

ICASSP 2016

Low-rank Matrix Recovery via Entropy Function



Dung N. Tran¹, Shuai Huang¹, Sang Chin² and Trac D. Tran¹

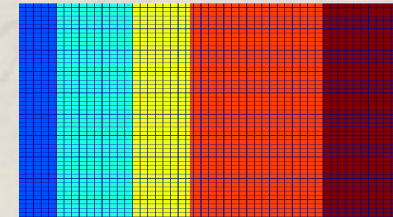
¹The Johns Hopkins University

²Draper Laboratory and Boston University

Low-rank matrix recovery

- ◆ **Main assumption:** input matrix is low-rank

$$\text{rank}(\mathbf{X}) = r \ll \min\{n_1, n_2\}$$



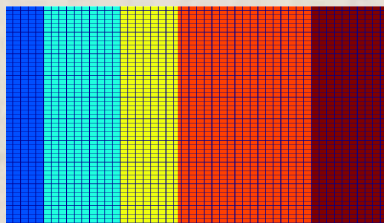
$$\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$$

- ◆ **Measurements:**

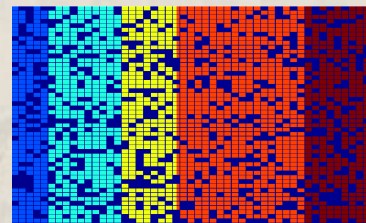
$$\mathbf{y} = \mathcal{A}(\mathbf{X}) + \boldsymbol{\epsilon}$$

- $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$: linear sampling operator
- $\mathbf{y} \in \mathbb{R}^m$: measurement vector
- $\boldsymbol{\epsilon} \in \mathbb{R}^m$: noise vector
- $m < n_1 n_2$

- ◆ **Reconstruction:** given $(\mathbf{y}, \mathcal{A})$, recover \mathbf{X} .



