# Bilinear Regression via Convex Programming without Lifting

**Sohail Bahmani**

Georgia Tech

# Bilinear regression

Given $\boldsymbol{a}_1, \ldots, \boldsymbol{a}_n \in \mathbb{R}^{d_1}$ and $\boldsymbol{a}'_1, \ldots, \boldsymbol{a}'_n \in \mathbb{R}^{d_2}$ we observe noisy *bilinear* measurements of unknown vectors $\boldsymbol{x}_\star \in \mathbb{R}^{d_1}$ and $\boldsymbol{x}'_\star \in \mathbb{R}^{d_2}$ as

$$y_1 = \boldsymbol{a}_1^\top \boldsymbol{x}_\star \, \boldsymbol{x}'^\top_\star \boldsymbol{a}'_1 + \xi_1$$
$$y_2 = \boldsymbol{a}_2^\top \boldsymbol{x}_\star \, \boldsymbol{x}'^\top_\star \boldsymbol{a}'_2 + \xi_2$$
$$\vdots \qquad \vdots$$
$$y_n = \boldsymbol{a}_n^\top \boldsymbol{x}_\star \, \boldsymbol{x}'^\top_\star \boldsymbol{a}'_n + \underbrace{\xi_n}_{\text{noise}} \, .$$

### Problem

*Estimate $\boldsymbol{x}_\star$ and $\boldsymbol{x}'_\star$ accurately, using the observations $(\boldsymbol{a}_i, \boldsymbol{a}'_i, y_i)$, $i = 1, \ldots, n$.*

A much more general **nonlinear regression** will be addressed towards the end.

## Some reminders

**Scaling ambiguity**:
$(x, x')$ is an accurate solution $\implies$ $(tx, t^{-1}x')$ is an accurate solution $\forall t \neq 0$.

**Computation**:
There are computationally hard instances. Randomness helps to avoid them.
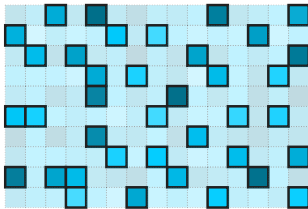
## Examples

In matrix notation,

$$y = Ax_\star \circ A'x'_\star + \xi,$$

where $A = [a_1 \; \cdots \; a_n]^\top$ and $A' = [a'_1 \; \cdots \; a'_n]^\top$.

**Matrix completion**:

$a_i = e_{r_i}$ and $a'_i = e_{c_i}$ are random coordinate indicator vectors



$$X = x_\star x'^\top_\star$$

$$y_i = X_{r_i, c_i} = a_{r_i}^\top X a'_{c_i}$$

# Examples

In matrix notation,

$$y = Ax_\star \circ A'x'_\star + \xi,$$

where $A = [a_1 \ \cdots \ a_n]^\top$ and $A' = [a'_1 \ \cdots \ a'_n]^\top$.

**Blind deconvolution**:

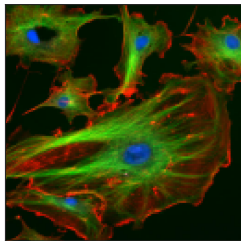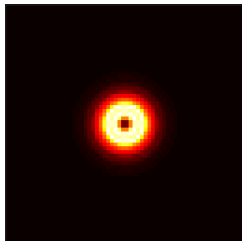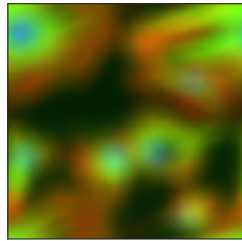$A$ and $A'$ interpreted as subspaces in the *Fourier* domain (ignoring $\mathbb{R}$ vs. $\mathbb{C}$)



image: $F^{-1}Ax_\star$     PSF: $F^{-1}A'x'_\star$     blurred image: $F^{-1}y$

## Related work ...

**SDP-relaxation**[*,†]: The conversion $xx'^{\mathsf{T}} \mapsto X$ maps the bilinear regression to a *matrix linear regression*

$$y_i = a_i^{\mathsf{T}} x_\star x_\star'^{\mathsf{T}} a_i' + \xi_i = a_i^{\mathsf{T}} X_\star a_i' + \xi_i \,.$$

Estimate the rank one solution through *nuclear norm* minimization

$$\underset{X}{\arg\min} \ \|X\|_*$$
$$\text{subject to } a_i^{\mathsf{T}} X a_i' = y_i \,, \quad i = 1, \ldots, n \,.$$

