# Intersecting singularities for multi-structured estimation

**Emile Richard**                                                                EMILE.RICHARD@MINES-PARISTECH.FR
CBIO Mines ParisTech, INSERM U900, Institut Curie

**Francis Bach**                                                                         FRANCIS.BACH@INRIA.FR
SIERRA project-team, INRIA - Département d'Informatique de l'Ecole Normale Supérieure, Paris

**Jean-Philippe Vert**                                                      JEAN-PHILIPPE.VERT@MINES-PARISTECH.FR
CBIO Mines ParisTech, INSERM U900, Institut Curie

## 1. Theoretical guarantees for lifted trace norm regularized estimation

We reformulated the "trace + 1" penalty using a linear mapping $\Pi$ and introduced a new penalty, the block norm, using $\Phi$. Using the general formalism of *lifted trace norms* we state theoretical results that help us better understand the behaviour of each of the two norms and compare them more easily. Due to space constraints, all proofs are postponed to appendices available as supplementary materials.

### 1.1. Lifted trace norms

We call *lifting* a linear mapping $\Lambda : \mathbb{R}^{n \times m} \to \mathbb{R}^{n' \times m'}$ and call the penalty induced by $\|\Lambda(X)\|_*$ on the matrix $X$ the $\Lambda$-trace or *lifted trace norm*. Such penalties have been used in compressed sensing (Hosseini Kamal & Vandergheynst, 2013), in statistics (Grave et al., 2011), and have similarities with fused sparsity inducing type of penalties $\|\Lambda(X)\|_1$ studied for instance by Dalalyan & Chen (2012); Vert & Bleakley (2010); Vaiter et al. (2012). Note that a lifted trace norm is not necessarily a norm. It verifies triangle inequality and positive homogeneity, but only separates points so becomes a norm if $\Lambda$ is injective (*i.e.*, $\Lambda(X) = 0 \Rightarrow X = 0$). We denote by $\|\Lambda\| = \max_{\|X\|_F \leq 1} \|\Lambda(X)\|_F$ the operator norm of the linear map $\Lambda$. The mapping $\Lambda^*$ denotes the adjoint operator of $\Lambda$. If $\Lambda(X) = U_{\Lambda(X)} \Sigma_{\Lambda(X)} V_{\Lambda(X)}^\top$ is the singular value decomposition of $\Lambda(X)$, the subgradient of the $\Lambda$-trace at $X$ is given by

$$\partial \|\Lambda(X)\|_* = \Big\{ \Lambda^* \left( U_{\Lambda(X)} V_{\Lambda(X)}^\top + \mathcal{P}_{\Lambda(X)}^\perp(Z) \right) \quad \text{where}$$
$$Z \in \mathbb{R}^{N \times M} \quad \text{and} \quad \|Z\|_{op} \leq 1 \Big\} \quad .$$

From this expression one can see that when $\Lambda(X)$ is rank deficient then $\|\Lambda(X)\|_*$ is nondifferentiable, in cases where the image of $\Lambda^*$ is the whole space $\mathbb{R}^{n \times m}$. This makes the rank of $\Lambda(X)$ a particularly interesting quantity in this context.

In the following $X^\star$ denotes the target matrix to be estimated and $\omega : \mathbb{R}^{n \times m} \to \mathbb{R}^d$ a set of linear measurements:

$$\omega(X) = \Big( \langle \Omega_1, X \rangle, \cdots, \langle \Omega_d, X \rangle \Big)^\top \quad .$$

We call the $\Omega_i$s design matrices and we will be interested in the estimation procedures (i) minimizing the least squares loss $\ell(X) = \frac{1}{d} \|\omega(X) - y\|_2^2$ penalized with lifted trace norm and (ii) minimizing the $\Lambda$-trace subject to $\omega(X) = \omega(X^\star)$.

### 1.2. Least squares regression with lifted trace-norm penalty

We consider linear regression and prove oracle inequalities for the estimation procedure using techniques introduced by Koltchinskii et al. (2011). That is, we consider the model

$$y = \omega(X^\star) + \epsilon \in \mathbb{R}^d$$

where $\epsilon \in \mathbb{R}^d$ having i.i.d zero mean entries.

**Assumption 1** *We assume that the lifting $\Lambda$ is orthogonal, that is $\Lambda^* \Lambda = \|\Lambda\|^2 Id$, which is for instance the case of $\Phi$ and $\Pi$.*

For the two orthogonal liftings of interest $\Pi$ and $\Phi$, the operator norms respectively are given by $\|\Pi\|^2 = (1-\beta)^2 + \beta^2$ and $\|\Phi\|^2 = (n+m)(1-\beta)^2 + \beta^2$.

**Definition 1** *The* cone of restriction $\mathcal{C}(X, \kappa, \Lambda)$ *is the*

set of matrices $B \in \mathbb{R}^{n \times m}$ satisfying

$$\|\mathcal{P}_{\Lambda(X)}^{\perp}(\Lambda(B))\|_* \leq \kappa \|\mathcal{P}_{\Lambda(X)}(\Lambda(B))\|_* \quad . \qquad (1)$$

*The* restricted eigenvalue *of* $\omega$ at $X$ is

$$\mu_{\kappa,\Lambda}(X) = \inf \left\{ \mu > 0 \quad such\ that \right.$$

$$\left. \|\mathcal{P}_{\Lambda(X)}(\Lambda(B))\|_F \leq \frac{\mu}{\sqrt{d}} \|\omega(B)\|_2 \ , \quad \forall B \in \mathcal{C}(X, \kappa, \Lambda) \right\} \quad .$$

Define the objective

$$\mathcal{L}(X) = \frac{1}{d}\|\omega(X) - y\|_2^2 + \lambda\|\Lambda(X)\|_* \ , \qquad (2)$$

and consider the following estimation procedure

$$\widehat{X} = \arg\min_{X \in \mathcal{S}} \mathcal{L}(X) \ , \qquad (3)$$

where $\mathcal{S} \subset \mathbb{R}^{n \times m}$ is the convex cone of admissible solutions. We can state the following oracle inequality on the estimate $\hat{X}$.