\mathbf{X}

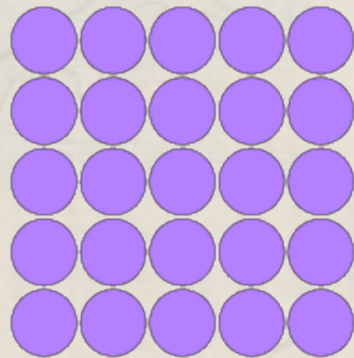


\mathbf{y}



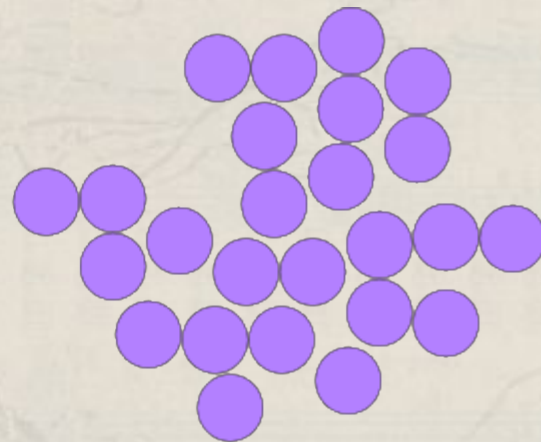
$\mathbf{X}?$

Entropy



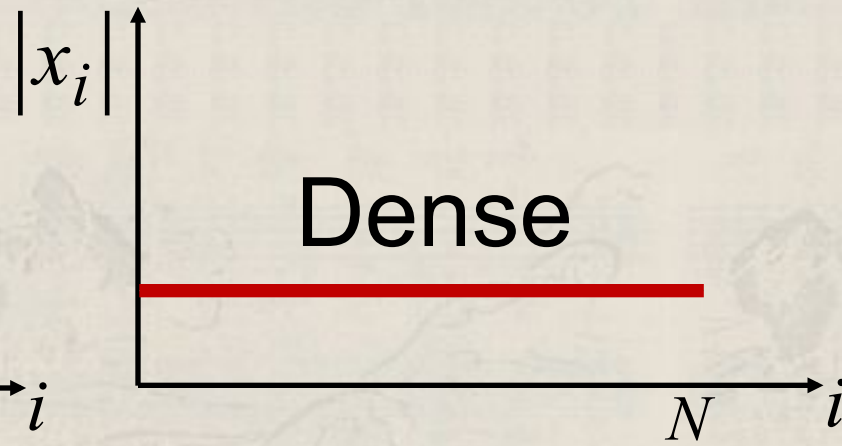
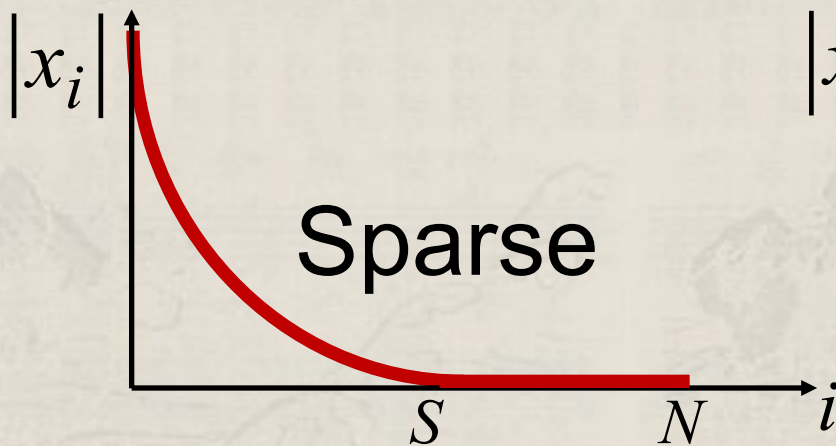
Low Entropy

Low Rank



High Entropy

High Rank



Entropy Minimization

- ◆ *Entropy Minimization for low-rank matrix recovery*

$$\min_{\mathbf{X}} h(\boldsymbol{\sigma}(\mathbf{X})) \quad \text{s.t.} \quad \mathcal{A}(\mathbf{X}) = \mathbf{y},$$

where

$h(\mathbf{x}) = -\sum_i \frac{|\mathbf{x}_i|}{\|\mathbf{x}\|_1} \log \frac{|\mathbf{x}_i|}{\|\mathbf{x}\|_1}$ is the entropy function,

$\boldsymbol{\sigma}(\mathbf{X}) = (\sigma_1(\mathbf{X}), \dots, \sigma_n(\mathbf{X}))$ is the *vector of singular values* of \mathbf{X} .

- ◆ *Questions of interest:*

- Why entropy minimization? → *Sparsity inducing property*
- How to solve it? → *ENM algorithm*
- What do we gain? → *Faster sampling rate (show empirically)*
- Why does it work theoretically? → *(Future work)*

Entropy Function Induces Sparsity

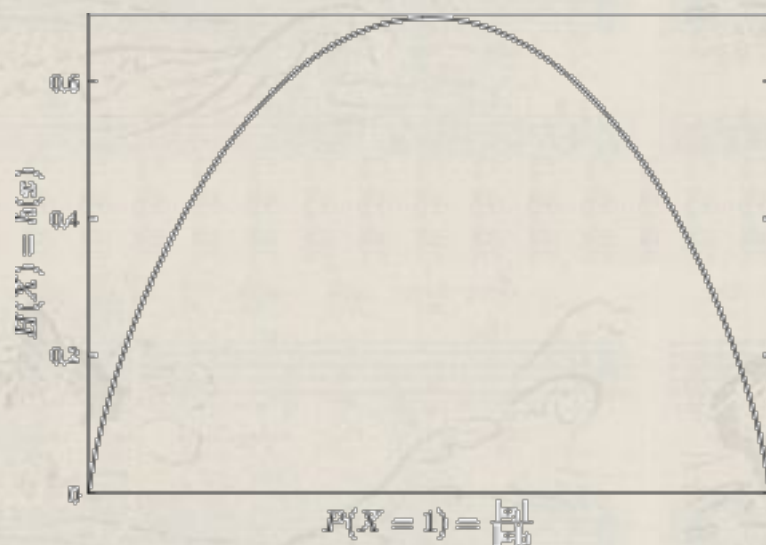
Recall $h(\mathbf{x}) = -\sum_i \frac{|x_i|}{\|\mathbf{x}\|_1} \log \frac{|x_i|}{\|\mathbf{x}\|_1}$ for $\mathbf{x} \in \mathbb{R}^n$