**Not scalable**, due to the prohibitive cost of SDP.

[*]Ahmed et al, "Blind deconvolution using convex programming," *IEEE Trans. Info. Theory*, 2014

[†]Cai & Zhang, "ROP: Matrix recovery via rank-one projections," Annals of Statistics, 2015

# Related work …

**SDP-relaxation**[*,†]: The conversion $\boldsymbol{x}\boldsymbol{x}'^{\mathsf{T}} \mapsto \boldsymbol{X}$ maps the bilinear regression to a *matrix linear regression*

$$y_i = \boldsymbol{a}_i^{\mathsf{T}} \boldsymbol{x}_\star \, \boldsymbol{x}_\star'^{\mathsf{T}} \boldsymbol{a}_i' + \xi_i = \boldsymbol{a}_i^{\mathsf{T}} \boldsymbol{X}_\star \boldsymbol{a}_i' + \xi_i \,.$$

Estimate the rank one solution through *nuclear norm* minimization

$$\underset{\boldsymbol{X}}{\operatorname{argmin}} \, \|\boldsymbol{X}\|_*$$
$$\text{subject to } \boldsymbol{a}_i^{\mathsf{T}} \boldsymbol{X} \boldsymbol{a}_i' = y_i \,, \quad i = 1, \ldots, n \,.$$

**Not scalable**, due to the prohibitive cost of SDP.

[*] Ahmed et al, "Blind deconvolution using convex programming," *IEEE Trans. Info. Theory*, 2014

[†] Cai & Zhang, "ROP: Matrix recovery via rank-one projections," Annals of Statistics, 2015

# Related work ...

**Nonconvex gradient descent**[*,†]: Similar to the *Wirtinger flow algorithm* for phase retrieval, run gradient descent on the residual error

$$f(\boldsymbol{x}, \boldsymbol{x}') = \frac{1}{n} \sum_{i=1}^{n} |\boldsymbol{a}_i^* \boldsymbol{x} \boldsymbol{x}'^* \boldsymbol{a}_i' - y_i|^2.$$

Regularized variants are proposed for blind deconvolution, blind calibration, & matrix completion.

**Light-tailed distribution** is needed for iteration analysis.
Analyses are often lengthy.

[*]Cambareri, Jacques, "Through the haze: A non-convex approach to blind gain calibration for linear random sensing models," Information & Inference, 2018.

[†]Li et al, "Rapid, robust, and reliable blind deconvolution via nonconvex optimization," *App. & Comp. Harmonic Analysis*, 2018.

# Related work …

**Nonconvex gradient descent**[*,†]: Similar to the *Wirtinger flow algorithm* for phase retrieval, run gradient descent on the residual error

$$f(\mathbf{x}, \mathbf{x}') = \frac{1}{n} \sum_{i=1}^{n} |\mathbf{a}_i^* \mathbf{x} \mathbf{x}'^* \mathbf{a}_i' - y_i|^2.$$

Regularized variants are proposed for blind deconvolution, blind calibration, & matrix completion.

**Light-tailed distribution** is needed for iteration analysis.

Analyses are often lengthy.

[*]Cambareri, Jacques, "Through the haze: A non-convex approach to blind gain calibration for linear random sensing models," Information & Inference, 2018.

[†]Li et al, "Rapid, robust, and reliable blind deconvolution via nonconvex optimization," *App. & Comp. Harmonic Analysis*, 2018.

## Related work

**BranchHull**[*]: With $s = \text{sgn}(Ax_\star)$ given, solve the **convex program**

$$\underset{x \in \mathbb{R}^{d_1}, x' \in \mathbb{R}^{d_2}}{\text{argmin}} \quad \|Ax\|_2^2 + \|A'x'\|_2^2$$

$$\text{subject to } s \circ (Ax) \circ (A'x') \geq |y|$$

$$s \circ (Ax) \geq 0 \,.$$

The sample complexity $n \gtrsim d_1 + d_2$ is shown for Gaussian $A$ and $A'$

Knowing $s$ is a **restrictive assumption**.
The relaxation seem **very sensitive to noise**.

[*]Aghasi et al, "BranchHull: Convex bilinear inversion from the entrywise product of signals with known signs,"
`arXiv:1702.04342`, 2017

# Related work

**BranchHull**[*]: With $s = \text{sgn}(Ax_\star)$ given, solve the **convex program**

$$\underset{x \in \mathbb{R}^{d_1}, x' \in \mathbb{R}^{d_2}}{\text{argmin}} \quad \|Ax\|_2^2 + \|A'x'\|_2^2$$

$$\text{subject to } s \circ (Ax) \circ (A'x') \geq |y|$$

$$s \circ (Ax) \geq 0 \,.$$

The sample complexity $n \gtrsim d_1 + d_2$ is shown for Gaussian $A$ and $A'$

Knowing $s$ is a **restrictive assumption**.

