# A. Supplementary Materials

In this section we provide the proofs of the results stated in the main paper. Throughout the section, we denote the objective function of $\ell_2$ and $\ell_1$ regularized problems with $J_\lambda(M)$ and $J_\mu(M)$ respectively, where

$$J_\mu(M) := -\mathbb{E}[x^\top M x] + \mu \mathrm{Tr}(M),$$

$$J_\lambda(M) := -\mathbb{E}[x^\top M x] + \frac{\lambda}{2}\|M\|_F^2.$$

Recall that the eigendecomposition of C is given by $C = \sum_{i=1}^d \lambda_i u_i u_i^\top$.

## A.1. Proofs of Section 2

*Proof of Lemma 2.2.* Note that the feasible set of Problem 6 is convex, and the objective $J_\lambda$ is $\lambda$-strongly convex. Hence, the optimum is unique. Since Slater condition is satisfied, strong duality holds and KKT conditions are necessary and sufficient for optimality. Let $(\lambda_i(M), v_i)$, $i \in [d]$ denote the eigenvalues and associated eigenvectors of M. KKT first-order optimality condition yields

$$0 = -C + \lambda M - \sum_{i=1}^d \gamma_i v_i v_i^\top + \sum_{i=1}^d \omega_i v_i v_i^\top + \beta I,$$

where $\gamma_i, \omega_i, \beta \geq 0$ are Lagrange multipliers for the constraints $M \succeq 0$, $M \preceq I$, and $\mathrm{Tr}(M) \leq k$ respectively. Note that except for C, every other term in the above equation has the same set of eigenvectors as M. We conclude that C and M have the same set of eigenvectors, i.e. $u_i = v_i$, $i \in [d]$. Rearranging, we get:

$$C = \sum_{i=1}^d (\lambda \lambda_i(M) - \gamma_i + \omega_i + \beta) u_i u_i^\top$$

Complementary slackness implies that $\gamma_i \lambda_i(M)$, $\omega_i(\lambda_i(M) - 1)$, and $\beta(\mathrm{Tr}(M) - k)$ all vanish. Now, for an admissible $\lambda$, we can verify that the following satisfy the KKT conditions:

$$M = \Pi_k(C)$$
$$\omega_i = (\lambda_i - \lambda_k) 1_{i \leq k}$$
$$\gamma_i = (\lambda_k - \lambda - \lambda_j) 1_{i > k}$$
$$\beta = \lambda_k - \lambda,$$

where $\Pi_k(C)$ returns the projection matrix corresponding to top-$k$ eigenspace of C and $1_{i \in \mathcal{A}}$ is the indicator of set $\mathcal{A}$. In particular, $\gamma_i = \lambda_i - \lambda_k \geq 0$ for $i \in [k]$. Also, since $\lambda < g_k$, for $i \in \{k+1, \ldots, d\}$ it holds that $\gamma_i = \lambda_k - \lambda - \lambda_j \geq 0$. It is easy to verify that the complementary slackness conditions are also satisfied for the above assignments of the primal and dual variables. $\square$

*Proof of Lemma 2.3.* As discussed in the proof of Lemma 2.2, since the feasible set is convex and the objective is strongly convex, the optimum is unique. Furthermore, Salter condition is satisfied. Therefore, strong duality holds and KKT conditions are necessary and sufficient for optimality. Following Lemma 2.2, KKT first-order optimality condition yields:

$$C = \sum_{i=1}^d (\lambda \lambda_i(M) - \gamma_i + \omega_i + \beta) u_i u_i^\top$$

