

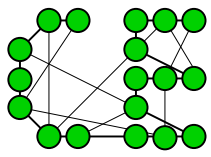
The Cross-Entropy Method

*A Unified Approach to Rare Event Simulation and
Stochastic Optimization*

*Dirk P. Kroese

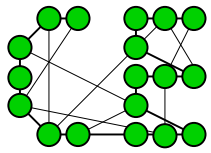
Reuven Y. Rubinstein

*Department of Mathematics, The University of Queensland, Australia
Faculty of Industrial Engineering and Management, Technion, Israel



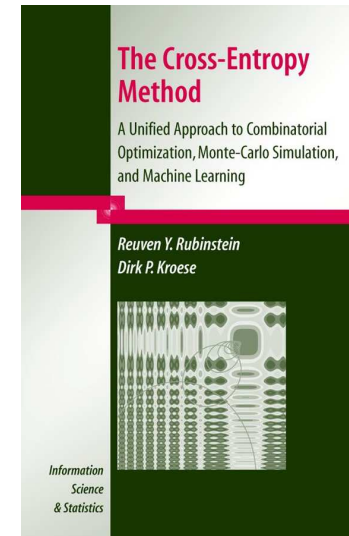
Contents

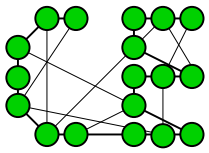
1. Introduction
2. CE Methodology
3. Application: Max-Cut Problem, etc.
4. Some Theory on CE
5. Conclusion



CE Matters

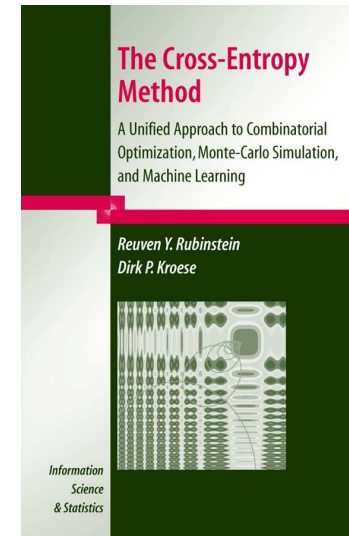
Book: R.Y. Rubinstein and D.P. Kroese.
The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization, Monte Carlo Simulation and Machine Learning, Springer-Verlag, New York, 2004.



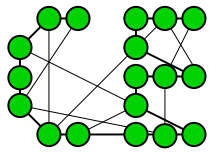


CE Matters

Book: R.Y. Rubinstein and D.P. Kroese.
The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization, Monte Carlo Simulation and Machine Learning, Springer-Verlag, New York, 2004.

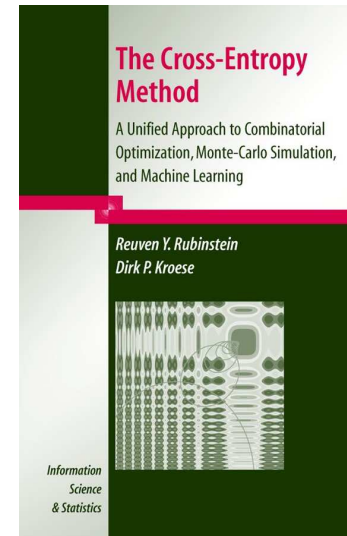


Special Issue: *Annals of Operations Research* (Jan 2005).



CE Matters

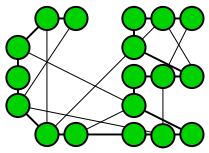
Book: R.Y. Rubinstein and D.P. Kroese.
The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization, Monte Carlo Simulation and Machine Learning, Springer-Verlag, New York, 2004.



Special Issue: *Annals of Operations Research* (Jan 2005).

The CE home page:

<http://www.cemethod.org>



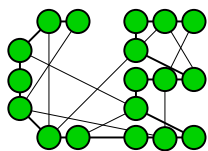
Introduction

The Cross-Entropy Method was originally developed as a simulation method for the estimation of *rare event* probabilities:

$$\text{Estimate } \mathbb{P}(S(\mathbf{X}) \geq \gamma)$$

\mathbf{X} : random vector/process taking values in some set \mathcal{X} .

S : real-values function on \mathcal{X} .



Introduction

The Cross-Entropy Method was originally developed as a simulation method for the estimation of *rare event* probabilities:

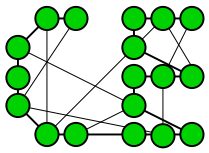
$$\text{Estimate } \mathbb{P}(S(\mathbf{X}) \geq \gamma)$$

\mathbf{X} : random vector/process taking values in some set \mathcal{X} .

S : real-values function on \mathcal{X} .

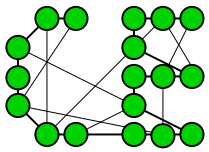
It was soon realised that the CE Method could also be used as an *optimization* method:

$$\text{Determine } \max_{\mathbf{x} \in \mathcal{X}} S(\mathbf{x})$$

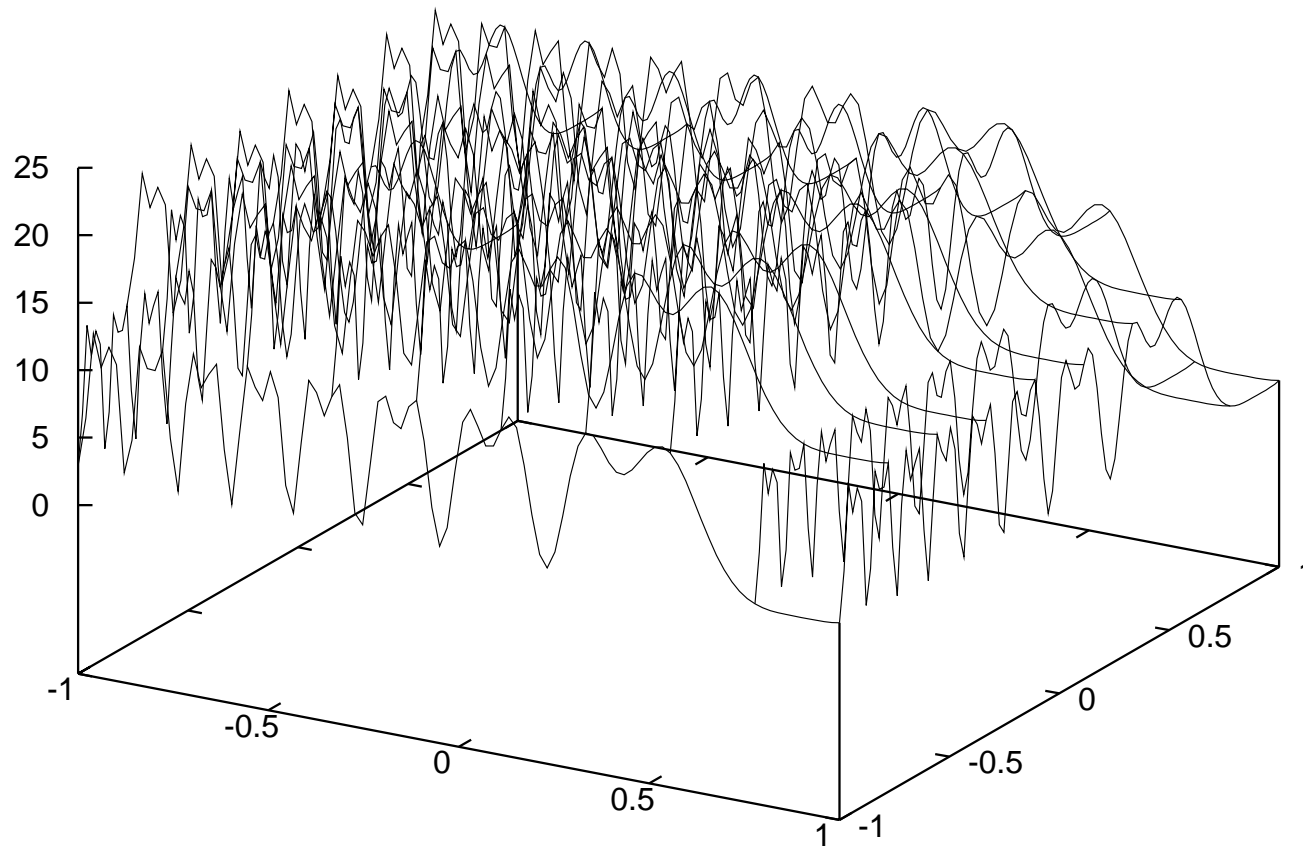


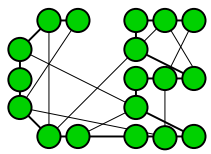
Some Applications

- **Combinatorial Optimization** (e.g., Travelling Salesman, Maximal Cut and Quadratic Assignment Problems)
- **Noisy Optimization** (e.g., Buffer Allocation, Financial Engineering)
- **Multi-Extremal Continuous Optimization**
- Pattern Recognition, Clustering and Image Analysis
- Production Lines and Project Management
- Network Reliability Estimation
- Vehicle Routing and Scheduling
- DNA Sequence Alignment



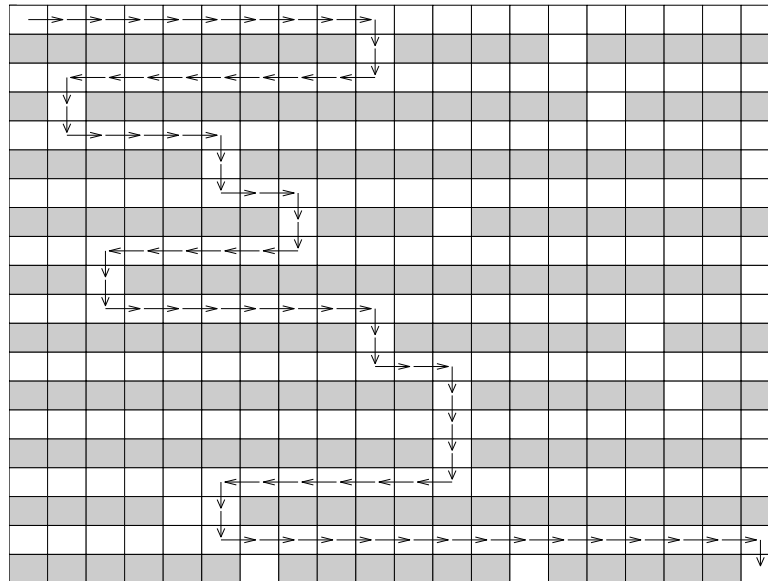
A Multi-extremal function

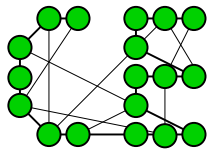




A Maze Problem

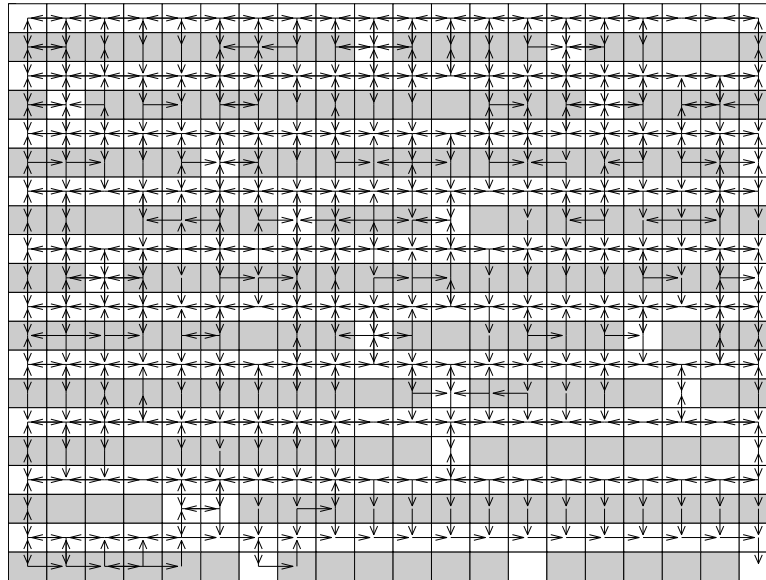
The Optimal Trajectory

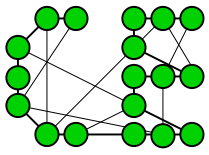




A Maze Problem

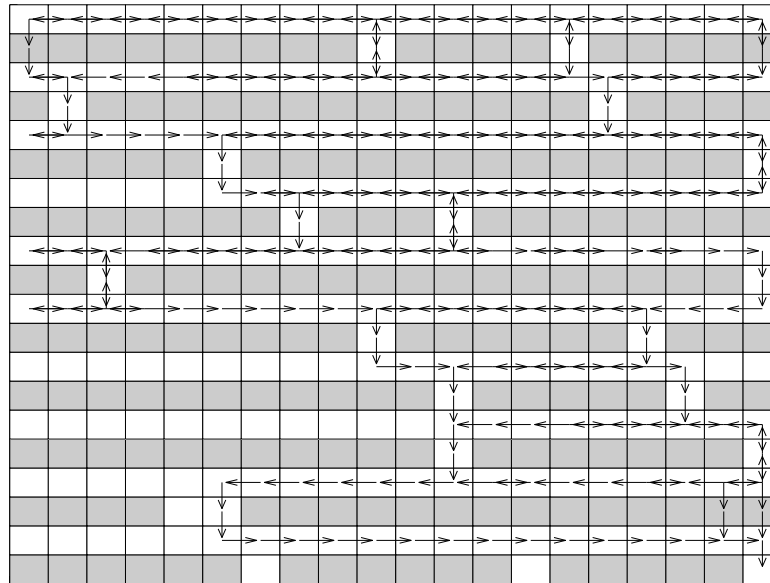
Iteration 1:

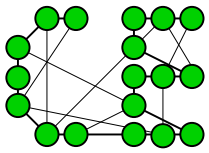




A Maze Problem

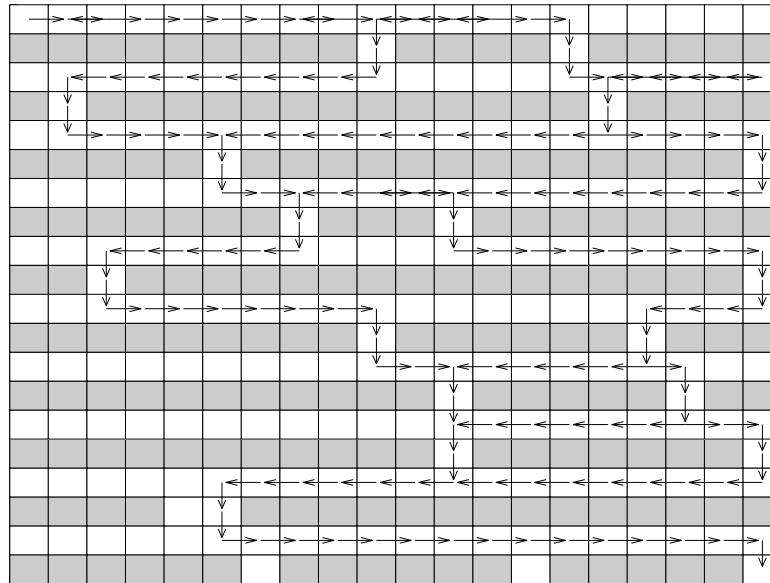
Iteration 2:

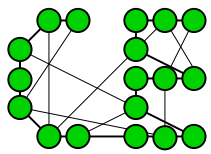




A Maze Problem

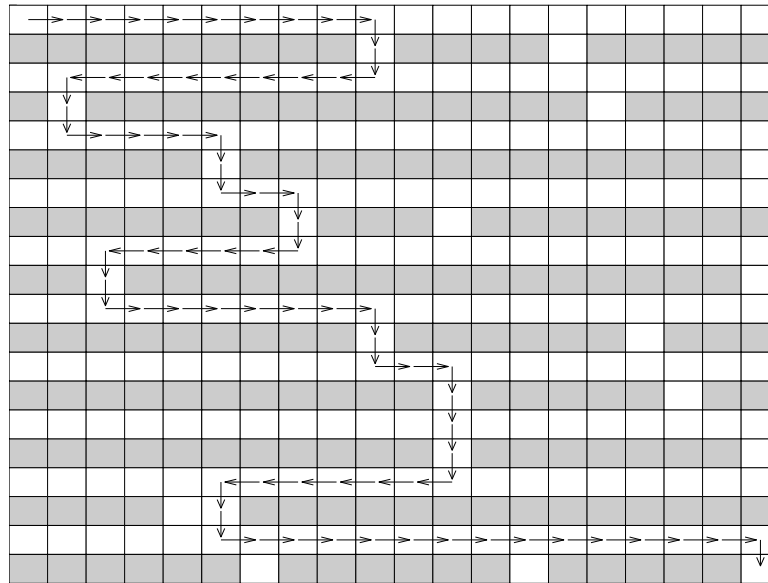
Iteration 3:

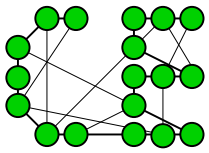




A Maze Problem

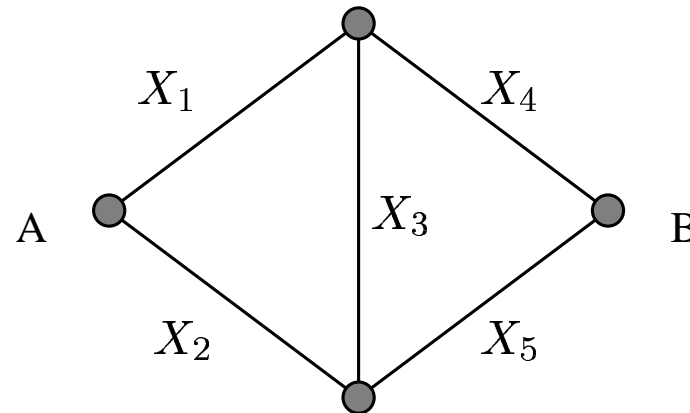
Iteration 4:



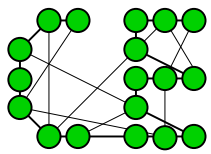


Example 1: Rare Event Simulation

Consider a randomly weighted graph:

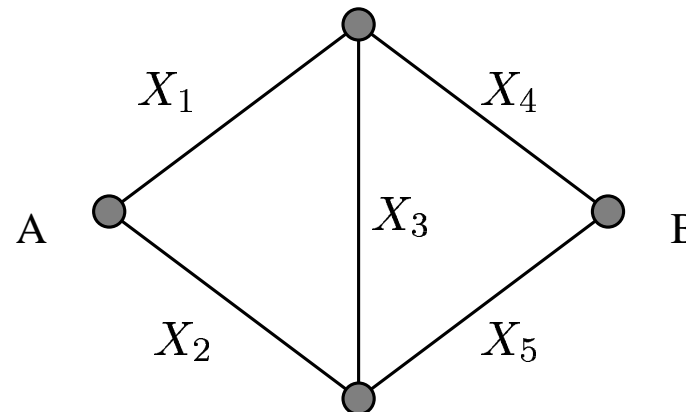


The random weights X_1, \dots, X_5 are independent and exponentially distributed with means u_1, \dots, u_5 .



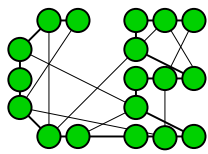
Example 1: Rare Event Simulation

Consider a randomly weighted graph:



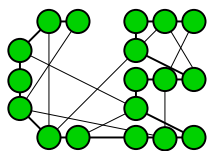
The random weights X_1, \dots, X_5 are independent and exponentially distributed with means u_1, \dots, u_5 .

Find the probability that the length of the shortest path from A to B is greater than or equal to γ .



Crude Monte Carlo (CMC)

Define $\mathbf{X} = (X_1, \dots, X_5)$ and $\mathbf{u} = (u_1, \dots, u_5)$. Let $S(\mathbf{X})$ be the **length of the shortest path** from node A to node B.

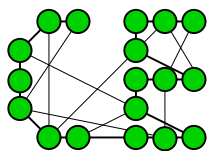


Crude Monte Carlo (CMC)

Define $\mathbf{X} = (X_1, \dots, X_5)$ and $\mathbf{u} = (u_1, \dots, u_5)$. Let $S(\mathbf{X})$ be the **length of the shortest path** from node A to node B.

We wish to estimate

$$\ell = \mathbb{P}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}I_{\{S(\mathbf{x}) \geq \gamma\}} .$$



Crude Monte Carlo (CMC)

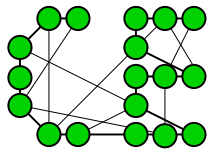
Define $\mathbf{X} = (X_1, \dots, X_5)$ and $\mathbf{u} = (u_1, \dots, u_5)$. Let $S(\mathbf{X})$ be the **length of the shortest path** from node A to node B.

We wish to estimate

$$\ell = \mathbb{P}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}I_{\{S(\mathbf{x}) \geq \gamma\}} .$$

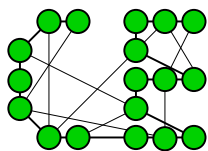
This can be done via **Crude Monte Carlo**: sample independent vectors from density $f(\mathbf{x}; \mathbf{u}) = \prod_{j=1}^5 \exp(-x_j/u_j)/u_j$, and estimate ℓ via

$$\frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} .$$



Importance Sampling (IS)

However, for small ℓ this requires a **very large simulation effort**.



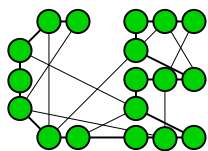
Importance Sampling (IS)

However, for small ℓ this requires a **very large simulation effort**.

A better way is to use **Importance Sampling**: draw $\mathbf{X}_1, \dots, \mathbf{X}_N$ from a **different** density g , and estimate ℓ via the estimator

$$\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} W(\mathbf{X}_i) ,$$

where $W(\mathbf{X}) = f(\mathbf{X})/g(\mathbf{X})$ is called the **likelihood ratio**.

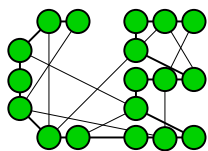


Which Change of Measure?

If we restrict ourselves to g such that X_1, \dots, X_5 are independent and exponentially distributed with means v_1, \dots, v_5 , then

$$W(\mathbf{x}; \mathbf{u}, \mathbf{v}) := \frac{f(\mathbf{x}; \mathbf{u})}{f(\mathbf{x}; \mathbf{v})} = \exp \left(- \sum_{j=1}^5 x_j \left(\frac{1}{u_j} - \frac{1}{v_j} \right) \right) \prod_{j=1}^5 \frac{v_j}{u_j}.$$

In this case the “change of measure” is determined by the **reference vector** $\mathbf{v} = (v_1, \dots, v_5)$.



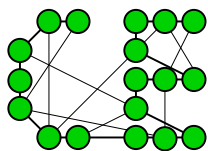
Which Change of Measure?

If we restrict ourselves to g such that X_1, \dots, X_5 are independent and exponentially distributed with means v_1, \dots, v_5 , then

$$W(\mathbf{x}; \mathbf{u}, \mathbf{v}) := \frac{f(\mathbf{x}; \mathbf{u})}{f(\mathbf{x}; \mathbf{v})} = \exp \left(- \sum_{j=1}^5 x_j \left(\frac{1}{u_j} - \frac{1}{v_j} \right) \right) \prod_{j=1}^5 \frac{v_j}{u_j}.$$

In this case the “change of measure” is determined by the **reference vector** $\mathbf{v} = (v_1, \dots, v_5)$.

Question: How do we find the optimal $\mathbf{v} = \mathbf{v}^*$?



Which Change of Measure?

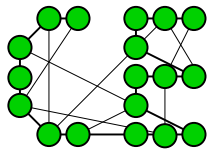
If we restrict ourselves to g such that X_1, \dots, X_5 are independent and exponentially distributed with means v_1, \dots, v_5 , then

$$W(\mathbf{x}; \mathbf{u}, \mathbf{v}) := \frac{f(\mathbf{x}; \mathbf{u})}{f(\mathbf{x}; \mathbf{v})} = \exp \left(- \sum_{j=1}^5 x_j \left(\frac{1}{u_j} - \frac{1}{v_j} \right) \right) \prod_{j=1}^5 \frac{v_j}{u_j}.$$

In this case the “change of measure” is determined by the **reference vector** $\mathbf{v} = (v_1, \dots, v_5)$.

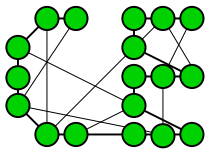
Question: How do we find the optimal $\mathbf{v} = \mathbf{v}^*$?

Answer: Let CE find it adaptively!



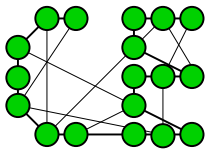
CE Algorithm

1 Define $\hat{v}_0 := u$. Set $t := 1$ (iteration counter).



CE Algorithm

- 1 Define $\hat{\mathbf{v}}_0 := \mathbf{u}$. Set $t := 1$ (iteration counter).
- 2 Update $\hat{\gamma}_t$: Generate $\mathbf{X}_1, \dots, \mathbf{X}_N$ according to $f(\cdot; \hat{\mathbf{v}}_{t-1})$. Let $\hat{\gamma}_t$ be the **worst of the $\rho \times N$ best** performances, provided this is less than γ . Else $\hat{\gamma}_t := \gamma$.

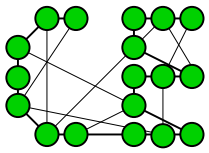


CE Algorithm

- 1 Define $\hat{\mathbf{v}}_0 := \mathbf{u}$. Set $t := 1$ (iteration counter).
- 2 Update $\hat{\gamma}_t$: Generate $\mathbf{X}_1, \dots, \mathbf{X}_N$ according to $f(\cdot; \hat{\mathbf{v}}_{t-1})$. Let $\hat{\gamma}_t$ be the **worst of the $\rho \times N$ best** performances, provided this is less than γ . Else $\hat{\gamma}_t := \gamma$.
- 3 Update $\hat{\mathbf{v}}_t$: Use the **same** sample to calculate, for $j = 1, \dots, n$,

$$\hat{v}_{t,j} = \frac{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \hat{\gamma}_t\}} W(\mathbf{X}_i; \mathbf{u}, \hat{\mathbf{v}}_{t-1}) X_{ij}}{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \hat{\gamma}_t\}} W(\mathbf{X}_i; \mathbf{u}, \hat{\mathbf{v}}_{t-1})}.$$

THIS UPDATING IS BASED ON CE.



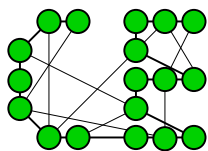
CE Algorithm

- 1 Define $\hat{\mathbf{v}}_0 := \mathbf{u}$. Set $t := 1$ (iteration counter).
- 2 Update $\hat{\gamma}_t$: Generate $\mathbf{X}_1, \dots, \mathbf{X}_N$ according to $f(\cdot; \hat{\mathbf{v}}_{t-1})$. Let $\hat{\gamma}_t$ be the **worst of the $\rho \times N$ best** performances, provided this is less than γ . Else $\hat{\gamma}_t := \gamma$.
- 3 Update $\hat{\mathbf{v}}_t$: Use the **same** sample to calculate, for $j = 1, \dots, n$,

$$\hat{v}_{t,j} = \frac{\sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \hat{\gamma}_t\}} W(\mathbf{X}_i; \mathbf{u}, \hat{\mathbf{v}}_{t-1}) X_{ij}}{\sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \hat{\gamma}_t\}} W(\mathbf{X}_i; \mathbf{u}, \hat{\mathbf{v}}_{t-1})}.$$

THIS UPDATING IS BASED ON CE.

- 4 If $\hat{\gamma}_t = \gamma$ then proceed to step 5; otherwise set $t := t + 1$ and reiterate from step 2.



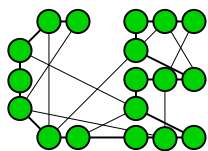
CE Algorithm

- 1 Define $\hat{\mathbf{v}}_0 := \mathbf{u}$. Set $t := 1$ (iteration counter).
- 2 Update $\hat{\gamma}_t$: Generate $\mathbf{X}_1, \dots, \mathbf{X}_N$ according to $f(\cdot; \hat{\mathbf{v}}_{t-1})$. Let $\hat{\gamma}_t$ be the **worst of the $\rho \times N$ best** performances, provided this is less than γ . Else $\hat{\gamma}_t := \gamma$.
- 3 Update $\hat{\mathbf{v}}_t$: Use the **same** sample to calculate, for $j = 1, \dots, n$,

$$\hat{v}_{t,j} = \frac{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \hat{\gamma}_t\}} W(\mathbf{X}_i; \mathbf{u}, \hat{\mathbf{v}}_{t-1}) X_{ij}}{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \hat{\gamma}_t\}} W(\mathbf{X}_i; \mathbf{u}, \hat{\mathbf{v}}_{t-1})}.$$

THIS UPDATING IS BASED ON CE.

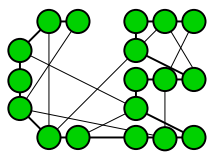
- 4 If $\hat{\gamma}_t = \gamma$ then proceed to step 5; otherwise set $t := t + 1$ and reiterate from step 2.
- 5 Estimate ℓ via the LR estimator, using the final $\hat{\mathbf{v}}_T$.