The relaxation seem **very sensitive to noise**.

[*]Aghasi et al, "BranchHull: Convex bilinear inversion from the entrywise product of signals with known signs," `arXiv:1702.04342`, 2017

# Our proposed convex relaxation

Given $\boldsymbol{x}_0 \approx \boldsymbol{x}_\star$ and $\boldsymbol{x}_0' \approx \boldsymbol{x}_\star'$ such that for a sufficiently small $\epsilon \geq 0$,

$$\left\| \begin{pmatrix} \boldsymbol{x}_0 - \boldsymbol{x}_\star \\ \boldsymbol{x}_0' - \boldsymbol{x}_\star' \end{pmatrix} \right\|_2 \leq \epsilon \left\| \begin{pmatrix} \boldsymbol{x}_\star \\ \boldsymbol{x}_\star' \end{pmatrix} \right\|_2,$$

we estimate $\boldsymbol{x}_\star$ and $\boldsymbol{x}_\star'$ through the convex program:

$$(\widehat{\boldsymbol{x}}, \widehat{\boldsymbol{x}}') = \operatorname*{argmax}_{\boldsymbol{x}, \boldsymbol{x}'} \boldsymbol{x}_0^\top \boldsymbol{x} + \boldsymbol{x}_0'^\top \boldsymbol{x}' - \frac{2}{n} \sum_{i=1}^n \max \left\{ \frac{1}{4}(\boldsymbol{a}_i^\top \boldsymbol{x} + \boldsymbol{a}_i'^\top \boldsymbol{x}')^2 - y_i, \frac{1}{4}(\boldsymbol{a}_i^\top \boldsymbol{x} - \boldsymbol{a}_i'^\top \boldsymbol{x}')^2 \right\}$$

# Our proposed convex relaxation

Given $\boldsymbol{x}_0 \approx \boldsymbol{x}_\star$ and $\boldsymbol{x}_0' \approx \boldsymbol{x}_\star'$ such that for a sufficiently small $\epsilon \geq 0$,

$$\left\| \begin{pmatrix} \boldsymbol{x}_0 - \boldsymbol{x}_\star \\ \boldsymbol{x}_0' - \boldsymbol{x}_\star' \end{pmatrix} \right\|_2 \leq \epsilon \left\| \begin{pmatrix} \boldsymbol{x}_\star \\ \boldsymbol{x}_\star' \end{pmatrix} \right\|_2,$$

we estimate $\boldsymbol{x}_\star$ and $\boldsymbol{x}_\star'$ through the convex program:

$$(\widehat{\boldsymbol{x}}, \widehat{\boldsymbol{x}}') = \operatorname*{argmax}_{\boldsymbol{x}, \boldsymbol{x}'} \boldsymbol{x}_0^\top \boldsymbol{x} + \boldsymbol{x}_0'^\top \boldsymbol{x}' - \frac{2}{n} \sum_{i=1}^n \max \left\{ \frac{1}{4} (\boldsymbol{a}_i^\top \boldsymbol{x} + \boldsymbol{a}_i'^\top \boldsymbol{x}')^2 - y_i, \frac{1}{4} (\boldsymbol{a}_i^\top \boldsymbol{x} - \boldsymbol{a}_i'^\top \boldsymbol{x}')^2 \right\}$$

**Intuition**: Using the identities

$$2 \max\{u, v\} = |u - v| + u + v, \qquad \begin{aligned} (u+v)^2 - (u-v)^2 &= 4uv \\ (u+v)^2 + (u-v)^2 &= 2(u^2 + v^2) \end{aligned},$$

the objective can be written as

$$\boldsymbol{x}_0^\top \boldsymbol{x} + \boldsymbol{x}_0'^\top \boldsymbol{x}' - \frac{1}{2n} \sum_{i=1}^n (\boldsymbol{a}_i^\top \boldsymbol{x})^2 + (\boldsymbol{a}_i'^\top \boldsymbol{x}')^2 - \frac{1}{n} \sum_{i=1}^n |\boldsymbol{a}_i^\top \boldsymbol{x} \, \boldsymbol{x}'^\top \boldsymbol{a}_i' - y_i|$$

