

Community Detection via Semidefinite Programming

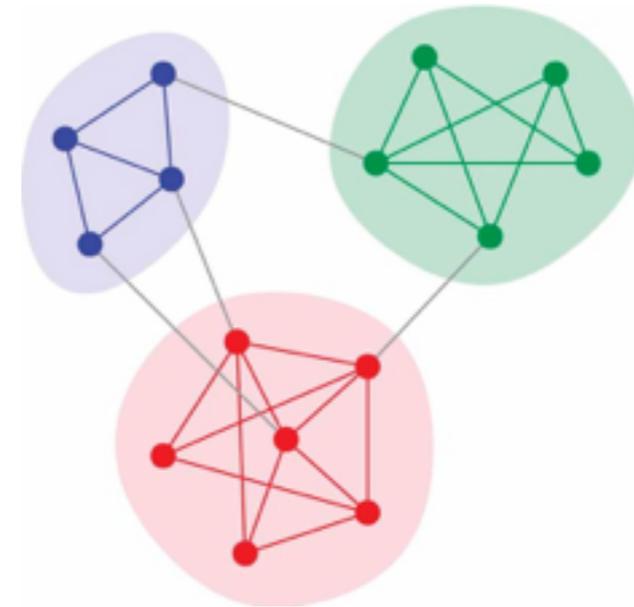
Federico Ricci-Tersenghi
(Sapienza University)

in collaboration with
Adel Javanmard and Andrea Montanari

PNAS 113, E2218 (2016)
J. Phys.: Conf. Ser. 699, 012015 (2016)

Communities detection problem

- Detecting communities/partitions/clusters in graphs is a widespread problem in many different disciplines



- We need fast (linear and scalable) algorithms
 - robust (real datasets are very noisy and not random)
 - close to optimal (on random ensemble benchmarks)

Benchmark for community detection

Hidden partition model or stochastic block model (SBM)

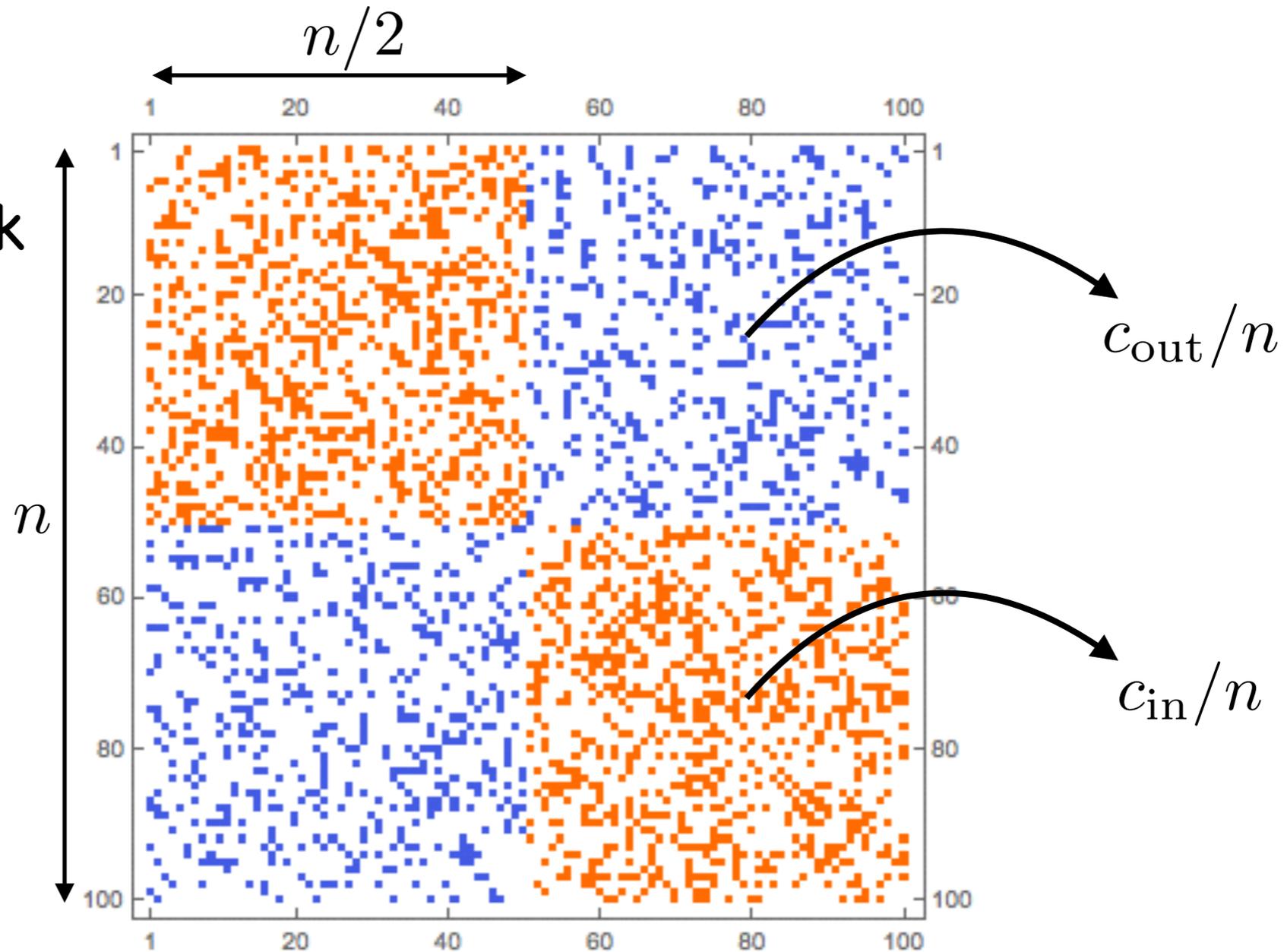
- Generate a partition of n nodes: e.g. q groups of size n/q
- Add independently edges between any pair of nodes according to the following probability

$$\mathbb{P}[(ij) \in E] = \begin{cases} c_{\text{in}}/n & \text{same group} \\ c_{\text{out}}/n & \text{different groups} \end{cases}$$

- Assortative model $c_{\text{in}} > c_{\text{out}}$
Disassortative model $c_{\text{in}} < c_{\text{out}}$

The hidden partition model

Stochastic block model (SBM) with $q = 2$

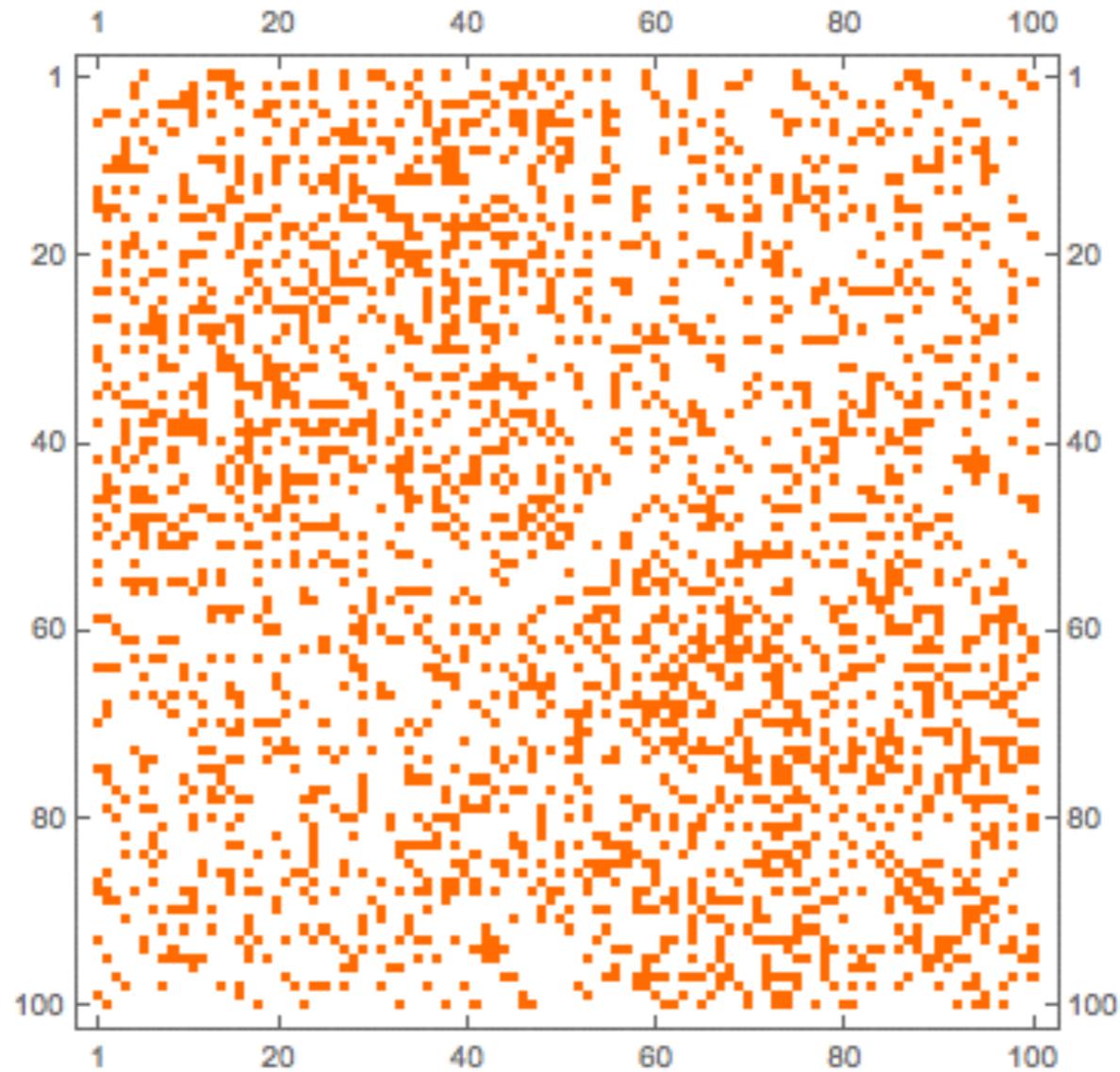


$$\mathbb{P}[(ij) \in E] = \begin{cases} c_{\text{in}}/n & \text{same group} \\ c_{\text{out}}/n & \text{different groups} \end{cases}$$

Assortative model:

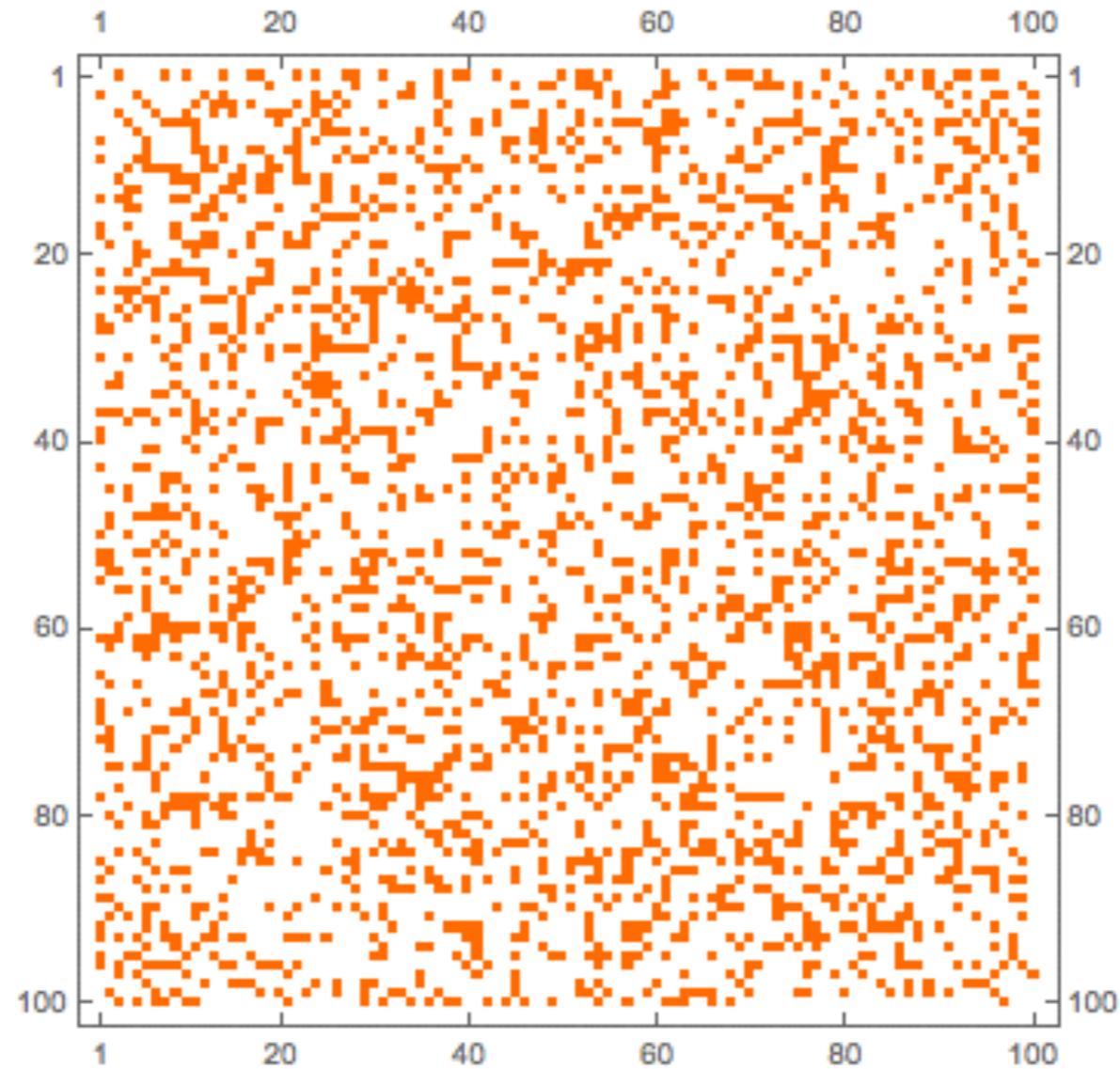
$$c_{\text{in}} > c_{\text{out}}$$

The hidden partition model



Colors are not provided !

