

Phase Transitions in Coding, Communications, and Inference

Andrea Montanari and Rüdiger Urbanke

Stanford University and EPFL, Lausanne

January 5, 2009

From coding to probabilistic inference

Binary symmetric channel



$$\underline{y} = \underline{x} \oplus \underline{z}, \quad \mathbb{H}_{\underline{x}} = 0.$$

$$\underline{z} = (z_1, z_2, \dots, z_n), \quad z_i \text{'s iid Bernoulli}(p).$$

Binary symmetric channel



$$\underline{y} = \underline{x} \oplus \underline{z}, \quad \mathbb{H}_{\underline{x}} = 0.$$

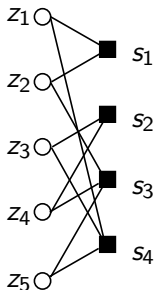
$$\underline{z} = (z_1, z_2, \dots, z_n), \quad z_i \text{'s iid Bernoulli}(p).$$

$$\underline{s} = \mathbb{H}\underline{y} = \mathbb{H}\underline{z}$$

$$\underline{s} = \mathbb{H}\underline{y} = \mathbb{H}\underline{z}$$

Syndrome decoding

$(z_1, \dots, z_n) \in \{0, 1\}^n$ is np sparse

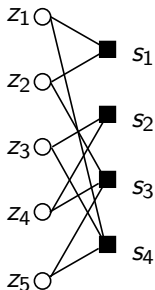


$$\underline{s} = \mathbb{H}\underline{z} \pmod{2}$$

s_1, \dots, s_m : linear observations of the noise vector.

Syndrome decoding

$(z_1, \dots, z_n) \in \{0, 1\}^n$ is np sparse



$$\underline{s} = \mathbb{H}\underline{z} \pmod{2}$$

s_1, \dots, s_m : linear observations of the noise vector.

Outline

- 1 Sparse vectors
- 2 Collaborative filtering
- 3 Conclusion

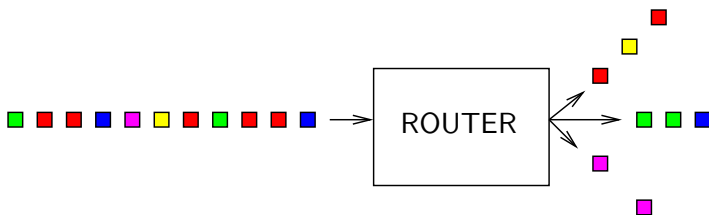
SPARSE VECTORS

Why sparse vectors? Network measurements

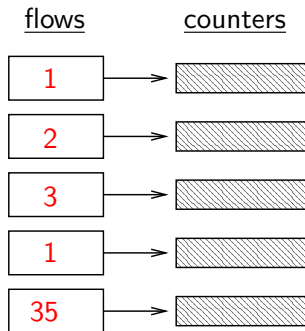
Network measurements

Internet core link :

- $10^{10} \div 10^{11}$ bits/sec
- Packet size 10^3 bits
- $10^6 \div 10^7$ flows per hour (mice... elephants)



A naive approach (1 flow \leftrightarrow 1 counter)



Problem 1: Memory

Processing time: 12 nanosec per packet

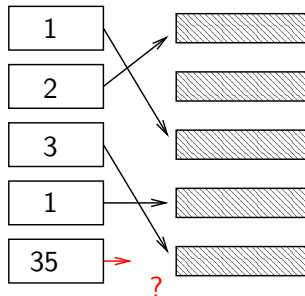
Space: 10^6 flows \times 64 bits counters = 8 MBytes

Problem 1: Memory

Processing time: 12 nanosec per packet

Space: 10^6 flows \times 64 bits counters = 8 MBytes

Problem 2: Flow-to-Counter association



Hybrid architectures

[Shah, Iyer, Prabhakar, McKeown 2002]

Sampling

[Estan, Varghese 2001; CISCO's NetFlow]

Compressed sensing

[Candes, Donoho, Romberg, Tao, Indyk, Gilbert, Tanner 2006-...]

Hybrid architectures

[Shah, Iyer, Prabhakar, McKeown 2002]

Sampling

[Estan, Varghese 2001; CISCO's NetFlow]

Compressed sensing

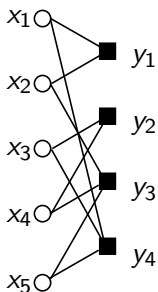
[Candes, Donoho, Romberg, Tao, Indyk, Gilbert, Tanner 2006-...]