**Proposition 1** *Under Assumption 1, for* $\lambda \geq \frac{3}{d}\|\Lambda(M)\|_{op}/\|\Lambda\|^2$, *where* $M = \sum_{i=1}^d \epsilon_i \Omega_i$, *the following holds:*

$$\|\omega(\widehat{X} - X^\star)\|_2^2 \leq$$

$$\inf_{X \in \mathcal{S}} \left\{ \|\omega(X - X^\star)\|_2^2 + \lambda^2 \mu_{5,\Lambda}(X)^2 \operatorname{rank}(\Lambda(X)) \right\} \quad .$$

Note that as (see the proof) $\widehat{X} - X^\star \in \mathcal{C}(X^\star, 5, \Lambda)$ and by orthogonality of $\Lambda$, we bound the estimation error by the prediction error $\|\widehat{X} - X^\star\|_F^2 \leq \frac{36\mu_{5,\Lambda}(X)^2 \operatorname{rank}(\Lambda(X^\star))}{\|\Lambda\|^2 d} \|\omega(\widehat{X} - X^\star)\|_2^2$ and hence the oracle inequality of Proposition 1 provides a abound on the estimation error.

We point out that using similar techniques, and under the stronger assumption called *Restricted Isometry Property* that assumes there exists $\mu > 0$ such that for any $X_1, X_2 \in \mathcal{S}$

$$\frac{1}{d}\|\omega(X_1 - X_2)\|_2^2 \geq \mu^{-2}\|X_1 - X_2\|_F^2 \ ,$$

one can state that for $\lambda \geq \frac{2}{d}\|\Lambda(M)\|_{op}/\|\Lambda\|^2$, we have

$$\mu^{-2}\|\widehat{X} - X^\star\|_F^2 \leq \|\omega(\widehat{X} - X^\star)\|_2^2 \leq$$

$$\inf_{X \in \mathcal{S}} \left\{ \|\omega(X - X^\star)\|_2^2 + \mu^2 c_0^2 \lambda^2 \operatorname{rank}(\Lambda(X)) \right\}$$

where $c_0 = \frac{\sqrt{2}+1}{2}$ and the first inequality being true if $X^\star \in \mathcal{S}$. In particular in the case of denoising

$\omega = id, y = X^\star + M$ considered for instance by <span style="color:blue">Chandrasekaran & Jordan</span> (2012), this proves that if $\lambda \geq \frac{2}{nm}\|\Lambda(M)\|_{op}/\|\Lambda\|^2$

$$\frac{1}{\sqrt{nm}}\|\widehat{X} - X^\star\|_F \leq c_0 \lambda \sqrt{\operatorname{rank}(\Lambda(X^\star))} \ .$$

### 1.3. Probabilistic results

The theoretical analysis of penalized estimation procedures by a norm highlights that when the dual norm of the noise is low the result is more attractive. This motivates us to understand the behavior of $\|\Lambda(G)\|_{op}$ where $G$ denotes the noise which we assume to Gaussian in this work. To this end let us first define the variance of a lifting using canonical matrices $E_{i,j}$ having 1 at the $(i,j)$ entry and 0 everywhere else as

$$v_\Lambda^2 = \|\sum_{i,j} \Lambda(E_{i,j})\Lambda(E_{i,j})^\top\|_{op} \vee \|\sum_{i,j} \Lambda(E_{i,j})^\top \Lambda(E_{i,j})\|_{op} \ .$$

Using results stated in (Tropp, 2010), we know that for a matrix $G$ having i.i.d. centered normal entries

$$\mathbb{E}\left[\|\Lambda(G)\|_{op}\right] \leq \sqrt{2v_\Lambda^2 \log(N + M)} \ ,$$

and we can control the deviation for $t > 0$ as

$$\mathbb{P}\left[\|\Lambda(G)\|_{op} \geq \sqrt{2v_\Lambda^2(\log(N + M) + t)}\right] \leq e^{-t} \ .$$

We can bound the $\Pi$s variance $v_\Pi^2(\beta) \leq \beta^2 \vee n(1 - \beta)^2$ and observe that by setting $\beta = \frac{\sqrt{n}}{1+\sqrt{n}}$ we get the upper bound on the expectation over standard normal matrices $G$

$$\mathbb{E}\|\Pi(G)\|_{op} \leq \sqrt{\frac{2n}{(1 + \sqrt{n})^2}\log(n + m + 2nm)} \ .$$

The variance of $\Phi$ can be controlled by $v_\Phi^2(\beta) \leq (1 + n)(1 - \beta)^2 + \beta^2$, which suggests to set $\beta = \frac{n+1}{n+2}$ in order to obtain

$$\mathbb{E}\|\Phi(G)\|_{op} \leq 2\sqrt{\frac{n + 1}{n + 2}\log(n + m)} \ .$$

We also define the observable variance under the linear map $\omega$ as

$$v_{\omega,\Lambda}^2 = \frac{1}{d}\left\|\sum_{i=1}^d \Lambda(\Omega_i)\Lambda(\Omega_i)^\top\right\|_{op} \vee \left\|\sum_{i=1}^d \Lambda(\Omega_i)^\top \Lambda(\Omega_i)\right\|_{op} \ ,$$

which is a function of $\beta$ for $\Pi$ and $\Phi$ and equal to $\frac{1}{nm}v_\Lambda^2$ in case of denoising $\omega = id$. We finally assume the noise vector elements $\epsilon_i$ are independently drawn from $\mathcal{N}(0, \sigma^2)$.