One can assert that the following set of primal and dual variables satisfy the KKT conditions:

$$M = \sum_{i=1}^p u_i u_i^\top + \frac{k-p}{q-p} \sum_{j=p+1}^q u_j u_j^\top$$

$$\omega_i = (\lambda_i - \lambda_{p+1} + \frac{k-q}{q-p}\lambda) 1_{i \leq p}$$

$$\gamma_i = (\lambda_{p+1} - \lambda_i - \frac{k-p}{q-p}\lambda) 1_{i > q}$$

$$\beta = \lambda_{p+1} - \frac{k-p}{q-p}\lambda,$$

where $1_{i \in \mathcal{A}}$ is the indicator of set $\mathcal{A}$. $\square$

*Proof of Theorem 2.4.* Let's denote the stochastic gradient at time $t$ by $\widehat{g}_t := -x_t x_t^\top + \lambda M_t$. Note that

$$\|\widehat{g}_t\|_F = \|-x_t x_t^\top + \lambda M_t\|_F \leq \|x_t\|^2 + \lambda \|M_t\|_F \leq 1 + \lambda\sqrt{k}.$$

Hence, $G^2 := (1 + \lambda\sqrt{k})^2 \geq \mathbb{E}[\|\widehat{g}_t\|^2]$ for all $t \in \{1, \ldots, T\}$. Let $M^*$ be the global optimum of Problem 6. Since $\lambda$ is admissible, it follows from Lemmas 2.2 and 2.3 that $M^*$ is also an optimum for Problem 4. Using the following Lemma from (Rakhlin et al., 2012), we bound the distance between $M_T$ and $M^*$.

**Lemma A.1** (Lemma 1 of (Rakhlin et al., 2012) ). Suppose $J$ is $\lambda$-strongly convex over a convex set $\mathcal{M}$, and that $\mathbb{E}[\|\widehat{g}_t\|_F^2] \leq G^2$. Then if we pick $\eta_t = \frac{1}{\lambda t}$, it holds for any $T$ that

$$\mathbb{E}[\|M_T - M^*\|_F^2] \leq \frac{4G^2}{\lambda^2 T}. \tag{14}$$

Now, we note that

$$\begin{aligned}
\|M^* - \tilde{M}\|_F &= \|M^* - M_T + M_T - \tilde{M}\|_F \\
&\leq \|M^* - M_T\|_F + \|M_T - \tilde{M}\|_F \\
&\leq 2\|M^* - M_T\|_F
\end{aligned}$$

where the second inequality holds because by definition of $\tilde{M}$, $\|M_T - \tilde{M}\|_F \leq \|M^* - M_T\|_F$. Plugging back the above into Equation (14) completes the proof. $\square$

*Proof of Theorem 2.5.* Let $M_* = \sum_{i=1}^{k} u_i u_i^\top$ and $\tilde{M} = \sum_{i=1}^{k} \tilde{u}_i \tilde{u}_i^\top$ be the eigendecompositions of $M_*$ and $\tilde{M}$ respectively. Denote the suboptimality gap by $\epsilon := \mathbb{E}[x^\top M_* x - x^\top \tilde{M} x]$, we have

$$
\begin{aligned}
\epsilon &= \sum_{i=1}^{k} u_i^\top \mathbb{E}[xx^\top] u_i - \mathbb{E}[\sum_{i=1}^{k} \tilde{u}_i^\top \mathbb{E}[xx^\top] \tilde{u}_i] \\
&= \sum_{i=1}^{k} \lambda_i - \mathbb{E}[\sum_{j=1}^{k} \tilde{u}_j^\top \sum_{i=1}^{d} \lambda_i u_i u_i^\top \tilde{u}_j] \\
&= \sum_{i=1}^{k} \lambda_i - \sum_{j=1}^{k} \sum_{i=1}^{d} \lambda_i \mathbb{E}[(u_i^\top \tilde{u}_j)^2]
\end{aligned}
$$

Noting that $\mathbb{E}[(u_i^\top \tilde{u}_j)^2] \geq 0$, we get

$$
\begin{aligned}
\epsilon &\leq \sum_{i=1}^{k} \lambda_i - \sum_{j=1}^{k} \sum_{i=1}^{k} \lambda_i \mathbb{E}[(u_i^\top \tilde{u}_j)^2] \\
&= \sum_{i=1}^{k} \lambda_i (1 - \sum_{j=1}^{k} \mathbb{E}[(u_i^\top \tilde{u}_j)^2])
\end{aligned}
$$