Let X : be a discrete random variable with possible values $\{1, \dots, n\}$:

$$P(X = i) = \frac{|x_i|}{\|\mathbf{x}\|_1}$$

$\Rightarrow \left\{ \frac{|x_1|}{\|\mathbf{x}\|_1}, \dots, \frac{|x_n|}{\|\mathbf{x}\|_1} \right\}$ is the distribution of X and $H(X) = h(\mathbf{x})$

Here, $H(\mathbf{X})$ is the Shannon entropy of \mathbf{X} .



Example: $\mathbf{x} \in \mathbb{R}^2$

$h(\mathbf{x}) = H(X)$ attains its maximum when $x_1 = x_2$
whereas its minima occur when \mathbf{x} is 1 sparse

Entropy Function Induces Sparsity

Consider a nonnegative diagonal matrix \mathbf{X}

Let $\mathbf{x} = \text{diag}(\mathbf{X})$, then $\mathbf{x} = \sigma(\mathbf{X})$. Assume \mathbf{x} is sparse.

$$\rightarrow \min_{\mathbf{x}} h(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}.$$

Lemma. *If there exists two solutions $\mathbf{x}_1 \neq \mathbf{x}_2$ to $\mathbf{A}\mathbf{x} = \mathbf{b}$, with $\mathbf{b} \neq \mathbf{0}$, in the same d -dimensional orthant ($d \leq n$), then there is at least one solutions \mathbf{x}' in some d' -dimensional orthant such that $d' < d$ and $h(\mathbf{x}') < \min\{h(\mathbf{x}_1), h(\mathbf{x}_2)\}$.*

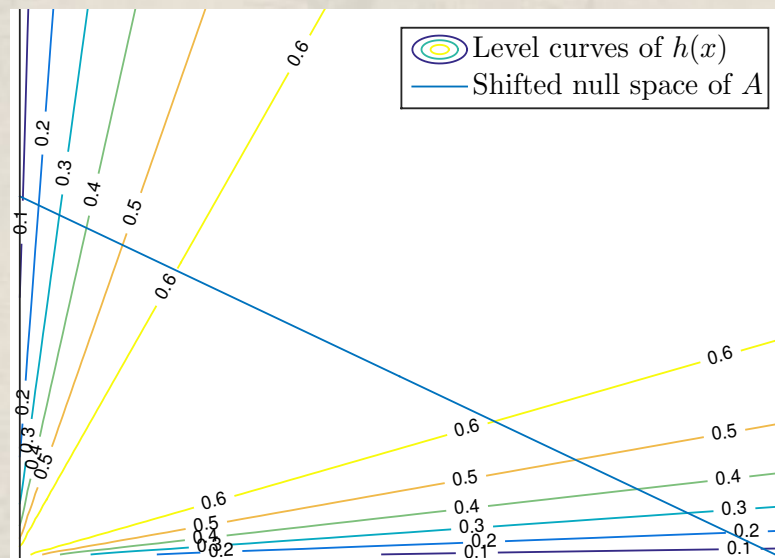


Fig. Illustration of Lemma 1: minimum entropy occurs at 1-sparse solution.

ENM Algorithm

- ◆ *Robust variant:*

$$\min_{\mathbf{X}} \lambda h(\boldsymbol{\sigma}(\mathbf{X})) + f(\mathbf{X}; \mathcal{A}, \mathbf{y}),$$

for some loss function $f : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}_+$ with Lipschitz continuous gradient.