Counter braids: Vanilla version

Counter braid: one layer

flows

counters



$$y = \mathbb{H}x$$

$$x = (x_1, \dots, x_n) \in \mathbb{N}^n$$

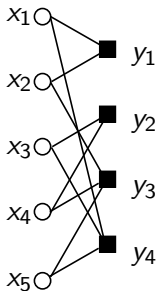
$$y = (y_1, \dots, y_m) \in \mathbb{N}^m$$

\mathbb{H} adjacency matrix

Decoding: probabilistic inference

flows

counters



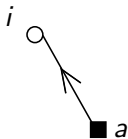
$$y = \mathbb{H}x$$

$$\mu(\underline{x}) = \frac{1}{Z} \prod_{a=1}^m \mathbb{I}(y_a = h_a^T x) \prod_{i=1}^n p_0(x_i)$$

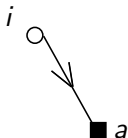
Message passing algorithm

BP messages $\rightarrow \nu_{i \rightarrow a}^{(t)}(x_i), \hat{\nu}_{a \rightarrow i}^{(t)}(x_i).$

Low complexity message passing algorithm



$\widehat{\nu}_{a \rightarrow i}^{(t)}$: counter a to flow i

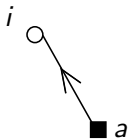


$\nu_{i \rightarrow a}^{(t)}$: flow i to counter a

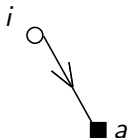
$$\nu_{i \rightarrow a}^{(t)}, \widehat{\nu}_{a \rightarrow i}^{(t)} \in \mathbb{N},$$

$$\nu_{i \rightarrow a}^{(0)} = 0$$

Low complexity message passing algorithm



$\widehat{\nu}_{a \rightarrow i}^{(t)}$: counter a to flow i

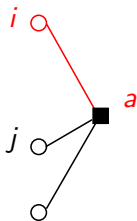


$\nu_{i \rightarrow a}^{(t)}$: flow i to counter a

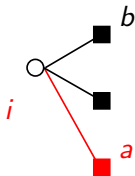
$$\nu_{i \rightarrow a}^{(t)}, \widehat{\nu}_{a \rightarrow i}^{(t)} \in \mathbb{N},$$

$$\nu_{i \rightarrow a}^{(0)} = 0$$

Decoding a counter braid



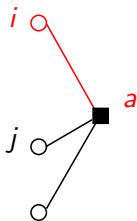
$$\hat{v}_{a \rightarrow i}^{(t)} = \left[y_a - \sum_{j \in \partial a \setminus i} v_{j \rightarrow a}^{(t)} \right]_+$$



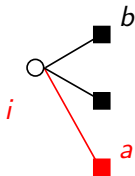
$$v_{i \rightarrow a}^{(t+1)} = \min \left\{ \hat{v}_{b \rightarrow i}^{(t)} : b \in \partial i \setminus a \right\} \quad \text{for } t \text{ even}$$

$$v_{i \rightarrow a}^{(t+1)} = \max \left\{ \hat{v}_{b \rightarrow i}^{(t)} : b \in \partial i \setminus a \right\} \quad \text{for } t \text{ odd}$$

Decoding a counter braid



$$\hat{v}_{a \rightarrow i}^{(t)} = \left[y_a - \sum_{j \in \partial a \setminus i} v_{j \rightarrow a}^{(t)} \right]_+$$



$$v_{i \rightarrow a}^{(t+1)} = \min \left\{ \hat{v}_{b \rightarrow i}^{(t)} : b \in \partial i \setminus a \right\} \quad \text{for } t \text{ even}$$

$$v_{i \rightarrow a}^{(t+1)} = \max \left\{ \hat{v}_{b \rightarrow i}^{(t)} : b \in \partial i \setminus a \right\} \quad \text{for } t \text{ odd}$$

Analysis and performances

The sandwich property

Proposition

$$\nu_i^{(0)} \leq \nu_i^{(2)} \leq \nu_i^{(4)} \leq \nu_i^{(6)} \leq \dots \leq x_i$$

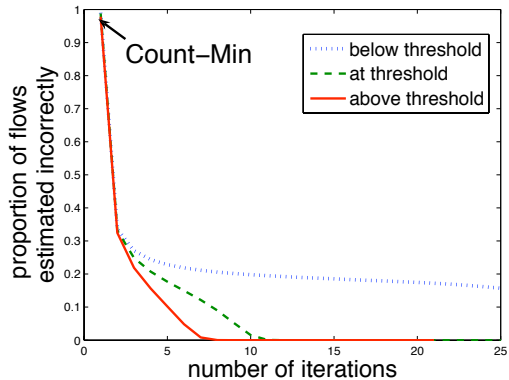
$$\nu_i^{(1)} \geq \nu_i^{(3)} \geq \nu_i^{(5)} \geq \nu_i^{(7)} \geq \dots \geq x_i$$

A statistical model?