# Our proposed convex relaxation

Given $x_0 \approx x_\star$ and $x_0' \approx x_\star'$ such that for a sufficiently small $\epsilon \geq 0$,

$$\left\| \begin{pmatrix} x_0 - x_\star \\ x_0' - x_\star' \end{pmatrix} \right\|_2 \leq \epsilon \left\| \begin{pmatrix} x_\star \\ x_\star' \end{pmatrix} \right\|_2,$$

we estimate $x_\star$ and $x_\star'$ through the convex program:

$$(\hat{x}, \hat{x}') = \operatorname*{argmax}_{x, x'} x_0^\top x + x_0'^\top x' - \frac{2}{n} \sum_{i=1}^n \max \left\{ \frac{1}{4}(a_i^\top x + a_i'^\top x')^2 - y_i, \frac{1}{4}(a_i^\top x - a_i'^\top x')^2 \right\}$$

**Equivalent folmulation**: Quadratically Constrained Linear Minimization

$$\operatorname*{argmax}_{x, x'} \max_w x_0^\top x + x_0'^\top x' - \frac{2}{n} \mathbf{1}^\top w$$

$$\text{subject to } \frac{1}{4}(a_i^\top x + a_i'^\top x')^2 - y_i \leq w_i, \qquad i = 1, \ldots, n$$

$$\frac{1}{4}(a_i^\top x - a_i'^\top x')^2 \leq w_i, \qquad i = 1, \ldots, n.$$

# Our proposed convex relaxation

Given $\boldsymbol{x}_0 \approx \boldsymbol{x}_\star$ and $\boldsymbol{x}_0' \approx \boldsymbol{x}_\star'$ such that for a sufficiently small $\epsilon \geq 0$,

$$\left\| \begin{pmatrix} \boldsymbol{x}_0 - \boldsymbol{x}_\star \\ \boldsymbol{x}_0' - \boldsymbol{x}_\star' \end{pmatrix} \right\|_2 \leq \epsilon \left\| \begin{pmatrix} \boldsymbol{x}_\star \\ \boldsymbol{x}_\star' \end{pmatrix} \right\|_2,$$

we estimate $\boldsymbol{x}_\star$ and $\boldsymbol{x}_\star'$ through the convex program:

$$(\widehat{\boldsymbol{x}}, \widehat{\boldsymbol{x}}') = \operatorname*{argmax}_{\boldsymbol{x}, \boldsymbol{x}'} \boldsymbol{x}_0^\top \boldsymbol{x} + \boldsymbol{x}_0'^\top \boldsymbol{x}' - \frac{2}{n} \sum_{i=1}^{n} \max\left\{ \frac{1}{4}(\boldsymbol{a}_i^\top \boldsymbol{x} + \boldsymbol{a}_i'^\top \boldsymbol{x}')^2 - y_i, \frac{1}{4}(\boldsymbol{a}_i^\top \boldsymbol{x} - \boldsymbol{a}_i'^\top \boldsymbol{x}')^2 \right\}$$

**Accelerated first-order methods**: We can use a *smoothed* formulation,

$$\operatorname*{argmax}_{\boldsymbol{x}, \boldsymbol{x}'} \boldsymbol{x}_0^\top \boldsymbol{x} + \boldsymbol{x}_0'^\top \boldsymbol{x}' - \frac{2}{\mu n} \sum_{i=1}^{n} \log\left( e^{\frac{\mu}{4}\left(\boldsymbol{a}_i^\top \boldsymbol{x} + \boldsymbol{a}_i'^\top \boldsymbol{x}'\right)^2 - \mu y_i} + e^{\frac{\mu}{4}\left(\boldsymbol{a}_i^\top \boldsymbol{x} - \boldsymbol{a}_i'^\top \boldsymbol{x}'\right)^2} \right)$$

# Theoretical guarantee[*]

We consider $a_i \overset{\text{i.i.d.}}{\sim} \text{Normal}(\mathbf{0}, \mathbf{I}_{d_1 \times d_1})$ and $a_i' \overset{\text{i.i.d.}}{\sim} \text{Normal}(\mathbf{0}, \mathbf{I}_{d_2 \times d_2})$, and $\boldsymbol{\xi} = \mathbf{0}$.

[*]Bahmani, "Estimation from non-linear observations via convex programming with application to bilinear regression," arXiv:1806.07307, 2018.