**Corollary 1 (Block norm)** *Consider the $\Phi$-trace penalty and calibrate for $t > 0$*

$$\lambda = \frac{6\sigma v_{\Phi,\omega}}{\beta^2 + (n+m)(1-\beta)^2}\sqrt{\frac{\log(n+m)+t}{d}} \ ,$$

*then with probability at least $1 - e^{-t}$,*

$$\|\omega(\widehat{X} - X^\star)\|_2^2 \leq \inf_{X \in \mathcal{S}} \left\{ \|\omega(X - X^\star)\|_2^2 \right.$$
$$\left. + c^2 \frac{\log(n+m)+t}{d} \mathbf{ranksity}(X) \right\} \ ,$$

*where $c = \frac{6\sigma\sigma v_{\Phi,\omega}\mu_{5,\Phi}(X)}{\beta^2+(n+m)(1-\beta)^2}$ depends on $\beta$.*

**Corollary 2 (Trace + 1)** *Consider the $\Pi$-trace penalty and calibrate for $t > 0$*

$$\lambda = \frac{3\sigma v_{\Pi,\omega}}{\beta^2 + (1-\beta)^2}\sqrt{\frac{2\log(n+m+2nm)+t}{d}} \ ,$$

*then with probability at least $1 - e^{-t}$,*

$$\|\omega(\widehat{X} - X^\star)\|_2^2 \leq \inf_{X \in \mathcal{S}} \left\{ \|\omega(X - X^\star)\|_2^2 \right.$$
$$\left. + c^2 \frac{\log(n+m+2nm)+t}{d}\Big(\mathrm{rank}(X) + \|X\|_0\Big) \right\} \ ,$$

*where $c = \frac{3\sqrt{2}\sigma v_{\Pi,\omega}\mu_{5,\Pi}(X)}{\beta^2+(1-\beta)^2}$.*

In both cases it is the minimizer of respectively $\beta \mapsto \frac{v_{\Phi,\omega}(\beta)}{\beta^2+(n+m)(1-\beta)^2}$ and $\frac{v_{\Pi,\omega}(\beta)}{\beta^2+(1-\beta)^2}$ that calibrates $\beta$. The two corollaries are interesting because they show that after a natural calibration of the tuning parameter $\lambda$, the convex estimation procedure (3) outputs the optimal estimators for the nonconvex penalties rank $+\ell_0$ and **ranksity**, respectively. In addition the multiplicative factor behind these estimators sharply reminds us of known optimal rates, such as $(\log n)/p$ for the Lasso.

### 1.4. Compressed sensing and exact recovery

Consider the constrained convex optimization problem

$$\min_X \|\Lambda(X)\|_* \ \text{ s.t. } \ \omega(X) = \omega(X^\star) \ , \qquad (4)$$

where the design matrices $\Omega_i$ are i.i.d. Gaussians. We have the following bound on the minimum required such observations for perfect recovery of $X^\star$.

**Proposition 2** *The minimum required number of Gaussian i.i.d. observations for achieving perfect recovery of $X^\star$ with overwhelming probability by solving (4) where $\Lambda$ is an orthogonal lifting is at most*

$$d_\Lambda = \mathbb{E}\left[\|\mathcal{P}^\perp_{\Lambda(X^\star)}(\Lambda(G))\|_{op}^2\right]\mathrm{rank}(\Lambda(X^\star)) + 1 \ ,$$

*the expectation being taken over the set of i.i.d. standard normal matrices $G$.*

In the case of the orthogonal lifting $\Phi$, the quantity $\|\mathcal{P}^\perp_{\Phi(X^\star)}(\Phi(G))\|_{op}$ can be naively bounded by $\|\Phi(G)\|_{op}^2$ for which we already have an upper bound.

**Corollary 3 (Block norm)** *For the $\Phi$-trace penalty, by taking $\beta = (n+1)/(n+2)$, $d_\Phi \leq 1 + 4\,\mathbf{ranksity}(X^\star)\,\log(n+m)$ i.i.d. Gaussian observations are enough to achieve with overwhelming probability perfect recovery of $X^\star$ by solving (4).*

For $\Pi$ the situation is simpler as we have a better understanding of the behavior of $\mathcal{P}^\perp_{\Pi(X^\star)}(\Pi(G))$. In fact

$$\|\mathcal{P}^\perp_{\Pi(X^\star)}(\Pi(G))\|_{op} =$$
$$\left\|\begin{pmatrix} (1-\beta)\mathcal{P}^\perp_{X^\star}(G) & 0 \\ 0 & \beta\,\mathrm{Diag}(\mathbf{vec}(\mathcal{Q}^\perp_{X^\star}(G))) \end{pmatrix}\right\|_{op} \ .$$

allows us to analyze the terms separately and state

**Corollary 4 (Trace + 1)** *In the case of $\Pi$-trace penalty, take $\beta = 1 - \frac{1}{\sqrt{n+m-2r}}$, and assume $r < m-2$, we have*

$$d_\Pi \leq 1 + c_1(r+s)\log\left(c_2 + \frac{nm-s}{2}\right)$$

*where $c_1 = \frac{8}{3}$ and $c_2 = 1 + e^{\frac{3}{4\beta^2}} \leq 2.3$.*

On a bi-clique of size $(k,l)$ we get $d_\Pi \leq c_1 kl\log(nm-s)$ and $d_\Phi \leq 4\{(n+m-1) + (k-1)(l-1)\}\log(n+m)$.

## References

Chandrasekaran, V. and Jordan, M. I. Computational and statistical tradeoffs via convex relaxation. *Preprint*, 2012.

Chandrasekaran, V., Recht, B., Parrilo, P.A., and Willsky, A.S. The convex geometry of linear inverse problems. *Foundations of Computational Mathematics*, 12, 2012.

Dalalyan, A. S. and Chen, Y. Fused sparsity and robust estimation for linear models with unknown variance. In *NIPS 2012*, 2012.

Grave, E., Obozinski, G., and Bach, F. Trace lasso: a trace norm regularization for correlated designs. In *Advances in Neural Information Processing Systems 24*, pp. 2187–2195, 2011.

Hosseini Kamal, M. and Vandergheynst, P. Joint low-rank and sparse light field modeling for dense multiview data compression. *ICASSP*, 2013.