X_1, X_2, \dots, X_n iid, $X_i \geq X_{\min}$

$$\mathbb{P}\{X_i > X_{\min}\} = \epsilon.$$

Typical runs



Density evolution

l : flows degree

r : counters degree

$$z_t \equiv \mathbb{P}\{\nu_{i \rightarrow a}^{(t)} \neq x_i\}$$

$$z_{t+1} = \begin{cases} (1 - (1 - z_t)^{r-1})^{l-1} & \text{for } t \text{ even,} \\ \epsilon(1 - (1 - z_t)^{r-1})^{l-1} & \text{for } t \text{ odd,} \end{cases}$$

Density evolution

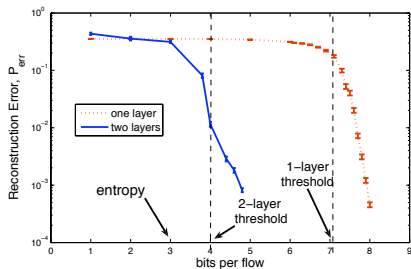
l : flows degree

r : counters degree

$$z_t \equiv \mathbb{P}\{\nu_{i \rightarrow a}^{(t)} \neq x_i\}$$

$$z_{t+1} = \begin{cases} (1 - (1 - z_t)^{r-1})^{l-1} & \text{for } t \text{ even,} \\ \epsilon(1 - (1 - z_t)^{r-1})^{l-1} & \text{for } t \text{ odd,} \end{cases}$$

Threshold in memory space



$$n = 1000, \mathbb{P}\{X_1 \geq x\} = x^{-3/2}$$

$$\gamma = \frac{\text{\#counters}}{\text{\#flows}} .$$

Optimal dimensionality reduction

Theorem (Donoho, Tanner, 2006)

Let $\gamma_{dens}(\epsilon)$ be the dimensionality reduction rate for ϵ -sparse sources, with Gaussian random matrices and LP decoding. Then

$$\gamma_{dens}(\epsilon) = 2 \cdot \epsilon \log(1/\epsilon) + O(\epsilon).$$

Theorem (Lu, M, Prabhakar, 2008)

Let $\gamma_{sparse}(\epsilon)$ be the dimensionality reduction rate for ϵ -sparse sources, sparse matrices and message passing decoding. Then

$$\gamma_{sparse}(\epsilon) \leq 2.09 \cdot \epsilon \log(1/\epsilon) + O(\epsilon).$$

Optimal dimensionality reduction

Theorem (Donoho, Tanner, 2006)

Let $\gamma_{dens}(\epsilon)$ be the dimensionality reduction rate for ϵ -sparse sources, with Gaussian random matrices and LP decoding. Then

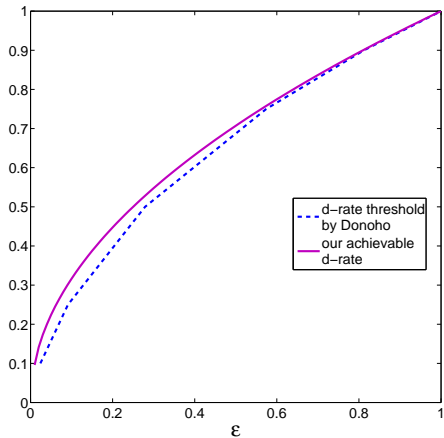
$$\gamma_{dens}(\epsilon) = 2 \cdot \epsilon \log(1/\epsilon) + O(\epsilon).$$

Theorem (Lu, M, Prabhakar, 2008)

Let $\gamma_{sparse}(\epsilon)$ be the dimensionality reduction rate for ϵ -sparse sources, sparse matrices and message passing decoding. Then

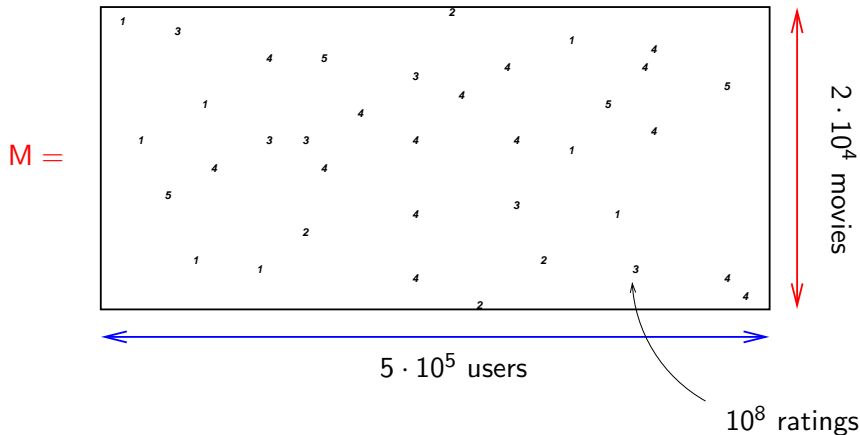
$$\gamma_{sparse}(\epsilon) \leq 2.09 \cdot \epsilon \log(1/\epsilon) + O(\epsilon).$$