Example

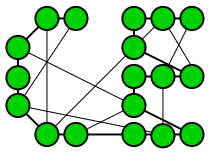
Level: $\gamma = 2$. Fraction of best performances: $\rho = 0.1$. Sample size in steps 2 – 4: $N = 1000$. Final sample size: $N_1 = 10^5$.

t	$\hat{\gamma}_t$	$\hat{\mathbf{v}}_t$				
0		0.250	0.400	0.100	0.300	0.200
1	0.575	0.513	0.718	0.122	0.474	0.335
2	1.032	0.873	1.057	0.120	0.550	0.436
3	1.502	1.221	1.419	0.121	0.707	0.533
4	1.917	1.681	1.803	0.132	0.638	0.523
5	2.000	1.692	1.901	0.129	0.712	0.564



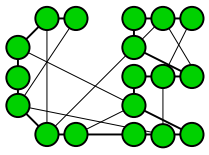
Example (cont.)

- The estimate was $1.34 \cdot 10^{-5}$,



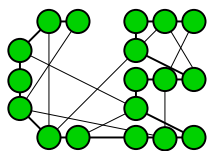
Example (cont.)

- The estimate was $1.34 \cdot 10^{-5}$,
- with an estimated **relative error** (that is, $\text{Std}(\hat{\ell})/\mathbb{E}\hat{\ell}$) of 0.03.



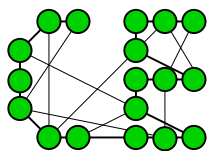
Example (cont.)

- The estimate was $1.34 \cdot 10^{-5}$,
- with an estimated **relative error** (that is, $\text{Std}(\hat{\ell})/\mathbb{E}\hat{\ell}$) of 0.03.
- The simulation time was only 3 seconds (1/2 second for table).



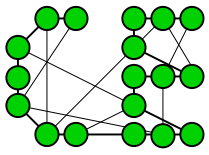
Example (cont.)

- The estimate was $1.34 \cdot 10^{-5}$,
- with an estimated **relative error** (that is, $\text{Std}(\hat{\ell})/\mathbb{E}\hat{\ell}$) of 0.03.
- The simulation time was only 3 seconds (1/2 second for table).
- CMC with $N_1 = 10^8$ samples gave an estimate $1.30 \cdot 10^{-5}$ with the same RE (0.03). The simulation time was 1875 seconds.



Example (cont.)

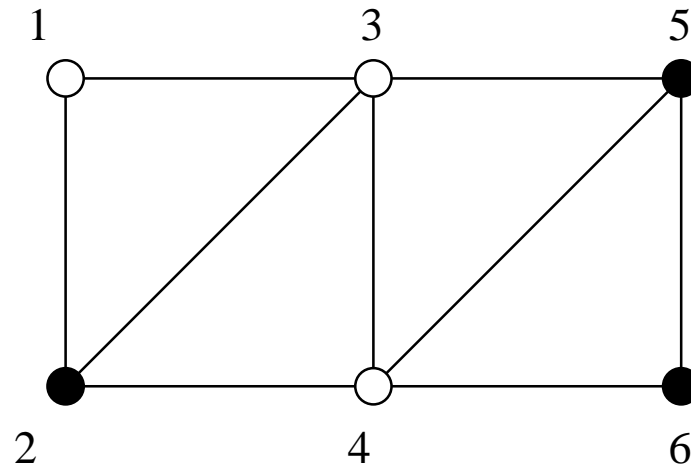
- The estimate was $1.34 \cdot 10^{-5}$,
- with an estimated **relative error** (that is, $\text{Std}(\hat{\ell})/\mathbb{E}\hat{\ell}$) of 0.03.
- The simulation time was only 3 seconds (1/2 second for table).
- CMC with $N_1 = 10^8$ samples gave an estimate $1.30 \cdot 10^{-5}$ with the same RE (0.03). The simulation time was 1875 seconds.
- With minimal effort we reduced our simulation time by a factor of 625.

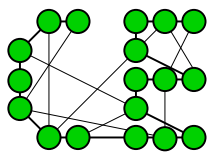


Example 2: The Max-Cut Problem

Consider a weighted graph G with node set $V = \{1, \dots, n\}$. Partition the nodes of the graph into two subsets V_1 and V_2 such that the sum of the weights of the edges going from one subset to the other is maximised.

Example



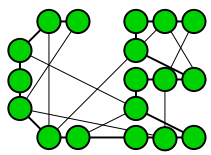


Cost matrix:

$$C = \begin{pmatrix} 0 & c_{12} & c_{13} & 0 & 0 & 0 \\ c_{21} & 0 & c_{23} & c_{24} & 0 & 0 \\ c_{31} & c_{32} & 0 & c_{34} & c_{35} & 0 \\ 0 & c_{42} & c_{43} & 0 & c_{45} & c_{46} \\ 0 & 0 & c_{53} & c_{54} & 0 & c_{56} \\ 0 & 0 & 0 & c_{64} & c_{65} & 0 \end{pmatrix}.$$

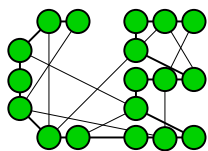
$\{V_1, V_2\} = \{\{1, 3, 4\}, \{2, 5, 6\}\}$ is a possible **cut**. The **cost** of the cut is

$$c_{12} + c_{32} + c_{35} + c_{42} + c_{45} + c_{46}.$$



Random Cut Vector

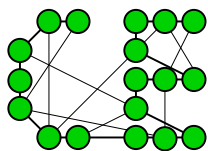
We can represent a cut via a **cut vector** $\mathbf{x} = (x_1, \dots, x_n)$, where $x_i = 1$ if node i belongs to same partition as 1, and 0 else.



Random Cut Vector

We can represent a cut via a **cut vector** $\mathbf{x} = (x_1, \dots, x_n)$, where $x_i = 1$ if node i belongs to same partition as 1, and 0 else.

For example, the cut $\{\{1, 3, 4\}, \{2, 5, 6\}\}$ can be represented via the cut vector $(1, 0, 1, 1, 0, 0)$.

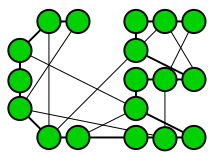


Random Cut Vector

We can represent a cut via a **cut vector** $\mathbf{x} = (x_1, \dots, x_n)$, where $x_i = 1$ if node i belongs to same partition as 1, and 0 else.

For example, the cut $\{\{1, 3, 4\}, \{2, 5, 6\}\}$ can be represented via the cut vector $(1, 0, 1, 1, 0, 0)$.

Let \mathcal{X} be the set of all cut vectors $\mathbf{x} = (1, x_2, \dots, x_n)$.



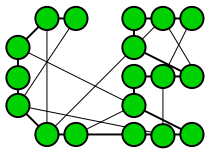
Random Cut Vector

We can represent a cut via a **cut vector** $\mathbf{x} = (x_1, \dots, x_n)$, where $x_i = 1$ if node i belongs to same partition as 1, and 0 else.

For example, the cut $\{\{1, 3, 4\}, \{2, 5, 6\}\}$ can be represented via the cut vector $(1, 0, 1, 1, 0, 0)$.

Let \mathcal{X} be the set of all cut vectors $\mathbf{x} = (1, x_2, \dots, x_n)$.

Let $S(\mathbf{x})$ be the corresponding cost of the cut.



Random Cut Vector

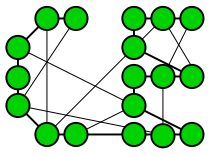
We can represent a cut via a **cut vector** $\mathbf{x} = (x_1, \dots, x_n)$, where $x_i = 1$ if node i belongs to same partition as 1, and 0 else.

For example, the cut $\{\{1, 3, 4\}, \{2, 5, 6\}\}$ can be represented via the cut vector $(1, 0, 1, 1, 0, 0)$.

Let \mathcal{X} be the set of all cut vectors $\mathbf{x} = (1, x_2, \dots, x_n)$.

Let $S(\mathbf{x})$ be the corresponding cost of the cut.

We wish to maximise $S(\mathbf{x})$ via the CE method.

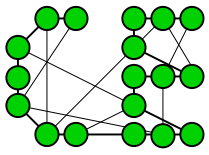


General CE Procedure

First, cast the original optimization problem of $S(\mathbf{x})$ into an associated rare-events estimation problem: the estimation of

$$\ell = \mathbb{P}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}I_{\{S(\mathbf{x}) \geq \gamma\}} .$$

Second, formulate a **parameterized random mechanism** to generate objects $\mathbf{X} \in \mathcal{X}$. Then, iterate the following steps:



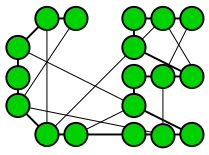
General CE Procedure

First, cast the original optimization problem of $S(\mathbf{x})$ into an associated rare-events estimation problem: the estimation of

$$\ell = \mathbb{P}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}I_{\{S(\mathbf{x}) \geq \gamma\}} .$$

Second, formulate a **parameterized random mechanism** to generate objects $\mathbf{X} \in \mathcal{X}$. Then, iterate the following steps:

- *Generate a random sample* of objects $\mathbf{X}_1, \dots, \mathbf{X}_N \in \mathcal{X}$ (e.g., cut vectors).



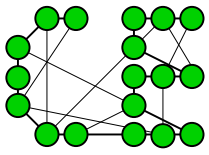
General CE Procedure

First, cast the original optimization problem of $S(\mathbf{x})$ into an associated rare-events estimation problem: the estimation of

$$\ell = \mathbb{P}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}I_{\{S(\mathbf{x}) \geq \gamma\}} .$$

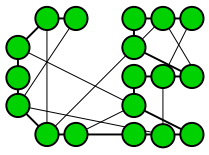
Second, formulate a **parameterized random mechanism** to generate objects $\mathbf{X} \in \mathcal{X}$. Then, iterate the following steps:

- *Generate a random sample* of objects $\mathbf{X}_1, \dots, \mathbf{X}_N \in \mathcal{X}$ (e.g., cut vectors).
- *Update the parameters* of the random mechanism (obtained via CE minimization), in order to produce a better sample in the next iteration.



Generation and Updating Formulas

Generation of cut vectors: The most natural and easiest way to generate the cut vectors is to let X_2, \dots, X_n be independent Bernoulli random variables with success probabilities p_2, \dots, p_n .

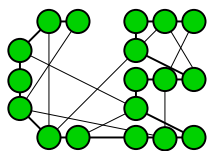


Generation and Updating Formulas

Generation of cut vectors: The most natural and easiest way to generate the cut vectors is to let X_2, \dots, X_n be independent Bernoulli random variables with success probabilities p_2, \dots, p_n .

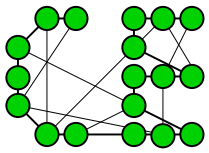
Updating formulas: From CE minimization: the updated probabilities are the maximum likelihood estimates of the ρN best samples:

$$\hat{p}_{t,j} = \frac{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \hat{\gamma}_t\}} I_{\{X_{ij}=1\}}}{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \hat{\gamma}_t\}}}, \quad j = 2, \dots, n .$$



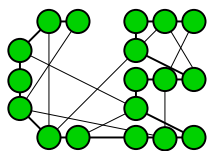
Algorithm

1 Start with $\hat{p}_0 = (1/2, \dots, 1/2)$. Let $t := 1$.



Algorithm

- 1 Start with $\hat{p}_0 = (1/2, \dots, 1/2)$. Let $t := 1$.
- 2 **Update $\hat{\gamma}_t$:** Draw X_1, \dots, X_N from $\text{Ber}(\hat{p}_t)$. Let $\hat{\gamma}_t$ be the worst performance of the $\rho \times 100\%$ best performances.

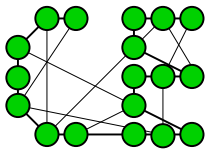


Algorithm

- 1 Start with $\hat{p}_0 = (1/2, \dots, 1/2)$. Let $t := 1$.
- 2 **Update $\hat{\gamma}_t$:** Draw $\mathbf{X}_1, \dots, \mathbf{X}_N$ from $\text{Ber}(\hat{p}_t)$. Let $\hat{\gamma}_t$ be the worst performance of the $\rho \times 100\%$ best performances.
- 3 **Update \hat{p}_t :** Use the same sample to calculate

$$\hat{p}_{t,j} = \frac{\sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \hat{\gamma}_t\}} I_{\{X_{ij}=1\}}}{\sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \hat{\gamma}_t\}}},$$

$j = 1, \dots, n$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{in})$, and increase t by 1.



Algorithm

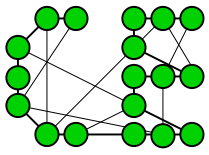
- 1 Start with $\hat{p}_0 = (1/2, \dots, 1/2)$. Let $t := 1$.
- 2 **Update $\hat{\gamma}_t$:** Draw $\mathbf{X}_1, \dots, \mathbf{X}_N$ from $\text{Ber}(\hat{p}_t)$. Let $\hat{\gamma}_t$ be the worst performance of the $\rho \times 100\%$ best performances.

- 3 **Update \hat{p}_t :** Use the same sample to calculate

$$\hat{p}_{t,j} = \frac{\sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \hat{\gamma}_t\}} I_{\{X_{ij}=1\}}}{\sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \hat{\gamma}_t\}}},$$

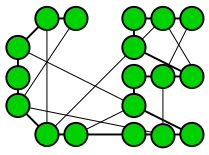
$j = 1, \dots, n$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{in})$, and increase t by 1.

- 4 If the stopping criterion is met, then stop; otherwise set $t := t + 1$ and reiterate from step 2.



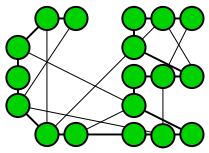
Example

- Results for the case with $n = 400$, $m = 200$ nodes are given next.



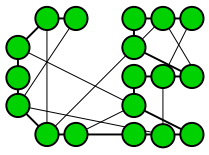
Example

- Results for the case with $n = 400, m = 200$ nodes are given next.
- Parameters: $\rho = 0.1, N = 1000$.



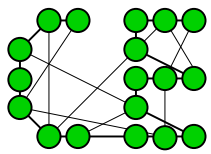
Example

- Results for the case with $n = 400$, $m = 200$ nodes are given next.
- Parameters: $\rho = 0.1$, $N = 1000$.
- The CPU time was only 100 seconds (Matlab, pentium III, 500 Mhz).

