

# Collective Matrix Completion

Mokhtar Z. Alaya

Joint work with:



Olga Klopp

ESSEC & CREST

MLMDA Seminar, Borelli Center  
June, 2022

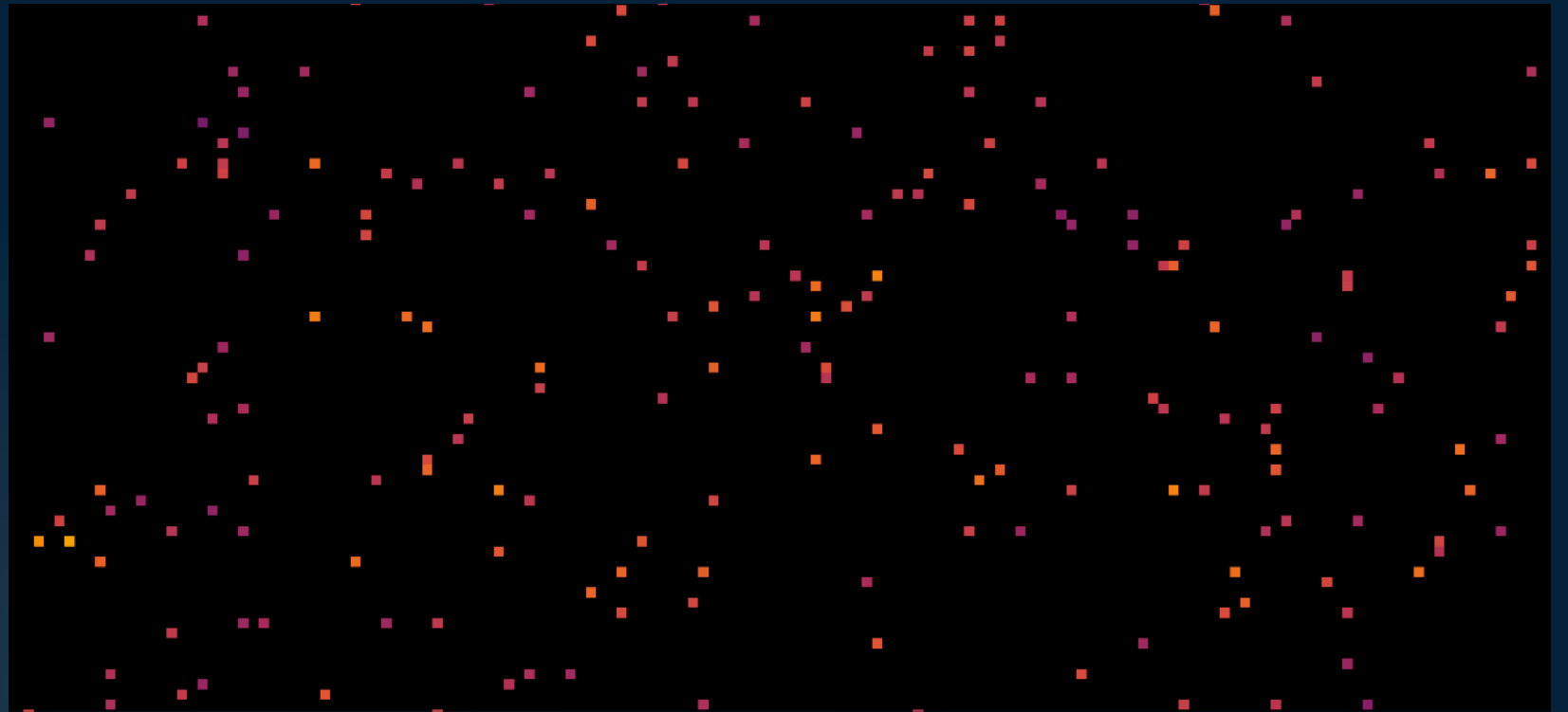


# Outline

1. Overview of matrix completion
2. Collective matrix completion
3. Numerical experiments

# Matrix completion is ...

$$\mathbf{X} =$$



- **Task:** given a partially observed data matrix  $\mathbf{X}$ , predict the unobserved entries.
- Application to recommender systems, system identification, image processing, microarray data, etc.

# Recommender systems, Netflix prize

- A popular example is the Netflix challenge (2006-2009).

$X =$

							...
	★★★★★	?	★★★★☆	?	?	?	...
	?	★★★★☆	?	?	★★★★☆	?	...
	?	?	?	★★★★☆	★★★★☆	?	...
	?	★★★★☆	★★★★☆	?	?	★★★★☆	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

- Dataset: 480K users, 18K movies, 100M ratings.
- Only 1,1% of the matrix is filled!



# Some issues ...

- In general, we cannot infer missing ratings without any other information.
- This problem is under-determined, more unknown than observations (100M  $\ll$  8.64M for Netflix).
- Low-rank assumption: fill matrix such that its rank is minimum.  
➔ A few factors explain most of the data.