The hidden partition model

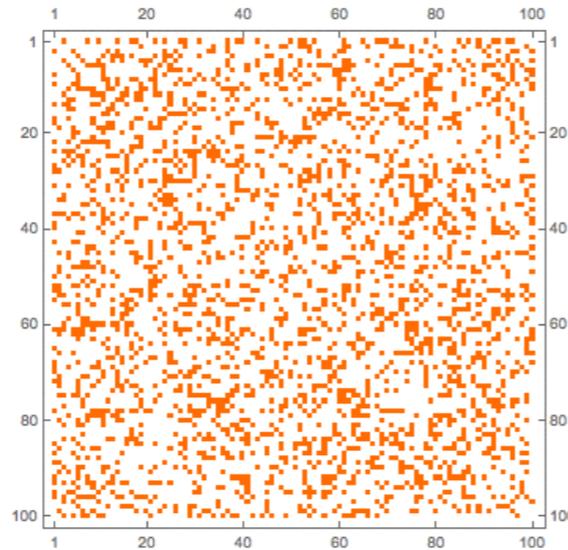


The right ordering neither !!

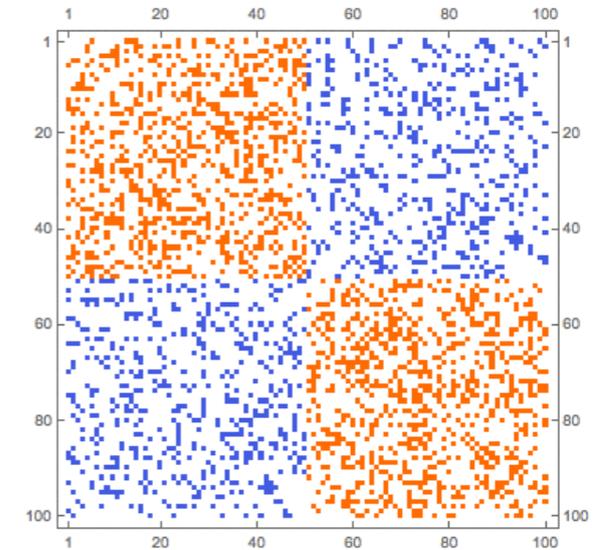
The hidden partition model

Given only the adjacency matrix

$$A_{ij} = A_{ji} = \mathbb{I}[(ij) \in E]$$



Infer the
hidden
partition



Hidden (true) partition $\rightarrow \mathbf{x}_0 \in \{+1, -1\}^n$

Estimated partition $\rightarrow \hat{\mathbf{x}}(G) \in \{+1, -1\}^n$

Quality of inference

via the overlap $\rightarrow Q = \frac{1}{n} |\langle \hat{\mathbf{x}}(G), \mathbf{x}_0 \rangle|$

Assortative SBM with 2 equal-size groups

Relevant parameters and threshold

- Mean degree $d = \frac{c_{\text{in}} + c_{\text{out}}}{2}$
- Signal-to-noise ratio $\lambda = \frac{c_{\text{in}} - c_{\text{out}}}{2\sqrt{d}}$
- Bayes optimal threshold $\lambda_c = 1$
 - Impossible detection for $\lambda < \lambda_c$
 - BP algorithm with $Q > 0$ for $\lambda > \lambda_c$

[Decelle, Krzakala, Moore, Zdeborova, 2011]
[Massoulié, 2013] [Mossel, Neeman, Sly, 2013]

Maximum Likelihood (ML)

- If no information on the generative model is given (apart being assortative and with 2 equal-size groups) a good choice is to maximize the likelihood

$$\text{maximize } \sum_{i,j} A_{i,j} x_i x_j$$

$$\text{subject to } x_i \in \{+1, -1\} \text{ and } \sum_i x_i = 0$$

- NP-hard problem

Spectral relaxation

- Relaxes the constraint $x \in \{+1, -1\}^n$
- Compute largest/smallest eigenvalues of a combination of adjacency (A) and degrees (D) matrices
Project the corresponding eigenvector to $\{+1, -1\}^n$

- Laplacian $L = D - A$

- Normalized Laplacian $D^{-1/2} L D^{-1/2}$

Eigenvector localization on high or low degree nodes

- Bethe Hessian $H(\lambda) = (\lambda^2 - 1)\mathbb{I} + D - \lambda A$
[Saade, Krzakala, Zdeborova, 2014]

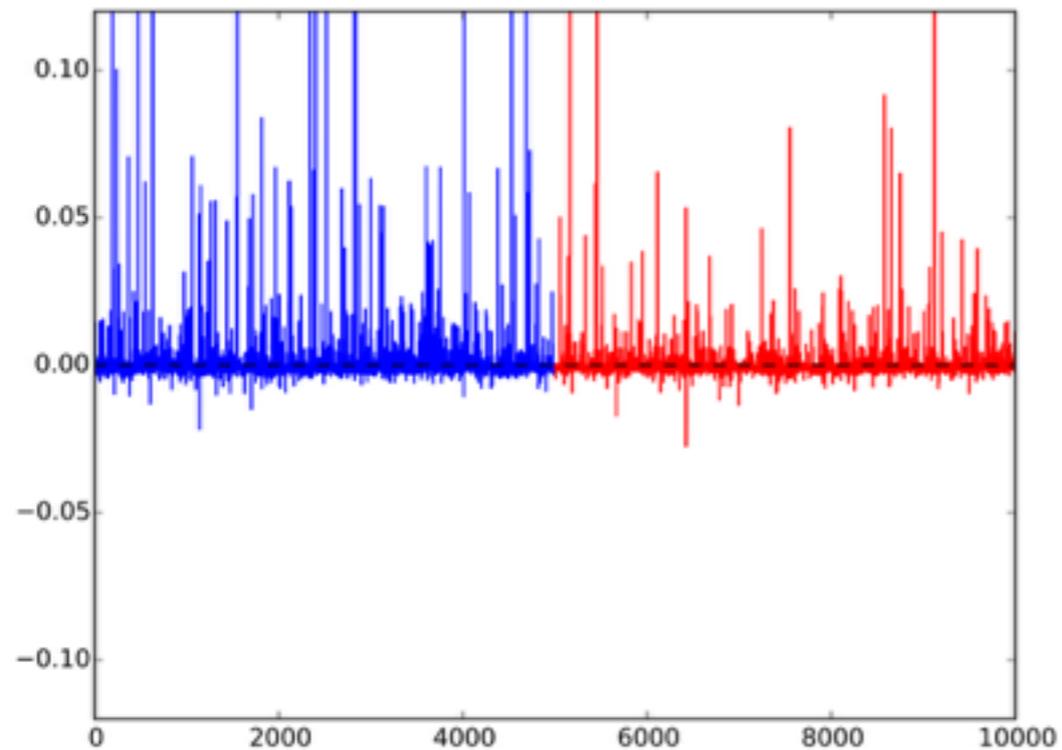
- z-Laplacian $L_z = zA - D$
[Banks, Moore, Newman, Zhang, 2014]

Eigenvector localization on cliques

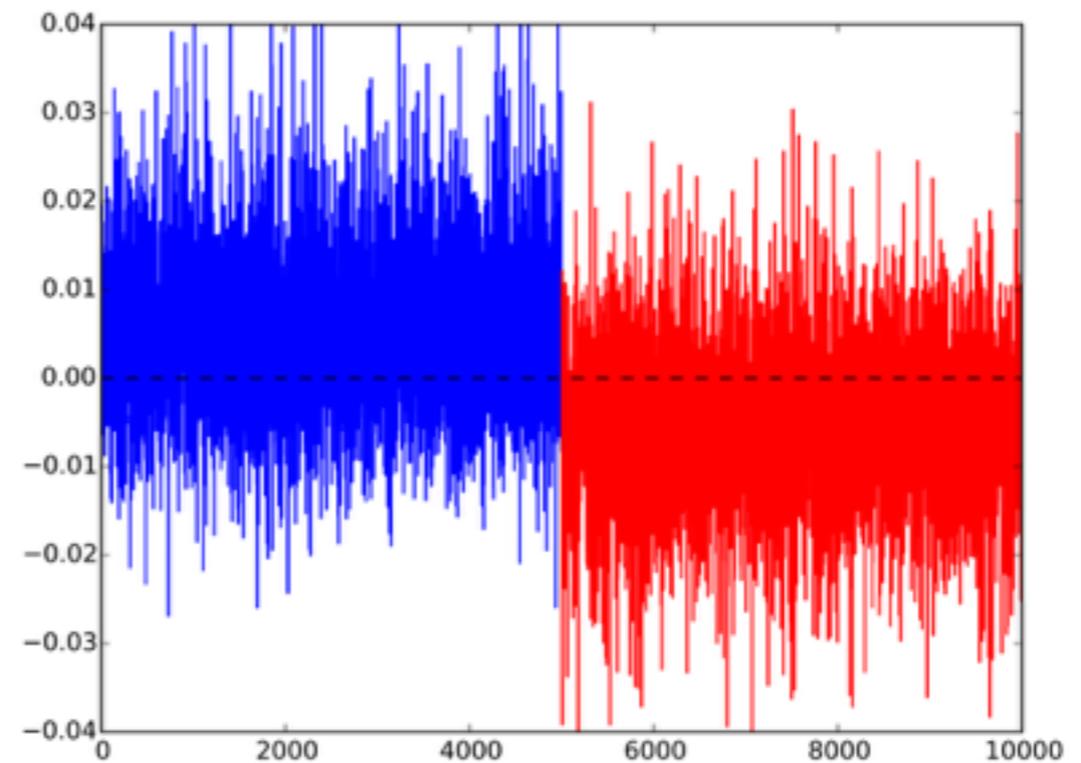
Spectral relaxation fails on sparse graphs

$$n = 10^4 \quad \lambda = 1.2$$

$$d = 3$$



$$d = 20$$



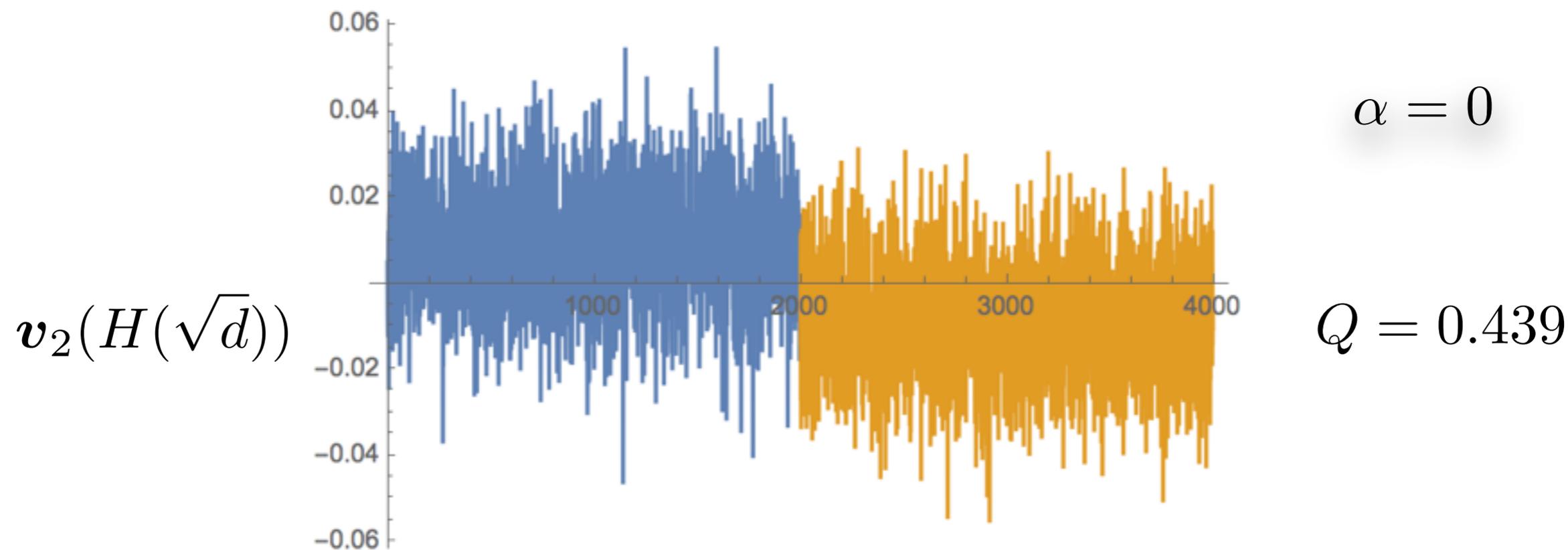
$$v_1(A^{\text{cen}})$$

Quasi-random graphs (SBM + random cliques)

- Generate a graph according to the SBM
- Choose a subset S of vertices of size $|S| = \alpha n$
- For each vertex in S connect all its neighbours
- The number of edges increases by $\sim \alpha d^2 n / 2$
i.e. by a fraction $\sim \alpha d$
- A robust inference method should work also for $\alpha > 0$
at least in the regime $\alpha \ll 1/d$