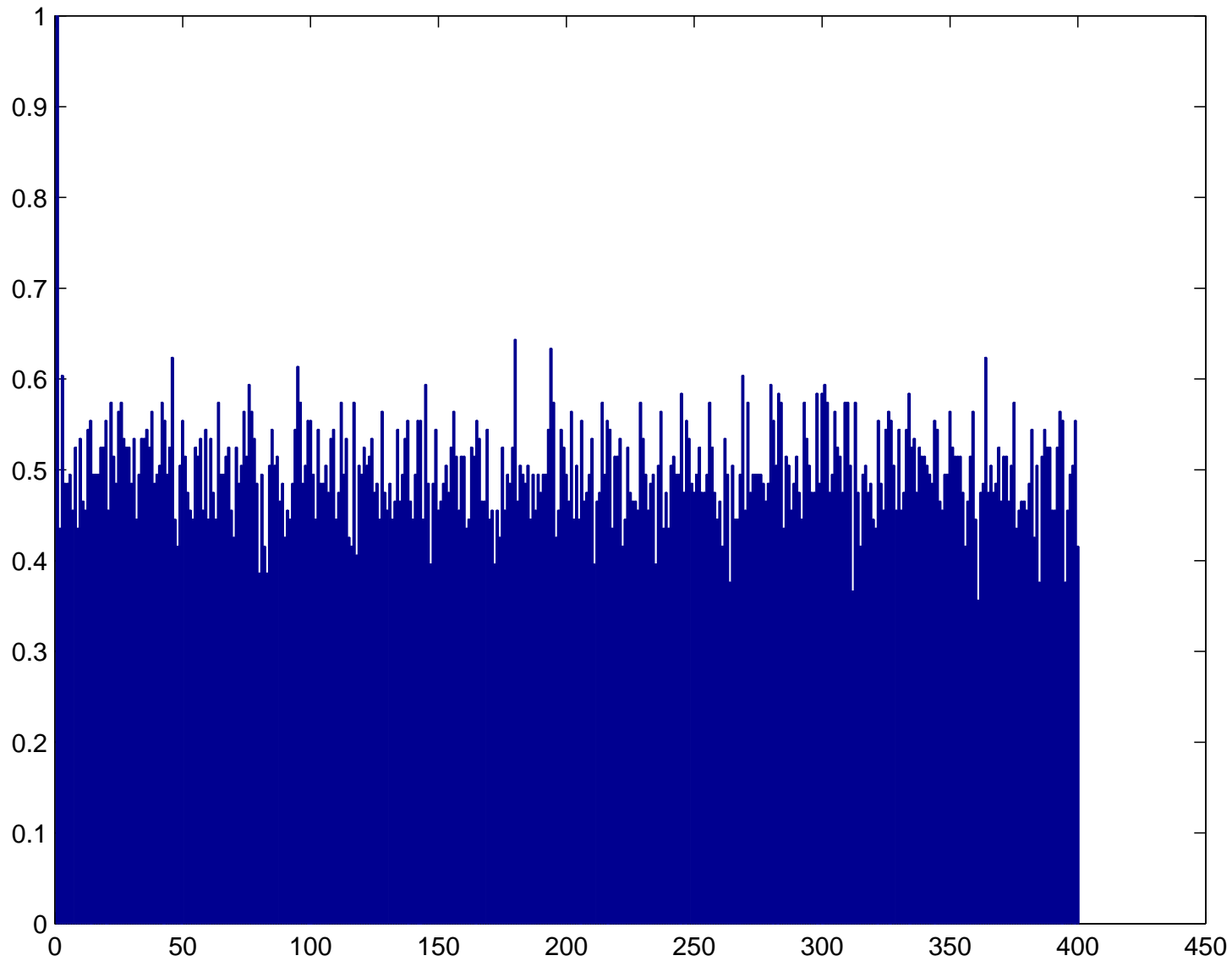


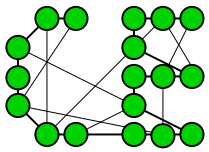
Example

- Results for the case with $n = 400$, $m = 200$ nodes are given next.
- Parameters: $\rho = 0.1$, $N = 1000$.
- The CPU time was only 100 seconds (Matlab, pentium III, 500 Mhz).
- The CE algorithm converges quickly, yielding the exact optimal solution 40000 in 22 iterations.