# Low rank minimization

- Denote by  $\Omega$  the set of entries of the matrix  $\mathbf{X}$  that have been observed: we know the values  $\mathbf{X}_{ij}$  for all  $(i, j) \in \Omega$ .

$$\underset{\mathbf{W}}{\text{minimize rank}}(\mathbf{W}) \quad \text{s. t.} \quad \mathbf{W}_{ij} = \underbrace{\mathbf{X}_{ij}}_{\text{observed entries}}, \quad \forall (i, j) \in \underbrace{\Omega}_{\text{sampling set}}.$$

- Or a slightly weaker version

$$\underset{\mathbf{W}}{\text{minimize rank}}(\mathbf{W}) \quad \text{s. t.} \quad \sum_{(i,j) \in \Omega} (\mathbf{X}_{ij} - \mathbf{W}_{ij})^2 \leq \delta.$$

- Or the regularization version

$$\underset{\mathbf{W}}{\text{minimize}} \sum_{(i,j) \in \Omega} (\mathbf{X}_{ij} - \mathbf{W}_{ij})^2 + \lambda \text{rank}(\mathbf{W}).$$

# Low rank minimization

- Non-convex problem and combinatorially NP-hard!

$$\text{rank}(\mathbf{X}) = \|\sigma(\mathbf{X})\|_0 = \sum_{i=1}^{\min \dim(\mathbf{X})} \underbrace{\mathbb{1}_{(\sigma_i(\mathbf{X}) > 0)}}_{i^{\text{th}} \text{ largest singular value}}.$$

- Replace the  $\ell_0$  pseudo-norm by the  $\ell_1$ -norm [Fazel ('02), Srebro et al ('05), Candès and Tao ('10), Negahban and Wainwright ('11), Davenport et al. ('14), Klopp ('14 and '15), ....].

**Nuclear / trace / 1-Schatten norm:**

$$\|\mathbf{X}\|_* = \|\sigma(\mathbf{X})\|_1 = \sum_{i=1}^{\min \dim(\mathbf{X})} \sigma_i(\mathbf{X}).$$

# Nuclear norm minimization

- Hence tempting to consider the nuclear norm minimization problem:

$$\underset{\mathbf{W}}{\text{minimize}} \|\mathbf{W}\|_* \quad \text{s. t.} \quad \sum_{(i,j) \in \Omega} (\mathbf{X}_{ij} - \mathbf{W}_{ij})^2 \leq \delta.$$

- Or equivalently the regularization / Lagrangian formulation:

$$\underset{\mathbf{W}}{\text{minimize}} \frac{1}{2} \sum_{(i,j) \in \Omega} (\mathbf{X}_{ij} - \mathbf{W}_{ij})^2 + \lambda \|\mathbf{W}\|_*.$$

- This is convex problem.

# Motivations of collective MC

- Data is often obtained from a collection of source matrices:

$$\mathcal{X} = (X^1, \dots, X^V)$$

$$\mathcal{X} = \left( \begin{array}{c|c|c} \begin{array}{cccccc} & & \text{dark green} & & & \\ & & \text{red} & & & \\ \text{orange} & & & & & \\ & & \text{yellow} & & \text{orange} & \\ & \text{dark green} & & & & \\ & & \text{dark green} & & & \\ & & & & & \text{red} \\ & & & & & \\ & & \text{red} & & & \end{array} & \begin{array}{cccccc} & & & \text{yellow} & & \\ & \text{red} & & & & \\ \text{dark green} & & & & & \\ & & \text{red} & & & \\ & & & & & \text{red} \\ & \text{orange} & & \text{red} & & \\ & & \text{yellow} & & & \\ & & & & & \end{array} & \dots & \begin{array}{cccccc} \text{red} & & & & \text{orange} & \\ & \text{yellow} & & & & \\ & & \text{dark green} & & & \\ \text{dark green} & & & & & \\ & & & & \text{orange} & \\ \text{red} & & & & & \text{dark green} \\ & \text{orange} & & & & \end{array} \end{array} \right)$$

$X^1 \qquad X^2 \qquad \qquad X^V$

- Cold-Start problem:** in recommender systems, when a new user has no rating it is impossible to predict his ratings.
- Shared structure among the sources can be useful to get better predictions.

# Collective MC: setup

- Each source view  $\mathbf{X}^v \in \mathbb{R}^{d_u \times d_v}$  and  $D = \sum_{v=1}^V d_v$ .
- **Model:** let  $B_{ij}^v$  be independent Bernoulli random variables and independent from  $\mathbf{X}_{ij}^v$  with parameter  $\pi_{ij}^v$ . Setting:

$$\mathbf{Y}^v = \mathbf{B}^v \odot \mathbf{X}^v \quad \text{that is} \quad Y_{ij}^v = B_{ij}^v X_{ij}^v.$$

# Collective MC: sampling scheme

- We consider **general sampling model** where we only assume that each entry is observed with a positive probability.

## Assumption 1

There exists a positive constant  $0 < p < 1$  such that

$$\min_{v \in [V]} \min_{(i,j) \in [d_u] \times [d_v]} \pi_{ij}^v \geq p.$$

[Klopp('15), Klopp et al. ('15), Lafond ('15), Cai and Zou ('16)]



# Collective MC: sampling scheme

- Let  $\pi_{i\bullet}^v$  (resp.  $\pi_{\bullet j}^v$ ) the probability of sampling a coefficient from  $i$ -th row (resp.  $j$ -th column) of  $\mathbf{X}^v$ . Namely:

$$\pi_{i\bullet}^v = \sum_{j \in [d_v]} \pi_{ij}^v \quad \text{and} \quad \pi_{\bullet j}^v = \sum_{i \in [d_u]} \pi_{ij}^v. \quad \text{Let } \pi_{i\bullet} = \sum_{v \in [V]} \pi_{i\bullet}^v.$$

## Assumption 2

There exists a positive constant  $\mu$  such that

$$\max_{v \in [V]} \max_{(i,j) \in [d_u] \times [d_v]} (\pi_{i\bullet}^v, \pi_{\bullet j}^v) \leq \mu.$$