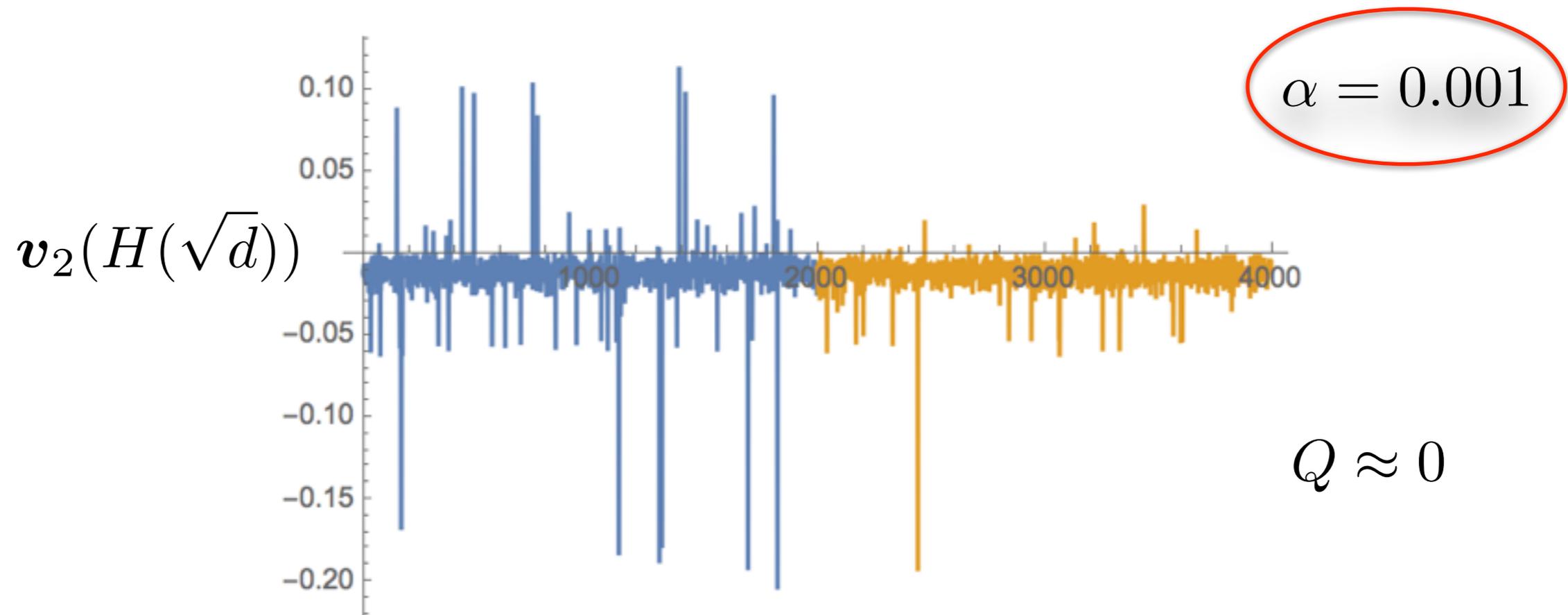
Improved spectral methods fail on quasi-random graphs

$$n = 4000 \quad d = 4 \quad \lambda = 1.1$$



Improved spectral methods fail on quasi-random graphs

$$n = 4000 \quad d = 4 \quad \lambda = 1.1$$

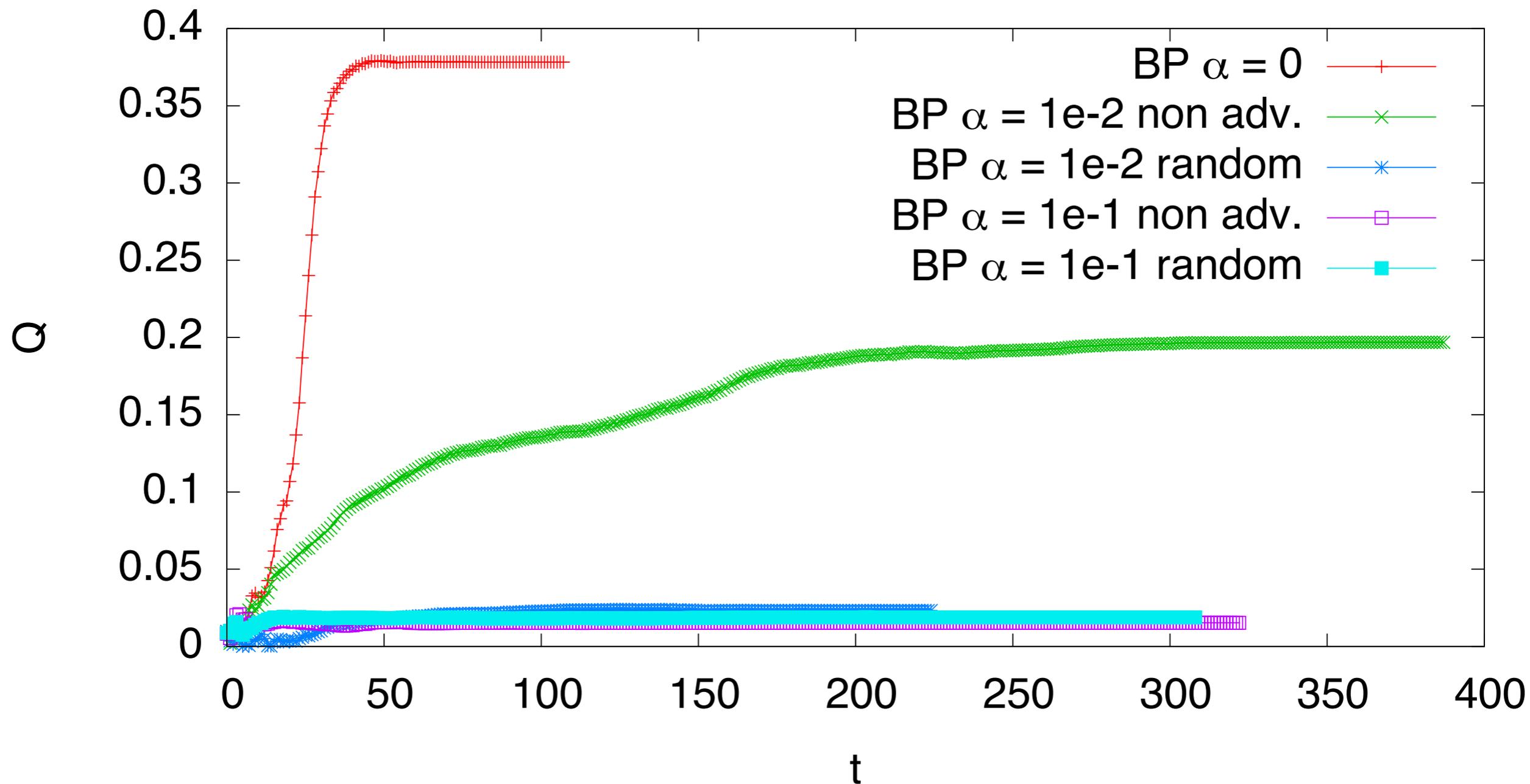


Quasi-random graphs (SBM + non adversarial cliques)

- Generate a graph according to the SBM
- Choose a subset S of vertices of size $|S| = \alpha n$
- For each vertex in S connect all its neighbours **belonging to the same community**
- Non adversarial cliques provide more information

SBM + cliques

$N=10^5$ $d=4$ $\lambda=1.1$



SDP: a better relaxation?

- Maximize $\sum_{i,j} A_{i,j} x_i x_j$ over $x \in \{+1, -1\}^n$

it is equivalent to maximize $\langle A, X \rangle \equiv \sum_{i,j} A_{ij} X_{ij}$

subject to $X \in \mathbb{R}^{n \times n}$, $X \succeq 0$ (i.e. all eigenvalues ≥ 0)

$X_{ii} = 1$ and X being of rank 1

- SDP relaxes the rank and maximizes $\langle A, X \rangle$ over the convex space of positive semidefinite matrices
- The maximizer is a matrix of rank $m \in [1, n]$ to be projected back on a rank 1 matrix...

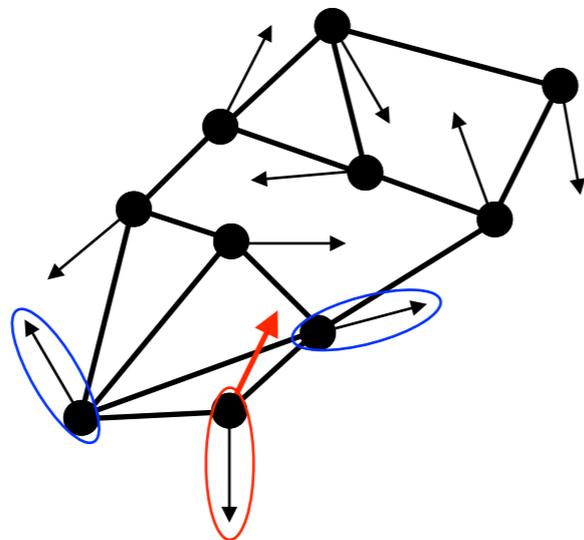
$$X^{\text{opt}} \longrightarrow \hat{x}^{\text{SDP}} (\hat{x}^{\text{SDP}})^{\text{T}}$$

SDP-based algorithm

- Maximize $\langle A, X \rangle$ over rank- m matrices = correlation matrices between m -components variables of unit norm

$$C_{ij} = \underline{x}_i \cdot \underline{x}_j, \quad \text{with } \underline{x}_i \in \mathbb{R}^m, \quad \|\underline{x}_i\|^2 = \underline{x}_i \cdot \underline{x}_i = 1$$

- Maximize $\sum_{(ij) \in E} \underline{x}_i \cdot \underline{x}_j$ subject to $\sum_i \underline{x}_i = \underline{0}$



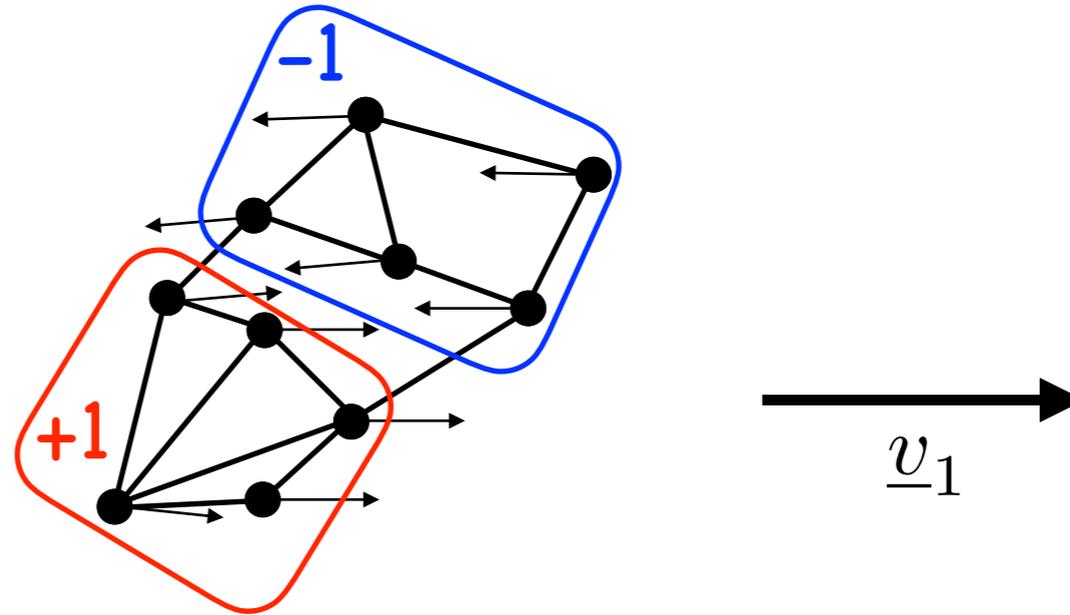
$$\underline{h}_i = \sum_{j \in \partial i} \underline{x}_j - \frac{d}{n} \sum_j \underline{x}_j$$

$$\underline{x}_i \leftarrow \frac{\underline{h}_i}{\|\underline{h}_i\|}$$

Greedy T=0 dynamics (very fast! no gradient used)

SDP-based algorithm

Once you reach a local maximum...



