proofs presented here by rescaling all bounds by a factor of  $k$ . We denote the error in the gradient at time  $t$  by  $\bar{\mathbf{E}}_t = \mathbf{C}_t - \tilde{\mathbf{C}}_t$  and  $\mathbf{E}_t = \mathbf{g}_t - \partial_t$ . We will need the following lemmas from [22].

**Lemma B.1** (Golden-Thompson inequality [9]). For arbitrary symmetric matrices  $A$  and  $B$ , it holds:

$$\text{Tr}(\exp(A+B)) \leq \text{Tr}(\exp(A)\exp(B)).$$

**Lemma B.2.** For any PSD matrix  $A$  and symmetric  $B, C$ ,  $B \preceq C$  implies  $\text{Tr}(AB) \leq \text{Tr}(AC)$ .

**Lemma B.3.** For any symmetric  $A$  such that  $0 \preceq A \preceq I$  and any  $\rho_1, \rho_2 \in \mathbb{R}$  the following holds

$$\exp(A\rho_1 + (I-A)\rho_2) \preceq A\exp(\rho_1) + (I-A)\exp(\rho_2).$$

We also need the following lemma.

**Lemma B.4.** For  $x = 1 + \sqrt{\frac{R}{L}}$  the following holds

$$\frac{-R + \log(x)L}{x-1} \geq L - 2\sqrt{RL}.$$

*Proof.* This is a simple consequence of the fact that  $\log(x) \geq (x-1) - (x-1)^2$  for  $x \geq 1$ .  $\square$

---

**Algorithm 2** Matrix Exponentiated Gradient for CCA (MEG-CCA)

---

**Input:** Training data  $\{(x_t, y_t)\}_{t=1}^T$ , step size  $\eta$ , auxiliary training data  $\{(x'_i, y'_i)\}_{t=i}^T$

**Output:**  $\tilde{M}$

```

Initialize:  $M_0 \leftarrow \frac{1}{d}I$ ,  $C_{x,0} \leftarrow \frac{1}{\tau} \sum_{i=1}^{\tau} x'_i x'^{\top}_i$ ,  $C_{y,0} \leftarrow \frac{1}{\tau} \sum_{i=1}^{\tau} y'_i y'^{\top}_i$ 
for  $t = 1$  to  $T$  do
   $C_{x,t} \leftarrow \frac{t+\tau-1}{t+\tau} C_{x,t-1} + \frac{1}{t+\tau} x_t x_t^{\top}$ ,  $W_{x,t} \leftarrow C_{x,t}^{-\frac{1}{2}}$ 
   $C_{y,t} \leftarrow \frac{t+\tau-1}{t+\tau} C_{y,t-1} + \frac{1}{t+\tau} y_t y_t^{\top}$ ,  $W_{y,t} \leftarrow C_{y,t}^{-\frac{1}{2}}$ 
   $\tilde{C}_t \leftarrow \begin{pmatrix} 0 & \partial_t \\ \partial_t^{\top} & 0 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} W_{x,t} x_t \\ W_{y,t} y_t \end{pmatrix} \begin{pmatrix} W_{x,t} x_t \\ W_{y,t} y_t \end{pmatrix}^{\top} - \frac{1}{2} \begin{pmatrix} W_{x,t} x_t \\ -W_{y,t} y_t \end{pmatrix} \begin{pmatrix} W_{x,t} x_t \\ -W_{y,t} y_t \end{pmatrix}^{\top}$ 
   $\hat{M}_t \leftarrow \frac{\exp(\log(M_{t-1}) + \eta \tilde{C}_t)}{\text{Tr}(\exp(\log(M_{t-1}) + \eta \tilde{C}_t))}$ 
   $M_t \leftarrow \mathcal{P}(\hat{M}_t)$  % projection is given by algorithm 4 in [27]
end for
 $\bar{M} = \frac{1}{T} \sum_{t=1}^T M_{t-1}$ 
 $\tilde{M} = \text{rounding}(\bar{M})$ 

```

---

**Lemma B.5.** Conditioned on the event  $\mathcal{A}_T$  occurring, after  $T$  iterations of Algorithm 2 with a step size  $\eta = \frac{1}{G} \log\left(1 + \sqrt{\frac{\log(d)}{GT}}\right)$ , where  $G = \frac{2B}{\sqrt{r_x r_y}}$  and  $M_0 = \frac{1}{d}I$  we have that,

$$\sum_{t=1}^T \text{Tr}(M_* \tilde{C}_t) - \sum_{t=1}^T \text{Tr}(M_{t-1} \tilde{C}_t) \leq 2\sqrt{G^2 T \log(d)}, \quad (18)$$

where  $M_*$  is an optimum of Problem (II).