[Klopp('15), Klopp et al. ('15), Lafond ('15), Cai and Zou ('16)]

# Case I: Exponential family noise

- We assume that the distribution of for each source  $\mathbf{X}^v$  depends on the matrix of parameters  $\mathbf{M}^v$  and satisfied a natural exponential family [Gunasekar et al. ('14); Cao and Xie ('16); Lafond ('15) ]

$$\mathbf{X}_{ij}^v | \mathbf{M}_{ij}^v \sim f_{h^v, G^v}(\mathbf{X}_{ij}^v | \mathbf{M}_{ij}^v) = h^v(\mathbf{X}_{ij}^v) \exp(\mathbf{X}_{ij}^v \mathbf{M}_{ij}^v - G^v(\mathbf{M}_{ij}^v)).$$

## Assumption 3

Assume that  $G^v(\cdot)$  is twice differentiable and there exists two constants  $L_\gamma^2, U_\gamma^2$

$$\sup_{\eta \in [-\gamma - \frac{1}{K}, \gamma + \frac{1}{K}]} (G^v)''(\eta) \leq U_\gamma^2 \quad \text{and} \quad \inf_{\eta \in [-\gamma - \frac{1}{K}, \gamma + \frac{1}{K}]} (G^v)''(\eta) \geq L_\gamma^2$$

for some  $K > 0$ .

# Exponential family noise: estimation procedure of $\mathcal{M} = (\mathbf{M}^1, \dots, \mathbf{M}^V)$

- Given the observations  $\mathcal{Y} = (\mathbf{Y}^1, \dots, \mathbf{Y}^V)$ , the normalized negative log-likelihood write as, for any  $\mathbf{W} = (\mathbf{W}^1, \dots, \mathbf{W}^V) \in \mathbb{R}^{d_u \times D}$ ,

$$\mathcal{L}_{\mathcal{Y}}(\mathbf{W}) = -\frac{1}{d_u D} \sum_{v \in [V]} \sum_{(i,j) \in [d_u] \times [d_v]} B_{ij}^v \left( \mathbf{Y}_{ij}^v \mathbf{W}_{ij}^v - G^v(\mathbf{W}_{ij}^v) \right)$$

- The nuclear norm penalized estimator  $\hat{\mathcal{M}}$  of  $\mathcal{M}$  is defined as:

$$\hat{\mathcal{M}} = (\hat{\mathbf{M}}^1, \dots, \hat{\mathbf{M}}^V) = \underset{\mathbf{W} \in \mathcal{C}_{\infty}(\gamma)}{\operatorname{argmin}} \mathcal{L}_{\mathcal{Y}}(\mathbf{W}) + \lambda \|\mathbf{W}\|_*$$

where  $\mathcal{C}_{\infty}(\gamma) = \{ \mathbf{W} \in \mathbb{R}^{d_u \times D} : \|\mathbf{W}\|_{\infty} \leq \gamma \}$ .

[Foygel et al. ('10), Salakhutdinov and Srebro ('10)]

# Exponential family noise: theoretical guarantee

- Upper bound on the rescaled Frobenius estimation risk:

## Theorem [A., Klopp 2019]

Assume that Assumptions 1, 2 and 3 hold and

$$\lambda = \mathcal{O}\left(\frac{(U_\gamma \vee K)(\sqrt{\mu} + (\log(d_u \vee D))^{3/2})}{d_u D}\right).$$

Then, with probability exceeding  $1 - 4/(d_u + D)$  one has,

$$\begin{aligned} \frac{1}{d_u D} \|\widehat{\mathcal{M}} - \mathcal{M}\|_F^2 &\lesssim \frac{\text{rank}(\mathcal{M})}{p^2 d_u D} \left( \gamma^2 + \frac{(U_\gamma \vee K)^2}{L_\gamma^4} \right) (\mu + \log^3(d_u \vee D)) \\ &\lesssim \frac{\text{rank}(\mathcal{M}) \mu}{p^2 d_u D}. \end{aligned}$$

# Exponential family noise: remarks

- For a close uniform sampling distribution, that is  $c_1 p \leq \pi_{ij}^v \leq c_2 p$

$$\frac{1}{d_u D} \|\widehat{\mathcal{M}} - \mathcal{M}\|_F^2 \lesssim \frac{\text{rank}(\mathcal{M})}{p(d_u \wedge D)}.$$

- Rate of convergence achieved by our estimator is faster compared to the penalization by the sum-nuclear-norm since

$$\text{rank}(\mathcal{M}) \leq \sum_{v=1}^V \text{rank}(\mathcal{M}^v).$$

- For small estimation error, one can choose  $p \geq \text{rank}(\mathcal{M}) / (d_u \wedge D)$ .

This implies

$$n \gtrsim \text{rank}(\mathcal{M})(d_u \vee D).$$

where  $n = \sum_{v \in [V]} \sum_{(i,j) \in [d_u] \times [d_v]} \pi_{ij}^v$  the expected number of observations

# Case 2: Distribution-free-setting

- We do not assume any specific model for the observations.
- We consider the risk of estimating  $\mathbf{X}^v$  with a loss function  $\ell^v(\cdot, \cdot)$ .