- Given the maximizer $\underline{x}^* = \{\underline{x}_1^*, \dots, \underline{x}_n^*\}$ compute the empirical covariance matrix (**m x m**)

$$\Sigma_{jk} = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i^*)_j (\underline{x}_i^*)_k$$

- Project on its principal eigenvector $\hat{x}_i^{\text{SDP}} = \text{sign}(\underline{x}_i \cdot \underline{v}_1)$

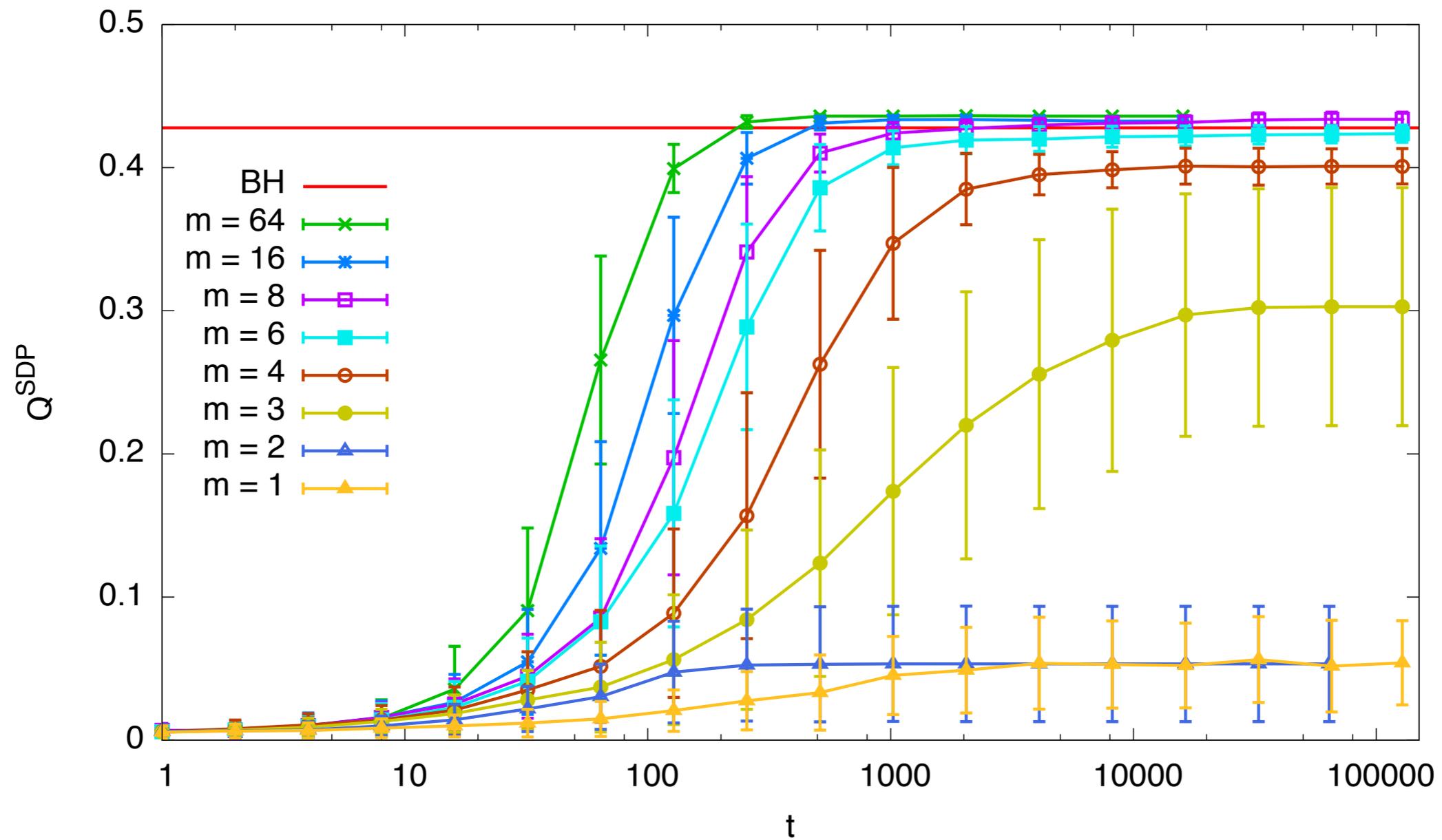
<http://web.stanford.edu/~montanar/SDPgraph/>

SDP-based algorithm

- Algorithm complexity $O(n m t_{\text{conv}})$ and quality of inference do depend on **m**
 - **m=1** \rightarrow ML, **very rough** objective function, NP-hard
 - **m=n** \rightarrow SDP, **convex** objective function
no local maxima for $m > \sqrt{2n}$ [Burer, Monteiro, 2003]
 - **m>1, but small** \rightarrow **smooth enough** objective function ?
local minima are "close enough" $O(m^{-1/2})$
to global minimum [Montanari, 2016]
- Running times grows very mildly with m and n
e.g. if stopping rule is max variation $< 10^{-3}$ $\rightarrow t_{\text{conv}} \propto n^{0.22}$

Small m values are fine!

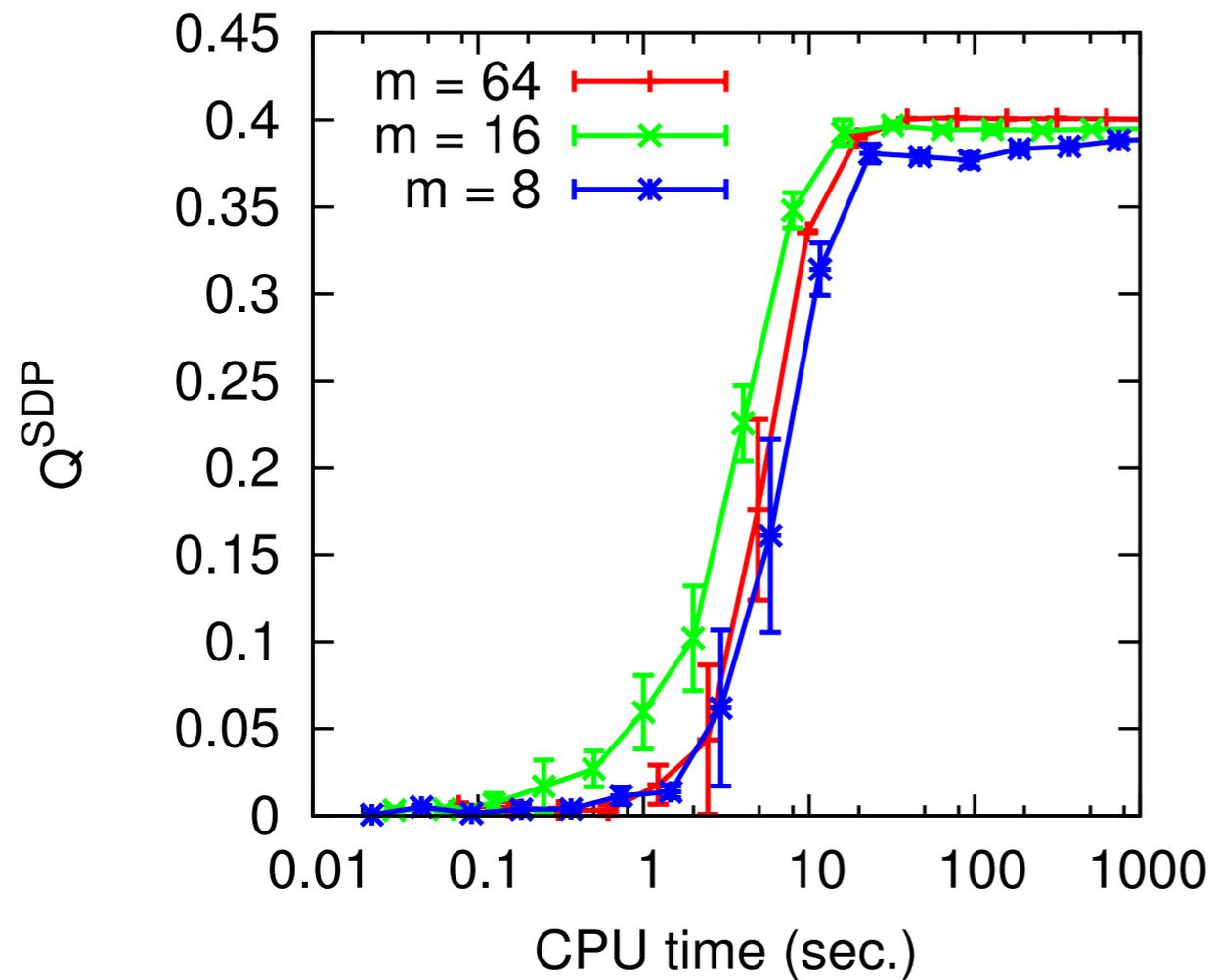
$$n = 4 \cdot 10^4 \quad d = 3 \quad \lambda = 1.1 \quad \alpha = 0.0$$



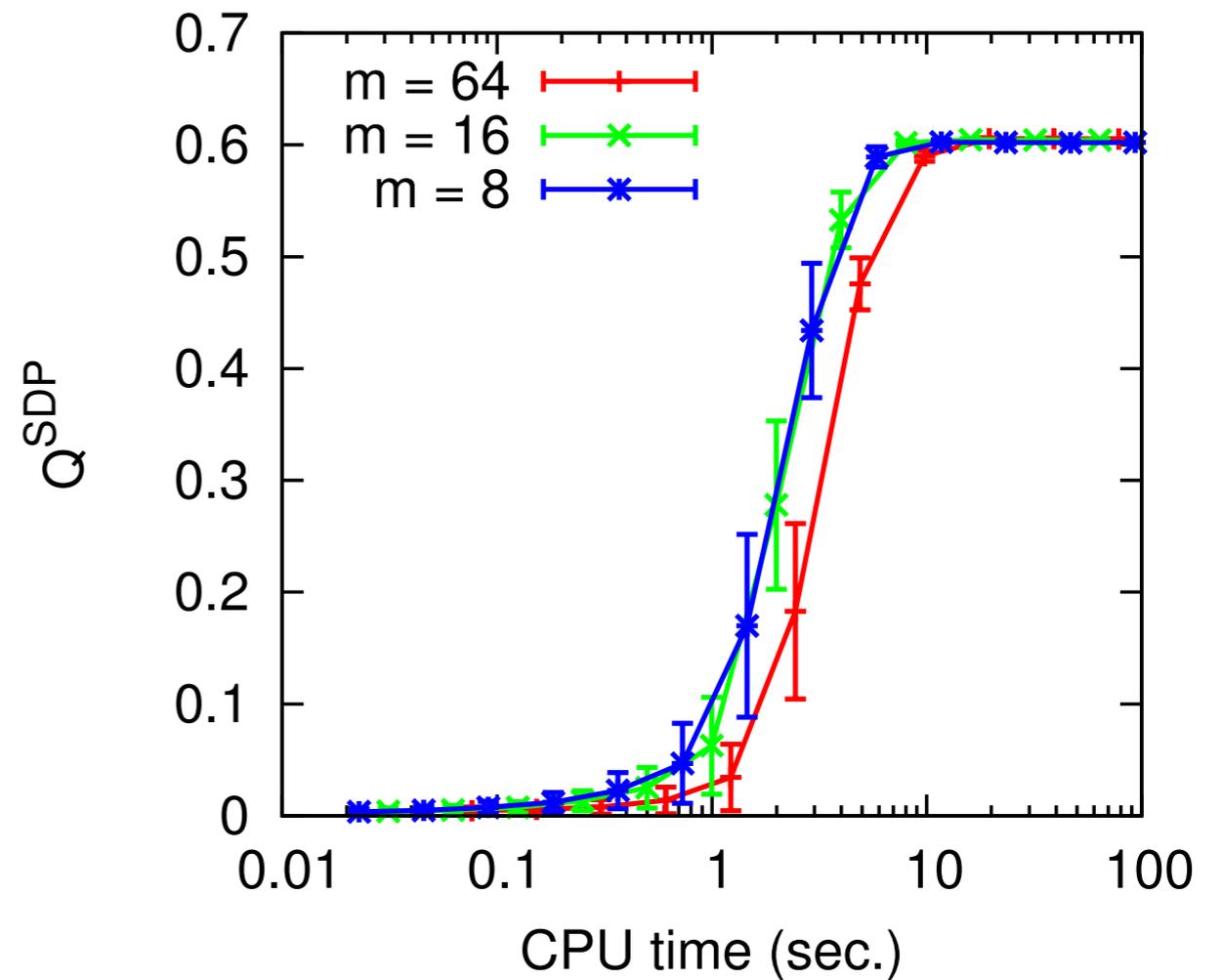
The algorithm is very fast!