Since $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_d$ and $1 - \sum_{j=1}^{k} \mathbb{E}[(u_i^\top \tilde{u}_j)^2] \geq 0$, we get that

$$
\epsilon \leq \lambda_1 \sum_{i=1}^{k} (1 - \mathbb{E}[\sum_{j=1}^{k} (u_i^\top \tilde{u}_j)^2]) = \lambda_1 (k - \mathbb{E}[\|U^\top \tilde{U}\|_F^2]) \tag{15}
$$

Now, we show that $k - \mathbb{E}[\|\tilde{U}^\top U_*\|_F^2] = \frac{1}{2}\mathbb{E}[\|M_* - \tilde{M}\|_F^2]$.

$$
\begin{aligned}
\|M_* - \tilde{M}\|_F^2 &= \|M_*\|_F^2 + \|\tilde{M}\|_F^2 - 2\langle M_*, \tilde{M}\rangle \\
&= 2k - 2\langle UU^\top, \tilde{U}\tilde{U}^\top\rangle \\
&= 2\left(k - \|\tilde{U}^\top U_*\|_F^2\right).
\end{aligned}
$$

Plugging back this result in Equation (15), we get

$$
\epsilon \leq \frac{\lambda_1}{2}\mathbb{E}[\|M_* - \tilde{M}\|_F^2] \leq \frac{8\lambda_1(1 + \lambda\sqrt{k})^2}{\lambda^2 T}
$$

where the last inequality follows from Equation (8). □

## A.2. Proofs of Section 3

*Proof of Lemma 3.1.* Since the problem is convex and Slater condition is satisfied, the KKT conditions are necessary and sufficient for optimality. Let $(\lambda_i(M), v_i)$, $i \in [d]$ denote the eigenvalues and associated eigenvectors of $M$. KKT first-order optimality condition yields

$$
0 = -C + \mu I - \sum_{i=1}^{d} \gamma_i v_i v_i^\top + \sum_{i=1}^{d} \omega_i v_i v_i^\top + \beta I,
$$

where $\gamma_i, \omega_i, \beta \geq 0$ are Lagrange multipliers for the constraints $M \succeq 0$, $M \preceq I$, and $\text{Tr}(M) \leq k$ respectively. Note that except for $C$, every other term in the above equation has the same set of eigenvectors as $M$. We conclude that $C$ and $M$ have the same set of eigenvectors, i.e. $u_i = v_i$, $i \in [d]$. Rearranging, we get:

$$
C = \sum_{i=1}^{d} (\mu - \gamma_i + \omega_i + \beta) u_i u_i^\top. \tag{16}
$$

Complementary slackness implies that $\gamma_i \lambda_i(M)$, $\omega_i(\lambda_i(M) - 1)$, and $\beta(\text{Tr}(M) - k)$ all vanish. Now, for an admissible $\mu$, we can assert that the following satisfy the KKT conditions:

$$
\begin{aligned}
M &= \Pi_k(C) \\
\omega_i &= (\lambda_i - \lambda_k) 1_{i \leq k} \\
\gamma_i &= (\lambda_k - \lambda_i) 1_{i > k} \\
\beta &= \lambda_k - \mu,
\end{aligned}
$$