Optimal dimensionality reduction

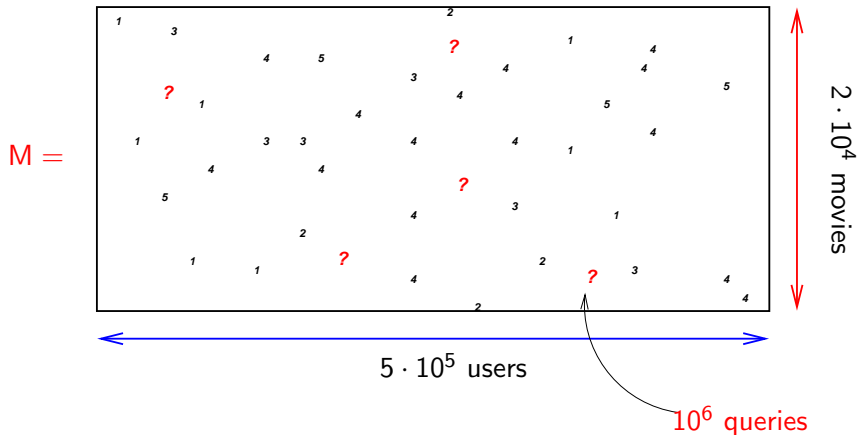


COLLABORATIVE FILTERING

Netflix dataset: A big (!) matrix



A big (!) matrix



You get a prize if...

RMSE < 0.8563 ; -)

Is this possible?

You get a prize if...

$\text{RMSE} < 0.8563$; -)

Is this possible?

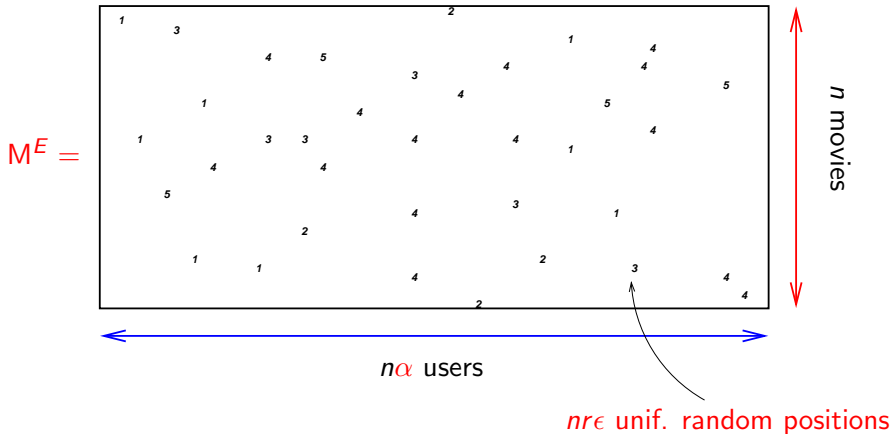
You get a prize if...

$\text{RMSE} < 0.8563$; -)

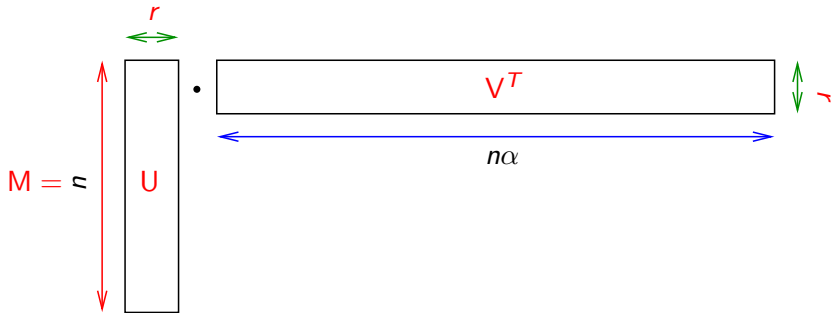
Is this possible?

A model: Random low-rank matrices

The observations

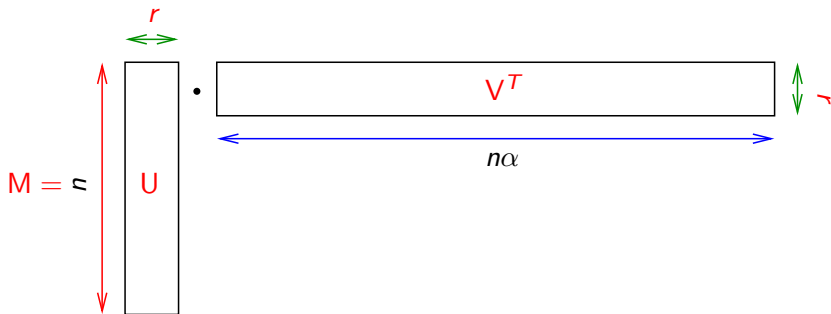


You need some structure!



$$r \leq n^{1/2}$$

You need some structure!



$$r \leq n^{1/2}$$

- U_{ik}, V_{ak} i.i.d.
- U, V random orthogonal.
- U, V 'incoherent'.