## Assumption 4

For every  $v$  the loss function  $\ell^v(y, \cdot)$  is  $\rho_v$ -Lipschitz in its second argument:

$$|\ell^v(y, x) - \ell^v(y, x')| \leq \rho_v |x - x'|.$$

# Distribution-free-setting: estimation procedure

- For any matrix  $\mathcal{Q} = (\mathcal{Q}^1, \dots, \mathcal{Q}^V)$ , we define the empirical risk as

$$R_{\mathcal{Y}}(\mathcal{Q}) = \frac{1}{d_u D} \sum_{v \in [V]} \sum_{(i,j) \in [d_u] \times [d_v]} B_{ij}^v \ell^v(\mathcal{Y}_{ij}^v, \mathcal{Q}_{ij}^v)$$

- We define the oracle as:

$$\hat{\mathcal{M}}^* = (\hat{\mathcal{M}}^{*1}, \dots, \hat{\mathcal{M}}^{*V}) = \operatorname{argmin}_{\mathcal{Q} \in \mathcal{C}_{\infty}(\gamma)} R(\mathcal{Q})$$

where  $R(\mathcal{Q}) = \mathbb{E}[R_{\mathcal{Y}}(\mathcal{Q})]$ .

- We consider excess risk  $R(\hat{\mathcal{M}}) - R(\hat{\mathcal{M}}^*)$ .



# Distribution-free-setting: estimation procedure

- For a tuning parameter  $\Lambda > 0$  the nuclear norm penalized estimator reads as

$$\hat{\mathcal{M}} \in \operatorname{argmin}_{\mathcal{Q} \in \mathcal{C}_\infty(\gamma)} \{R(\mathcal{Q}) + \Lambda \|\mathcal{Q}\|_*\}$$

## Assumption 5

There exists a constant  $\varsigma > 0$  such that for every  $\mathcal{Q} \in \mathcal{C}_\infty(\gamma)$ , one has

$$R(\mathcal{Q}) - R(\mathcal{M}^\star) \geq \frac{\varsigma}{pd_u D} \|\mathcal{Q} - \mathcal{M}^\star\|_F^2$$

- Assumption 4 is called “Bernstein” condition [Mendelson (2008); Bartlett et al., (2004); Alquier et al., (2017); Elsen and van de Geer, (2018)].

# Distribution-free-setting: theoretical guarantee

## Theorem [A., Klopp 2019]

Assume that Assumptions 1, 2, 4 and 5 hold and set  $\rho = \max_{v \in [V]} \rho_v$ .  
Let

$$\Lambda = \mathcal{O} \left( \frac{\rho(\sqrt{\mu} + \sqrt{\log(d_u \vee D)})}{d_u D} \right).$$

Then, with probability exceeding  $1 - 4/(d_u + D)$  one has,

$$R(\hat{\mathcal{M}}) - R(\mathcal{M}^\star) \lesssim \frac{\text{rank}(\mathcal{M}^\star)}{p} \frac{(\rho^2 + \rho^{3/2} \sqrt{\gamma/\varsigma})(\mu + \log(d_u \vee D))}{d_u D}$$

# 3. Numerical Experiments

# Optimization of $\hat{\mathcal{M}} = (\hat{\mathcal{M}}^1, \dots, \hat{\mathcal{M}}^V) = \underset{\mathcal{W} \in \mathcal{C}_\infty(\gamma)}{\operatorname{argmin}} \mathcal{L}_{\mathcal{Y}}(\mathcal{W}) + \lambda \|\mathcal{W}\|_*$

- Proximal gradient (PG): [Beck and Teboulle ('09), Cai et al. ('09), Mazumder et al., ('10); Yao and Kwok ('15)]
- The PG generates a sequence of estimates

$$\mathcal{W}_{t+1} = \operatorname{prox}_{\frac{\lambda}{L} \|\cdot\|_*}(\mathcal{Z}_t), \text{ where } \mathcal{Z}_t = \mathcal{W}_t - \frac{1}{L} \nabla \mathcal{L}_{\mathcal{Y}}(\mathcal{W}_t)$$

- Assume a singular value decomposition  $\mathcal{W} = \mathcal{U} \Sigma \mathcal{V}^\top$ , then one has

$$\operatorname{prox}_{\frac{\lambda}{L} \|\cdot\|_*}(\mathcal{W}) = \operatorname{SVT}_{\lambda/L}(\mathcal{W}) = \mathcal{U} \operatorname{diag}((\sigma_1 - \lambda/L)_+, \dots, (\sigma_r - \lambda/L)_+) \mathcal{V}^\top$$

[Cai et al. ('10)]

# Power method to reduce complexity

- To compute  $\mathcal{W}_{t+1}$  we need to perform an SVD of  $\mathcal{Z}_t$   $\mathcal{O}((d_u \wedge D)d_u D)$
- We do not require to do a full SVD only a few  $k_t$  singular values of  $\mathcal{Z}_t$  which are large than  $\lambda/L$ .
- As  $\mathcal{W}_t$  converges to a low rank solution then  $k_t$  will be small during iterations.
- [Yao and Kwok ('15)] showed the following result:

$$\text{SVT}_{\lambda/L}(\mathcal{Z}_t) = \mathcal{Q} \text{SVT}_{\lambda/L}(\mathcal{Q}^\top \mathcal{Z}_t) \quad \mathcal{O}(k_t d_u D)$$