where $\Pi_k(C)$ returns the projection matrix corresponding to top-$k$ eigenspace of $C$ and $1_{i \in \mathscr{A}}$ is the indicator of set $\mathscr{A}$. We shall now prove that $M$ is the unique global optimum. Assume $M_0 \neq M$ achieves the optimum as well. Observe that all extreme points of the feasible set are rank-$k$ projection matrices. Since the objective is linear, it must be the case that $M_0$ is a convex combination of some extreme points of the feasible set, i.e. $M_0 = \sum_{i=1}^{r} \alpha_i M_i$ where $M_1, \ldots, M_r$ are all rank-$k$ projection matrices and $\sum_{i=1}^{r} \alpha_i = 1$ and $\alpha_i \geq 0, \forall i \in 1, \ldots, r$. Hence, $\text{Tr}(M_0) = \sum_{i=1}^{r} \alpha_i \text{Tr}(M_i) = k$. Since $-\mathbb{E}_x[x^\top Mx] + \mu\text{Tr}(M) = -\mathbb{E}_x[x^\top M_0 x] + \mu\text{Tr}(M_0)$ and $\text{Tr}(M) = \text{Tr}(M_0) = k$, it should be the case that $-\mathbb{E}_x[x^\top Mx] = -\mathbb{E}_x[x^\top M_0 x]$, which contradicts with the eigengap assumption ($g_k > 0$). □

**Key insights on the admissibility condition** $(\mu \leq \lambda_k)$**:** The stationary condition in Equation (16) gives a system of linear equations with $d$ equations:

$$
\lambda_i = \mu - \gamma_i + \omega_i + \beta, \forall i \in [d]. \tag{17}
$$

If $\mu > \lambda_k$, since $\beta \geq 0$, it should hold for all $j \in \{k, \ldots, d\}$ that $-\gamma_j + \omega_j < 0$. Since $\omega_j \geq 0$, we should have $\gamma_j > 0$. By complementary slackness, we conclude that $\lambda_j(M^*) = 0$. This implies that $\text{Tr}(M^*) < k$, because at least $d - k + 1$ eigenvalues of $M^*$ are equal to zero, and the rest are at most one. In this case, $M^*$ cannot be a solution of Problem 4. Hence, the admissibility condition in the statement of Lemma 3.1.

*Proof of Theorem 3.2.* Lets denote the stochastic gradient at time $t$ by $\hat{g}_t := -x_t x_t^\top + \mu I$. Note that

$$
\|\hat{g}_t\|_F = \| - x_t x_t^\top + \mu I\|_F \leq \|x_t\|^2 + \mu\|I\|_F \leq 1 + \mu\sqrt{d}.
$$

Hence, $\mathbb{E}[\|\widehat{g}_t\|^2] \leq (1 + \mu\sqrt{d})^2 =: G^2$ for all iterates. Also note that the diameter of the feasible set is at most $2\sqrt{k}$, that is, $\|M - M'\|_F \leq 2\sqrt{k} =: D$ for all $M$ and $M'$ in the feasible set. By analysis of SGD for convex smooth problems (see Theorem 2 of (Shamir & Zhang, 2013)), we get that for $\eta_t = \frac{c}{\sqrt{t}}$,

$$\mathbb{E}[J_\mu(M_T)] - J_\mu(M^*) \leq (\frac{D^2}{c} + cG^2)\frac{2 + \log T}{\sqrt{T}}$$
$$\leq (\frac{4k}{c} + c(1 + \mu\sqrt{d})^2)\frac{2 + \log T}{\sqrt{T}}$$

For $T > 1$, we have $2 + \log T \leq 4 \log T$. Choosing $c = \frac{\sqrt{4k}}{1 + \mu\sqrt{d}}$ gives

$$\mathbb{E}[J_\mu(M_T)] - J_\mu(M_*) \leq \frac{16\sqrt{k}(1 + \mu\sqrt{d})\log T}{\sqrt{T}}. \quad (18)$$

To prove the claim of the theorem, we start from its left hand side where we have for any $T > 1$:

$$\mathbb{E}[x^\top M^* x] - \mathbb{E}[x^\top M_T x] = \mathbb{E}[J_\mu(M_T)] - J_\mu(M^*) \\ + \mu\mathrm{Tr}\,(M^* - \mathbb{E}[M_T]) \quad (19)$$