*Proof of Lemma B.5* The proof closely follows proof of Lemma 3.1 in [22], however, we provide it for completeness. Lemma 2.1 implies that  $W_{x,t} \preceq \sqrt{\frac{2}{r_x}}I$  and  $W_{y,t} \preceq \sqrt{\frac{2}{r_y}}I$  with probability  $1 - \delta$ . Let  $F(W) = \text{Tr}(W \log(W) - W)$  be the von Neumann entropy and denote by  $\Delta(A, B)$ , the von Neumann divergence induced by  $F(\cdot)$ . More precisely,  $\Delta(A, B) = \text{Tr}(A \log(A) - A \log(B) - A + B)$ . First we note that the update step (13) (after substituting  $C_t$  with  $\tilde{C}_t$ ) is invariant under perturbing the  $\tilde{C}_t$ 's by a multiple of the identity [25], so we can assume that each  $\tilde{C}_t \succeq 0$ . Since  $\max(\|x_t\|^2, \|y_t\|^2) \leq B$ ,  $W_{x,t} \preceq \sqrt{\frac{2}{r_x}}I$  and  $W_{y,t} \preceq \sqrt{\frac{2}{r_y}}I$ , we see that  $G = \frac{2B}{\sqrt{r_x r_y}}$  is such that  $\tilde{C}_t - \lambda_{\min}(\tilde{C}_t)I \preceq GI$ . Also it holds that  $\text{Tr}(M^* \tilde{C}_t) \leq \|M^*\|_2 \|\tilde{C}_t\|_* \leq 2 \|\tilde{C}_t\|_2 \leq G$ ,

where the second to last inequality holds because  $\tilde{C}_t$  is a rank-2 matrix for all  $t$ . We begin by considering the difference  $\Delta(M, M_{t-1}) - \Delta(M, \hat{M}_t)$  for any  $M$  in the feasible set of (12).

$$\begin{aligned}\Delta(M, M_{t-1}) - \Delta(M, \hat{M}_t) &= \text{Tr}(M(\log(M) - \log(M_{t-1}))) - \text{Tr}(M(\log(M) - \log(\hat{M}_t))) \\ &= -\text{Tr}\left(M\left(\log(M_{t-1}) - \log\left(\frac{\exp(\log(M_{t-1}) + \eta\tilde{C}_t)}{\text{Tr}(\exp(\log(M_{t-1}) + \eta\tilde{C}_t))}\right)\right)\right) \\ &= \eta\text{Tr}(M\tilde{C}_t) - \log(\text{Tr}(\exp(\log(M_{t-1}) + \eta\tilde{C}_t))),\end{aligned}$$

where the first equality holds by the fact  $\text{Tr}(M) = \text{Tr}(M_{t-1})$  and the second inequality holds by expanding  $\hat{M}_t$ , according to (13). We now bound  $\log(\text{Tr}(\exp(\log(M_{t-1}) + \eta\tilde{C}_t)))$ . By Golden-Thompson's inequality B.1, we have

$$\text{Tr}(\exp(\log(M_{t-1}) + \eta\tilde{C}_t)) \leq \text{Tr}(M_{t-1} \exp(\eta\tilde{C}_t)).$$

Next, since  $0 \preceq \frac{\tilde{C}_t}{G} \preceq I$ , we use Lemma B.3 on  $\exp(\eta\tilde{C}_t)$  with  $\rho_0 = 0$  and  $\rho_1 = G\eta$  to get  $\exp(\eta\tilde{C}_t) \preceq \frac{\tilde{C}_t}{G}(\exp(G\eta) - 1) + I$ . By Lemma B.2, we now have

$$\text{Tr}(M_{t-1} \exp(\eta\tilde{C}_t)) \leq \text{Tr}\left(M_{t-1} + M_{t-1} \frac{\tilde{C}_t}{G}(\exp(G\eta) - 1)\right),$$

which implies

$$\log(\text{Tr}(M_{t-1} \exp(\eta\tilde{C}_t))) \leq \log\left(1 + \frac{\text{Tr}(M_{t-1}\tilde{C}_t)}{G}(\exp(G\eta) - 1)\right) \leq \frac{\text{Tr}(M_{t-1}\tilde{C}_t)}{G}(\exp(G\eta) - 1),$$

where last inequality holds since  $\log(1+x) \leq x$ . Thus

$$\Delta(M, M_{t-1}) - \Delta(M, \hat{M}_t) \geq \eta\text{Tr}(M\tilde{C}_t) - (\exp(G\eta) - 1) \frac{\text{Tr}(M_{t-1}\tilde{C}_t)}{G}.$$