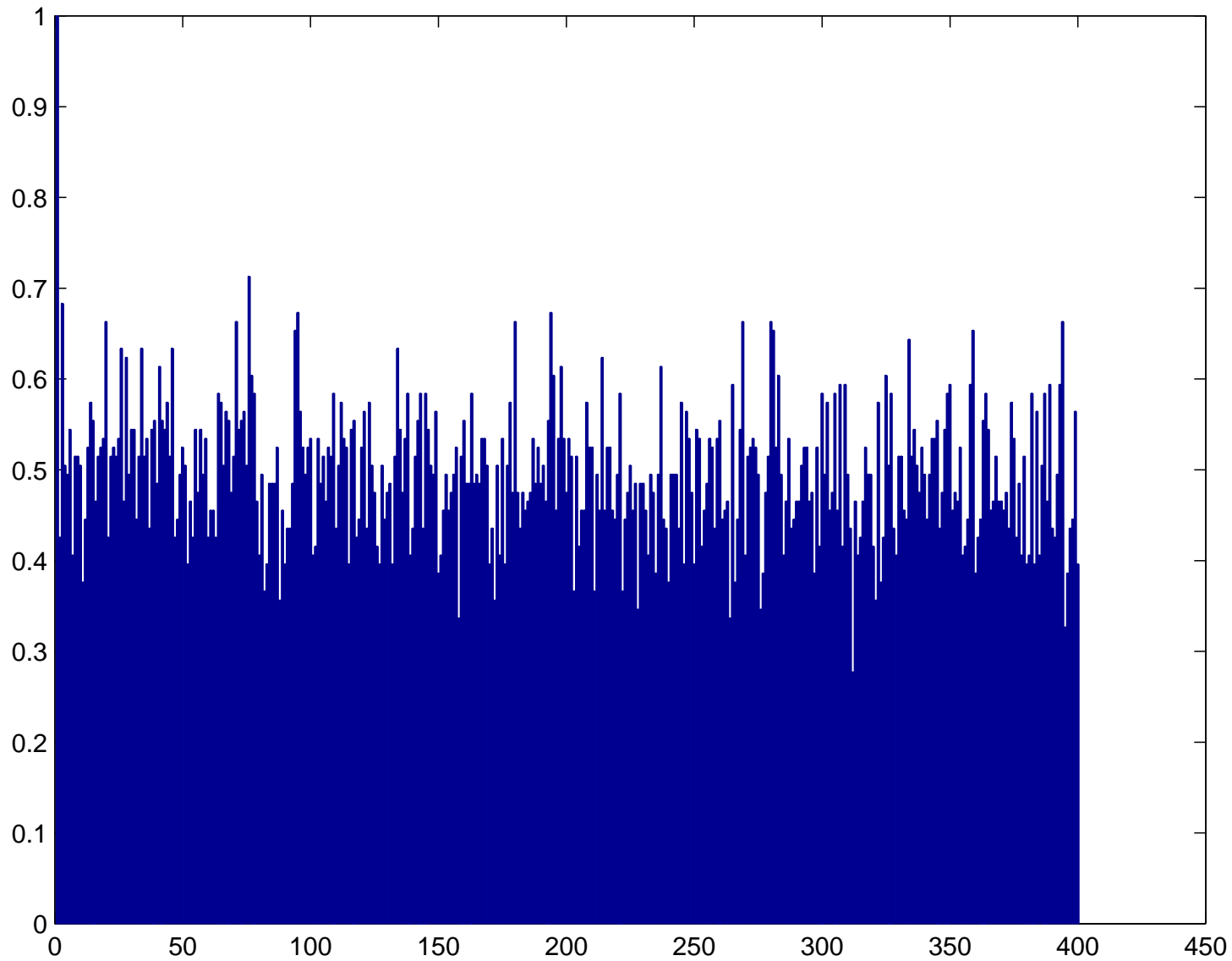


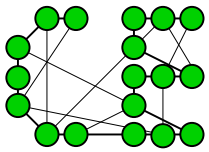
Max-Cut



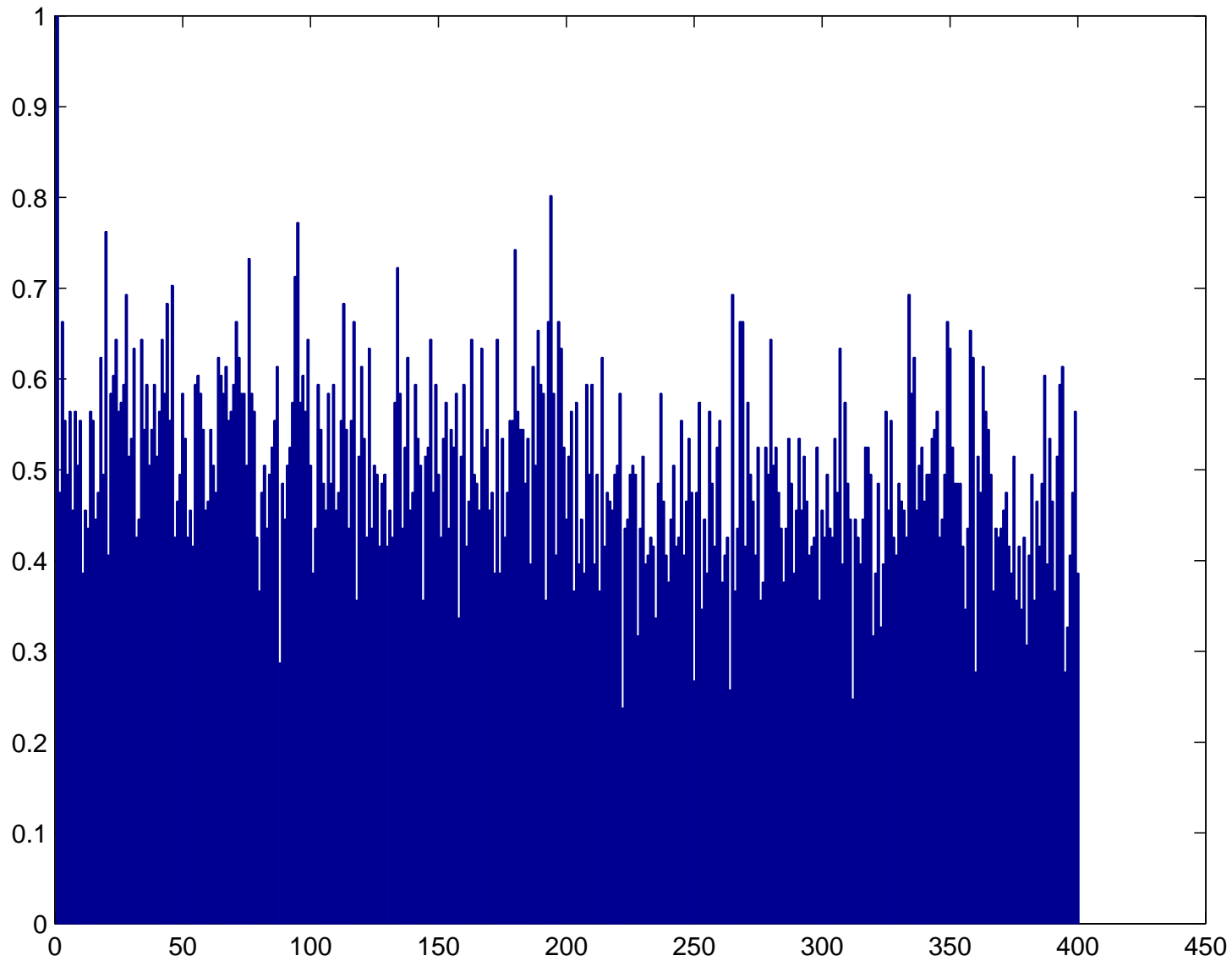


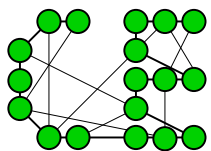
Max-Cut



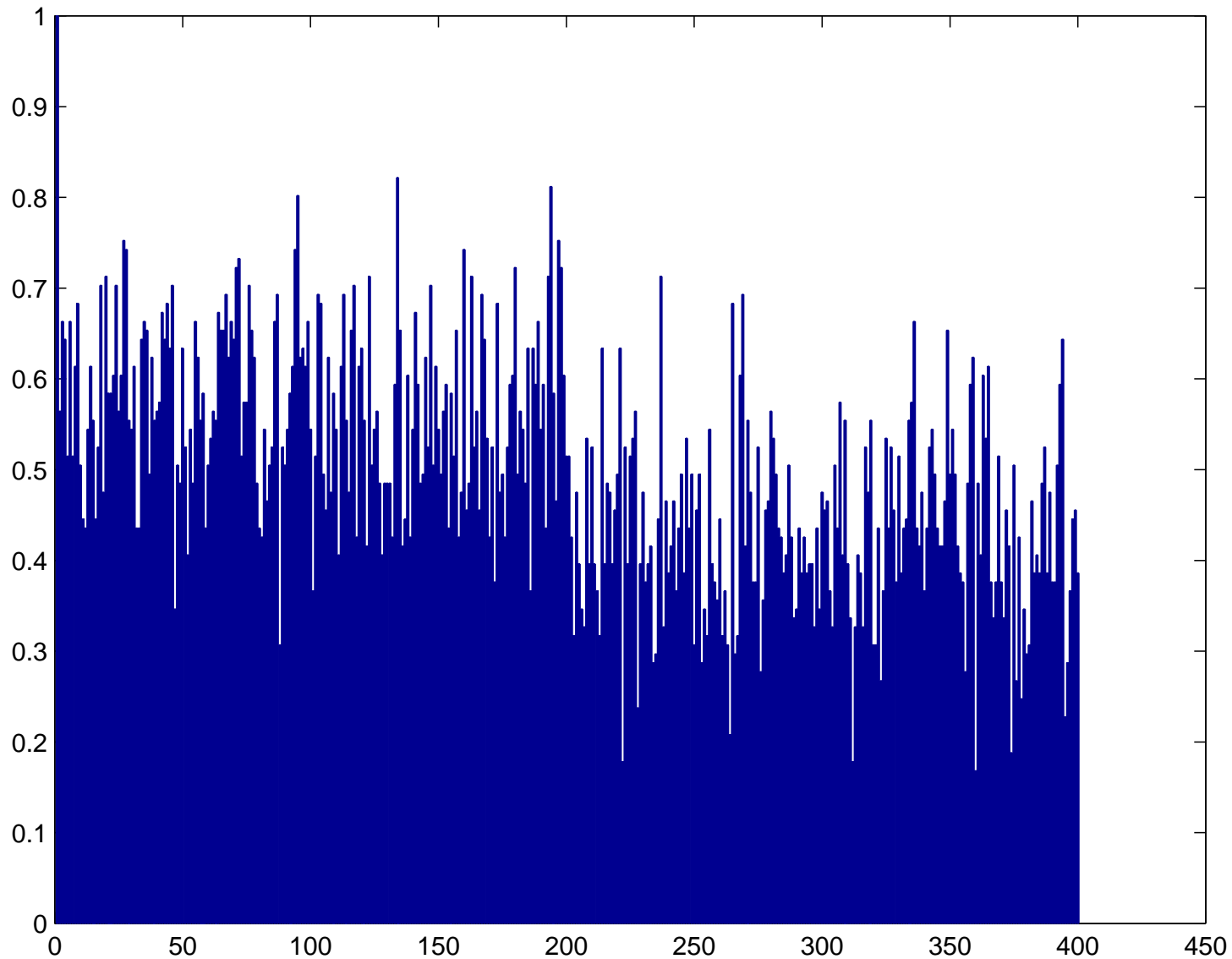


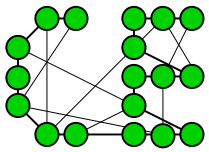
Max-Cut



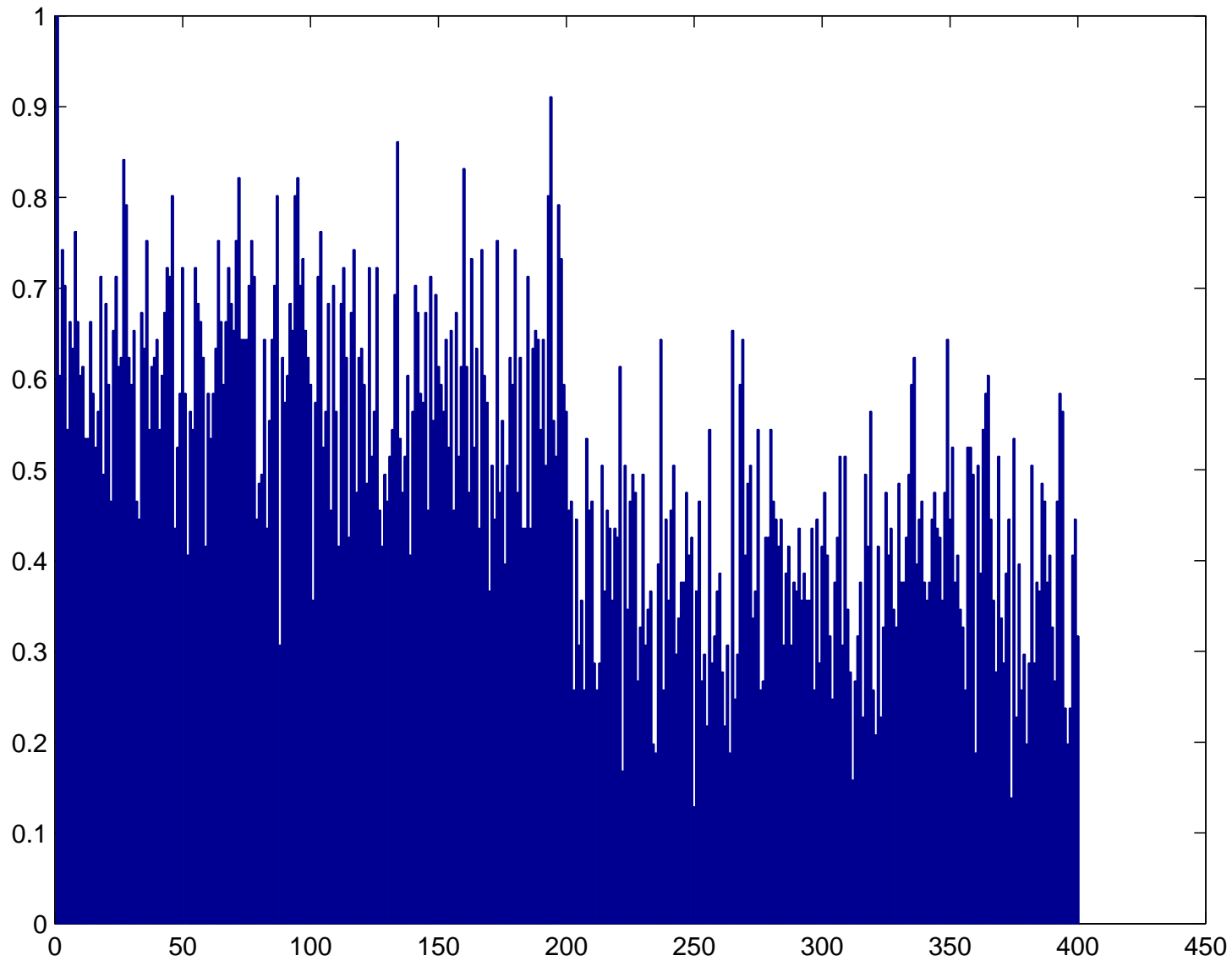


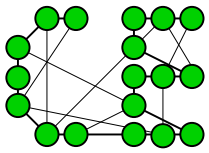
Max-Cut



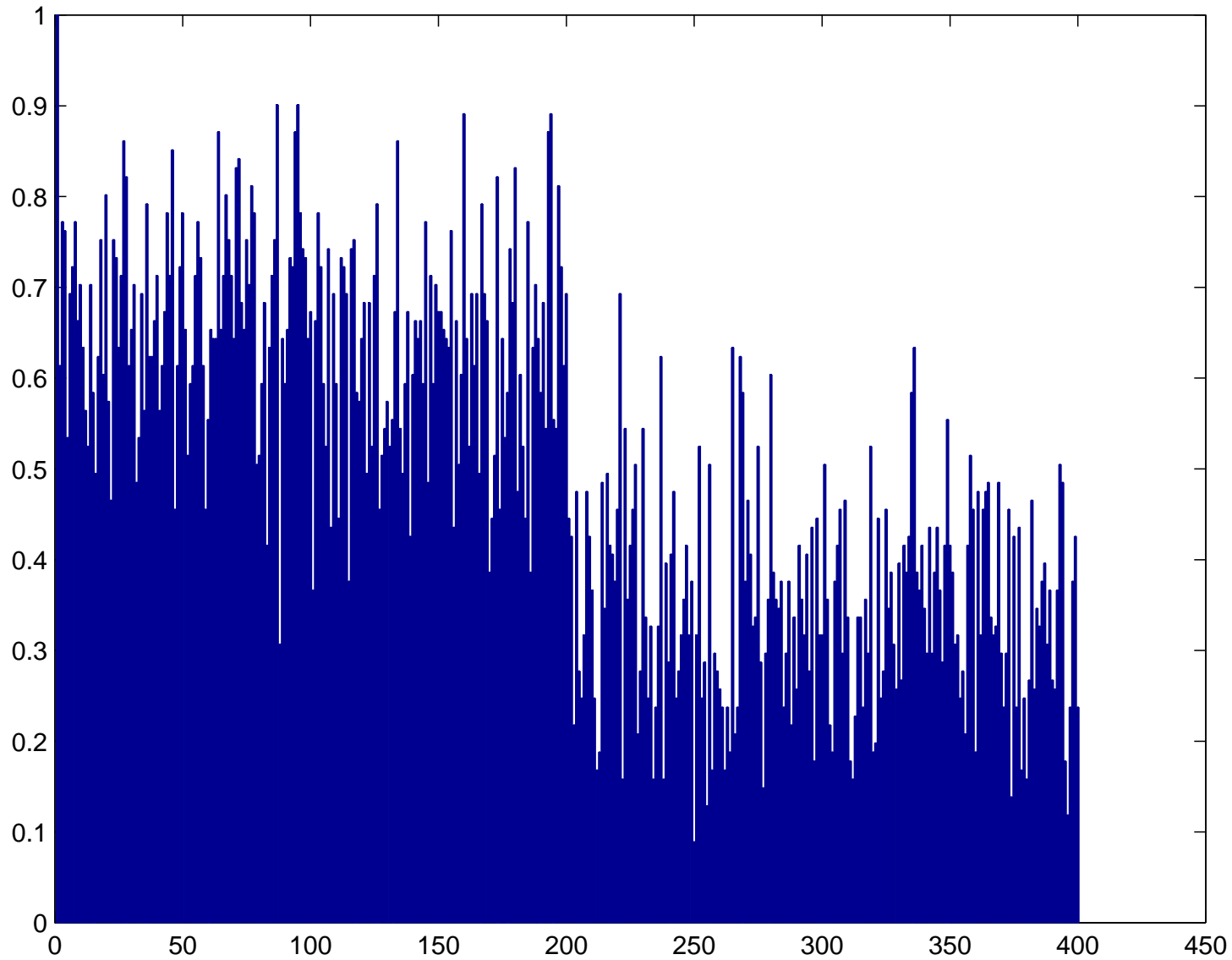