$$\underline{n = 10^5} \quad d = 3$$

$$\lambda = 1.1$$

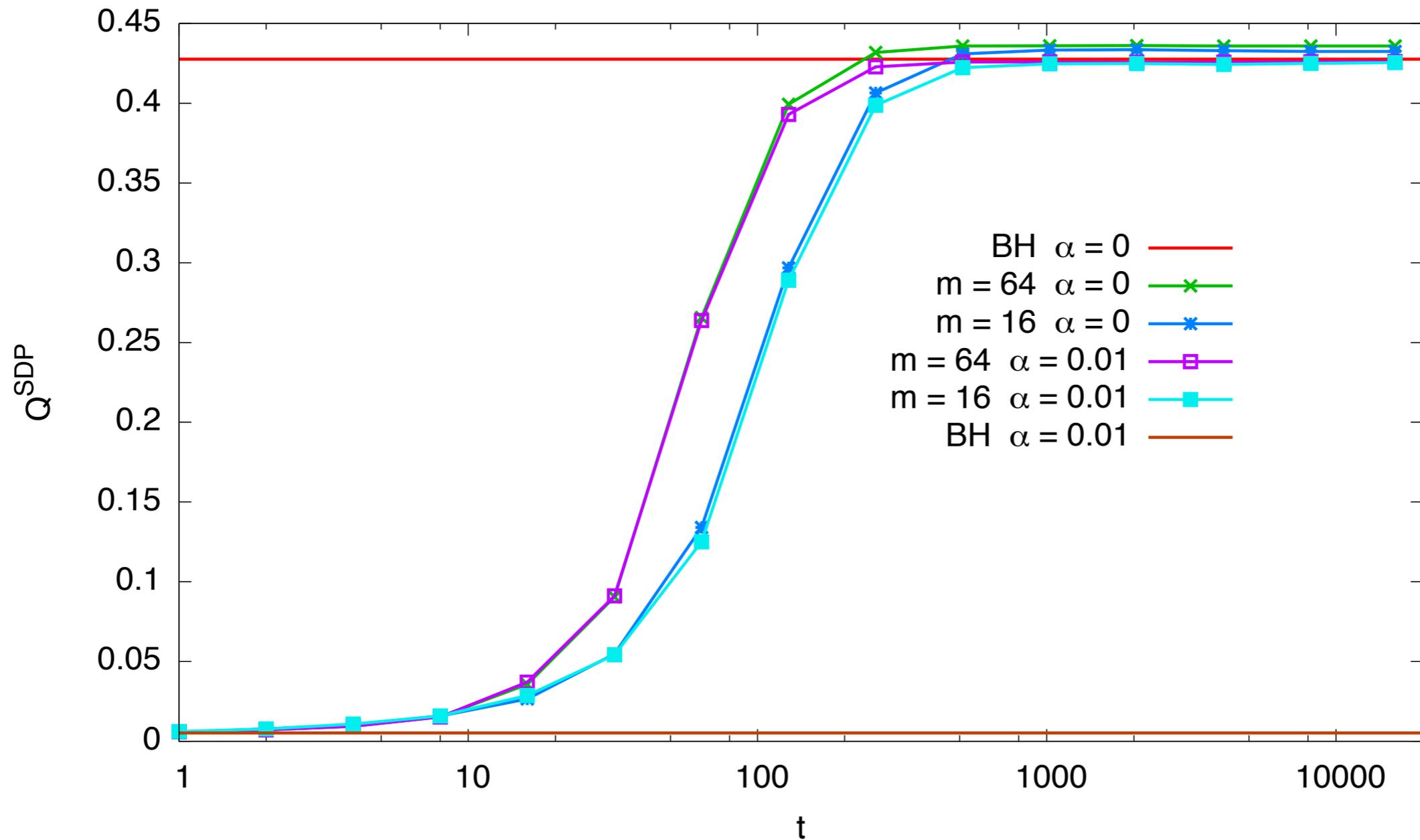


$$\lambda = 1.2$$



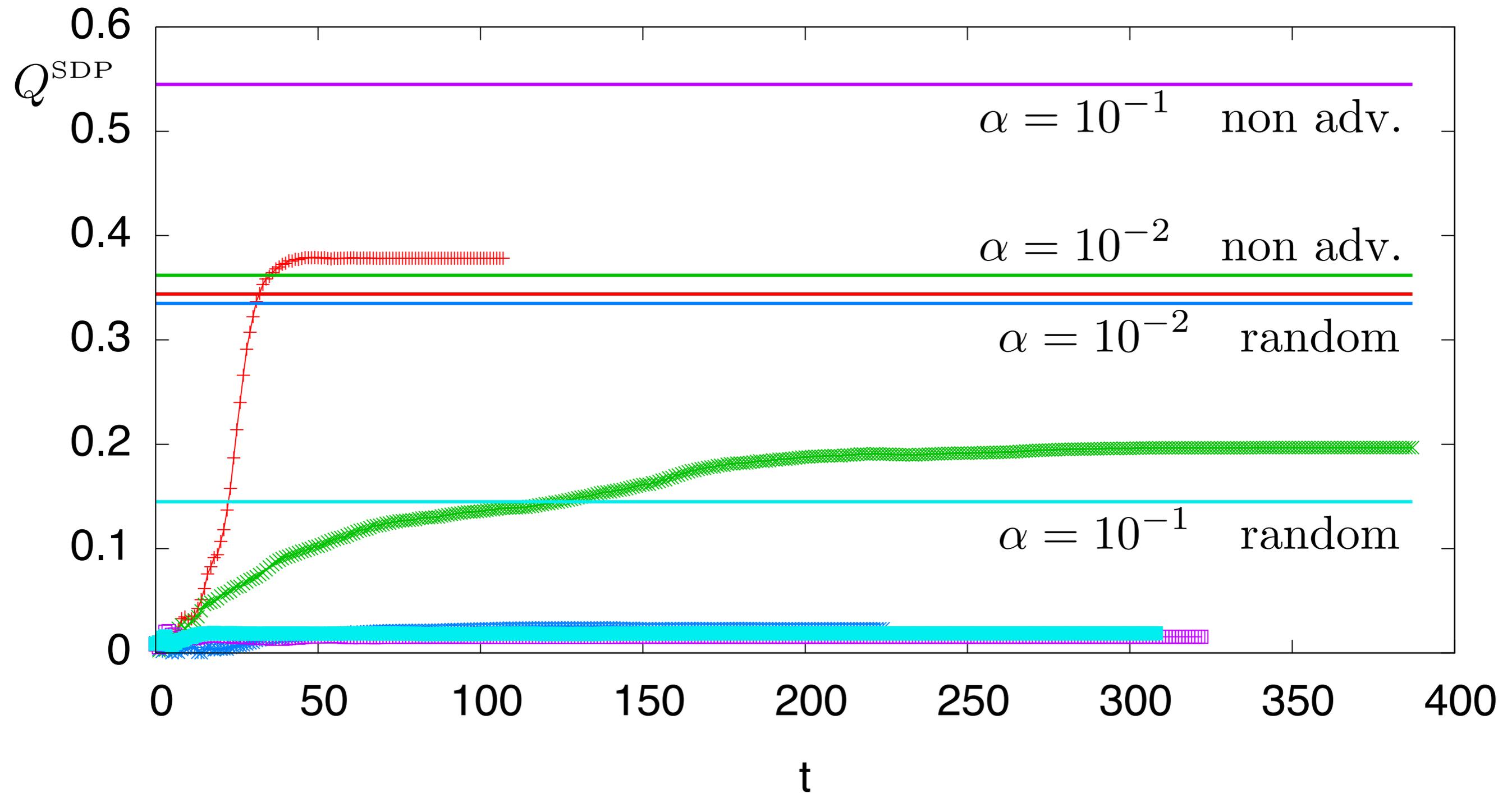
The algorithm is very robust!

$$n = 4 \cdot 10^4 \quad d = 3 \quad \lambda = 1.1$$



The algorithm is very robust!

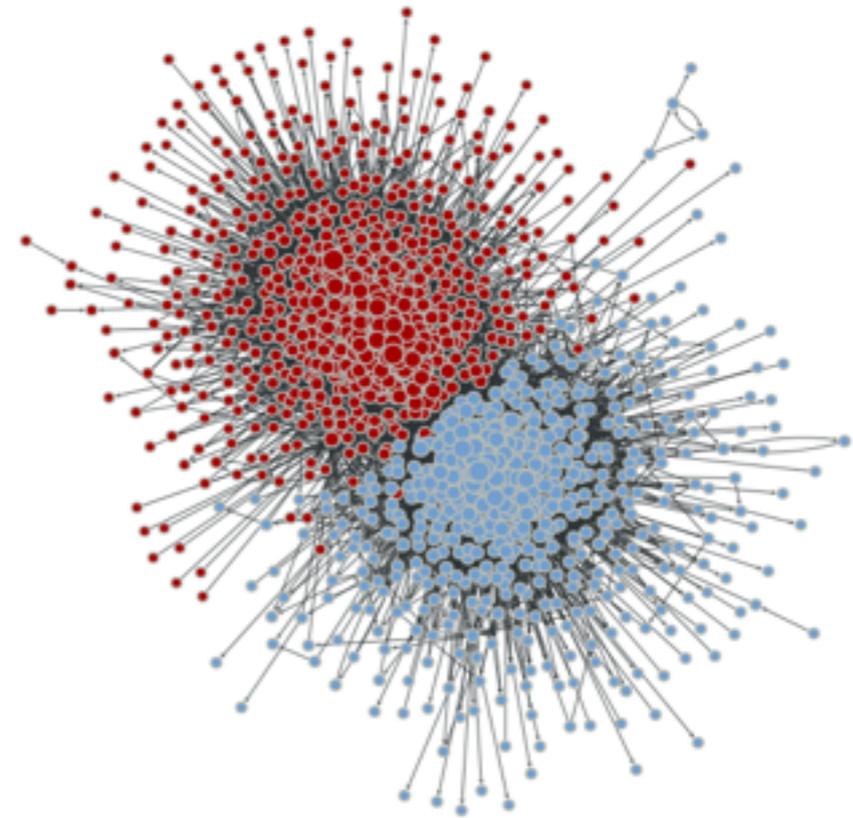
$N=10^5$ $d=4$ $\lambda=1.1$



A real-world network (political blogs)

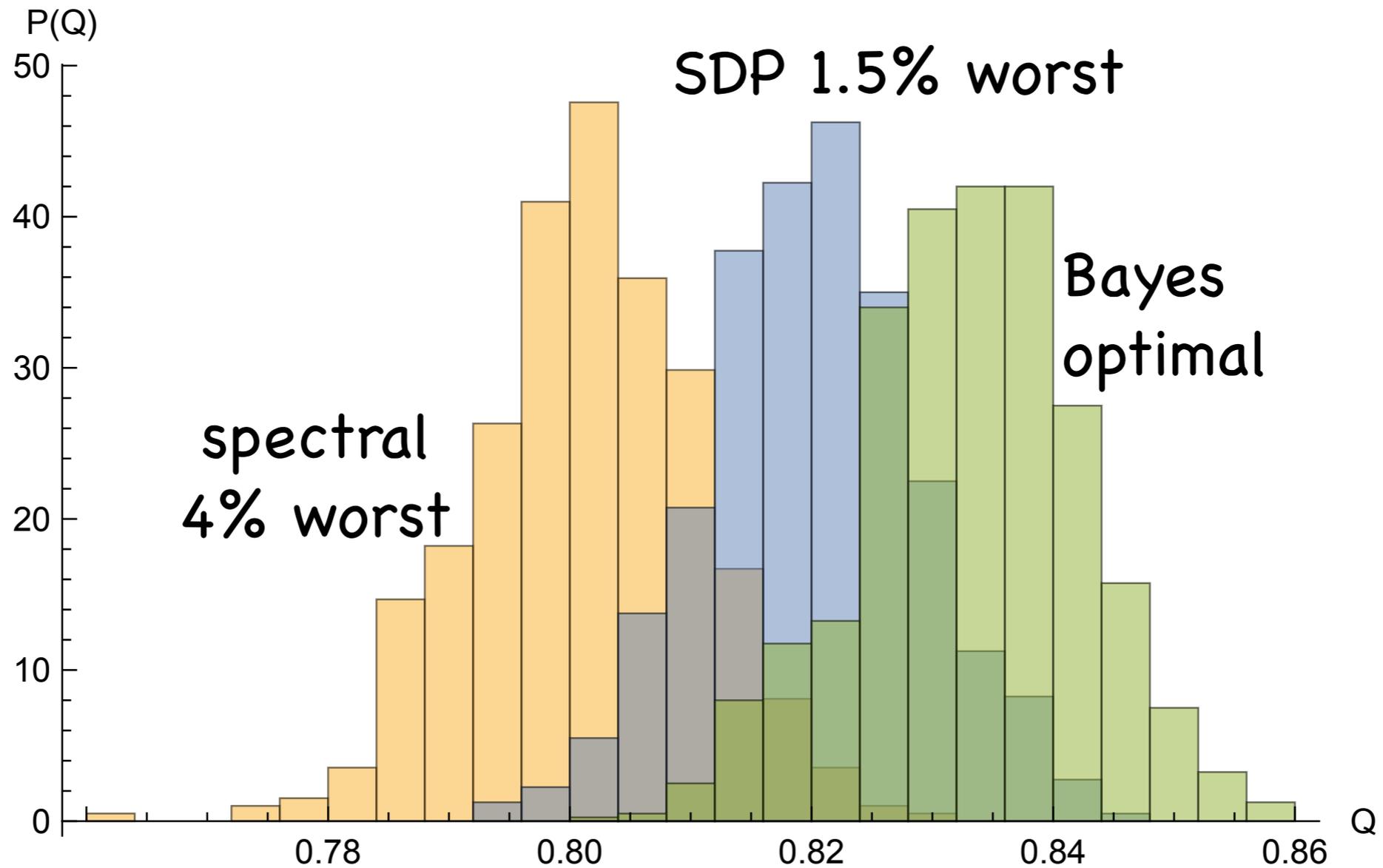
1222 nodes, 16714 edges

	overlap	cut size
Bethe Hessian z-Laplacian	0.865794	1271
Adjacency	0.86743	1268
X-Laplacian	0.918167	1250
Low rank SDP	0.903437	1221
“ground truth”	1.0	1575