## Algorithm 2: Power Method: PowerMethod( $\mathcal{Z}, \mathcal{R}, \epsilon$ )

1. **input:**  $\mathcal{Z} \in \mathbb{R}^{d_u \times D}$ , initial  $\mathcal{R} \in \mathbb{R}^{D \times k}$  for warm-start, tolerance  $\delta$ ;
2. **initialize**  $\mathcal{W}_1 = \mathcal{Z}\mathcal{R}$ ;
3. **for**  $t = 1, 2, \dots$ , **do**
4.      $\mathcal{Q}_{t+1} = \text{QR}(\mathcal{W}_t)$ ; // QR denotes the QR factorization
5.      $\mathcal{W}_{t+1} = \mathcal{Z}(\mathcal{Z}^\top \mathcal{Q}_{t+1})$ ;
6.     **if**  $\|\mathcal{Q}_{t+1}\mathcal{Q}_{t+1}^\top - \mathcal{Q}_t\mathcal{Q}_t^\top\|_F \leq \delta$  **then**  
        └ break;
7. **return**  $\mathcal{Q}_{t+1}$ .

[Halko et al ('11)]

# Approximate SVT based on power method

---

**Algorithm 3:** Approximate SVT:  $\text{Approx-SVT}(\mathcal{Z}, \mathcal{R}, \lambda, \delta)$

---

1. **input:**  $\mathcal{Z} \in \mathbb{R}^{d_u \times D}$ ,  $\mathcal{R} \in \mathbb{R}^{D \times k}$ , thresholds  $\lambda$  and  $\delta$ ;
  2.  $\mathcal{Q} = \text{PowerMethod}(\mathcal{Z}, \mathcal{R}, \delta)$ ; // Approximate the top  $k_t$  left singular values
  3.  $[\mathcal{U}, \Sigma, \mathcal{V}] = \text{SVD}(\mathcal{Q}^\top \mathcal{Z})$ ;
  4.  $\mathcal{U} = \{u_i | \sigma_i > \lambda\}$ ;
  5.  $\mathcal{V} = \{v_i | \sigma_i > \lambda\}$ ;
  6.  $\Sigma = \max(\Sigma - \lambda \mathcal{I}, \mathbf{0})$ ; // ( $\mathcal{I}$  denotes the identity matrix)
  7. **return**  $\mathcal{Q}\mathcal{U}, \Sigma, \mathcal{V}$ .
- [Yao and Kwok ('15)]

**Algorithm 4: PLAIS-Impute for Collective Matrix Completion** $\mathcal{O}(1/T^2)$ 

1. **input:** observed collective matrix  $\mathbf{Y}$ , parameter  $\lambda$ , decay parameter  $\nu \in (0, 1)$ , tolerance  $\varepsilon$ ;
2.  $[\mathbf{U}_0, \lambda_0, \mathbf{V}_0] = \text{rank-1 SVD}(\mathbf{Y})$ ;
3. **initialize**  $c = 1$ ,  $\delta_0 = \|\mathbf{Y}\|_F$ ,  $\mathbf{W}_0 = \mathbf{W}_1 = \lambda_0 \mathbf{U}_0 \mathbf{V}_0^\top$ ;
4. **for**  $t = 1, \dots, T$  **do**
  5.  $\delta_t = \nu^t \delta_0$ ; // Regularization is dynamically reduced by continuation strategy
  6.  $\lambda_t = \nu^t(\lambda_0 - \lambda) + \lambda$ ;
  7.  $\theta_t = (c - 1)/(c + 2)$ ;
  8.  $\mathbf{Q}_t = (1 + \theta_t)\mathbf{W}_t - \theta_t\mathbf{W}_{t-1}$ ; //Acceleration (FISTA)
  9.  $\mathbf{Z}_t = \nabla \mathcal{L}_{\mathbf{Y}}(\mathbf{Q}_t)$ ;
  10.  $\mathbf{V}_{t-1} = \mathbf{V}_{t-1} - \mathbf{V}_t(\mathbf{V}_t^\top \mathbf{V}_{t-1})$ ; //Warm-start
  11.  $\mathbf{R}_t = \text{QR}([\mathbf{V}_t, \mathbf{V}_{t-1}])$ ;
  12.  $[\mathbf{U}_{t+1}, \mathbf{\Sigma}_{t+1}, \mathbf{V}_{t+1}] = \text{Approx-SVT}(\mathbf{Z}_t, \mathbf{R}_t, \lambda_t, \delta_t)$ ; //Approximate SVT
  13. **if**  $\mathcal{F}_\lambda(\mathbf{U}_{t+1}\mathbf{\Sigma}_{t+1}\mathbf{V}_{t+1}^\top) > \mathcal{F}_\lambda(\mathbf{U}_t\mathbf{\Sigma}_t\mathbf{V}_t^\top)$  **then** //Restart the algorithm if the objective function increases  
     $c = 1$ ;
  14. **else**  
     $c = c + 1$ ;
  15. **if**  $|\mathcal{F}_\lambda(\mathbf{U}_{t+1}\mathbf{\Sigma}_{t+1}\mathbf{V}_{t+1}^\top) - \mathcal{F}_\lambda(\mathbf{U}_t\mathbf{\Sigma}_t\mathbf{V}_t^\top)| \leq \varepsilon$  **then**  
    **break**;
16. **return**  $\mathbf{W}_{T+1}$ .