- ◆ *Technique: linearization*

At current estimate $\boldsymbol{\sigma}^t$:

$$h(\boldsymbol{\sigma}) \approx h(\boldsymbol{\sigma}^t) + \nabla h(\boldsymbol{\sigma}^t)^T (\boldsymbol{\sigma} - \boldsymbol{\sigma}^t)$$

$$f(\mathbf{X}) \approx f(\mathbf{X}^t) + \nabla f(\mathbf{X}^t)^T (\mathbf{X} - \mathbf{X}^t) + \frac{\rho}{2} \|\mathbf{X} - \mathbf{X}^t\|_F^2,$$

$$\Rightarrow \mathbf{X}^{t+1} = \operatorname{argmin}_{\mathbf{X}} \lambda \nabla h(\boldsymbol{\sigma}^t)^T \boldsymbol{\sigma} + f(\mathbf{X}^t)$$

$$+ \nabla f(\mathbf{X}^t)^T (\mathbf{X} - \mathbf{X}^t) + \frac{\rho}{2} \|\mathbf{X} - \mathbf{X}^t\|_F^2$$

$$= \operatorname{argmin}_{\mathbf{X}} \lambda \nabla h(\boldsymbol{\sigma}^t)^T \boldsymbol{\sigma} + \frac{\rho}{2} \left\| \mathbf{X} - \left(\mathbf{X}^t - \frac{1}{\rho} \nabla f(\mathbf{X}^t) \right) \right\|_F^2.$$

ENM Algorithm

$$\mathbf{X}^{t+1} = \underset{\mathbf{X}}{\operatorname{argmin}} \lambda \nabla h(\boldsymbol{\sigma}^t)^T \boldsymbol{\sigma} + \frac{\rho}{2} \left\| \mathbf{X} - \left(\mathbf{X}^t - \frac{1}{\rho} \nabla f(\mathbf{X}^t) \right) \right\|_F^2.$$

Lemma. Let h be the entropy function, and let $\boldsymbol{\sigma}$ be a positive vector, then

$$\frac{\partial h(\boldsymbol{\sigma})}{\partial \sigma_i} = -\frac{\log \sigma_i}{\|\boldsymbol{\sigma}\|_1} + \frac{\sum_j \sigma_j \log \sigma_j}{\|\boldsymbol{\sigma}\|_1^2}.$$

Lemma. If $\sigma_1^t \geq \sigma_2^t \geq \dots \geq \sigma_n^t \geq 0$, then $0 \leq \frac{\partial h(\boldsymbol{\sigma}^t)}{\partial \sigma_1} \leq \frac{\partial h(\boldsymbol{\sigma}^t)}{\partial \sigma_2} \leq \dots \leq \frac{\partial h(\boldsymbol{\sigma}^t)}{\partial \sigma_n}$.

Lemma. Let $\lambda > 0$, $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$, and $0 \leq w_1 \leq w_2 \leq \dots \leq w_n$, where $n = \min\{n_1, n_2\}$. Let \mathbf{X}^* be the optimal solution of the minimization problem

$$\min_{\mathbf{X}} \lambda \sum_i w_i \sigma_i(\mathbf{X}) + \frac{1}{2} \|\mathbf{X} - \mathbf{Z}\|_F^2, \quad (1)$$

then

$$\mathbf{X}^* = \mathbf{U} \mathcal{D}_{\lambda \mathbf{w}}(\boldsymbol{\Sigma}) \mathbf{V}^T, \quad (2)$$

where $\mathbf{Z} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$ and $\mathcal{D}_{\lambda \mathbf{w}}(\boldsymbol{\Sigma}) = \operatorname{diag}\{(\sigma_i - \lambda w_i)_+\}$ is the singular value shrinkage operator.

ENM Algorithm

Algorithm 1 ENtropy-Minimization (ENM)

input: measurements $(\mathcal{A}, \mathbf{y})$, $\lambda > 0$, and $\rho > L_f$.

initialization: \mathbf{X}^0 .

while not converged **do**

 Update the weights:

$$w_i^t = -\frac{\log \sigma_i^t}{\|\boldsymbol{\sigma}^t\|_1} + \frac{\sum_j \sigma_j^t \log \sigma_j^t}{\|\boldsymbol{\sigma}^t\|_1^2}, \quad i = 1, \dots, n \quad (1)$$

 Update the estimate:

$$\mathbf{X}^{t+1} = U\mathcal{D}_{(\lambda/\rho)\mathbf{w}}(\boldsymbol{\Sigma})\mathbf{V}^T, \quad (2)$$

 where $\mathbf{X}^t - \frac{1}{\rho}\nabla f(\mathbf{X}^t) = U\boldsymbol{\Sigma}\mathbf{V}^T$.

end while

output: Estimated solution $\hat{\mathbf{X}}$.