Koltchinskii, V., Lounici, K., and Tsybakov, A. Nuclear norm penalization and optimal rates for noisy matrix completion. *Annals of Statistics*, 2011.

Tropp, J. A. User-friendly tail bounds for sums of random matrices. *ArXiv e-prints*, April 2010.

Vaiter, S., Peyré, G., Dossal, C., and Fadili, J. Robust sparse analysis regularization. *to appear in IEEE Transactions on Information Theory*, 2012.

Vert, J.-P and Bleakley, K. Fast detection of multiple change-points shared by many signals using group lars. *Advances in Neural Information Processing Systems 23 (NIPS)*, pp. 2343–2351, 2010.

**Proof of Proposition 1.** Pick $X \in \mathcal{S}$, in the convex cone of admissible solutions. Let $\mathcal{P}_{\Lambda(X)}$ denote the projector onto $\mathbf{span}(\Lambda(X))$. We start by setting some technical lemmas.

**Lemma 1** *For all $M \in \mathbb{R}^{n \times n}$, we have*

$$\langle M, \hat{X} - X \rangle \leq$$
$$\|\mathcal{P}_{\Lambda(X)}(\Lambda(M))\|_* \|\mathcal{P}_{\Lambda(X)}(\Lambda(\hat{X} - X))\|_{op}/\|\Lambda\|^2$$
$$+ \|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(M))\|_{op} \|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X} - X))\|_*/\|\Lambda\|^2$$

*and*

$$\langle M, \hat{X} - X \rangle \leq \tag{5}$$
$$\sqrt{2 \ \mathrm{rank}(\Lambda(M))} \|\mathcal{P}_{\Lambda(X)}(\Lambda(M))\|_{op} \tag{6}$$
$$\|\mathcal{P}_{\Lambda(X)}(\Lambda(\hat{X} - X))\|_F/\|\Lambda\|^2 \tag{7}$$
$$+ \|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(M))\|_{op} \tag{8}$$
$$\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X} - X))\|_*/\|\Lambda\|^2 \tag{9}$$

**Lemma 2** *There exists $Z \in \partial\|\Lambda(X)\|_*$ such that*

$$-\langle Z, \hat{X} - X \rangle \leq$$
$$\sqrt{\mathrm{rank}(\Lambda(X))}\|\hat{X} - X\|_F\|\Lambda\| - \|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X}))\|_*$$

*and*

$$-\langle Z, \hat{X} - X \rangle \leq \|\mathcal{P}_{\Lambda(X)}(\Lambda(\hat{X} - X))\|_* - \|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X}))\|_* . \tag{10}$$

**Lemma 3** *Let $M = \sum_{i=1}^d \epsilon_i \Omega_i$, we have*

$$\nabla\|\omega(\hat{X}) - y\|_2^2 = 2\langle \omega(\hat{X} - X^\star), \omega(\hat{X} - X)\rangle - 2\langle M, \hat{X} - X\rangle . \tag{11}$$

By optimality, an element of the subgradient of $\mathcal{L}$ at $\hat{X}$ belongs to the normal cone of $\mathcal{S}$ at $\hat{X}$. We have $\langle \partial\mathcal{L}(\hat{X}), \hat{X} - X\rangle \leq 0$. On the other hand, by the monotonicity of the subgradient of the convex function $\|\Lambda(\cdot)\|_*$ we have $\langle \hat{X} - X, \hat{Z} - Z\rangle \geq 0$. Therefore we can deduce by using Lemma 3, that for $M = \sum_{i=1}^d \epsilon_i \Omega_i$,

$$\langle \partial\mathcal{L}(\hat{X}), \hat{X} - X\rangle - \lambda\langle \hat{Z} - Z, \hat{X} - X\rangle \leq 0 \tag{12}$$
$$\Leftrightarrow \quad \langle \frac{1}{d}\nabla\|\omega(\hat{X}) - y\|_2^2 + \lambda Z, \hat{X} - X\rangle \leq 0 \tag{13}$$
$$\Leftrightarrow \quad \frac{2}{d}\langle \omega(\hat{X} - X^\star), \omega(\hat{X} - X)\rangle \leq \tag{14}$$
$$\frac{2}{d}\langle M, \hat{X} - X\rangle - \lambda\langle Z, \hat{X} - X\rangle . \tag{15}$$

We recall the identity

$$2\langle \omega(\hat{X} - X^\star), \omega(\hat{X} - X)\rangle =$$
$$\|\omega(\hat{X} - X^\star)\|_2^2 + \|\omega(\hat{X} - X)\|_2^2 - \|\omega(X - X^\star)\|_2^2 .$$

It shows that if $\langle \omega(\hat{X} - X^\star), \omega(\hat{X} - X)\rangle \leq 0$, then the bound trivially holds. So lets assume $\langle \omega(\hat{X} - X^\star), \omega(\hat{X} - X)\rangle > 0$.

In this case the bound (10) in Lemma 2 and equation (15) imply

$$\lambda\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X}))\|_* \leq \frac{2}{d}\langle M, \hat{X} - X\rangle + \lambda\|\mathcal{P}_{\Lambda(X)}(\Lambda(\hat{X} - X))\|_* . \tag{16}$$

By using Lemma 1 , first inequality (5), we have

$$(\lambda - \frac{2}{d}\frac{\|\Lambda(M)\|_{op}}{\|\Lambda\|^2})\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X} - X))\|_*$$
$$\leq (\lambda + \frac{2}{d}\frac{\|\Lambda(M)\|_{op}}{\|\Lambda\|^2})\|\mathcal{P}_{\Lambda(X)}(\Lambda(\hat{X} - X))\|_* .$$

This shows that for $\lambda \geq \frac{3}{d}\|\Lambda(M)\|_{op}/\|\Lambda\|^2$, by using the fact that for $x \geq 3, \frac{x-2}{x+2} \geq \frac{1}{5}$, the following holds true

$$\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X} - X))\|_* \leq 5\|\mathcal{P}_{\Lambda(X)}(\Lambda(\hat{X} - X))\|_* .$$