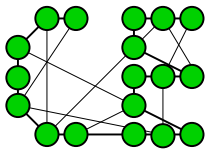
Max-Cut



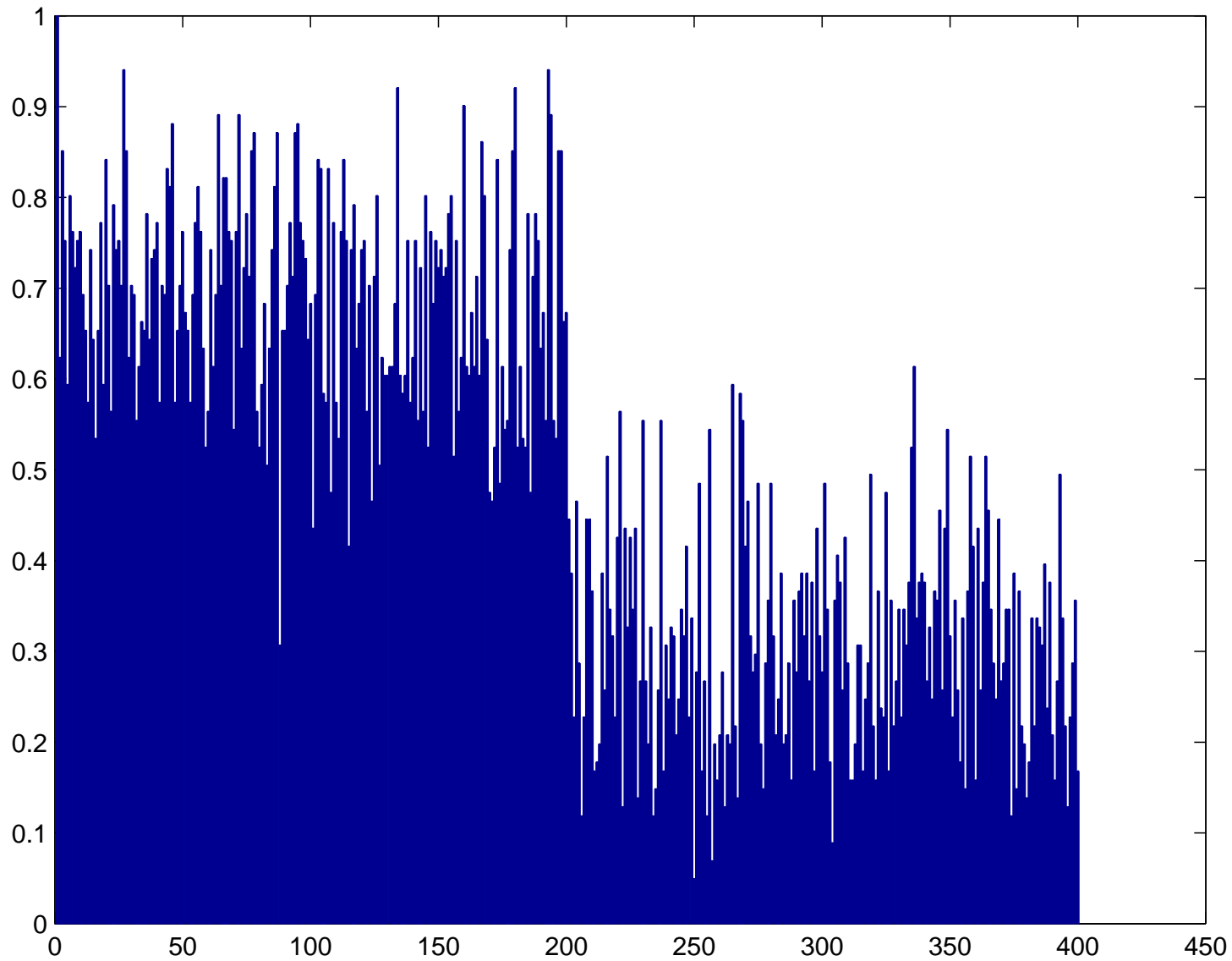


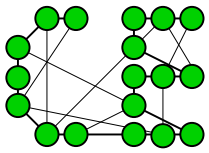
Max-Cut



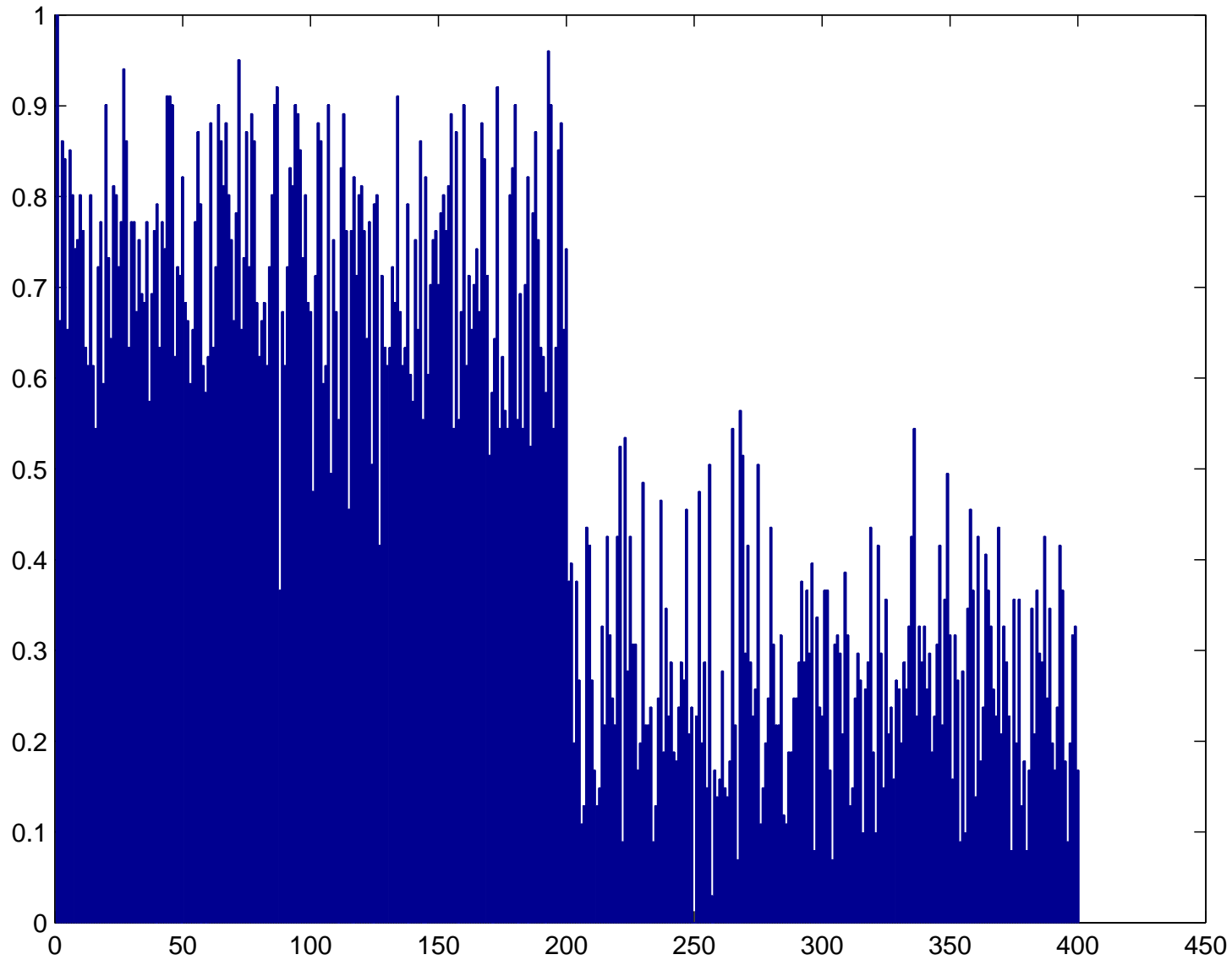


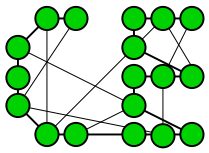
Max-Cut



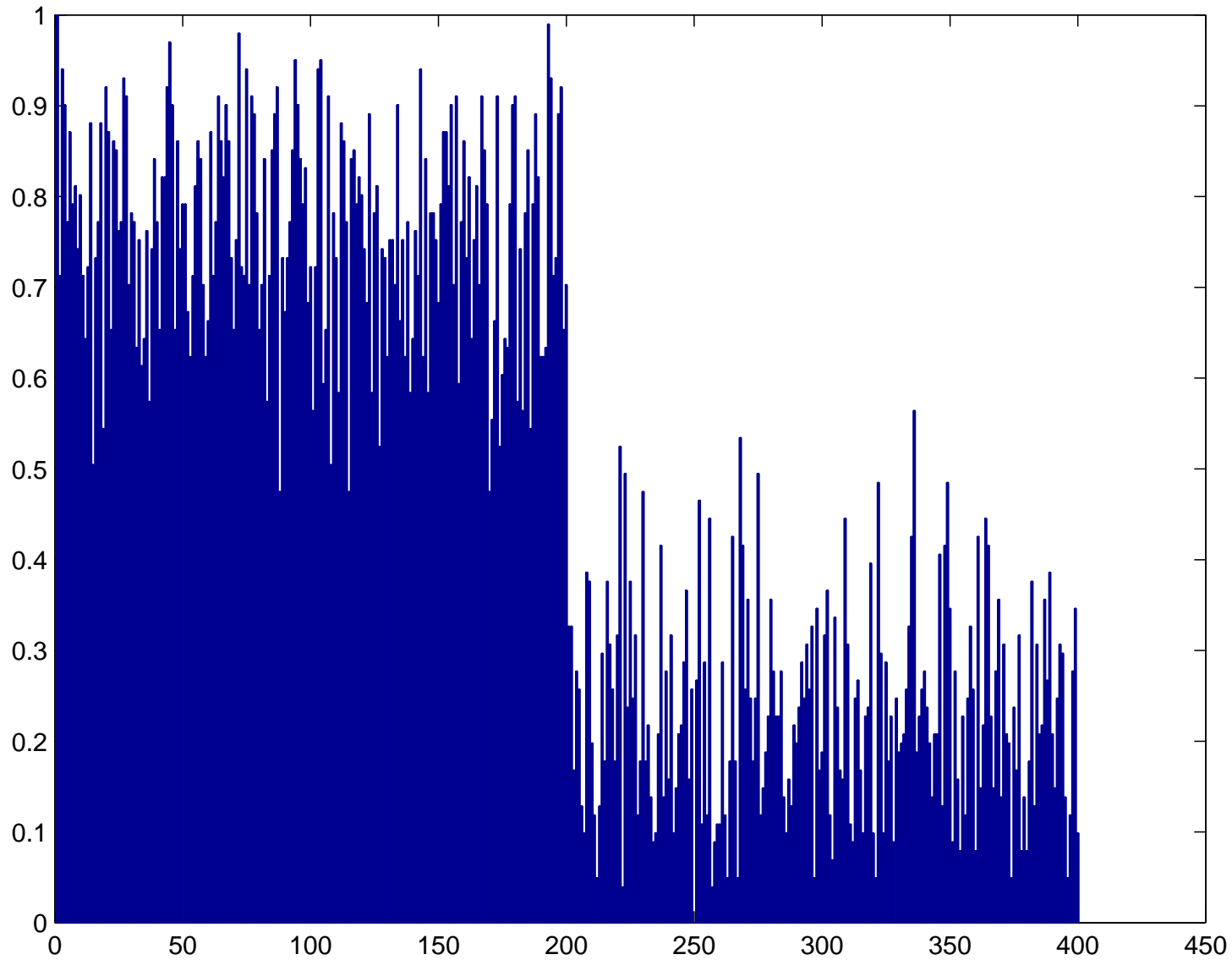


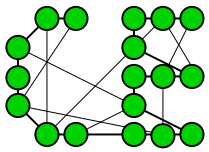
Max-Cut



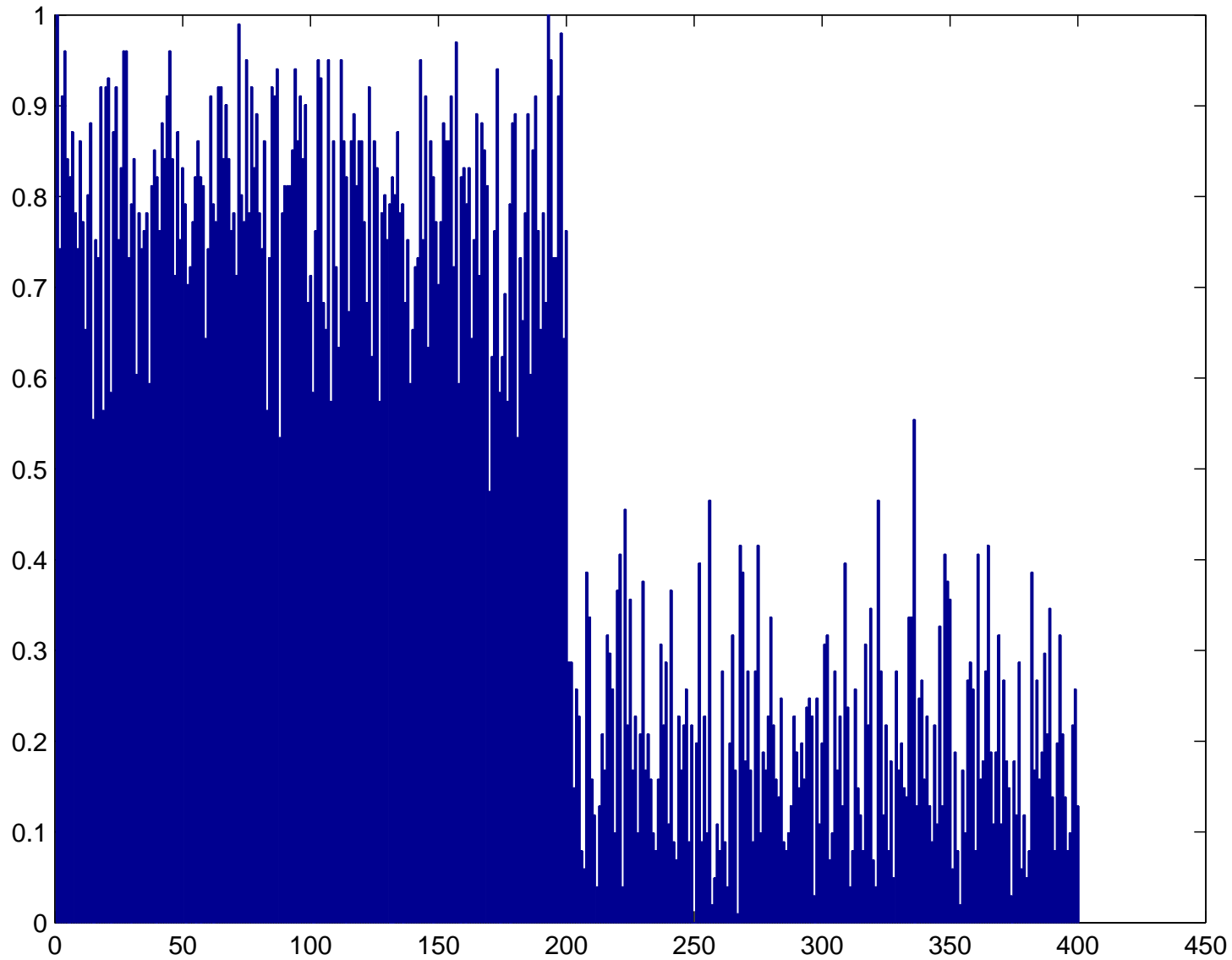


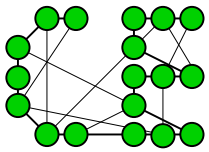
Max-Cut