- U_{ik}, V_{ak} i.i.d.
- U, V random orthogonal.
- U, V 'incoherent'.

- U_{ik}, V_{ak} i.i.d.
- U, V random orthogonal.
- U, V 'incoherent'.

$$D(M, \hat{M}) \equiv \left\{ \frac{1}{n^2 \alpha} \sum_{i,a} |M_{ia} - \hat{M}_{ia}|^2 \right\}^{1/2}$$

Theorem (Candés, Recht, 2008)

If

$$\epsilon \geq C n^{1/5} \log n$$

then whp

1. *M is unique given the observed entries.*
2. *M is the unique minimum of a SDP.*

cf. also [Recht, Fazel, Parrilo 2007]

Theorem (Candés, Recht, 2008)

If

$$\epsilon \geq C n^{1/5} \log n$$

then whp

1. *M is unique given the observed entries.*
2. *M is the unique minimum of a SDP.*

cf. also [Recht, Fazel, Parrilo 2007]

Theorem (Candés, Recht, 2008)

If

$$\epsilon \geq C n^{1/5} \log n$$

then whp

1. *M is unique given the observed entries.*
2. *M is the unique minimum of a SDP.*

cf. also [Recht, Fazel, Parrilo 2007]

Theorem (Candés, Recht, 2008)

If

$$\epsilon \geq C n^{1/5} \log n$$

then whp

1. *M is unique given the observed entries.*
2. *M is the unique minimum of a SDP.*

cf. also [Recht, Fazel, Parrilo 2007]

Theorem (Candés, Recht, 2008)

If

$$\epsilon \geq C n^{1/5} \log n$$

then whp

1. *M is unique given the observed entries.*
2. *M is the unique minimum of a SDP.*

cf. also [Recht, Fazel, Parrilo 2007]

Great, but...

1. $n^{1/5}$ observations for 1 bit of information?
2. RMSE = 0?
3. SDP = $O(n^{4...6})$. Substitute $n = 10^5$...

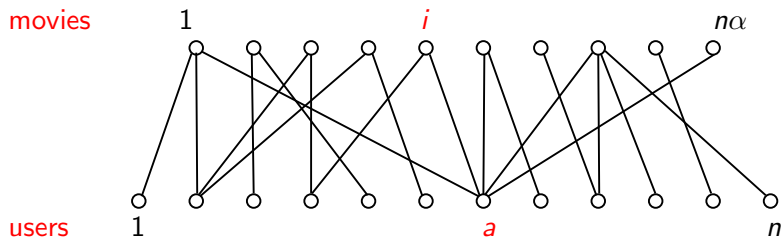
1. $n^{1/5}$ observations for 1 bit of information?
2. RMSE = 0?
3. SDP = $O(n^{4\dots 6})$. Substitute $n = 10^5 \dots$

1. $n^{1/5}$ observations for 1 bit of information?
2. RMSE = 0?
3. SDP = $O(n^{4...6})$. Substitute $n = 10^5$...

1. $n^{1/5}$ observations for 1 bit of information?
2. RMSE = 0?
3. SDP = $O(n^{4\dots 6})$. Substitute $n = 10^5 \dots$

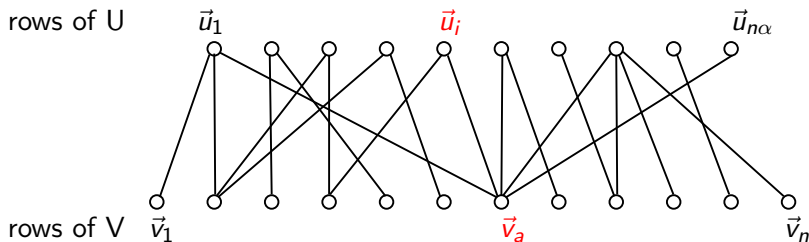
A graphical model

The graph



$(i, a) \in E \Leftrightarrow$ User a rated movie i .

The graphical model



$$\mu(\{\vec{u}_i\}, \{\vec{v}_a\}) = \frac{1}{Z} \prod_{(i,a) \in E} \mathbb{I}(\vec{u}_i \cdot \vec{v}_a = M_{ia}) \prod_{i=1}^{n\alpha} p_0(\vec{u}_i) \prod_{a=1}^n p_0(\vec{v}_a).$$

Messages $\nu_{i \rightarrow a}(\vec{u}_i)$, $\nu_{a \rightarrow i}(\vec{v}_a)$.