We first bound the quantity $\mathrm{Tr}\,(M^* - M')$ in terms of the difference in objectives achieved by $M^*$ and $M'$. Eigendecomposition of $C - \mu I$ is given by $C - \mu I = \sum_{i=1}^d (\lambda_i - \mu)u_i u_i^\top$. By definition of the principal subspace we have $M^* = \sum_{i=1}^k u_i u_i^\top$. Since $u_1, \ldots, u_d$ form an orthogonal basis for $\mathbb{R}^d$, we have the following equalities for the trace of the difference:

$$\mathrm{Tr}\,(M^* - M') = \sum_{i=1}^d u_i^\top(M^* - M')u_i$$
$$= \sum_{i=1}^k \underbrace{u_i^\top(M^* - M')u_i}_{\geq 0} + \sum_{i=k+1}^d \underbrace{u_i^\top(M^* - M')u_i}_{\leq 0}$$

Note that by definition of $M^*$, and the fact that $\|M'\|_2 \leq 1$, we get that $u_i^\top(M^* - M')u_i \geq 0$ for all $i \in \{1, \ldots, k\}$. On the other hand, since $\lambda_i M^* = 0$ for $i \in \{k + 1, \ldots, d\}$, we have that $u_i^\top(M^* - M')u_i \leq 0$. We use this property to upper bound the trace by scaling up the $\geq 0$ part and scaling

down the $\leq 0$ part:

$$\mathrm{Tr}\,(M^* - M') \leq \sum_{i=1}^k \frac{\lambda_i - \mu}{\lambda_k - \mu}u_i^\top(M^* - M')u_i$$
$$+ \sum_{i=k+1}^d \frac{\lambda_i - \mu}{\lambda_k - \mu}u_i^\top(M^* - M')u_i$$
$$= \frac{1}{\lambda_k - \mu}\sum_{i=1}^d (\lambda_i - \mu)u_i^\top(M^* - M')u_i$$
$$= \frac{1}{\lambda_k - \mu}\langle C - \mu I, M^* - M'\rangle$$
$$= \frac{1}{\lambda_k - \mu}(J_\mu(M') - J_\mu(M^*))$$

where the first inequality holds because for $i \in \{1, \ldots, k\}$, $0 < \lambda_k - \mu \leq \lambda_i - \mu$ so that $\frac{\lambda_i - \mu}{\lambda_k - \mu} \geq 1$. Furthermore, for $i \in \{k+1, \ldots, d\}$, $0 < \lambda_k - \mu \geq \lambda_i - \mu$ so that $\frac{\lambda_i - \mu}{\lambda_k - \mu} \leq 1$. Observe that $\mu\mathrm{Tr}\,(M^* - M') \leq \frac{\mu}{\lambda_k - \mu}(J_\mu(M') - J_\mu(M^*))$ whenever $\mu$ is an admissible regularization parameter, which together with Equation 19 imply

$$\mathbb{E}[x^\top M^* x] - \mathbb{E}[x^\top M_T x] \leq (\mathbb{E}[J_\mu(M_T)] - J_\mu(M^*)) \\ + \frac{\mu}{\lambda_k - \mu}(\mathbb{E}[J_\mu(M_T)] - J_\mu(M^*)) \\ \leq \frac{\lambda_k}{\lambda_k - \mu}(\mathbb{E}[J_\mu(M_T)] - J_\mu(M^*))$$

Plugging back the above in Equaition (18), we get

$$\mathbb{E}[x^\top M_* x] - \mathbb{E}[x^\top M_T x] \leq \frac{16\lambda_k\sqrt{k}(1 + \mu\sqrt{d})\log T}{(\lambda_k - \mu)\sqrt{T}}.$$

When $\mu \leq \min\{\frac{\lambda_k}{2}, \frac{1}{\sqrt{d}}\}$, it is easy to see that the right hand side above is bounded by $\frac{64\sqrt{k}\log T}{\sqrt{T}}$, which completes the proof.