Experiment Results

- ◆ *Random subsampling (matrix completion):*

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* \quad \text{s.t.} \quad X_{ij} = M_{ij}, \quad (i, j) \in \Omega$$

- ◆ *Synthetic data*

$$\mathbf{X} \in \mathbb{R}^{100 \times 100}$$

$$m = 0.5n_1n_2 = 5000 \text{ samples}$$

$$r = \text{rank}(\mathbf{X}) \text{ varies}$$

$$\frac{\|\hat{\mathbf{X}} - \mathbf{M}\|_F^2}{\|\mathbf{M}\|_F^2} < 10^{-3} \rightarrow \text{successful}$$

100 trials for each r

Algorithms:

- ENtropy Minimization (ENM)
- Singular Value Thresholding (SVT) [Cai *et al.*]
- Augmented Lagrange Multiplier (ALM) [Lin *et al.*]

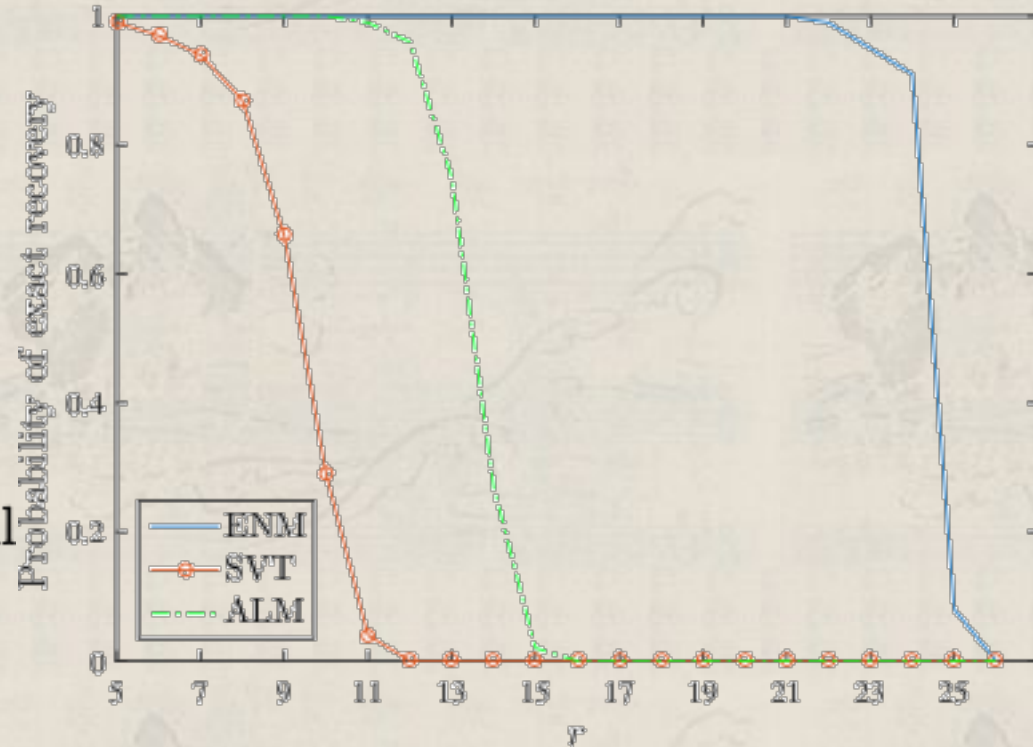


Fig. Probability of exact recovery on synthetic data.

Experiment Results

- ◆ *Random subsampling (matrix completion):*

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* \quad \text{s.t.} \quad X_{ij} = M_{ij}, \quad (i, j) \in \Omega$$