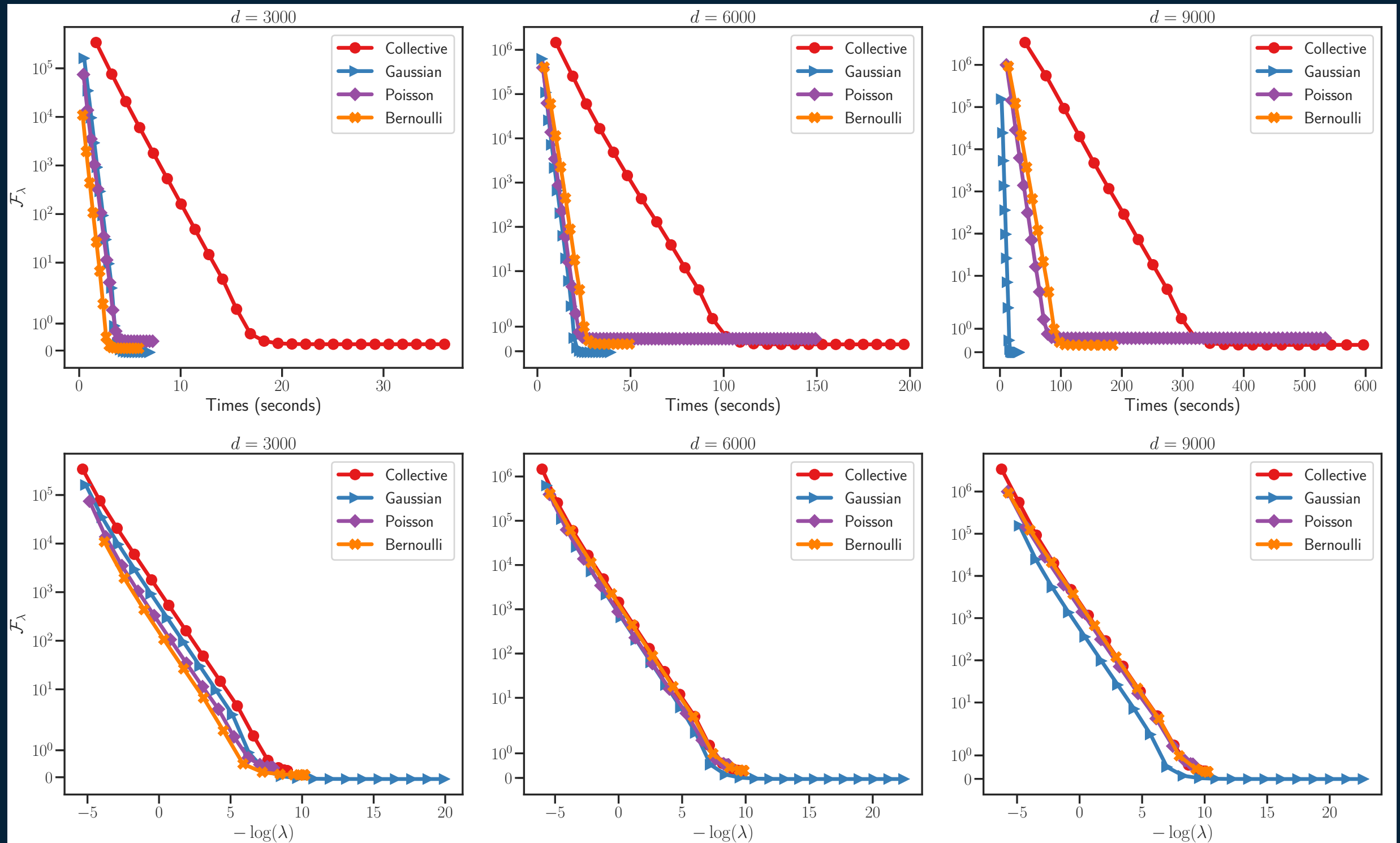
# Experimental results on synthetic data

Each source matrix  $\mathbf{M}^v$  is constructed as  $\mathbf{M}^v = \mathbf{L}^v \mathbf{R}^{v\top}$  where  $\mathbf{L}^v \in \mathbb{R}^{d \times r_v}$  and  $\mathbf{R}^v \in \mathbb{R}^{d_v \times r_v}$

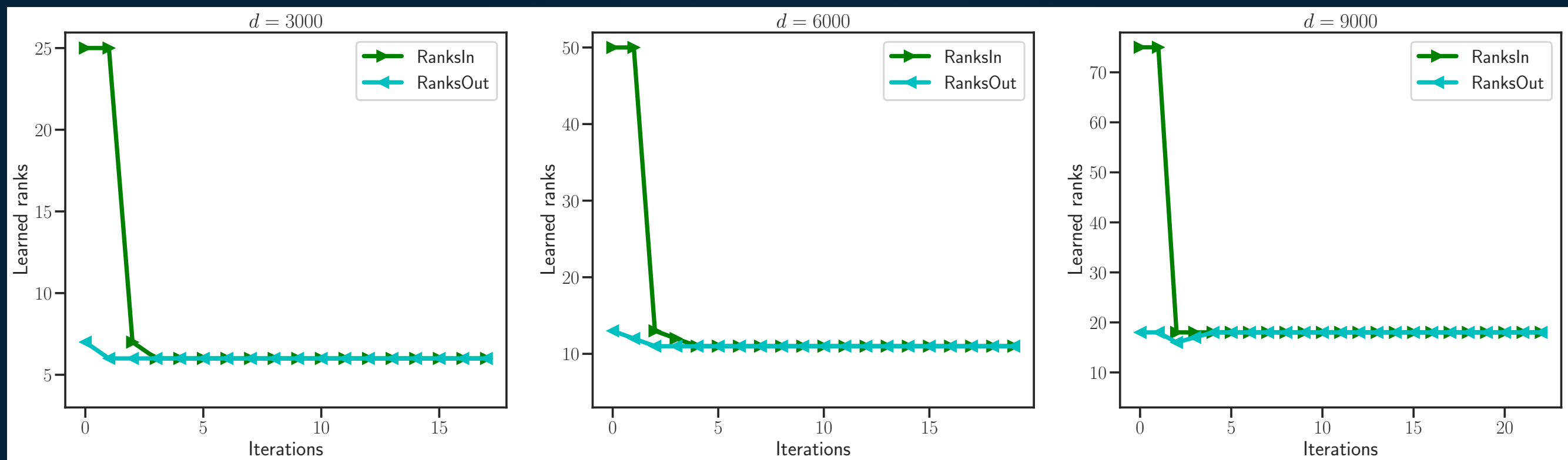
A fraction of the entries of  $\mathbf{M}^v$  is removed uniformly at random with probability  $p \in [0, 1]$ .