$\square$

*Proof of Lemma 3.3.* Let denote the tail at $t$-th iterate by $\delta_t = \sum_{i=k+1}^d \lambda_i(M_t)$. Expand the right hand side to get

$$\langle C - \mu I, M_* - M_t\rangle = \sum_{i=1}^k \lambda_i - \langle C, M_t\rangle \\ - \mu(\mathrm{Tr}\,(M_*) - \mathrm{Tr}\,(M_t))$$

By Von Neumann's trace inequality, we know $\langle C, M_t\rangle \leq \sum_{i=1}^d \lambda_i\lambda_i(M_t)$. Substituting in the above equality and

expanding, we get

$$\langle \mathbf{C} - \mu\mathbf{I}, \mathbf{M}_* - \mathbf{M}_t \rangle$$

$$\geq \sum_{i=1}^{k} \lambda_i - \sum_{i=1}^{d} \lambda_i \lambda_i(\mathbf{M}_t) - \mu(\sum_{i=1}^{k} 1 - \sum_{i=1}^{d} \lambda_i(\mathbf{M}_t))$$

$$= \sum_{i=1}^{k} \lambda_i(1 - \lambda_i(\mathbf{M}_t)) - \sum_{i=k+1}^{d} \lambda_i \lambda_i(\mathbf{M}_t)$$

$$- \mu \sum_{i=1}^{k} (1 - \lambda_i(\mathbf{M}_t)) + \mu \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t)$$

$$\geq \lambda_k \sum_{i=1}^{k} (1 - \lambda_i(\mathbf{M}_t)) - \lambda_{k+1} \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t)$$

$$- \mu \sum_{i=1}^{k} (1 - \lambda_i(\mathbf{M}_t)) + \mu \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t)$$

$$= (\lambda_k - \mu) \sum_{i=1}^{k} (1 - \lambda_i(\mathbf{M}_t)) - (\lambda_{k+1} - \mu) \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t)$$

$$= (\lambda_k - \mu) \left( k - \sum_{i=1}^{k} \lambda_i(\mathbf{M}_t) \right) - (\lambda_{k+1} - \mu)\delta_t$$

$$\geq (\lambda_k - \mu)\delta_t - (\lambda_{k+1} - \mu)\delta_t = (\lambda_k - \lambda_{k+1})\delta_t$$

where the second inequality holds because $\lambda_i$'s are sorted in descending order. $\qquad\square$

*Proof of Lemma 3.4.* Simply follows from the following inequalities:

$$\sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t + \eta_t \mathbf{x}_t \mathbf{x}_t^\top) \leq \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t) + \sum_{i=1}^{d-k} \lambda_i(\eta_t \mathbf{x}_t \mathbf{x}_t^\top)$$

$$= \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t) + \eta_t \|\mathbf{x}_t\|^2$$

where the first inequality is due to Lidskii's (see, e.g. (Tao, 2012)). Taking expectation of both sides

$$\mathbb{E}[\sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t + \eta_t \mathbf{x}_t \mathbf{x}_t^\top)] \leq \mathbb{E}[\sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t) + \eta_t \|\mathbf{x}_t\|^2]$$

$$\leq \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t) + \eta_t$$

where the last inequality holds because $\mathbb{E}[\|\mathbf{x}_t\|^2] \leq 1$ by assumptions of Theorem 3.2. $\qquad\square$

*Proof of Theorem 3.5.* As a consequence of Lemma 3.3 and

Lemma 3.4, we have that

$$\tilde{\delta}_t := \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t + \eta_t \mathbf{x}_t \mathbf{x}_t^\top) \leq \sum_{i=k+1}^{d} \lambda_i(\mathbf{M}_t) + \eta_t$$

$$\leq \frac{1}{g_k} \langle \mathbf{C} - \mu\mathbf{I}, \mathbf{M}_* - \mathbf{M}_t \rangle + \eta_t$$

$$\leq \frac{16\sqrt{k}(1 + \mu\sqrt{d})\log t}{g_k \sqrt{t}} + \frac{4\sqrt{k}}{(1 + \mu\sqrt{d})\sqrt{t}}$$

where the last inequality follows from Equation (18). After shrinking the spectrum by $-\eta_t\mu\mathbf{I}$, at most $\frac{\tilde{\delta}_t}{\mu\eta_t}$ eigenvalues in the tail will not be eliminated.