# Theoretical guarantee[*]

We consider $a_i \overset{\text{i.i.d.}}{\sim} \text{Normal}(\mathbf{0}, I_{d_1 \times d_1})$ and $a_i' \overset{\text{i.i.d.}}{\sim} \text{Normal}(\mathbf{0}, I_{d_2 \times d_2})$, and $\boldsymbol{\xi} = \mathbf{0}$.

> **"Spectral" method of choosing $x_0$ and $x_0'$:**
> Let $S_n = n^{-1} \sum_{i=1}^{n} y_i a_i a_i'^\top$ which satisfies $\mathbb{E} S_n = x_\star x_\star'^\top$. Then we choose,
> $$x_0 = \|S_n\|^{1/2} u_{\max}(S_n) , \qquad\qquad x_0' = \|S_n\|^{1/2} v_{\max}(S_n) .$$
> that meet the required $\epsilon$ relative error if $n \overset{\epsilon}{\gtrsim} (d_1 + d_2) \log(d_1 + d_2)$.

[*]Bahmani, "Estimation from non-linear observations via convex programming with application to bilinear regression," arXiv:1806.07307, 2018.

# Theoretical guarantee[*]

We consider $a_i \overset{\text{i.i.d.}}{\sim} \text{Normal}(\mathbf{0}, \mathbf{I}_{d_1 \times d_1})$ and $a_i' \overset{\text{i.i.d.}}{\sim} \text{Normal}(\mathbf{0}, \mathbf{I}_{d_2 \times d_2})$, and $\boldsymbol{\xi} = \mathbf{0}$.

> **Theorem**
>
> *For a sufficiently large absolute constant $C > 0$, if*
> $$n \geq C \max \left\{ d_1 + d_2, \log\left(8/\delta\right) \right\},$$
> *then with probability $\geq 1 - \delta$ the estimates are exact up to the scaling ambiguity, i.e., for some $t \neq 0$, we have $\hat{x} = t x_\star$ and $\hat{x}' = t^{-1} x_\star'$.*
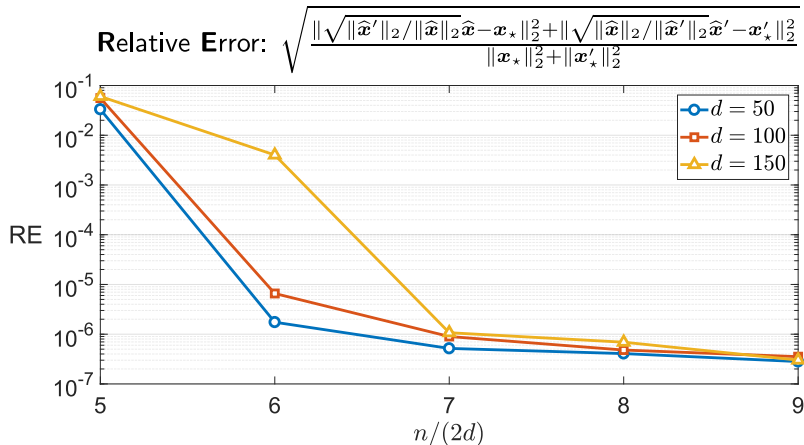
[*]Bahmani, "Estimation from non-linear observations via convex programming with application to bilinear regression," arXiv:1806.07307, 2018.

# Simulation

**Setup**: $d_1 = d_2 = d$. Measurement vectors $\boldsymbol{a}_i, \boldsymbol{a}_i' \overset{\text{i.i.d.}}{\sim} \text{Normal}(\boldsymbol{0}, \boldsymbol{I}_{d \times d})$.

**Solver**: Gurobi for QCLM formulation

Plots show the **median** of the relative error over **100 trials**.

$$\text{R}\text{elative } \text{E}\text{rror}: \sqrt{\frac{\|\sqrt{\|\widehat{\boldsymbol{x}}'\|_2/\|\widehat{\boldsymbol{x}}\|_2}\widehat{\boldsymbol{x}} - \boldsymbol{x}_\star\|_2^2 + \|\sqrt{\|\widehat{\boldsymbol{x}}\|_2/\|\widehat{\boldsymbol{x}}'\|_2}\widehat{\boldsymbol{x}}' - \boldsymbol{x}_\star'\|_2^2}{\|\boldsymbol{x}_\star\|_2^2 + \|\boldsymbol{x}_\star'\|_2^2}}$$

# Nonlinear parametric regression[*]      <span>■</span>

**Difference of convex regression**:

For a random data point **a**, the observation function is **given** in the DC form as

$$f_a(\cdot) = f_a^+(\cdot) - f_a^-(\cdot),$$

where the functions $f_a^+$ and $f_a^-$ are both convex.