As a consequence, $\hat{X} - X \in \mathcal{C}(X, 5, \Lambda)$. On the other hand, by using Lemma 1, second inequality (9) and (15) we have

$$\frac{1}{d}\Big(\|\omega(\widehat{X}-X^\star)\|_2^2+\|\omega(\widehat{X}-X)\|_2^2-\|\omega(X-X^\star)\|_2^2\Big)$$

$$\leq \frac{2}{d}\Big(\sqrt{2\,\mathrm{rank}(\Lambda(X))}\frac{\|\Lambda(M)\|_{op}}{\|\Lambda\|^2}\|\mathcal{P}_{\Lambda(X)}(\Lambda(\widehat{X}-X))\|_F$$

$$+\frac{\|\Lambda(M)\|_{op}}{\|\Lambda\|^2}\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\widehat{X}))\|_*\Big)$$

$$+\lambda\sqrt{\mathrm{rank}(\Lambda(X))}\|\mathcal{P}_{\Lambda(X)}(\Lambda(\widehat{X}-X))\|_F-\lambda\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\widehat{X}))\|_* \ .$$

By using the definition of the restricted eigenvalue $\mu(X)=\mu_{5,\Lambda}(X)$, given that $\widehat{X}-X\in\mathcal{C}(X,5,\Lambda)$,

$$\|\mathcal{P}_{\Lambda(X)}(\Lambda(\widehat{X}-X))\|_F\leq\frac{\mu(X)}{\sqrt{d}}\|\omega(\hat{X}-X)\|_2$$

so we can write, again thanks to $\lambda\geq\frac{3}{d}\|\Lambda(M)\|_{op}/\|\Lambda\|^2$,

$$\frac{1}{d}\Big(\|\omega(\widehat{X}-X^\star)\|_2^2+\|\omega(\widehat{X}-X)\|_2^2-\|\omega(X-X^\star)\|_2^2\Big)$$

$$\leq\frac{\mu(X)}{\sqrt{d}}\lambda\sqrt{\mathrm{rank}(\Lambda(X))}\Big(1+\frac{2\sqrt{2}}{3}\Big)\|\omega(\widehat{X}-X)\|_F \ .$$

So by $bx-x^2\leq\big(\frac{b}{2}\big)^2$ we finally get

$$\frac{1}{d}\|\omega(\widehat{X}-X^\star)\|_2^2\leq$$

$$\frac{1}{d}\|\omega(X-X^\star)\|_2^2+\lambda^2\frac{\mu(X)^2}{4d}(1+\frac{2\sqrt{2}}{3})^2\,\mathrm{rank}(\Lambda(X)) \ . \ \square$$

**Proof of Lemma 1** Let us decompose $\Lambda(M)$ onto the direct sum formed by the span of $\Lambda(X)$ and the orthogonal space:

$$\Lambda(M)=\mathcal{P}_{\Lambda(X)}(M)+\mathcal{P}_{\Lambda(X)}^\perp(M) \ .$$

By using assumption 1 and Holder's inequality twice

$$\langle M,\hat{X}-X\rangle=\langle\Lambda(M),\Lambda(\hat{X}-X)\rangle/\|\Lambda\|^2\leq$$

$$\|\mathcal{P}_{\Lambda(X)}(\Lambda(M))\|_*\|\mathcal{P}_{\Lambda(X)}(\Lambda(\hat{X}-X))\|_{op}/\|\Lambda\|^2$$

$$+\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(M))\|_{op}\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X}))\|_*/\|\Lambda\|^2$$

The other bound is obtained in a similar fashion by using Cauchy-Schwarz on the first term and also the fact that $\|\mathcal{P}_{\Lambda(X)}(M)\|_F\leq\sqrt{2\ \mathrm{rank}(\Lambda(X))}\|M\|_F$ since we can write $\mathcal{P}_{\Lambda(X)}(M)=(I-UU^\top)MVV^\top+UU^\top M$ for $U$ and $V$ singular vectors of $\Lambda(X)$.

**Proof of Lemma 2.** Let

$$Z=\Lambda^*\Big(U_{\Lambda(X)}V_{\Lambda(X)}^\top+\mathcal{P}_{\Lambda(X)}^\perp(W)\Big)$$

denote an element of the subgradient of $\|\Lambda(\cdot)\|_*$, where $\|W\|_{op}\leq 1$ . Take $W=-UV^\top$ where $U\Sigma V^\top=\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X}))$ is a singular value decomposition, then $\|W\|_{op}=1$ and

$$\langle\Lambda^*(\mathcal{P}_{\Lambda(X)}^\perp(W)),\hat{X}-X\rangle=$$

$$\langle\mathcal{P}_{\Lambda(X)}^\perp(W),\Lambda(\hat{X}-X)\rangle=$$

$$-\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\hat{X}))\|_*$$

(17)

so we can write

$$-\langle Z,\widehat{X}-X\rangle=$$

$$-\langle\Lambda^*\Big(U_{\Lambda(X)}V_{\Lambda(X)}^\top\Big),\hat{X}-X\rangle$$

$$+\langle\Lambda^*(\mathcal{P}_{\Lambda(X)}^\perp(W)),\hat{X}-X\rangle=$$

$$-\langle U_{\Lambda(X)}V_{\Lambda(X)}^\top,\Lambda(\hat{X}-X)\rangle$$

$$+\langle\mathcal{P}_{\Lambda(X)}^\perp(W),\Lambda(\hat{X}-X)\rangle=$$

$$-\langle U_{\Lambda(X)}V_{\Lambda(X)}^\top,\Lambda(\hat{X}-X)\rangle$$

$$-\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\widehat{X}))\|_* \ .$$

We know that $\|U_{\Lambda(X)}V_{\Lambda(X)}^\top\|_F^2\leq\mathrm{rank}(\Lambda(X))$. By Cauchy-Schwarz

$$-\langle Z,\widehat{X}-X\rangle\leq\sqrt{\mathrm{rank}(\Lambda(X))}\|\widehat{X}-X\|_F\|\Lambda\|-\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\widehat{X}))\|_* \ .$$