Equivalently,

$$\text{Tr}(M_{t-1}\tilde{C}_t) \geq G \frac{\Delta(M, \hat{M}_t) - \Delta(M, M_{t-1}) + \eta\text{Tr}(M\tilde{C}_t)}{\exp(G\eta) - 1}.$$

By Generalized Pythagorean Theorem

$$\text{Tr}(M_{t-1}\tilde{C}_t) \geq G \frac{\Delta(M, M_t) - \Delta(M, M_{t-1}) + \eta\text{Tr}(M\tilde{C}_t)}{\exp(G\eta) - 1}.$$

Summing from  $t = 1$  to  $T$  and using the fact the Bregman divergence is positive we have

$$\sum_{t=1}^T \text{Tr}(M_{t-1}\tilde{C}_t) \geq G \frac{-\Delta(M, M_0) + \eta \sum_{t=1}^T \text{Tr}(M_*\tilde{C}_t)}{\exp(G\eta) - 1}.$$

To complete the proof notice that  $\Delta(M, M_0) \leq \log(d)$  and apply lemma B.4 with  $\eta = \frac{1}{G} \log\left(1 + \sqrt{\frac{G \log(d)}{GT}}\right)$ .  $\square$

**Lemma B.6.** Assume that the event  $\mathcal{A}_t$  occurs and that  $E_t$  has no repeated singular values. It holds that

$$-\frac{\kappa}{\sqrt{t}}I \preceq \mathbb{E}_{x_t, y_t} [\bar{E}_t | \mathcal{A}_t] \preceq \frac{\kappa}{\sqrt{t}}I.$$

*Proof.* By the properties of self-adjoint dilation, we have  $\mathbb{E}_{x_t, y_t} [\|E_t\|_2 | \mathcal{A}_t] = \mathbb{E}_{x_t, y_t} [\|\bar{E}_t\|_2 | \mathcal{A}_t]$ . By Jensen's inequality and Lemma 2.2 we have  $\|\mathbb{E}_{x_t, y_t} [\bar{E}_t | \mathcal{A}_t]\|_2 \leq \mathbb{E}_{x_t, y_t} [\|\bar{E}_t\|_2 | \mathcal{A}_t] = \mathbb{E}_{x_t, y_t} [\|E_t\|_2 | \mathcal{A}_t] \leq \frac{\kappa}{\sqrt{t}}$  and thus  $-\frac{\kappa}{\sqrt{t}}I \preceq \mathbb{E}_{x_t, y_t} [\bar{E}_t | \mathcal{A}_t] \preceq \frac{\kappa}{\sqrt{t}}I$ .  $\square$



*Proof of Lemma 3.1* Since  $M_{t-1}$  is independent of  $(x_t, y_t)$  we have  $\mathbb{E}_{x_t, y_t} [\text{Tr}(M_{t-1} \bar{E}_t) | \mathcal{A}_t] = \text{Tr}(M_{t-1} \mathbb{E}_{x_t, y_t} [\bar{E}_t | \mathcal{A}_t])$ . From Lemma B.6, we know that  $\mathbb{E}_{x_t, y_t} [\bar{E}_t | \mathcal{A}_t] \preceq \frac{\kappa}{\sqrt{t}} \mathbf{I}$  and since  $M_{t-1} \succeq 0$ , Lemma B.2 implies  $\text{Tr}(M_{t-1} \mathbb{E}_{x_t, y_t} [\bar{E}_t | \mathcal{A}_t]) \leq \frac{\kappa}{\sqrt{t}} \text{Tr}(M_{t-1}) = \frac{\kappa}{\sqrt{t}}$ . Similarly using that  $-\frac{\kappa}{\sqrt{t}} \mathbf{I} \preceq \mathbb{E}_{x_t, y_t} [\bar{E}_t | \mathcal{A}_t]$ , we have  $\mathbb{E}_{x_t, y_t} [\text{Tr}(M_* \bar{E}_t) | \mathcal{A}_t] \geq -\frac{\kappa}{\sqrt{t}}$  and the result follows.  $\square$