		$i.i.d. \mathcal{N}(0.5, 1)$	$i.i.d. \mathcal{P}(0.5)$	$i.i.d. \mathcal{B}(0.5)$	
		$\mathbf{M}^1$	$\mathbf{M}^2$	$\mathbf{M}^3$	$\mathcal{M}$
		(Gaussian)	(Poisson)	(Bernoulli)	(Collective)
exp. 1	dimension	$3000 \times 1000$	$3000 \times 1000$	$3000 \times 1000$	$3000 \times 3000$
	rank	5	5	5	unknown
exp. 2	dimension	$6000 \times 2000$	$6000 \times 2000$	$6000 \times 2000$	$6000 \times 6000$
	rank	10	10	10	unknown
exp. 3	dimension	$9000 \times 3000$	$9000 \times 3000$	$9000 \times 3000$	$9000 \times 9000$
	rank	15	15	15	unknown

# Experimental results on synthetic data: convergence of the objective functions



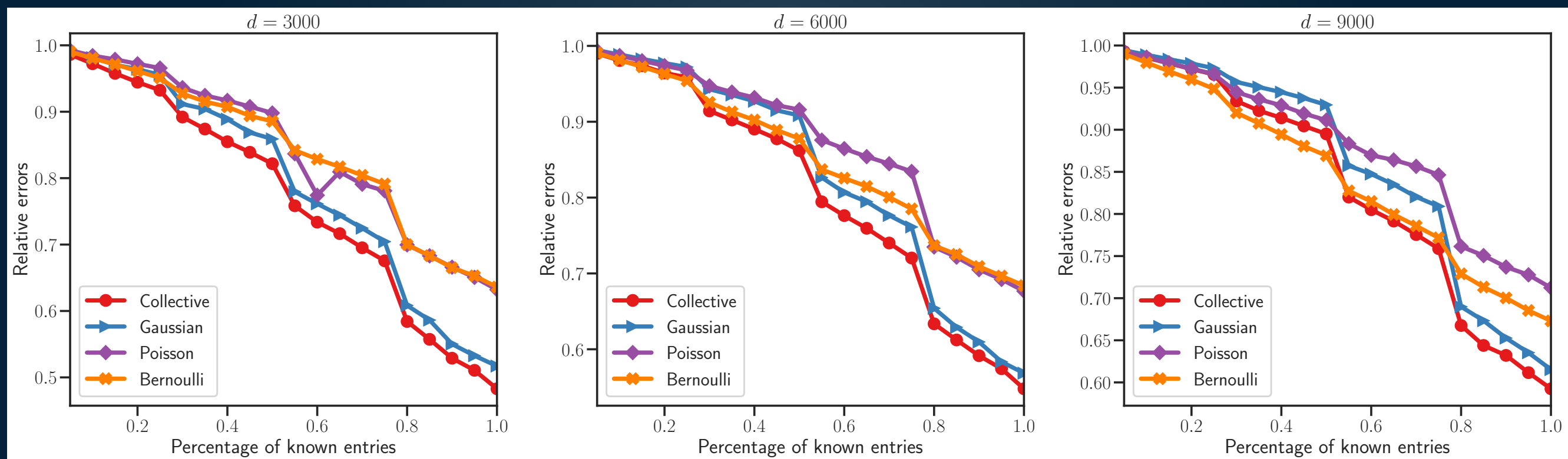
# Experimental results on synthetic data: Learning ranks curve



# Experimental results on synthetic data: evaluation of the estimator

- Our metric matrix completion is defined by the relative error, [Cai. et al. ('10); Davenport et al. ('14); Cai and Zhou ('13)],

$$\text{RE}(\widehat{\mathbf{W}}, \mathbf{W}) = \frac{\|\widehat{\mathbf{W}} - \mathbf{W}^o\|_F}{\|\mathbf{W}^o\|_F}$$