[*]Bahmani, "Estimation from non-linear observations via convex programming with application to bilinear regression," arXiv:1806.07307, 2018.

# Nonlinear parametric regression[*]

**Difference of convex regression**:

For a random data point $a$, the observation function is **given** in the DC form as

$$f_a(\cdot) = f_a^+(\cdot) - f_a^-(\cdot),$$

where the functions $f_a^+$ and $f_a^-$ are both convex.

Estimate the parameter $x_\star$, from observations at i.i.d. data points $a_1, \ldots, a_n$, i.e.,

$$y_i = f_{a_i}^+(x_\star) - f_{a_i}^-(x_\star) + \xi_i, \qquad\qquad i = 1, \ldots, n.$$

[*]Bahmani, "Estimation from non-linear observations via convex programming with application to bilinear regression," arXiv:1806.07307, 2018.

# Nonlinear parametric regression[*]

**Difference of convex regression**:

For a random data point $\boldsymbol{a}$, the observation function is **given** in the DC form as

$$f_{\boldsymbol{a}}(\cdot) = f_{\boldsymbol{a}}^+(\cdot) - f_{\boldsymbol{a}}^-(\cdot),$$

where the functions $f_{\boldsymbol{a}}^+$ and $f_{\boldsymbol{a}}^-$ are both convex.

Estimate the parameter $\boldsymbol{x}_\star$, from observations at i.i.d. data points $\boldsymbol{a}_1, \ldots, \boldsymbol{a}_n$, i.e.,

$$y_i = f_{\boldsymbol{a}_i}^+(\boldsymbol{x}_\star) - f_{\boldsymbol{a}_i}^-(\boldsymbol{x}_\star) + \xi_i, \qquad\qquad i = 1, \ldots, n.$$

**Estimator**

Given $\boldsymbol{x}_0 \approx \frac{1}{2n} \sum_{i=1}^n \nabla f_{\boldsymbol{a}_i}^+(\boldsymbol{x}_\star) + \nabla f_{\boldsymbol{a}_i}^-(\boldsymbol{x}_\star)$, we formulate the estimator as

$$\hat{\boldsymbol{x}} = \operatorname*{argmax}_{\boldsymbol{x}} \ \boldsymbol{x}_0^\top \boldsymbol{x} - \frac{1}{n} \sum_{i=1}^n \max\{f_{\boldsymbol{a}_i}^+(\boldsymbol{x}) - y_i, f_{\boldsymbol{a}_i}^-(\boldsymbol{x})\}$$

[*]Bahmani, "Estimation from non-linear observations via convex programming with application to bilinear regression," `arXiv:1806.07307`, 2018.

# Nonlinear parametric regression[*]

**Difference of convex regression**:

For a random data point $\boldsymbol{a}$, the observation function is **given** in the DC form as

$$f_{\boldsymbol{a}}(\cdot) = f_{\boldsymbol{a}}^+(\cdot) - f_{\boldsymbol{a}}^-(\cdot),$$

where the functions $f_{\boldsymbol{a}}^+$ and $f_{\boldsymbol{a}}^-$ are both convex.

**Theorem (simplified)**

Let $\Lambda = \sup_{\boldsymbol{h} \in \mathbb{S}^{d-1}} \mathbb{E} |\boldsymbol{h}^\mathsf{T} \nabla f_{\boldsymbol{a}}(\boldsymbol{x}_\star)|$ and $\lambda = \inf_{\boldsymbol{h} \in \mathbb{S}^{d-1}} \mathbb{E} |\boldsymbol{h}^\mathsf{T} \nabla f_{\boldsymbol{a}}(\boldsymbol{x}_\star)|$. Then, for a sufficiently accurate $\boldsymbol{x}_0$, with probability $\geq 1 - \delta$, having

$$n \gtrsim \max \left\{ \frac{\Lambda^2}{\lambda^2} \log \left( \frac{2}{\delta} \right), \frac{\Lambda^3}{\lambda^3} d \right\},$$
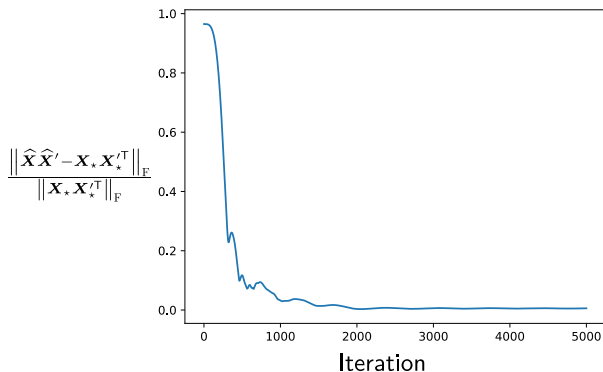
guarantees

$$\|\widehat{\boldsymbol{x}} - \boldsymbol{x}_\star\|_2 \lesssim \frac{\|\boldsymbol{\xi}\|_1}{\lambda n}.$$