Similarly if we use Holder's instead of Cauchy-Schwarz, and thanks to $\|U_{\Lambda(X)}V_{\Lambda(X)}^\top\|_{op}=1$ ,

$$-\langle Z,\widehat{X}-X\rangle\leq$$

$$\|\mathcal{P}_{\Lambda(X)}(\Lambda(\widehat{X}-X))\|_*-\|\mathcal{P}_{\Lambda(X)}^\perp(\Lambda(\widehat{X}))\|_* \ . \ \square$$

**Proof of Lemma 3.**

Given that $\nabla\|\omega(\widehat{X})-y\|_2^2=2\sum_{i=1}^d\Omega_i\langle\Omega_i,\widehat{X}\rangle-y_i\Omega_i$, we obtain

$$\langle\nabla\|\omega(\widehat{X})-y\|_2^2,\widehat{X}-X\rangle$$

$$=2\sum_{i=1}^d\langle(\langle\Omega_i,\widehat{X}\rangle-y_i)\Omega_i,\widehat{X}-X\rangle$$

$$=2\sum_{i=1}^d(\langle\Omega_i,\widehat{X}\rangle-y_i)\langle\Omega_i,\widehat{X}-X\rangle$$

$$=2\langle\omega(\hat{X})-y,\omega(\widehat{X}-X)\rangle$$

$$=2\langle\omega(\hat{X}-X^\star)+\omega(X^\star)-y,\omega(\widehat{X}-X)\rangle$$

$$=2\langle\omega(\widehat{X}-X^\star),\omega(\widehat{X}-X)\rangle-2\langle\epsilon,\omega(\widehat{X}-X)\rangle$$

$$=2\langle\omega(\widehat{X}-X^\star),\omega(\widehat{X}-X)\rangle-2\langle M,\widehat{X}-X\rangle \ . \ \square$$

**Proof of Proposition 2.** By orthogonality of $\Lambda$ we have

$$\|\Lambda\|^2 G = \Lambda^* \Lambda(G) = \Lambda^* \left( \mathcal{P}_{\Lambda(X^\star)}(\Lambda(G)) + \mathcal{P}^\perp_{\Lambda(X^\star)}(\Lambda(G)) \right)$$

Lets build an appropriate element of the normal cone of the $\Lambda$-trace at $X^\star$

$$Z(G) = \frac{1}{\|\Lambda\|^2} \Lambda^*(\mathcal{P}^\perp_{\Lambda(X^\star)}(\Lambda(G))) +$$
$$\frac{\|\mathcal{P}^\perp_{\Lambda(X^\star)}(\Lambda(G))\|_{op}}{\|\Lambda\|^2} \Lambda^* \left( U_{\Lambda(X^\star)} V^\perp_{\Lambda(X^\star)} \right) \quad,$$

and get by Cauchy-Schwarz inequality

$$\|Z(G) - G\|_F^2 = \frac{\|\mathcal{P}^\perp_{\Lambda(X^\star)}(\Lambda(G))\|_{op}^2}{\|\Lambda\|^2} \|\Lambda^* U_{\Lambda(X^\star)} V^\perp_{\Lambda(X^\star)}\|_F^2$$
$$\leq \|\mathcal{P}^\perp_{\Lambda(X^\star)}(\Lambda(G))\|_{op}^2 \operatorname{rank}(\Lambda(X^\star)) \quad.$$

By Lemma 2.7 in (Chandrasekaran et al., 2012) this bounds the squared gaussian width of the tangent cone to $\|\Lambda(\cdot)\|_*$ at $X^\star$ intersected with the unit sphere. We conclude by using Corollary 3.3 from the same paper. $\square$

**Proof of Corollary 4**

Let $s = \|X^\star\|_0$ and $r = \operatorname{rank}(X^\star)$. First lets show that for any $G \in \mathbb{R}^{n \times m}$

$$\|\mathcal{P}^\perp_{\Pi(X^\star)}(\Pi(G))\|_{op} =$$
$$\left\| \begin{pmatrix} (1-\beta)\mathcal{P}^\perp_{X^\star}(G) & 0 \\ 0 & \beta \operatorname{Diag}(\mathbf{vec}(\mathcal{Q}^\perp_{X^\star}(G))) \end{pmatrix} \right\|_{op} .$$

In fact as the singular value decomposition of $\Pi(X^\star)$ can be written (up to permutations of rows and columns) using the matrices

$$U_{\Pi(X^\star)} = \begin{pmatrix} U_{X^\star} & 0 \\ 0 & \operatorname{Diag}(\mathbf{vec}(\operatorname{sgn}(X^\star))) \end{pmatrix}$$

and

$$V_{\Pi(X^\star)} = \begin{pmatrix} V_{X^\star} & 0 \\ 0 & \operatorname{Diag}(\mathbf{vec}(|\operatorname{sgn}(X^\star)|)) \end{pmatrix}$$

the formula $\mathcal{P}^\perp(Z) = (I - UU^\top)Z(I - VV^\top)$ implies the result. Since the gaussian distribution is isotropic

we know that $\|\mathcal{P}^\perp_{\Pi(X^\star)}(G)\|_{op}$ is distributed as the operator norm of a $(n-r) \times (m-r)$ gaussian matrix and that $\|\mathcal{Q}^\perp_{\Pi(X^\star)}(G)\|_\infty$ is distributed as the $\ell_\infty$ norm of a vector of length $nm - s$ having iid standard normal entries.