A small simulation

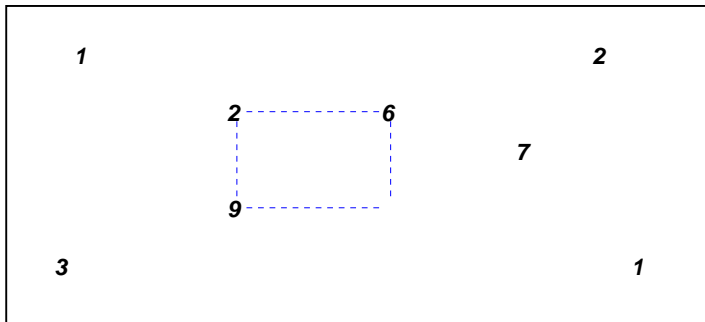
Messages $\nu_{i \rightarrow a}(\vec{u}_i)$, $\nu_{a \rightarrow i}(\vec{v}_a)$.

A small simulation

$O(n)$ entries are enough (practice)

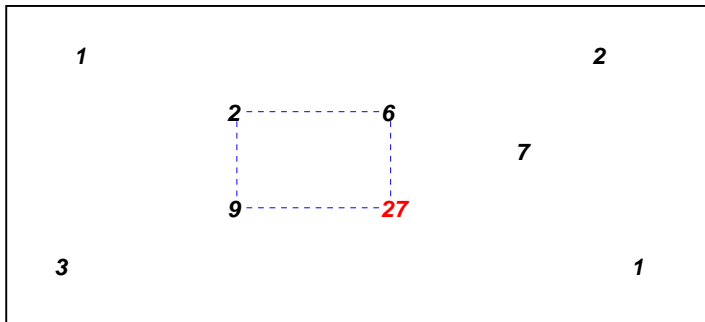
Rank = 1: an easy trick

M =

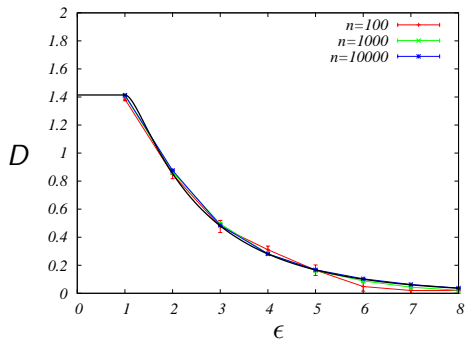


Rank = 1: an easy trick

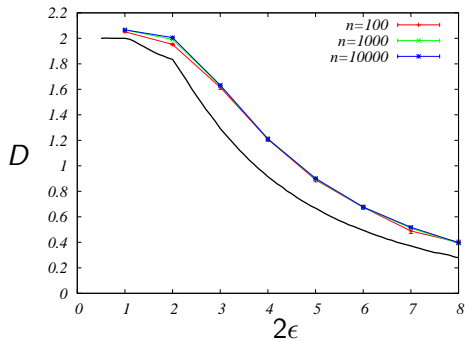
M =



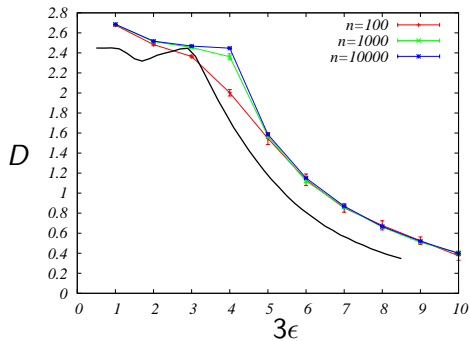
Rank = 1: Trick vs. Belief Propagation



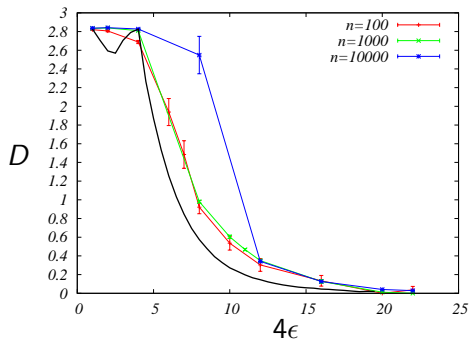
Rank = 2: Belief Propagation



Rank = 3: Belief Propagation



Rank = 4: Belief Propagation



$O(n)$ entries are enough (theory)

Can you prove anything about BP?

Assume $p_0(\cdot)$ uniform over $\{-1, +1\}^r$

Uniform fixed point $\nu_{i \rightarrow a}^*(\cdot) = \nu_{a \rightarrow i}^*(\cdot) = p_0(\cdot)$

- Becomes unstable for $r\epsilon > 1$.
- Two new $+/-$ symmetric fixed points.
- Become unstable for $\epsilon > \text{const.}$:-)

Can you prove anything about BP?