[*]Bahmani, "Estimation from non-linear observations via convex programming with application to bilinear regression," arXiv:1806.07307, 2018.

# Simulation: rank $> 1$

**Setup**: For $d = 128, r = 3$ the signal is $\boldsymbol{X}_\star \boldsymbol{X}_\star'^\mathsf{T}$ with $\boldsymbol{X}_\star, \boldsymbol{X}_\star' \in \mathbb{R}^{d \times r}$ and the measurement functions are $f_{\boldsymbol{a}_i, \boldsymbol{a}_i'}^\pm (\boldsymbol{X}, \boldsymbol{X}') = \frac{1}{4} \| \boldsymbol{X}^\mathsf{T} \boldsymbol{a}_i \pm \boldsymbol{X}'^\mathsf{T} \boldsymbol{a}_i' \|_\mathsf{F}^2$ for vectors $\boldsymbol{a}_i$, $\boldsymbol{a}_i' \overset{\text{i.i.d.}}{\sim} \text{Normal}(\boldsymbol{0}, \boldsymbol{I}_{d \times d})$. The observations are then $y_i = \boldsymbol{a}_i^\mathsf{T} \boldsymbol{X}\, \boldsymbol{X}'^\mathsf{T} \boldsymbol{a}_i'$

**Solver**: Nesterov's accelerated gradient method for the smoothed variant

## Proof skecth

By convexity of $f_{a_i}^{\pm}$ we have

$$\max\{f_{a_i}^+ (x_\star + h) - y_i, f_{a_i}^- (x_\star + h)\}$$
$$\geq \max\{h^\top \nabla f_{a_i}^+ (x_\star), h^\top \nabla f_{a_i}^- (x_\star)\} + f_{a_i}^- (x_\star) - (\xi_i)_+ .$$

It suffices to show that

$$\frac{1}{2n} \sum_{i=1}^n \left| \left( \nabla f_{a_i}^+ (x_\star) - \nabla f_{a_i}^- (x_\star) \right)^\top h \right| \geq \left\| x_0 - \frac{1}{2n} \sum_{i=1}^n \nabla f_{a_i}^+ (x_\star) + \nabla f_{a_i}^- (x_\star) \right\|_2 \|h\|_2 .$$

Using a *PAC-Bayesian* argument (à la Catoni*), we show

$$\frac{1}{2n} \sum_{i=1}^n \left| \left( \nabla f_{a_i}^+ (x_\star) - \nabla f_{a_i}^- (x_\star) \right)^\top h \right| \gtrsim \lambda \|h\|_2 ,$$

with high probability.

*Catoni and Giulini, "Dimension-free PAC-Bayesian bounds for matrices, vectors, and linear least squares regression," arXiv:1712.02747.

## Proof skecth

By convexity of $f_{a_i}^{\pm}$ we have

$$\max\{f_{a_i}^+\left(\boldsymbol{x}_\star + \boldsymbol{h}\right) - y_i, f_{a_i}^-\left(\boldsymbol{x}_\star + \boldsymbol{h}\right)\}$$
$$\geq \max\{\boldsymbol{h}^\mathsf{T}\nabla f_{a_i}^+\left(\boldsymbol{x}_\star\right), \boldsymbol{h}^\mathsf{T}\nabla f_{a_i}^-\left(\boldsymbol{x}_\star\right)\} + f_{a_i}^-\left(\boldsymbol{x}_\star\right) - \left(\xi_i\right)_+ \,.$$

It suffices to show that

$$\frac{1}{2n}\sum_{i=1}^{n}\left|\left(\nabla f_{a_i}^+\left(\boldsymbol{x}_\star\right) - \nabla f_{a_i}^-\left(\boldsymbol{x}_\star\right)\right)^\mathsf{T}\boldsymbol{h}\right| \geq \left\|\boldsymbol{x}_0 - \frac{1}{2n}\sum_{i=1}^{n}\nabla f_{a_i}^+\left(\boldsymbol{x}_\star\right) + \nabla f_{a_i}^-\left(\boldsymbol{x}_\star\right)\right\|_2 \|\boldsymbol{h}\|_2 \,.$$