Let $J = \mathcal{Q}^\perp_{X^\star}(G)$ and $H = \mathcal{P}^\perp_{X^\star}(G)$ and

$$z = \max \left\{ (1-\beta)^2 \|H\|_{op}^2 \ , \ \beta^2 \|J\|_\infty^2 \right\} \quad,$$

and notice that by Jensen inequality, for all $t > 0$

$$\exp\left( t \, \mathbb{E}[z] \right) \leq \mathbb{E} \exp(tz)$$
$$\leq \mathbb{E} \exp(t(1-\beta)^2 \|H\|_{op}^2) + \sum_{i=1}^{nm-s} \mathbb{E} \exp(t\beta^2 J_i)$$
$$= \mathbb{E} \exp(t(1-\beta)^2 \|H\|_{op}^2) + \frac{nm - s}{\sqrt{1 - 2t\beta^2}} \quad,$$

where $J_i$s are iid $\chi^2$ variables and the last relation being the moment generating function of $\chi^2$. For bounding the term $\mathbb{E} \exp(t(1-\beta)^2 \|H\|_{op}^2)$, let us recall

$$\mathbb{P}[\|H\|_{op} > \sqrt{n-r} + \sqrt{m-r} + s] \leq \exp(-s^2/2)$$

and introduce $f(x) = \exp(t(1-\beta)^2 x^2)$. We have $f^{-1}(z) = \frac{1}{1-\beta} \sqrt{\frac{\log(z)}{t}}$ strictly increasing $[1; \infty) \to \mathbb{R}$. Denoting $R = \sqrt{n-r} + \sqrt{m-r}$ we have the sequence of inequalities

$$\mathbb{E} \exp(t(1-\beta)^2 \|H\|_{op}^2) \tag{18}$$
$$= \mathbb{E} f(\|H\|_{op}) \tag{19}$$
$$= \int_1^\infty \mathbb{P}[f(\|H\|_{op}) > h] \ dh \tag{20}$$
$$\leq \int_1^{1+f(R)} 1 \ dh \tag{21}$$
$$+ \int_{1+f(R)}^\infty \mathbb{P}[f(\|H\|_{op}) > h] dh \tag{22}$$
$$= f(R) \tag{23}$$
$$+ \int_0^\infty \mathbb{P}[\|H\|_{op} > f^{-1}(f(R) + 1 + \zeta)] d\zeta \tag{24}$$
$$\leq f(R) \tag{25}$$
$$+ \int_0^\infty \mathbb{P}[\|H\|_{op} > R + f^{-1}(1 + \zeta)] d\zeta \tag{26}$$
$$\leq f(R) \tag{27}$$
$$+ \int_0^\infty 2ts(1-\beta)^2 \ \exp\left(-s^2/2 + ts^2(1-\beta)^2\right) ds \tag{28}$$
$$\leq f(R) + 1 \tag{29}$$

where (26) is due to the sublinearity of $f^{-1}(z) = \frac{1}{(1-\beta)}\sqrt{\frac{\log(z)}{t}}$ :

$$f^{-1}(z + z') \le f^{-1}(z) + f^{-1}(z')$$

and (29) is true for any $t < \frac{1}{2(1-\beta)^2}$. We have for $t < \frac{1}{2}\min\left(\frac{1}{(1-\beta)^2}, \frac{1}{\beta^2}\right)$,

$$\mathbb{E}[z] \le$$
$$\frac{1}{t}\log\left\{1 + \exp[2t(1-\beta)^2(n + m - 2r)] + \frac{nm - s}{\sqrt{1 - 2t\beta^2}}\right\}$$

By taking $t = \frac{3}{8\beta^2}$ and $(1-\beta)^2 = \frac{1}{n+m-2r}$ the latter expression gives

$$\mathbb{E}[z] \le$$
$$\frac{8\beta^2}{3}\log\left\{1 + e^{\frac{3}{4\beta^2}} + \frac{nm - s}{2}\right\}.$$

The bound in Proposition 4 (skippin 1+)becomes

$$(r + s)\frac{8\beta^2}{3}\log\left\{1 + e^{\frac{3}{4\beta^2}} + \frac{nm - s}{2}\right\}$$
$$\le c_1(r + s)\log\left\{c_2 + \frac{nm - s}{2}\right\}$$

where $c_1 = \frac{8}{3}$ and $c_2 = 1 + e^{\frac{3}{4\beta^2}} \le 2.3$.

**Lemma 4** *The variance (see (Tropp, 2010)) of the set of $\Phi(E_{i,j})$s where $1 \le i \le n$, $1 \le j \le m$ is bounded by*

$$\sigma^2 =$$
$$\|\sum_{i=1}^{n}\sum_{j=1}^{m}\Phi(E_{i,j})\Phi(E_{i,j})^\top\|_{op}$$
$$\vee \|\sum_{i=1}^{n}\sum_{j=1}^{m}\Phi(E_{i,j})^\top\Phi(E_{i,j})\|_{op}$$
$$\le (1 + (n \vee m))(1 - \beta)^2 + \beta^2$$

**Proof of Lemma 4.** Lets recall for $E_{i_1,j_1}^{n_1,m_1}$ and $E_{i_2,j_2}^{n_2,m_2}$ denoting canonical elements of size $n_1 \times m_1$ and $n_2 \times m_2$, the Kronecker product expression:

$$E_{i_1,j_1}^{n_1,m_1} \otimes E_{i_2,j_2}^{n_2,m_2} = E_{(i_1-1)n_2+i_2,(j_1-1)m_2+j_2}^{n_1 n_2,m_1 m_2}.$$

Using this and by expressing $I_n = \sum_{i=1}^{n}E_{i,i}^{n,n}$, after some algebra we get

$$\Phi(E_{i,j})\Phi(E_{i,j})^\top =$$
$$\left((1 - \beta)^2 E_{j,j}^{m,m} \otimes I_n\right.$$
$$\left. + (1 - \beta)^2 I_m \otimes E_{i,i}^{n,n} + \beta^2 E_{i+n(j-1),i+n(j-1)}^{nm,nm}\right).$$