$$\frac{\tilde{\delta}_t}{\mu\eta_t} = \frac{\frac{16\sqrt{k}(1+\mu\sqrt{d})\log t}{g_k\sqrt{t}} + \frac{4\sqrt{k}}{(1+\mu\sqrt{d})\sqrt{t}}}{\mu\frac{2}{1+\mu\sqrt{d}}\sqrt{\frac{k}{t}}}$$

$$\leq \frac{8\log t(1 + \mu\sqrt{d})^2}{\mu g_k} + \frac{2}{\mu} \leq \frac{33\log t}{\mu g_k}.$$

where the last inequality holds since $\mu \leq 1/\sqrt{d}$. $\qquad\square$

### A.3. Proofs of Section 4

*Proof of Lemma 4.1.* Note that the feasible set of Problem 12 is convex, and the objective function is $\lambda$-strongly convex. Hence, the optimum is unique. Since Slater condition is satisfied, strong duality holds and KKT conditions are necessary and sufficient for optimality. Let $(\lambda_i(\mathbf{M}), \mathbf{v}_i)$, $i \in [d]$ denote the eigenvalues and associated eigenvectors of M. KKT first-order optimality condition yields

$$0 = -\mathbf{C} + \mu\mathbf{I} + \lambda\mathbf{M} - \sum_{i=1}^{d} \gamma_i \mathbf{v}_i \mathbf{v}_i^\top + \sum_{i=1}^{d} \omega_i \mathbf{v}_i \mathbf{v}_i^\top + \beta\mathbf{I},$$

where $\gamma_i, \omega_i, \beta \geq 0$ are Lagrange multipliers for the constraints $\mathbf{M} \succeq 0$, $\mathbf{M} \preceq \mathbf{I}$, and $\mathrm{Tr}(\mathbf{M}) \leq k$ respectively. Note that except for C, every other term in the above equation has the same set of eigenvectors as M. We conclude that C and M have the same set of eigenvectors, i.e. $\mathbf{u}_i = \mathbf{v}_i$, $i \in [d]$. Rearranging, we get:

$$\mathbf{C} = \sum_{i=1}^{d} (\mu + \lambda\lambda_i(\mathbf{M}) - \gamma_i + \omega_i + \beta)\mathbf{u}_i \mathbf{u}_i^\top$$

Complementary slackness implies that $\gamma_i \lambda_i(\mathbf{M})$, $\omega_i(\lambda_i(\mathbf{M}) - 1)$, and $\beta(\mathrm{Tr}(\mathbf{M}) - k)$ all vanish. Now, If $(\lambda, \mu)$ is an admissible regularization pair, we can verify that the following satisfy the KKT conditions:

$$\mathbf{M} = \Pi_k(\mathbf{C})$$
$$\omega_i = (\lambda_i - \lambda_k)1_{i \leq k}$$
$$\gamma_i = (\lambda_k - \lambda - \lambda_i)1_{i > k}$$
$$\beta = \lambda_k - \mu - \lambda,$$

where $\Pi_k(\mathbf{C})$ returns the projection matrix corresponding to top-$k$ eigenspace of $\mathbf{C}$ and $1_{i \in \mathscr{A}}$ is the indicator of set $\mathscr{A}$. In particular, $\gamma_i = \lambda_i - \lambda_k \geq 0$ for $i \in [k]$. Also, since $\lambda < g_k$, for $i \in \{k+1, \ldots, d\}$ it holds that $\gamma_i = \lambda_k - \lambda - \lambda_j \geq 0$. It is easy to verify that the complementary slackness conditions are also satisfied for the above assignments of the primal and dual variables. $\qquad\square$