- ◆ *Real data: MIT Logo*

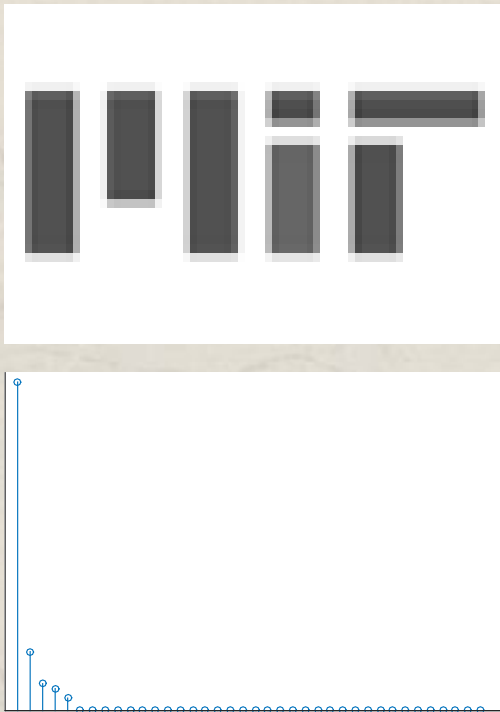


Fig. MIT logo and its singular values

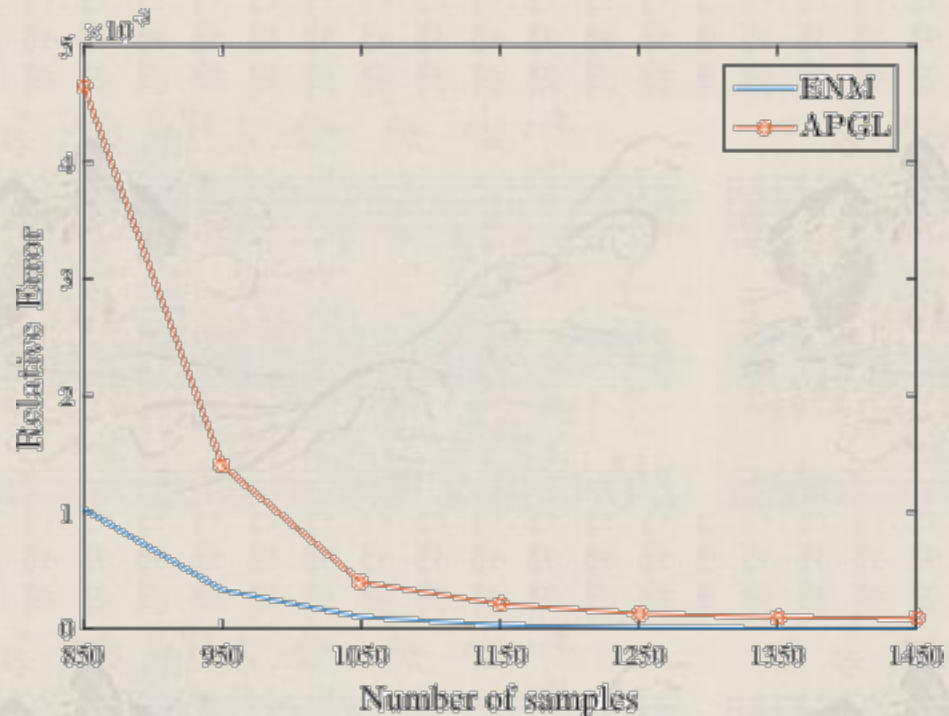


Fig. Low-rank matrix completion on MIT logo.

Algorithms:

- ENM vs. Accelerated Proximal Gradient with Line Search (APGL) [Toh & Yun]

What About Re-weighted ℓ_1 ?

- ◆ ENM

$$\min_{\mathbf{x}} h(\mathbf{x}) \text{ s.t. } \mathbf{Ax} = \mathbf{y}$$

Not Separable in Its Parameters

- ◆ Reweighted ℓ_1

$$\min_{\mathbf{x}} \|\log(\mathbf{x})\|_1 \text{ s.t. } \mathbf{Ax} = \mathbf{y}$$

- ◆ Common optimization problem

$$\min_{\mathbf{w}, \mathbf{x}} \|\mathbf{w} \circ \mathbf{x}\|_1 + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2$$

Key Difference

Discussion: ENM Weights

- ◆ Closed-form thresholds

$$w_i^t = -\frac{\log \sigma_i^t}{\|\boldsymbol{\sigma}^t\|_1} + \frac{\sum_j \sigma_j^t \log \sigma_j^t}{\|\boldsymbol{\sigma}^t\|_1^2}$$

- ◆ Threshold proven to be larger when associated singular value gets smaller → lower-rank solution from shrinkage operation!
- ◆ ENM encourages singular values to have a Laplacian distribution

Final Discussion

- ◆ Both ENM and Re-weighted ℓ_1 can be solved with the same strategy
- ◆ ENM seems to empirically offer a better weighting scheme than Re-weighted ℓ_1
- ◆ Critical Questions:
 - How can we theoretically justify that?
 - Any other interesting applications of information theory tools/concepts to sparse problems?