Adding up the terms results in a very simple object:

$$\sum_{i=1}^{n}\sum_{j=1}^{m}\Phi(E_{i,j})\Phi(E_{i,j})^\top$$
$$= (1 - \beta)^2 I_m \otimes I_n + (1 - \beta)^2 I_m \otimes I_n + \beta^2 I_{nm}$$
$$= \left(2(1 - \beta)^2 + \beta^2\right) I_{nm}.$$

The second term is also quite friendly, in fact

$\Phi(E_{i,j})^\top \Phi(E_{i,j}) =$

$$
\begin{pmatrix}
(1-\beta)^2 E_{i,i}^{n,n} \otimes I_n & (1-\beta)^2 E_{in,jm}^{n^2,m^2} & \beta(1-\beta) E_{ni,n(j-1)+i}^{n^2,nm} \\
(1-\beta)^2 E_{jm,in}^{m^2,n^2} & (1-\beta)^2 I_m \otimes E_{j,j}^{m,m} & \beta(1-\beta) E_{mj,n(j-1)+i}^{m^2,nm} \\
\beta(1-\beta) E_{i+n(j-1),ni}^{nm,n^2} & \beta(1-\beta) E_{i+n(j-1),mj}^{nm,m^2} & \beta^2 E_{i+n(j-1),i+n(j-1)}^{nm,nm}
\end{pmatrix}
$$

$$
= \begin{pmatrix}
(1-\beta)^2 \sum_{k=1}^n E_{n(i-1)+k,n(i-1)+k}^{n^2,n^2} & (1-\beta)^2 E_{in,jm}^{n^2,m^2} & \beta(1-\beta) E_{ni,n(j-1)+i}^{n^2,nm} \\
(1-\beta)^2 E_{jm,in}^{m^2,n^2} & (1-\beta)^2 \sum_{k=1}^m E_{j+(k-1)m,j+(k-1)m}^{m^2,m^2} & \beta(1-\beta) E_{mj,n(j-1)+i}^{m^2,nm} \\
\beta(1-\beta) E_{i+n(j-1),ni}^{nm,n^2} & \beta(1-\beta) E_{i+n(j-1),mj}^{nm,m^2} & \beta^2 E_{i+n(j-1),i+n(j-1)}^{nm,nm}
\end{pmatrix}
$$

$$
= \begin{pmatrix}
(1-\beta)^2 E_{ni,ni}^{n^2,n^2} & (1-\beta)^2 E_{in,jm}^{n^2,m^2} & \beta(1-\beta) E_{ni,n(j-1)+i}^{n^2,nm} \\
(1-\beta)^2 E_{jm,in}^{m^2,n^2} & (1-\beta)^2 E_{mj,mj}^{m^2,m^2} & \beta(1-\beta) E_{mj,n(j-1)+i}^{m^2,nm} \\
\beta(1-\beta) E_{i+n(j-1),ni}^{nm,n^2} & \beta(1-\beta) E_{i+n(j-1),mj}^{nm,m^2} & \beta^2 E_{i+n(j-1),i+n(j-1)}^{nm,nm}
\end{pmatrix}
$$

$$
+ \begin{pmatrix}
(1-\beta)^2 \sum_{k\neq i}^n E_{n(i-1)+k,n(i-1)+k}^{n^2,n^2} & 0_{n^2,m^2} & 0_{n^2,nm} \\
0_{m^2,n^2} & (1-\beta)^2 \sum_{k\neq j}^m E_{j+(k-1)m,j+(k-1)m}^{m^2,m^2} & 0_{m^2,nm} \\
0_{nm,n^2} & 0_{nm,m^2} & 0_{nm,nm}
\end{pmatrix}
$$

Adding up the terms we get on the one hand matrices having only diagonal terms (from the second term of the last equality) and on the other hand (first term) pairwise orthogonal matrices which are also orthogonal to the diagonal terms. The second bunch of matrices that can be written, up to row and column permutations, as the following matrix

$$
\begin{pmatrix}
(1-\beta)^2 m & (1-\beta)^2 & \beta(1-\beta) \\
(1-\beta)^2 & (1-\beta)^2 n & \beta(1-\beta) \\
\beta(1-\beta) & \beta(1-\beta) & \beta^2
\end{pmatrix}
= \begin{pmatrix}
1-\beta & 0 & 0 \\
0 & 1-\beta & 0 \\
0 & 0 & \beta
\end{pmatrix}
\begin{pmatrix}
m & 1 & 1 \\
1 & n & 1 \\
1 & 1 & 1
\end{pmatrix}
\begin{pmatrix}
1-\beta & 0 & 0 \\
0 & 1-\beta & 0 \\
0 & 0 & \beta
\end{pmatrix} .
$$

Using triangle inequality

$$
\left\| \begin{pmatrix}
(1-\beta)^2 m & (1-\beta)^2 & \beta(1-\beta) \\
(1-\beta)^2 & (1-\beta)^2 n & \beta(1-\beta) \\
\beta(1-\beta) & \beta(1-\beta) & \beta^2
\end{pmatrix} \right\|_{op}
$$

$$
= \left\| \begin{pmatrix}
1-\beta & 0 & 0 \\
0 & 1-\beta & 0 \\
0 & 0 & \beta
\end{pmatrix}
\left\{ \begin{pmatrix}
m-1 & 0 & 0 \\
0 & n-1 & 0 \\
0 & 0 & 0
\end{pmatrix} + \begin{pmatrix}
1 & 1 & 1 \\
1 & 1 & 1 \\
1 & 1 & 1
\end{pmatrix} \right\}
\begin{pmatrix}
1-\beta & 0 & 0 \\
0 & 1-\beta & 0 \\
0 & 0 & \beta
\end{pmatrix} \right\|_{op}
$$

$$
\leq (1-\beta)^2 (1 + (n \vee m)) + \beta^2 \quad,
$$

so

$$
\sum_{i=1}^n \sum_{j=1}^m \Phi(E_{i,j})^\top \Phi(E_{i,j}) \leq (1-\beta)^2 (1 + (n \vee m)) + \beta^2 \quad . \square
$$