Using a *PAC-Bayesian* argument (à la Catoni[*]), we show

$$\frac{1}{2n}\sum_{i=1}^{n}\left|\left(\nabla f_{a_i}^+\left(\boldsymbol{x}_\star\right) - \nabla f_{a_i}^-\left(\boldsymbol{x}_\star\right)\right)^\mathsf{T}\boldsymbol{h}\right| \gtrsim \lambda\|\boldsymbol{h}\|_2 \,,$$

with high probability.

## Proof skecth

By convexity of $f_{a_i}^{\pm}$ we have

$$\max\{f_{a_i}^+ (x_\star + h) - y_i, f_{a_i}^- (x_\star + h)\}$$
$$\geq \max\{h^\top \nabla f_{a_i}^+ (x_\star), h^\top \nabla f_{a_i}^- (x_\star)\} + f_{a_i}^- (x_\star) - (\xi_i)_+ .$$

It suffices to show that

$$\frac{1}{2n} \sum_{i=1}^n \left| \left( \nabla f_{a_i}^+ (x_\star) - \nabla f_{a_i}^- (x_\star) \right)^\top h \right| \geq \left\| x_0 - \frac{1}{2n} \sum_{i=1}^n \nabla f_{a_i}^+ (x_\star) + \nabla f_{a_i}^- (x_\star) \right\|_2 \|h\|_2 .$$

Using a *PAC-Bayesian* argument (à la Catoni[*]), we show

$$\frac{1}{2n} \sum_{i=1}^n \left| \left( \nabla f_{a_i}^+ (x_\star) - \nabla f_{a_i}^- (x_\star) \right)^\top h \right| \gtrsim \lambda \|h\|_2 ,$$

with high probability.

[*]Catoni and Giulini, "Dimension-free PAC-Bayesian bounds for matrices, vectors, and linear least squares regression," `arXiv:1712.02747`.

# Final remarks

Apply the general result in special cases (e.g., matrix completion, blind deconvolution, …). Adding regularization could be necessary.

Choice of the *anchor* vector $x_0$ in special problems

Other applications where the approach applies (e.g, machine learning)

Analyzing the *iterated* method: $\cdots \to x_0^{(t)} \to \hat{x}^{(t)} \to x_0^{(t+1)} \to \hat{x}^{(t+1)} \to \cdots$

# Final remarks

Apply the general result in special cases (e.g., matrix completion, blind deconvolution, …). Adding regularization could be necessary.

Choice of the *anchor* vector $x_0$ in special problems

Other applications where the approach applies (e.g, machine learning)

Analyzing the *iterated* method: $\cdots \to x_0^{(t)} \to \hat{x}^{(t)} \to x_0^{(t+1)} \to \hat{x}^{(t+1)} \to \cdots$

# Final remarks

Apply the general result in special cases (e.g., matrix completion, blind deconvolution, ...). Adding regularization could be necessary.

Choice of the *anchor* vector $x_0$ in special problems

Other applications where the approach applies (e.g, machine learning)

Analyzing the *iterated* method: $\cdots \to x_0^{(t)} \to \widehat{x}^{(t)} \to x_0^{(t+1)} \to \widehat{x}^{(t+1)} \to \cdots$

# Final remarks

Apply the general result in special cases (e.g., matrix completion, blind deconvolution, …). Adding regularization could be necessary.

Choice of the *anchor* vector $x_0$ in special problems

Other applications where the approach applies (e.g, machine learning)

Analyzing the *iterated* method: $\cdots \rightarrow x_0^{(t)} \rightarrow \widehat{x}^{(t)} \rightarrow x_0^{(t+1)} \rightarrow \widehat{x}^{(t+1)} \rightarrow \cdots$

# Final remarks

Apply the general result in special cases (e.g., matrix completion, blind deconvolution, …). Adding regularization could be necessary.

Choice of the *anchor* vector $x_0$ in special problems

Other applications where the approach applies (e.g, machine learning)

Analyzing the *iterated* method: $\cdots \to x_0^{(t)} \to \widehat{x}^{(t)} \to x_0^{(t+1)} \to \widehat{x}^{(t+1)} \to \cdots$

## Thank you.

preprint @ `arXiv:1806.07307`