Assume $p_0(\cdot)$ uniform over $\{-1, +1\}^r$

Uniform fixed point $\nu_{i \rightarrow a}^*(\cdot) = \nu_{a \rightarrow i}^*(\cdot) = p_0(\cdot)$

- Becomes unstable for $r\epsilon > 1$.
- Two new $+/-$ symmetric fixed points.
- Become unstable for $\epsilon > \text{const.}$:-)

Can you prove anything about BP?

Assume $p_0(\cdot)$ uniform over $\{-1, +1\}^r$

Uniform fixed point $\nu_{i \rightarrow a}^*(\cdot) = \nu_{a \rightarrow i}^*(\cdot) = p_0(\cdot)$

- Becomes unstable for $r\epsilon > 1$.
- Two new $+/-$ symmetric fixed points.
- Become unstable for $\epsilon > \text{const.}$:-)

Can you prove anything about BP?

Assume $p_0(\cdot)$ uniform over $\{-1, +1\}^r$

Uniform fixed point $\nu_{i \rightarrow a}^*(\cdot) = \nu_{a \rightarrow i}^*(\cdot) = p_0(\cdot)$

- Becomes unstable for $r\epsilon > 1$.
- Two new $+/-$ symmetric fixed points.
- Become unstable for $\epsilon > \text{const.}$:-)

Naive spectral algorithm

$$M_{ia}^E = \begin{cases} M_{ia} & \text{if } (i, a) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

Projection

$$M^E = \sum_{i=1}^n \sigma_i x_i y_i^T, \quad \sigma_1 \geq \sigma_2 \geq \dots$$

$$\text{Tr}_r(M^E) = \frac{n\sqrt{\alpha}}{\epsilon} \sum_{i=1}^r \sigma_i x_i y_i^T.$$

$$M_{ia}^E = \begin{cases} M_{ia} & \text{if } (i, a) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

Projection

$$M^E = \sum_{i=1}^n \sigma_i x_i y_i^T, \quad \sigma_1 \geq \sigma_2 \geq \dots$$

$$\text{Tr}_r(M_E) = \frac{n\sqrt{\alpha}}{\epsilon} \sum_{i=1}^r \sigma_i x_i y_i^T.$$

If $\epsilon = O(1)$, 'spurious' singular values $\Omega(\sqrt{\log n / (\log \log n)})$.

Trimming

$$\tilde{M}_{ia}^E = \begin{cases} M_{ia}^E & \text{if } \deg(i) \leq 2 \mathbb{E} \deg(i), \quad \deg(a) \leq 2 \mathbb{E} \deg(a), \\ 0 & \text{otherwise.} \end{cases}$$

If $\epsilon = O(1)$, 'spurious' singular values $\Omega(\sqrt{\log n / (\log \log n)})$.

Trimming

$$\tilde{M}_{ia}^E = \begin{cases} M_{ia}^E & \text{if } \deg(i) \leq 2 \mathbb{E} \deg(i), \quad \deg(a) \leq 2 \mathbb{E} \deg(a), \\ 0 & \text{otherwise.} \end{cases}$$

SPECTRAL MATRIX COMPLETION(matrix M^E)

- 1: Trim M^E , and let \tilde{M}^E be the output;
 - 2: Project \tilde{M}^E to $T_r(\tilde{M}^E)$;
 - 3: Clean residual errors by coordinate descent in the factors.
-

SVD of \tilde{M}^E

- Standard algorithms $\rightarrow O(n^3)$
- Iterative $\rightarrow O(nr\epsilon \log n)$

Theorem (Keshavan, Oh, M, 2009)

For each $\delta > 0$, if $\epsilon \geq C(\alpha, \delta)$, then with high probability

$$\|\mathbf{M} - \mathbf{T}_r(\tilde{\mathbf{M}}^E)\|_{\mathbf{F}}^2 \leq n^2 r \delta.$$

Theorem (Keshavan, Oh, M, 2009)

If $\epsilon \geq C'(\alpha) \log n$, then SPECTRAL MATRIX COMPLETION returns, with high probability, the matrix \mathbf{M} .

Theorem (Keshavan, Oh, M, 2009)

For each $\delta > 0$, if $\epsilon \geq C(\alpha, \delta)$, then with high probability

$$\|M - \text{Tr}_r(\tilde{M}^E)\|_F^2 \leq n^2 r \delta.$$

Theorem (Keshavan, Oh, M, 2009)

If $\epsilon \geq C'(\alpha) \log n$, then SPECTRAL MATRIX COMPLETION returns, with high probability, the matrix M .