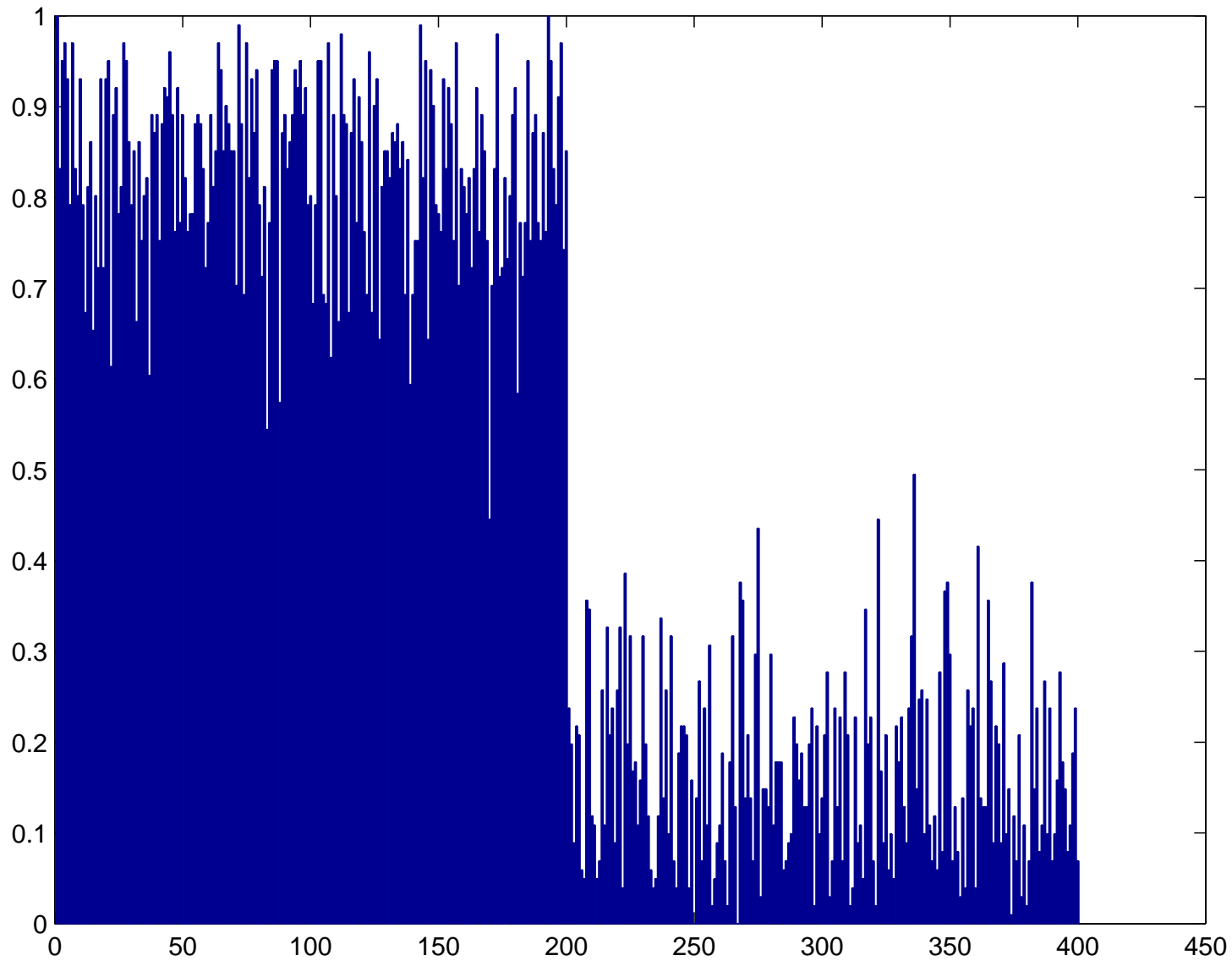


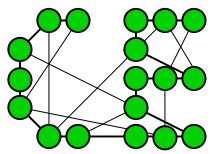
Max-Cut



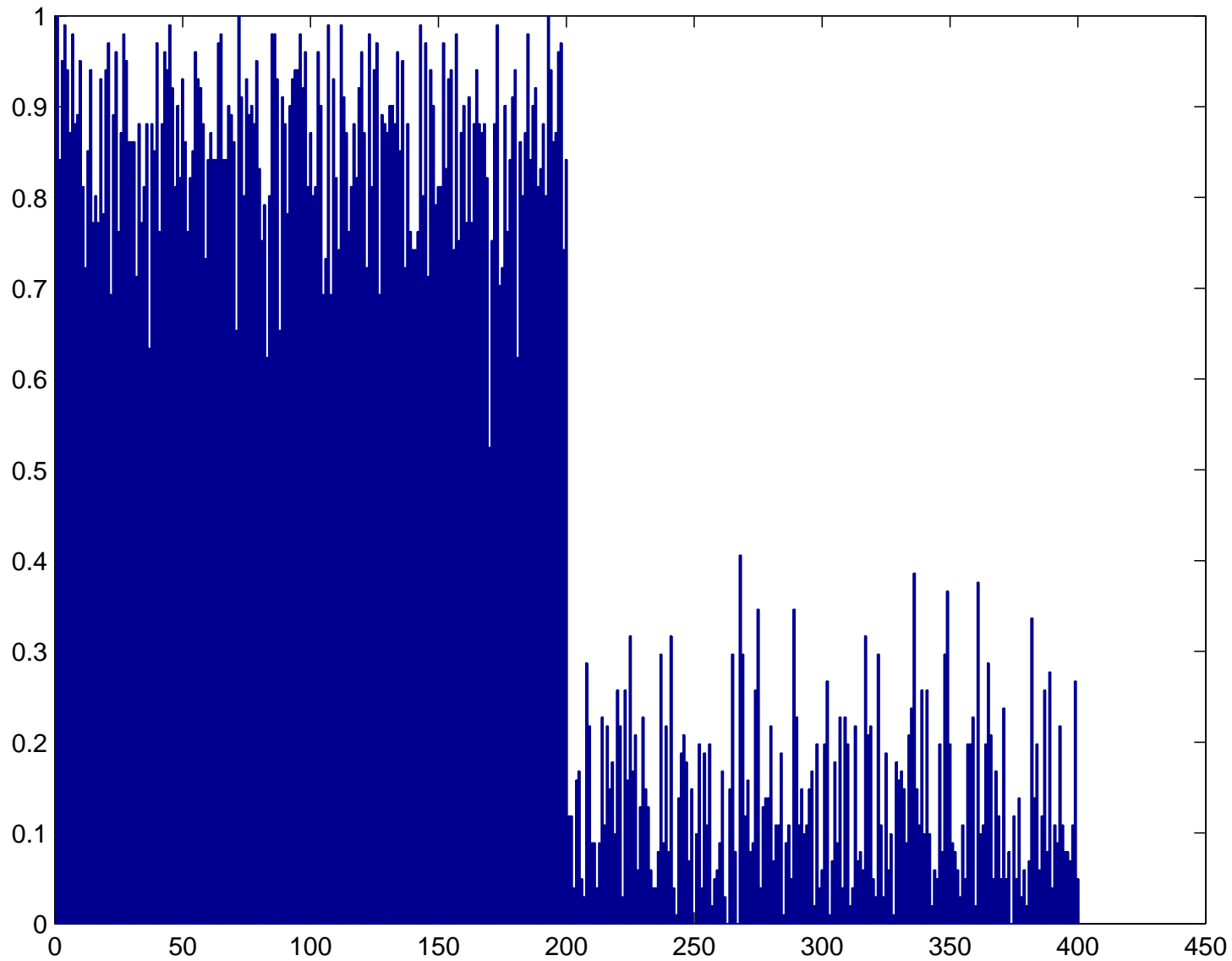


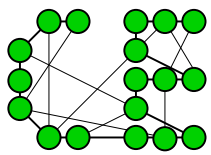
Max-Cut



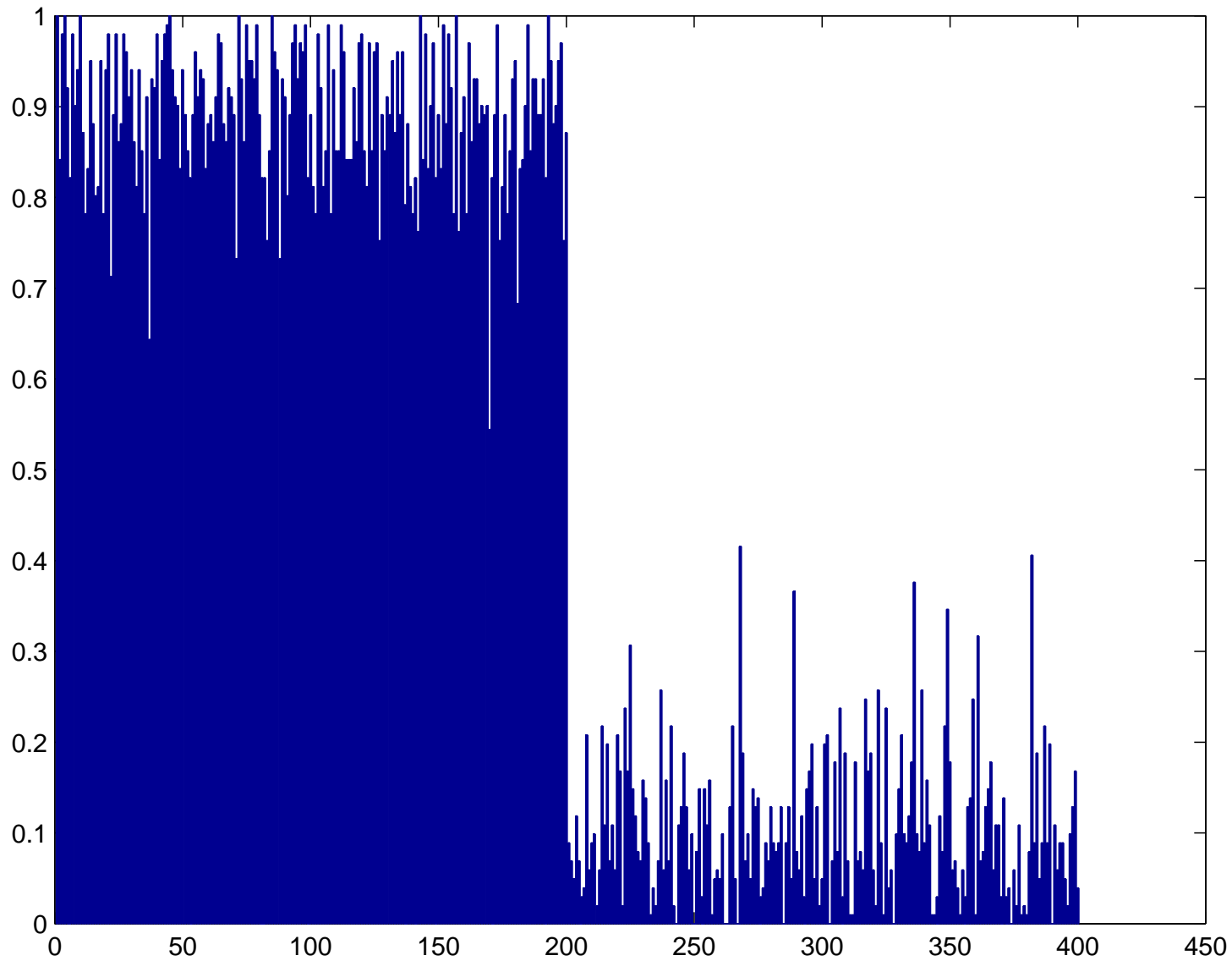


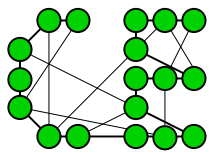
Max-Cut



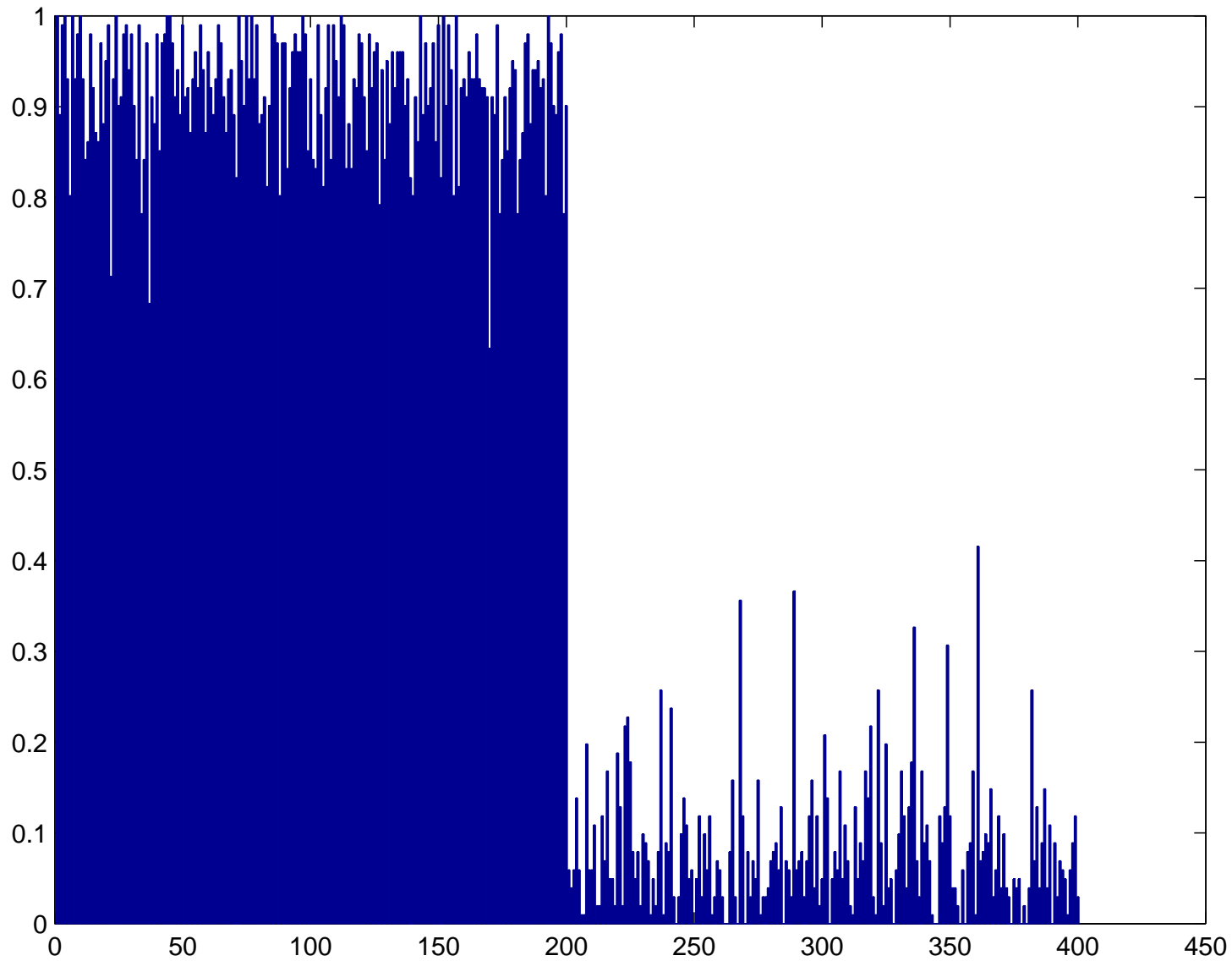


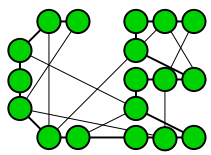
Max-Cut



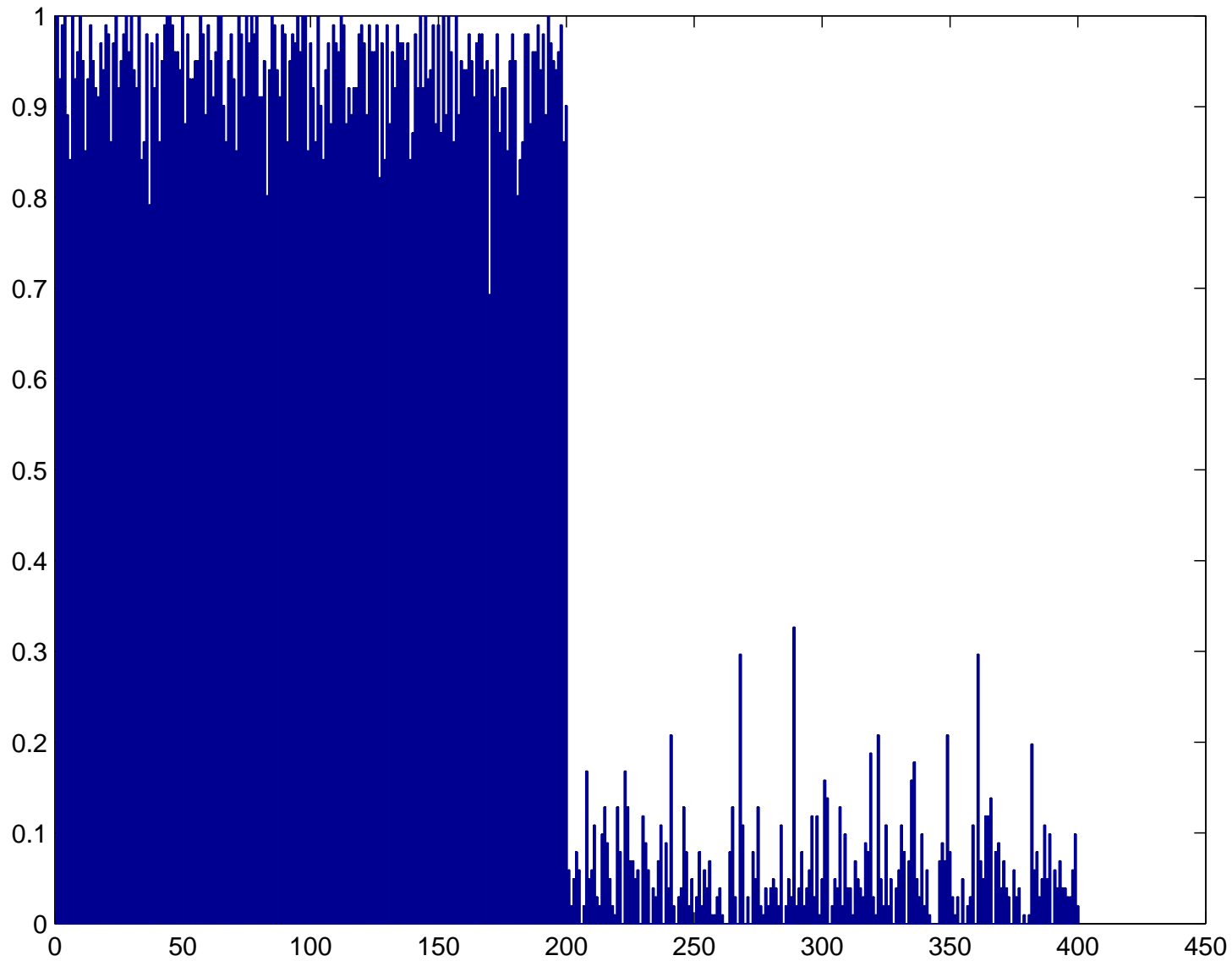