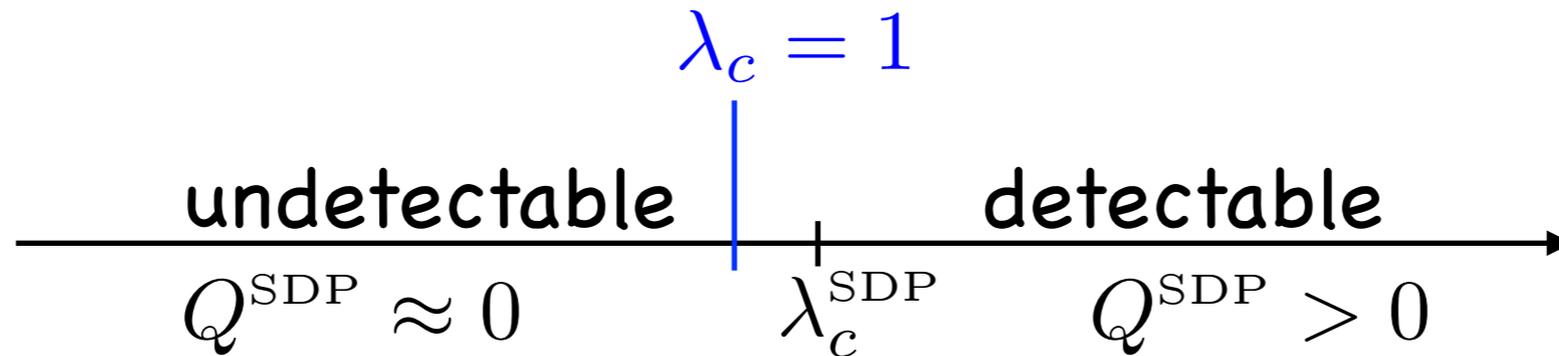


A quantitative comparison

BH SDP BP $d = 4$ $\lambda = 1.5$



SDP quasi-optimality

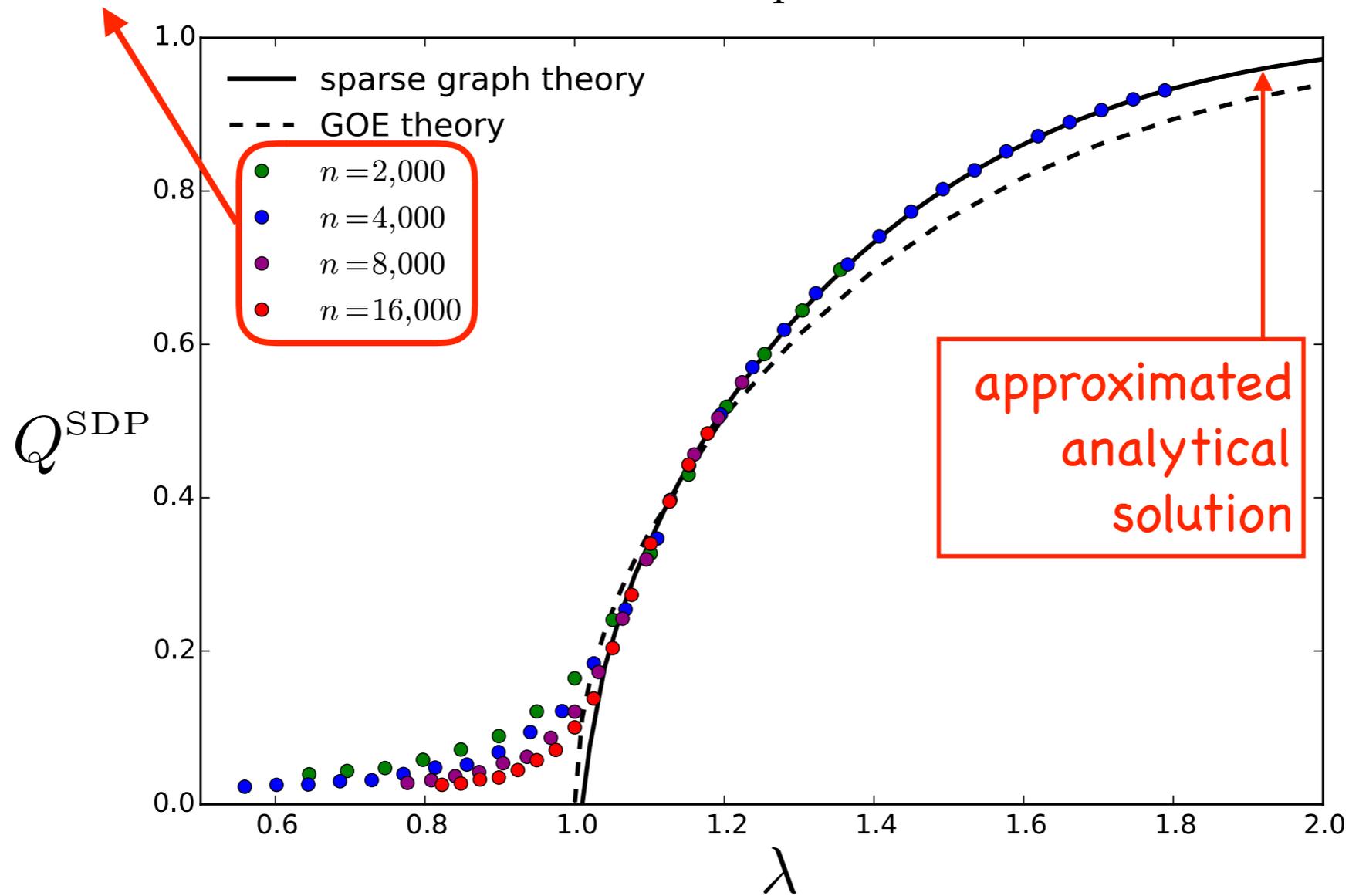


- We estimate λ_c^{SDP} by solving the statistical physics of models with m -component spin variables, in $m \rightarrow \infty$ limit
- Running the SDP-based algorithm for very large m values (= solve the exact cavity equations)
- Solving analytically via an approximate ansatz

SDP quasi-optimality

SDP-based algorithm
(very large m values)

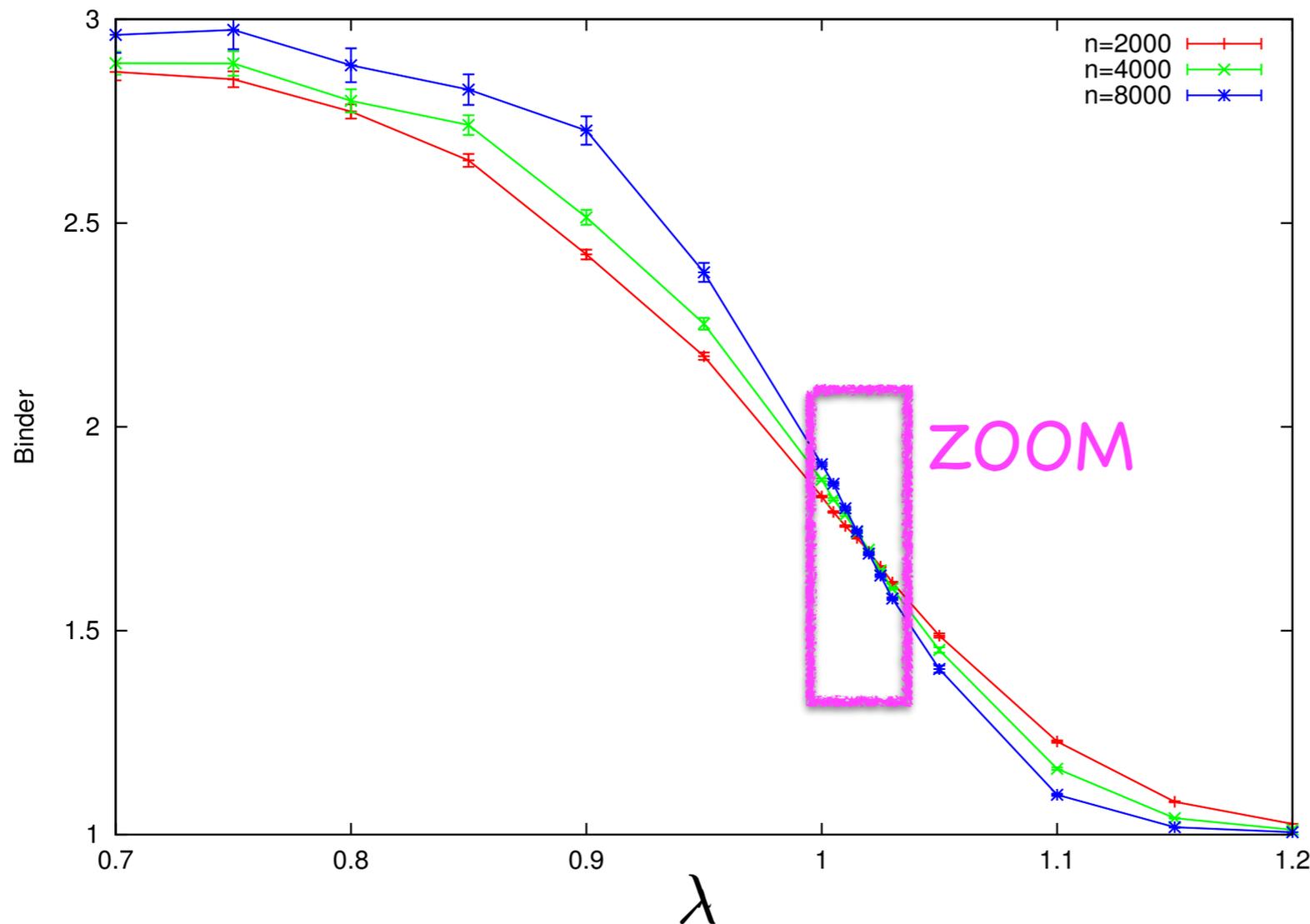
$$d = 5 \quad N_{\text{samples}} = 500$$



$$Q^{\text{SDP}} = \frac{1}{n} \mathbb{E} \{ |\langle \mathbf{x}^{\text{SDP}}, \mathbf{x}_0 \rangle| \}$$

Computing the threshold

- Crossing of the Binder cumulants to locate exactly λ_c^{SDP}



$$d = 5$$

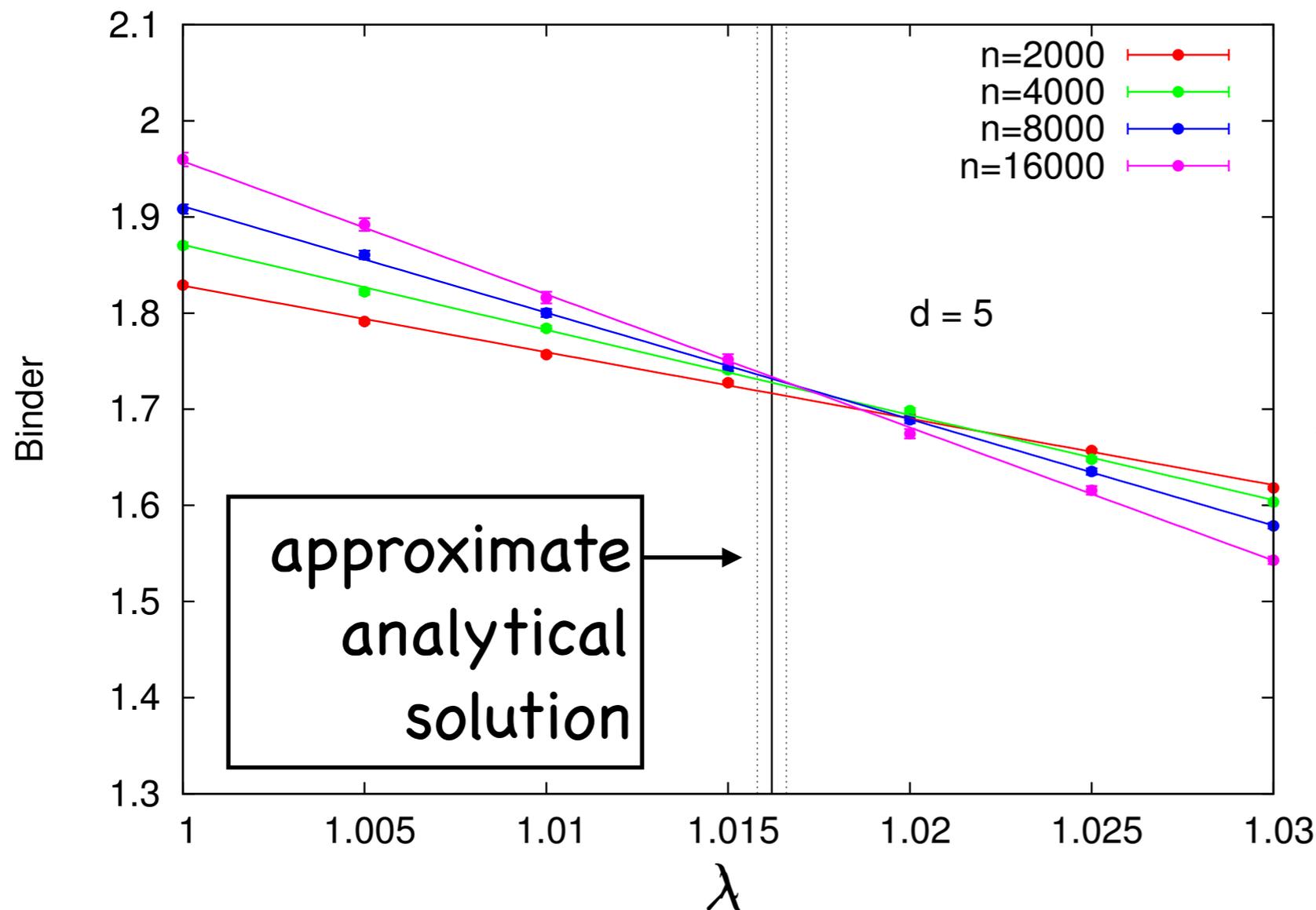
$$N_{\text{samples}} \geq 10^5$$

$$B = \frac{n \mathbb{E}\{\langle \mathbf{x}^{\text{SDP}}, \mathbf{x}_0 \rangle^4\}}{(\mathbb{E}\{\langle \mathbf{x}^{\text{SDP}}, \mathbf{x}_0 \rangle^2\})^2}$$

SDP-based algorithm
for very large n values

Computing the threshold

- Crossing of the Binder cumulants to locate exactly λ_c^{SDP}



$$d = 5$$

$$N_{\text{samples}} \geq 10^5$$

$$B = \frac{n \mathbb{E}\{\langle \mathbf{x}^{\text{SDP}}, \mathbf{x}_0 \rangle^4\}}{(\mathbb{E}\{\langle \mathbf{x}^{\text{SDP}}, \mathbf{x}_0 \rangle^2\})^2}$$

SDP-based algorithm
for very large n values

Statistical physics analytical approach

- Unified framework: statistical physics models with m -component variables: $\underline{x}_i \in \mathbb{R}^m$, $\|\underline{x}_i\| = 1$

$$P(\underline{x}) = \frac{1}{Z} \exp \left[\beta \sum_{(ij) \in E} \underline{x}_i \cdot \underline{x}_j \right]$$

- **Bayes**: $m = 1$, $\tanh(\beta) = \lambda/\sqrt{d}$
- **ML**: $m = 1$, $\beta \rightarrow \infty$
- **SDP**: $m \rightarrow \infty$, $\beta \rightarrow \infty$

Statistical physics analytical approach

- Ansatz for the marginals in m-component dense models

$$P_i(\underline{x}_i) = \frac{1}{Z_i} \exp \left[2m\beta (\underline{\xi}_i^\top \underline{x}_i + \underline{x}_i^\top \mathbf{C}_i \underline{x}_i) \right]$$

$$\underline{x}_i \in \mathbb{F}^m, \|\underline{x}_i\| = 1 \quad \underline{\xi}_i \sim \mathcal{N}(\underline{\mu}, \mathbf{Q}) \quad \mathbf{C}_i = \mathbf{C}$$