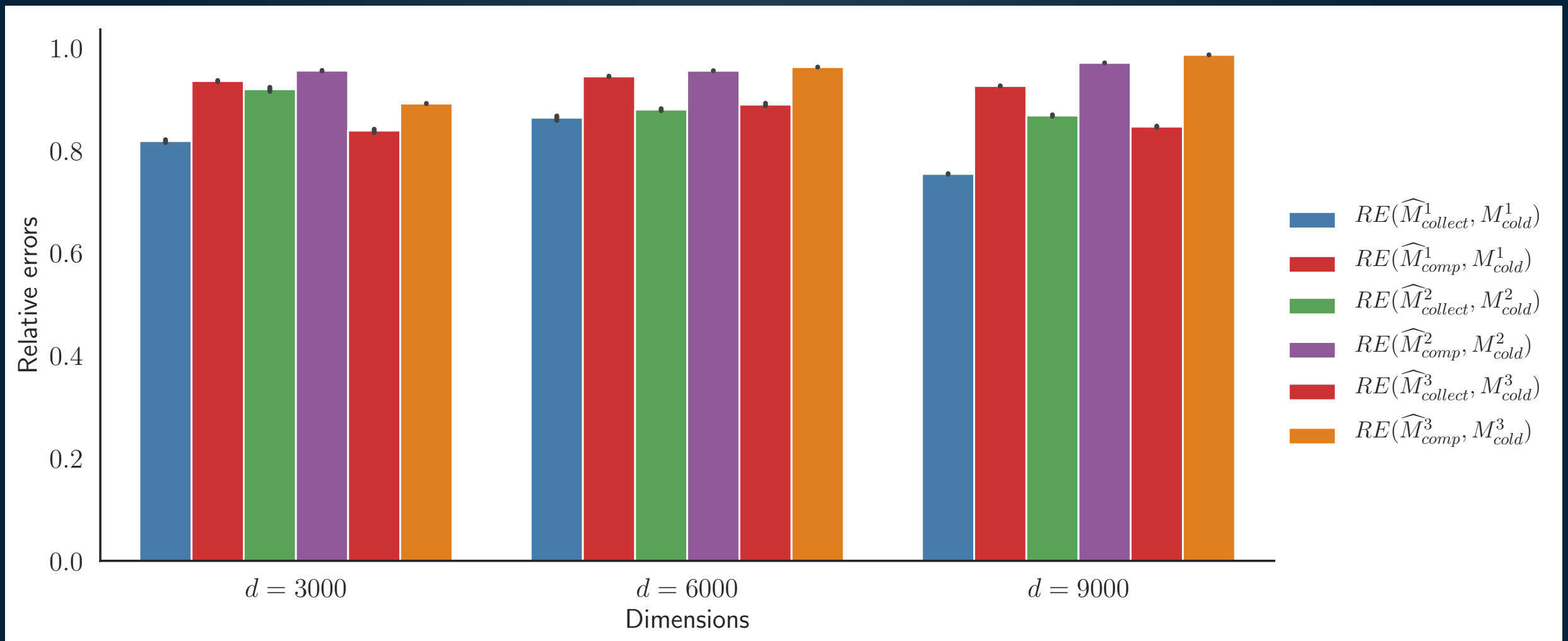
Performance on the synthetic data in terms of relative errors between the target and the estimator matrices

# Experimental results on synthetic data:

## Cold-Start problem

- We construct the "cold" collective matrices: we extract vector of known entries of the chosen matrix and we set the first 1/5 fraction of its entries to be equal to zero.

$$\mathcal{M}_{\text{cold}}^1 = (M_{\text{cold}}^1, M^2, M^3), \mathcal{M}_{\text{cold}}^2 = (M^1, M_{\text{cold}}^2, M^3), \text{ and } \mathcal{M}_{\text{cold}}^3 = (M^1, M^2, M_{\text{cold}}^3).$$



# Take home message

- Recovering a low-rank matrix when the data are collected from multiple and heterogeneous source matrices.
- Estimators are based on minimizing the sum of a goodness-of-fit term and the nuclear norm penalization of the whole collective matrix.
- Upper bounds on the prediction risk of the estimators.
- Empirical evidence of the efficiency of the collective matrix completion approach in the case of joint low-rank structure compared to estimate each source matrices separately.

# References

- Alaya, Mokhtar Z., and Olga Klopp. 2019. “Collective Matrix Completion.” *Journal of Machine Learning Research* 20:1–43.
- T. Cai and W. X. Zhou. A max-norm constrained minimization approach to 1-bit matrix completion. *J. Mach. Learn. Res.*, 14(1):3619–3647, 2013.
- E. J. Candes and T. Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.
- M. Fazel. *Matrix Rank Minimization with Applications*. PhD thesis, Stanford University, 2002.
- O. Klopp. Noisy low-rank matrix completion with general sampling distribution. *Bernoulli*, 20(1):282–303, 2014.
- O. Klopp. Matrix completion by singular value thresholding: Sharp bounds. *Electron. J. Statist.*, 9(2):2348–2369, 2015.
- Q. Yao and J. T. Kwok. Accelerated inexact soft-impute for fast large-scale matrix completion. In *Proceedings of the 24th International Conference on Artificial Intelligence*, IJCAI’15, pages 4002–4008. AAAI Press, 2015.



*Thank you!*

# Distribution-free-setting: remarks

- In 1-bit matrix completion with logistic (resp. hinge) loss, the Bernstein assumption is satisfied with  $\varsigma = 1/(4e^{2\gamma})$  (resp.  $\varsigma = 2\tau$ , such that  $|\dot{M}_{ij}^v - 1/2| \geq \tau, \forall v \in [V], (i, j) \in [d_u] \times [d_v]$ ) [Alquier et al. (2017)].
- The excess risk with respect to these two losses under the uniform sampling is obtained without a logarithmic factor [Alquier et al. (2017)],

$$R(\hat{\mathcal{M}}) - R(\dot{\mathcal{M}}) \lesssim \frac{\text{rank}(\dot{\mathcal{M}})}{p(d_u \wedge D)}.$$