Key technical result

$$M = \sum_{i=1}^r \Sigma_i \underline{u}_i \underline{v}_i^T,$$

$$\|\underline{u}_i\| = \sqrt{n}, \quad \underline{u}_i^T \underline{u}_j = 0, \quad \|\underline{v}_i\| = \sqrt{n\alpha}, \quad \underline{v}_i^T \underline{v}_j = 0.$$

Theorem

If $\{\underline{u}_i\}, \{\underline{v}_i\}$ are incoherent, then, w.h.p.

$$\begin{aligned} |\sigma_q - \epsilon r \Sigma_q| &\leq C r \sqrt{\epsilon} \log \epsilon && \text{for } q \leq r, \\ \sigma_q &\leq C r \sqrt{\epsilon} && \text{for } q > r. \end{aligned}$$

Key technical result

$$M = \sum_{i=1}^r \Sigma_i \underline{u}_i \underline{v}_i^T,$$

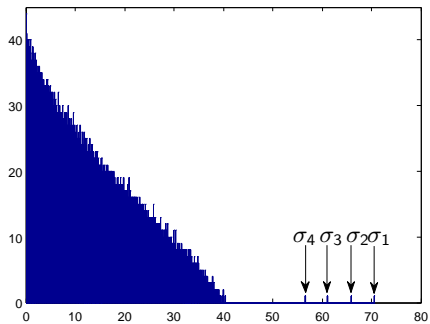
$$\|\underline{u}_i\| = \sqrt{n}, \quad \underline{u}_i^T \underline{u}_j = 0, \quad \|\underline{v}_i\| = \sqrt{n\alpha}, \quad \underline{v}_i^T \underline{v}_j = 0.$$

Theorem

If $\{\underline{u}_i\}, \{\underline{v}_i\}$ are incoherent, then, w.h.p.

$$\begin{aligned} |\sigma_q - \epsilon r \Sigma_q| &\leq C r \sqrt{\epsilon} \log \epsilon && \text{for } q \leq r, \\ \sigma_q &\leq C r \sqrt{\epsilon} && \text{for } q > r. \end{aligned}$$

$n = 10000, r = 4, \epsilon = 12.5$



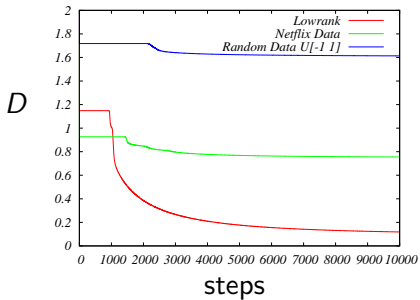
Back to the data

Is Netflix a random low-rank matrix?

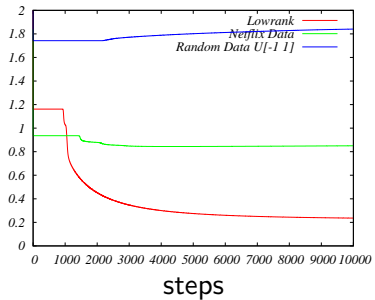
Compare for coordinate descent (SimonFunk).

$$(n = 5 \cdot 10^3, \alpha = 1)$$

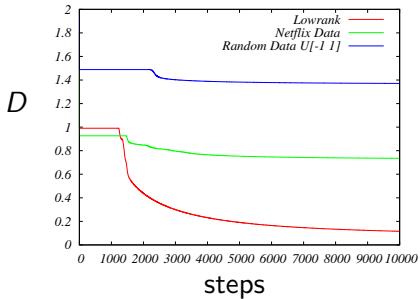
fit error



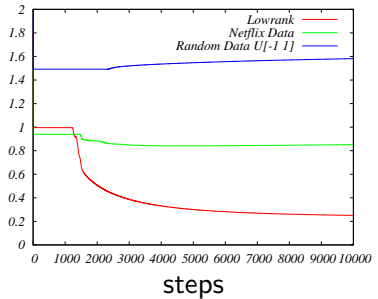
pred. error



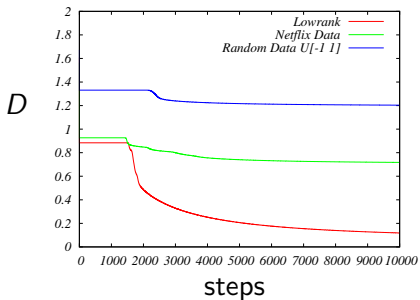
fit error



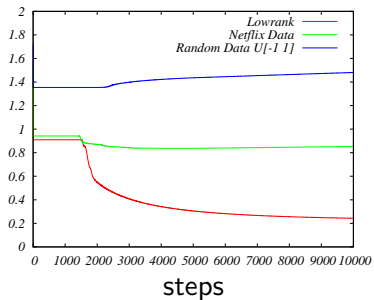
pred. error



fit error



pred. error



CONCLUSION

Enough information, measurements, . . . \Rightarrow Threshold

Engineering phase transitions.

Enough information, measurements, . . . \Rightarrow Threshold

Engineering phase transitions.