- Self consistency equations in the dense case

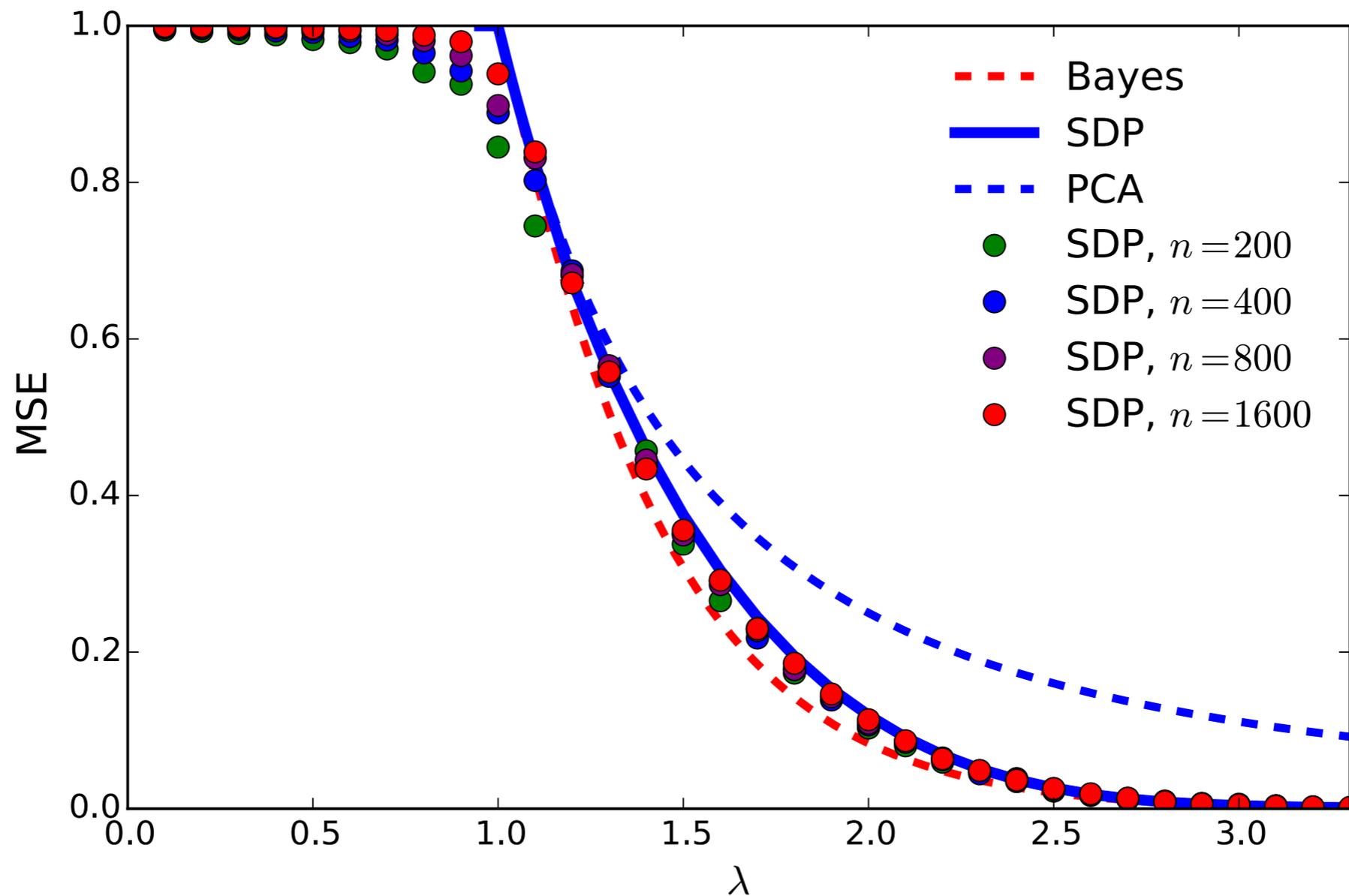
$$\underline{\mu} = \lambda \mathbb{E}[\langle \underline{x} \rangle]$$

$$\mathbf{Q} = \mathbb{E}[\langle \underline{x} \rangle \langle \underline{x}^\top \rangle]$$

$$\mathbf{C} = \beta m \mathbb{E}[\langle \underline{x} \underline{x}^\top \rangle - \langle \underline{x} \rangle \langle \underline{x}^\top \rangle]$$

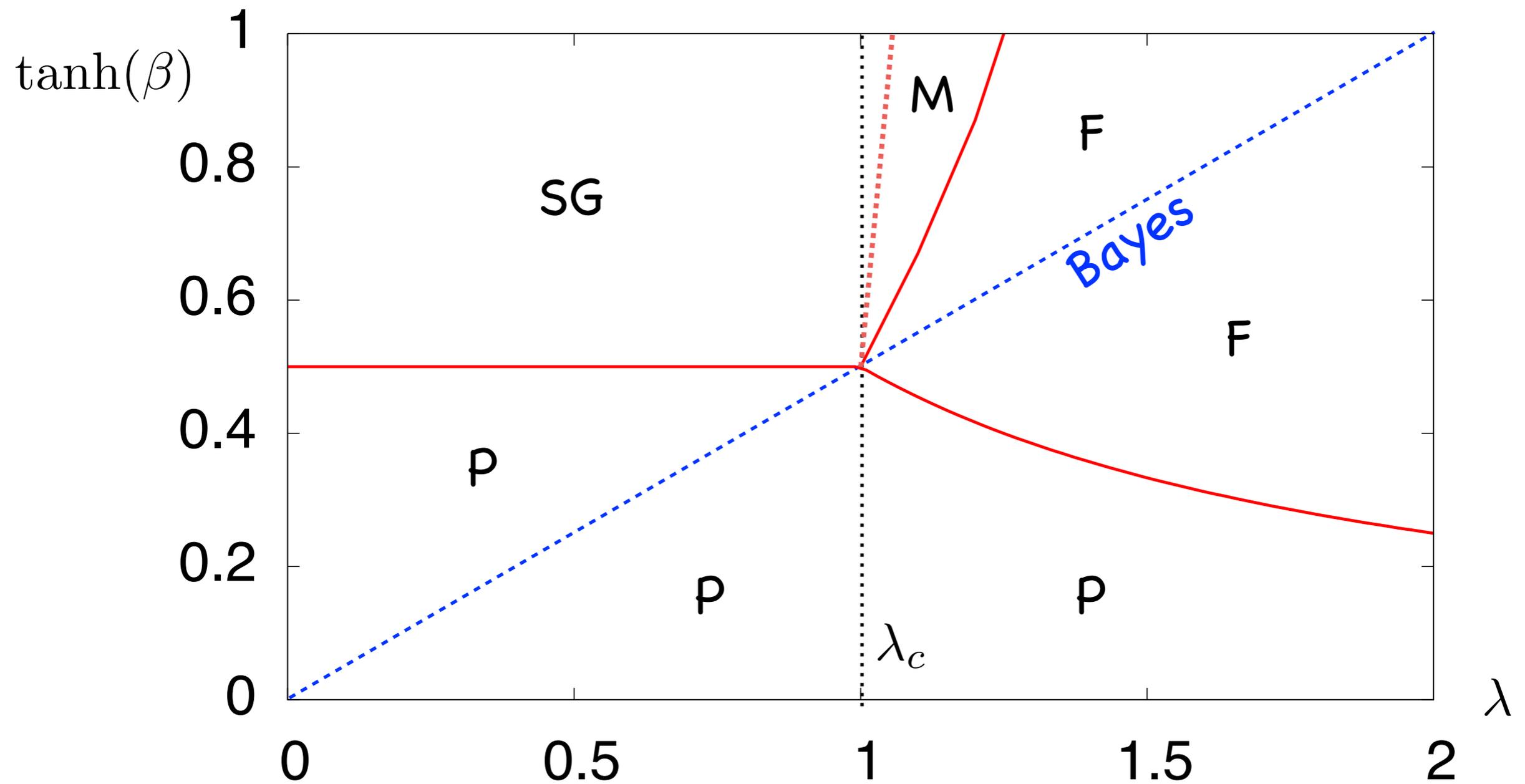
Analytical solution: dense real case

$$\text{MSE}_n(\hat{\boldsymbol{x}}) \equiv \frac{1}{n} \mathbb{E} \left\{ \min_{s \in \{+1, -1\}} \|\hat{\boldsymbol{x}}(\mathbf{Y}) - s \boldsymbol{x}_0\|_2^2 \right\}$$



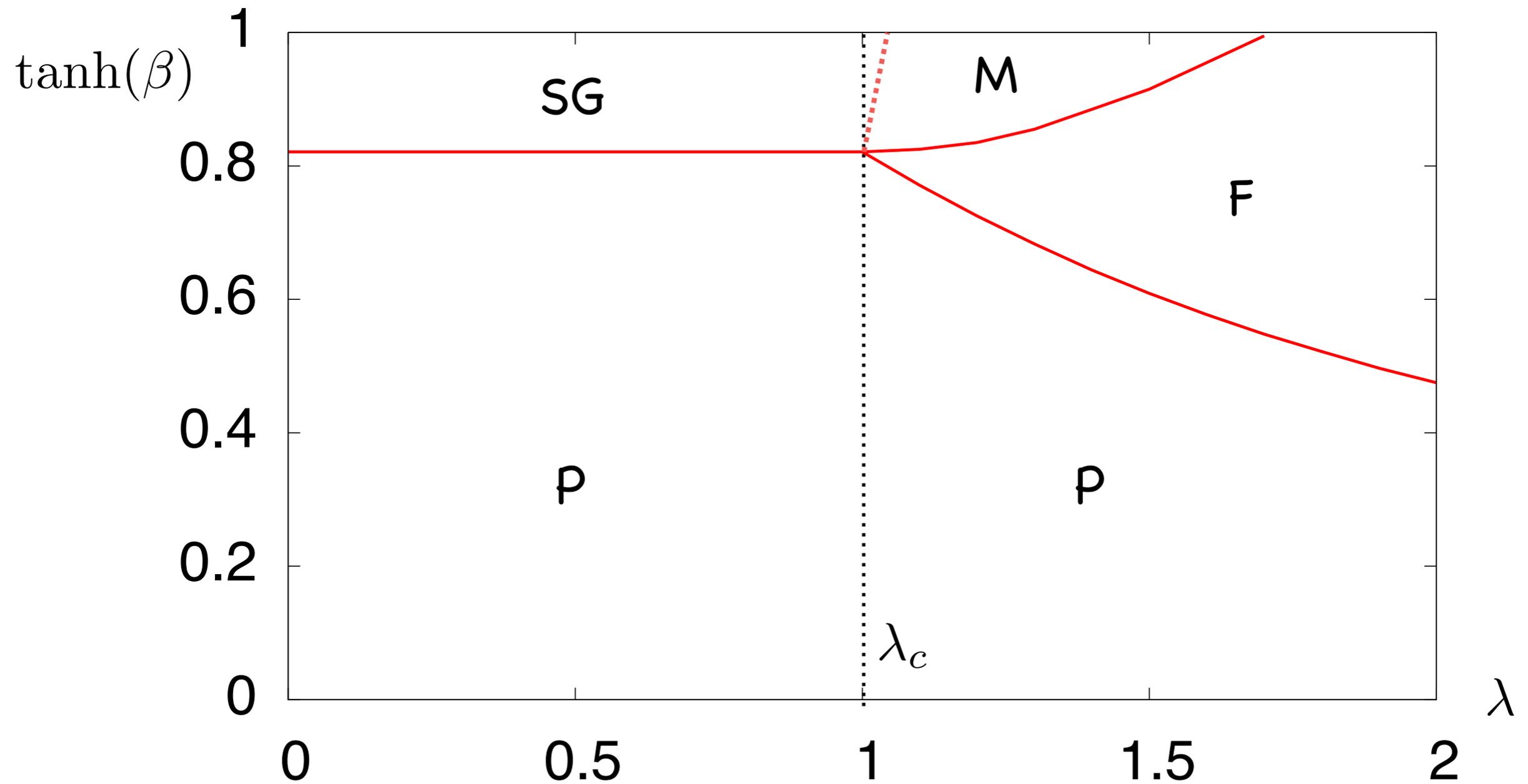
Phase diagrams in the sparse case (SBM $d=4$)

Ising ($m=1$)



Phase diagrams in the sparse case (SBM $d=4$)

XY model ($m=2$)



Approximate analytical solution (SBM)

- In the recovery phase we assume the $O(m)$ symmetry to break along the first component, while preserving $O(m-1)$

$$\underline{x}_i = (s_i, \boldsymbol{\tau}_i), \quad s_i \in \mathbb{R}, \quad \boldsymbol{\tau}_i \in \mathbb{R}^{m-1}$$

- We write the marginal for \underline{x}_i as

$$\exp \left\{ 2\beta\sqrt{mc_i} \langle \mathbf{z}_i, \boldsymbol{\tau}_i \rangle + 2\beta mh_i s_i - \beta mr_i s_i^2 + O_m(1) \right\} \delta \left(s_i^2 + \|\boldsymbol{\tau}_i\|_2^2 - 1 \right)$$

with $\mathbf{z}_i \sim N(0, \mathbf{I}_{m-1})$

- Approximate because the \mathbf{z}_i are correlated
- It should be valid in the limits $d \rightarrow 1$ and $d \rightarrow \infty$

Approximate analytical solution (SBM)

$$\exp \left\{ 2\beta \sqrt{m c_i} \langle \mathbf{z}_i, \boldsymbol{\tau}_i \rangle + 2\beta m h_i s_i - \beta m r_i s_i^2 + O_m(1) \right\} \delta \left(s_i^2 + \|\boldsymbol{\tau}_i\|_2^2 - 1 \right)$$