Entropy Function Induces Sparsity

Lemma. *If there exists two solutions $\mathbf{x}_1 \neq \mathbf{x}_2$ to $A\mathbf{x} = \mathbf{b}$, with $\mathbf{b} \neq \mathbf{0}$, in the same d -dimensional orthant ($d \leq n$), then there is at least one solutions \mathbf{x}' in some d' -dimensional orthant such that $d' < d$ and $h(\mathbf{x}') < \min\{h(\mathbf{x}_1), h(\mathbf{x}_2)\}$.*

Proof.(sketch) For any $\mathbf{x} \in \mathbb{R}^n$, define $\hat{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|_1}$, then $h(\mathbf{x}) = h(\hat{\mathbf{x}})$.

Let $\mathbf{x}' = (1 - \lambda)\mathbf{x}_1 + \lambda\mathbf{x}_2$, then there is an one-to one mapping between \mathbf{x}' and $\hat{\mathbf{x}}'$.

Furthermore, $\hat{\mathbf{x}}' = (1 - \hat{\lambda})\hat{\mathbf{x}}_1 + \hat{\lambda}\hat{\mathbf{x}}_2$, for some constant $\hat{\lambda}$.

As $h(\hat{\mathbf{x}}')$ is concave, its minima are archived at the boundaries of the current orthant.

This also true for $h(\mathbf{x}')$ by the one-to-one mapping between \mathbf{x}' and $\hat{\mathbf{x}}'$.

Therefore, we can tune λ so that \mathbf{x}' lies at a boundary of the current orthant and minimum entropy is archived.

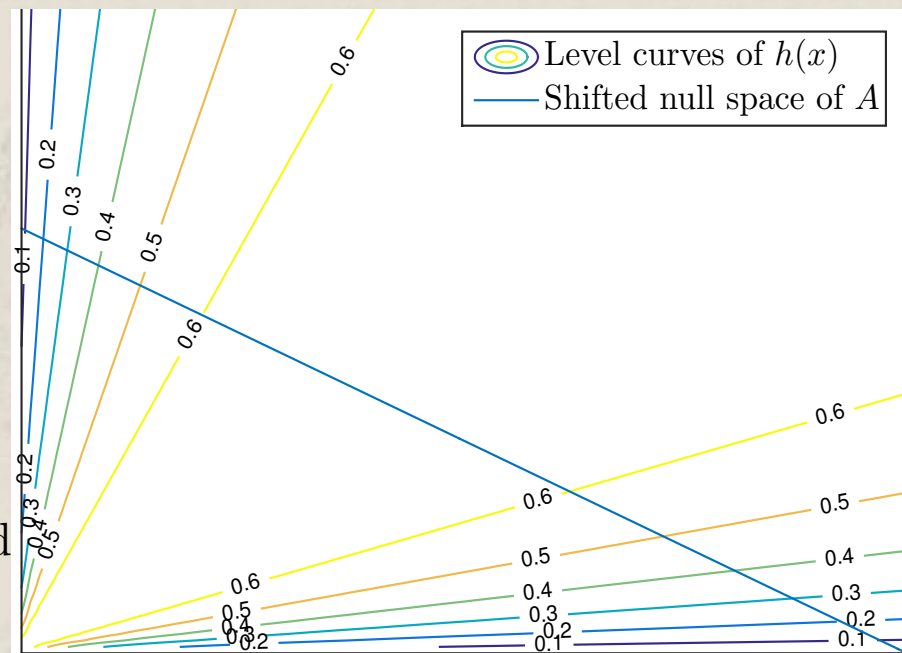


Fig. Illustration of Lemma 1: minimum entropy occurs at 1-sparse solution.