*Proof of Theorem 3.2* By Lemma B.5 we have

$$\sum_{t=1}^T \text{Tr}(M_* (C_t - \bar{E}_t)) - \sum_{t=1}^T \text{Tr}(M_{t-1} (C_t - \bar{E}_t)) \leq 2\sqrt{G^2 T \log(d)}. \quad (19)$$

Let  $\mathbb{E}_\tau[\cdot]$  denote the expectation w.r.t.  $(x_t, y_t)_{t=1}^T$ . We now compute the expectations of the two terms on the left hand side of (19)

$$\sum_{t=1}^T \mathbb{E} [\text{Tr}(M_* (C_t - \bar{E}_t)) | \mathcal{A}_t] = T \text{Tr}(M_* C) - \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_* \bar{E}_t) | \mathcal{A}_t] \quad (20)$$

The second term expands as follows

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} (C_t - \bar{E}_t)) | \mathcal{A}_t] &= \sum_{t=1}^T \text{Tr}(\mathbb{E}_t [M_{t-1} (C_t - \bar{E}_t) | \mathcal{A}_t]) \\ &= \sum_{t=1}^T \text{Tr}(\mathbb{E}_{t-1} [\mathbb{E}_t [M_{t-1} (C_t - \bar{E}_t) | (x_i, y_i)_{i=1}^{t-1}, \mathcal{A}_t]]) \\ &= \sum_{t=1}^T \text{Tr}(\mathbb{E}_{t-1} [M_{t-1} | \mathcal{A}_t] \mathbb{E}_t [(C_t - \bar{E}_t) | \mathcal{A}_t]) \\ &= \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C) | \mathcal{A}_t] - \sum_{t=1}^T \text{Tr}(\mathbb{E}_{t-1} [M_{t-1} | \mathcal{A}_t] \mathbb{E}_t [\bar{E}_t | \mathcal{A}_t]) \\ &= \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C) | \mathcal{A}_t] - \sum_{t=1}^T \mathbb{E}_{t-1} [\mathbb{E}_{x_t, y_t} [\text{Tr}(M_{t-1} \bar{E}_t) | \mathcal{A}_t]] , \end{aligned} \quad (21)$$

where the second equality holds by smoothing property of expectation and the third equality holds because  $M_{t-1}$  is conditionally independent of  $C_t$  and  $\bar{E}_t$ . Putting together (19), (20), (21) we have

$$\begin{aligned} T \text{Tr}(M_* C) - \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C) | \mathcal{A}_t] &\leq 2\sqrt{G^2 T \log(d)} + \sum_{t=1}^T [\mathbb{E}_{t-1} [\mathbb{E}_{x_t, y_t} [\text{Tr}(M_{t-1} \bar{E}_t) | \mathcal{A}_t]] - \mathbb{E} [\text{Tr}(M_* \bar{E}_t) | \mathcal{A}_t]] \\ &= 2\sqrt{G^2 T \log(d)} + \sum_{t=1}^T \mathbb{E}_{t-1} [\mathbb{E}_{x_t, y_t} [\text{Tr}(M_{t-1} \bar{E}_t) - \text{Tr}(M_* \bar{E}_t) | \mathcal{A}_t]] \\ &\leq 2\sqrt{G^2 T \log(d)} + \sum_{t=1}^T \frac{\kappa}{\sqrt{t}} \leq 2\sqrt{G^2 T \log(d)} + 2\sqrt{T} \kappa \end{aligned} \quad (22)$$

where the second inequality follows from Lemma 3.1 and the last inequality follows from Lemma A.6

Next we bound:

$$\begin{aligned} T \text{Tr}(M_* C) - \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C)] &= T \text{Tr}(M_* C) - \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C) | \mathcal{A}_t] (1 - \delta) - \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C) | \bar{\mathcal{A}}_t] \delta \\ &\leq T \text{Tr}(M_* C) - \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C) | \mathcal{A}_t] - \sum_{t=1}^T \mathbb{E} [\text{Tr}(M_{t-1} C) | \bar{\mathcal{A}}_t] \delta \\ &\leq 2\sqrt{G^2 T \log(d)} + 2\sqrt{T} \kappa. \end{aligned}$$

To finish the proof we only need to divide both sides by  $T$ .  $\square$