- Cavity method \rightarrow self consistency equation for marginals

$$c_0 = \sum_{i=1}^k \frac{c_i}{\rho_i^2},$$

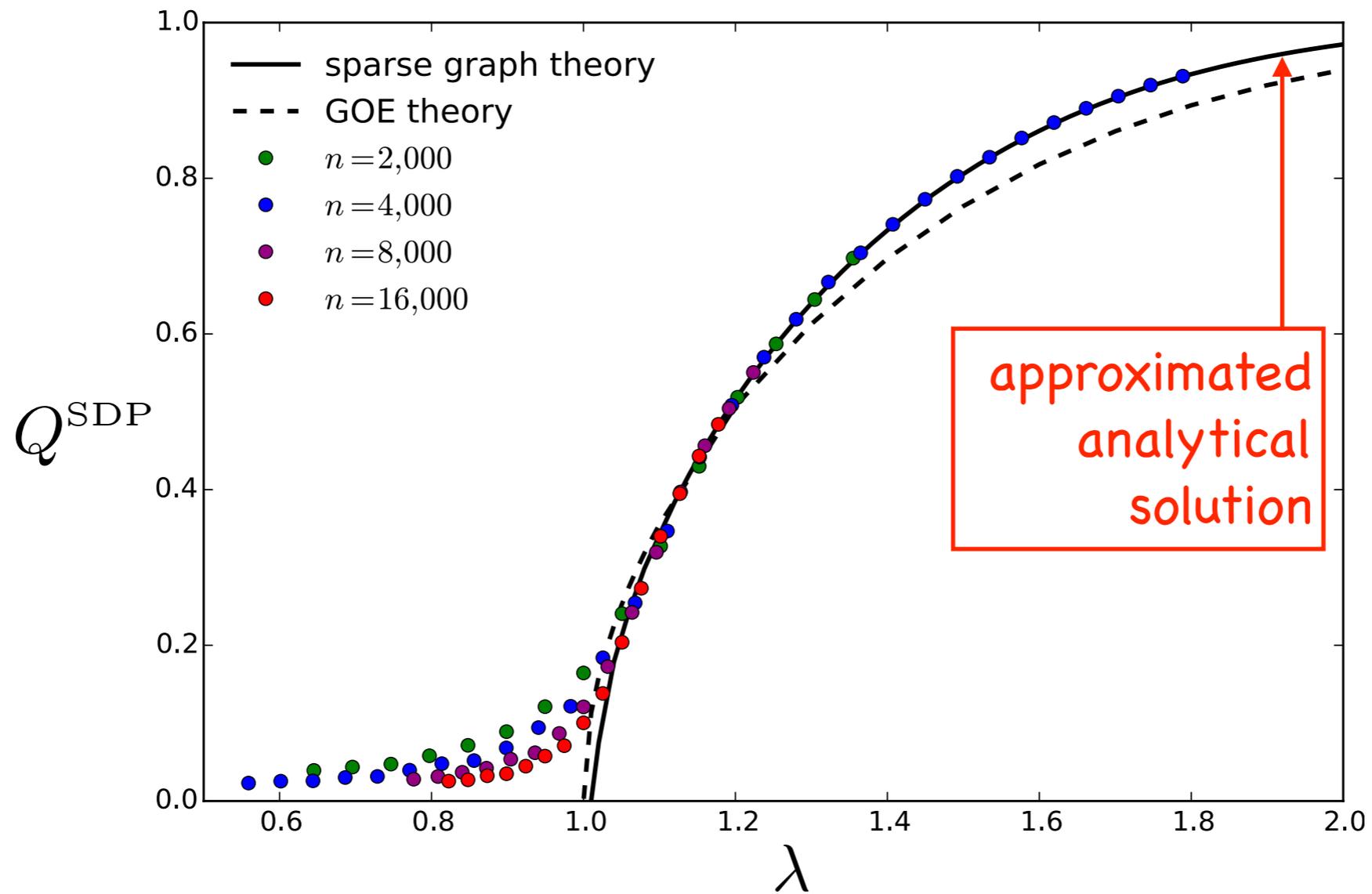
$$h_0 = \sum_{i=1}^k \frac{h_i}{\rho_i + r_i},$$

$$r_0 = \sum_{i=1}^k \left\{ \frac{1}{\rho_i} - \frac{1}{\rho_i + r_i} + \left(1 + \frac{(1 + c_i)r_i}{\rho_i^3} \right)^{-1} \frac{h_i^2}{(\rho_i + r_i)^3} \right\}$$

$$1 = \frac{h_i^2}{(\rho_i + r_i)^2} + \frac{1 + c_i}{\rho_i^2}$$

- Solve by population dynamics
- At the fixed point $Q^{\text{SDP}} = \mathbb{E}[\text{sign}(h^*)]$

Approximate analytical solution (SBM)



Approximate analytical solution (SBM)

- Linearize the cavity equations to locate the threshold
- To linear order in $h \implies r_i = 0$

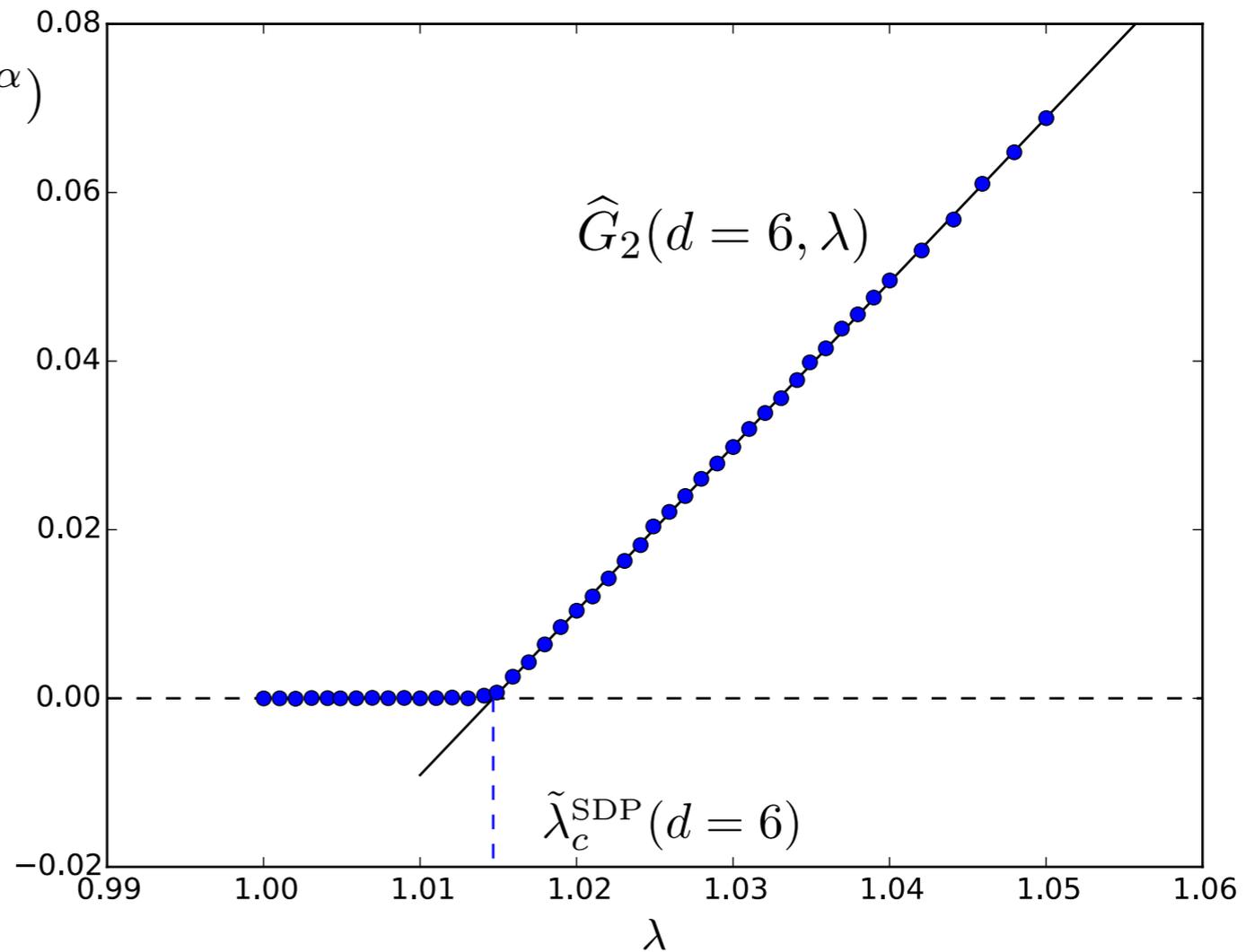
$$c_0 = \sum_{i=1}^k \frac{c_i}{\rho_i^2},$$
$$h_0 = \sum_{i=1}^k \frac{h_i}{\rho_i + r_i},$$
$$r_0 = \sum_{i=1}^k \left\{ \frac{1}{\rho_i} - \frac{1}{\rho_i + r_i} + \left(1 + \frac{(1 + c_i)r_i}{\rho_i^3} \right)^{-1} \frac{h_i^2}{(\rho_i + r_i)^3} \right\}$$
$$1 = \frac{h_i^2}{(\rho_i + r_i)^2} + \frac{1 + c_i}{\rho_i^2}$$



$$c_0 = \sum_{i=1}^k \frac{c_i}{1 + c_i},$$
$$h_0 = \sum_{i=1}^k \frac{h_i}{\sqrt{1 + c_i}}$$

Approximate analytical solution (SBM)

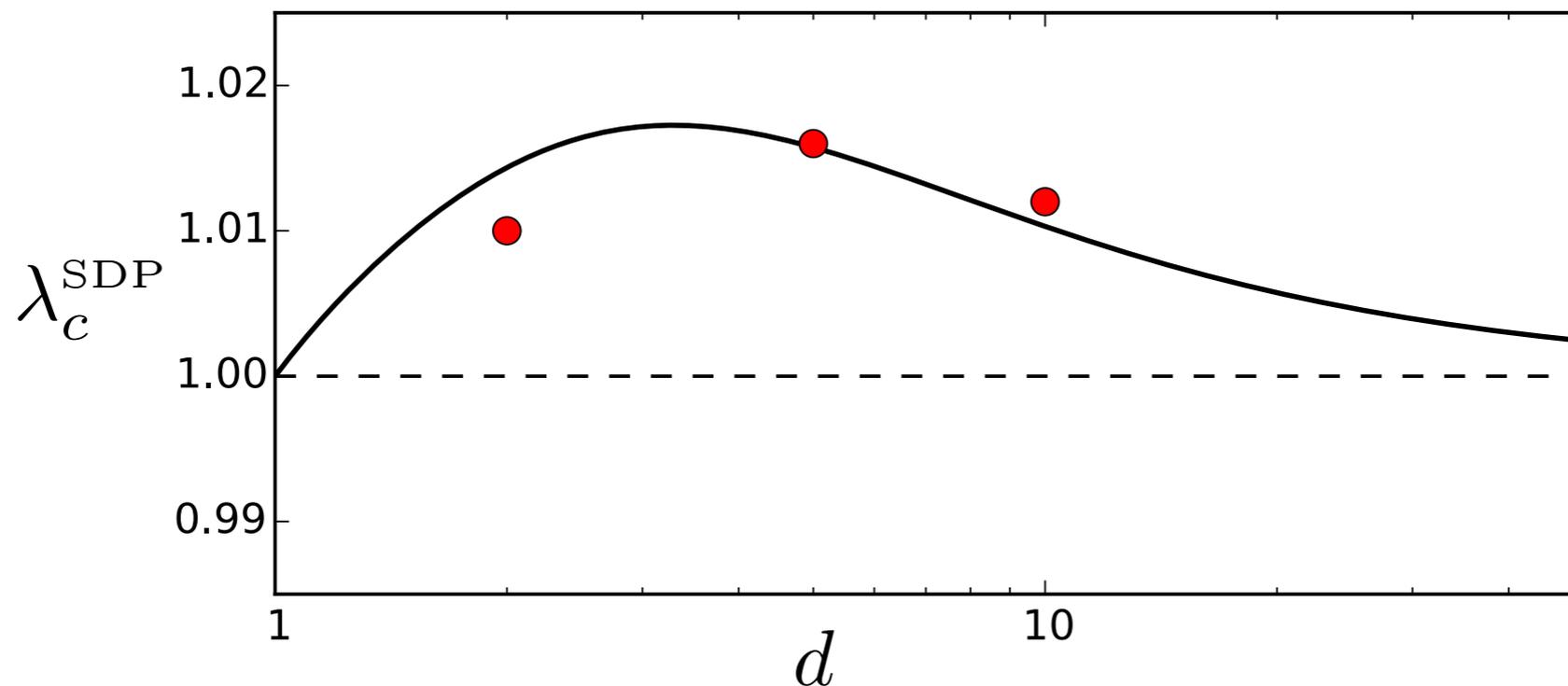
$$G_\alpha(d, \lambda) \equiv \liminf_{t \rightarrow \infty} \frac{1}{t\alpha} \log \mathbb{E}(|h^t|^\alpha)$$



$$(c^{t+1}; h^{t+1}) \stackrel{\text{d}}{=} \left(\sum_{i=1}^{L_+ + L_-} \frac{c_i^t}{1 + c_i^t}; \sum_{i=1}^{L_+ + L_-} \frac{s_i h_i^t}{\sqrt{1 + c_i^t}} \right)$$

Analytical solution: sparse case (SBM)

- SDP at most 2% sub-optimal!



- **Red points:** numerical solution of the replica/cavity equations (crossing of Binder cumulants)
- **Black line:** approximated analytical solution

Some conclusions...

- SDP relaxations are very effective:
 - robust and quasi-optimal
 - may outperform spectral relaxations
- Better than SDP are SDP-inspired algorithms (small m)
<http://web.stanford.edu/~montanar/SDPgraph/>
- It is worth studying the statistical physics of models with m -component variables:
 - unifying framework to study and solve several estimators in statistical inference
 - different physics, better algorithms