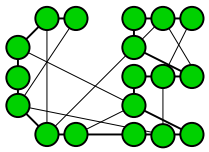
Max-Cut



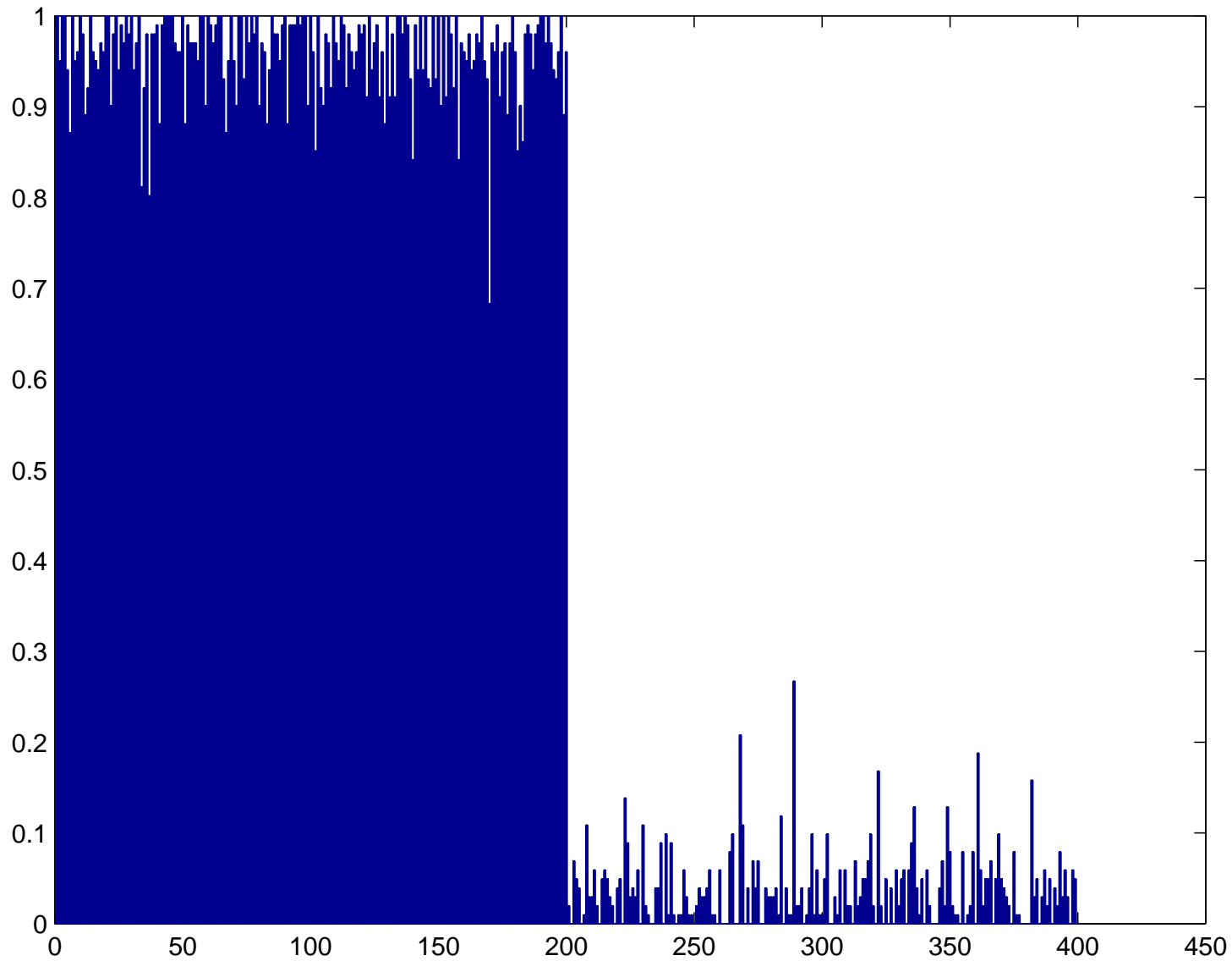


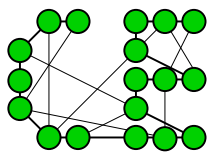
Max-Cut



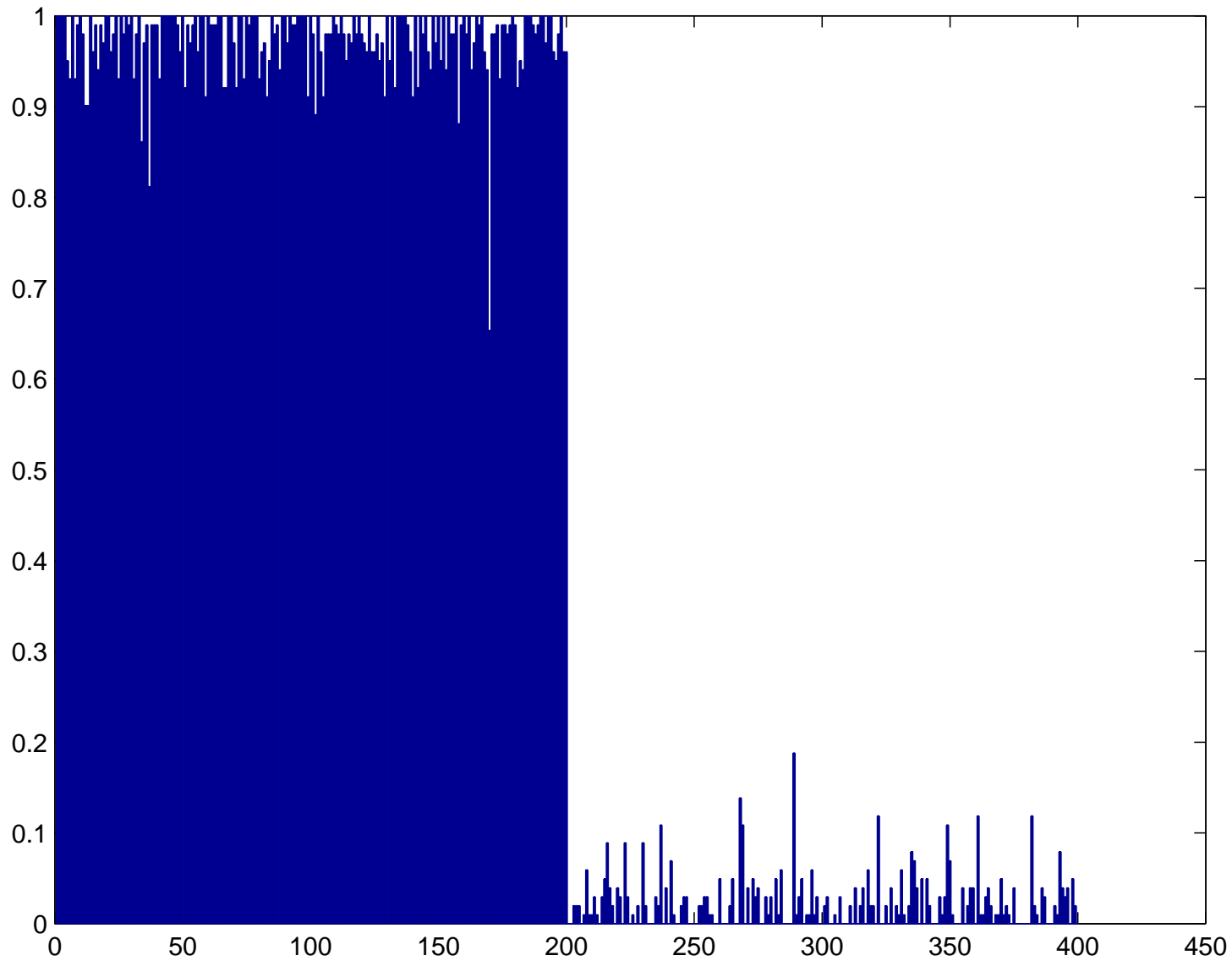


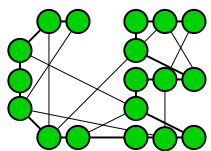
Max-Cut



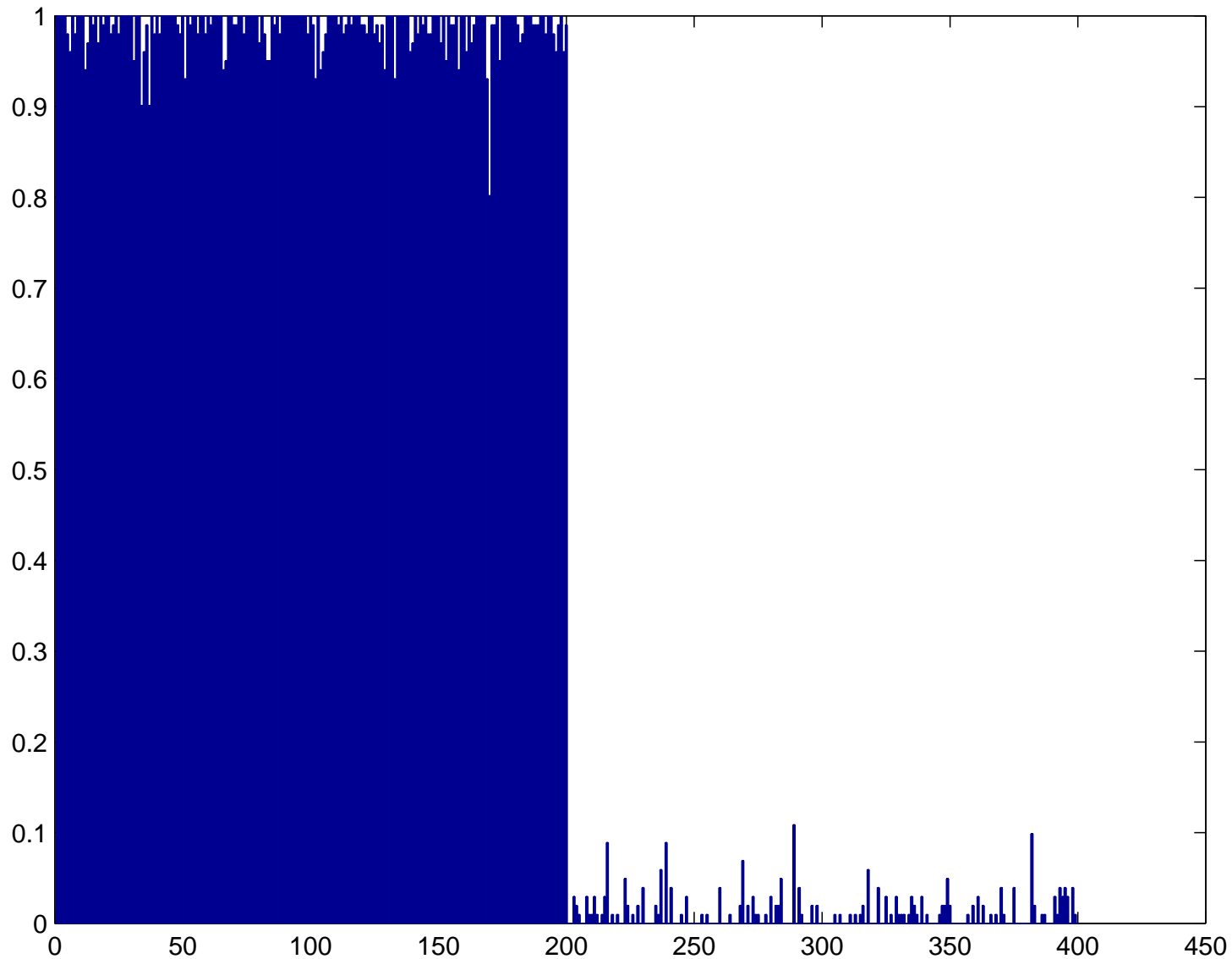


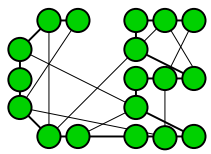
Max-Cut



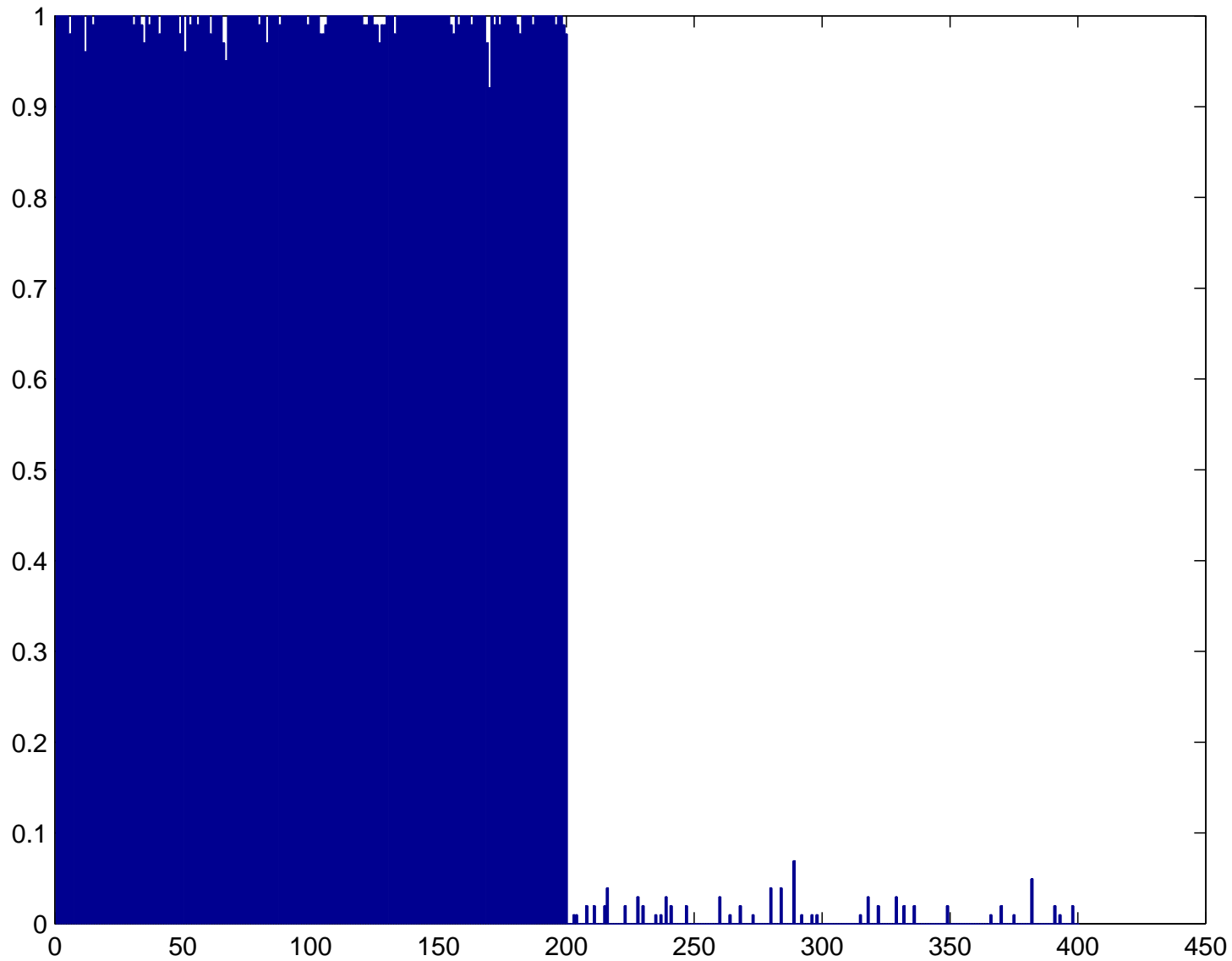


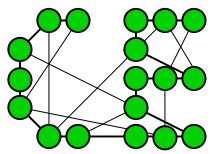
Max-Cut



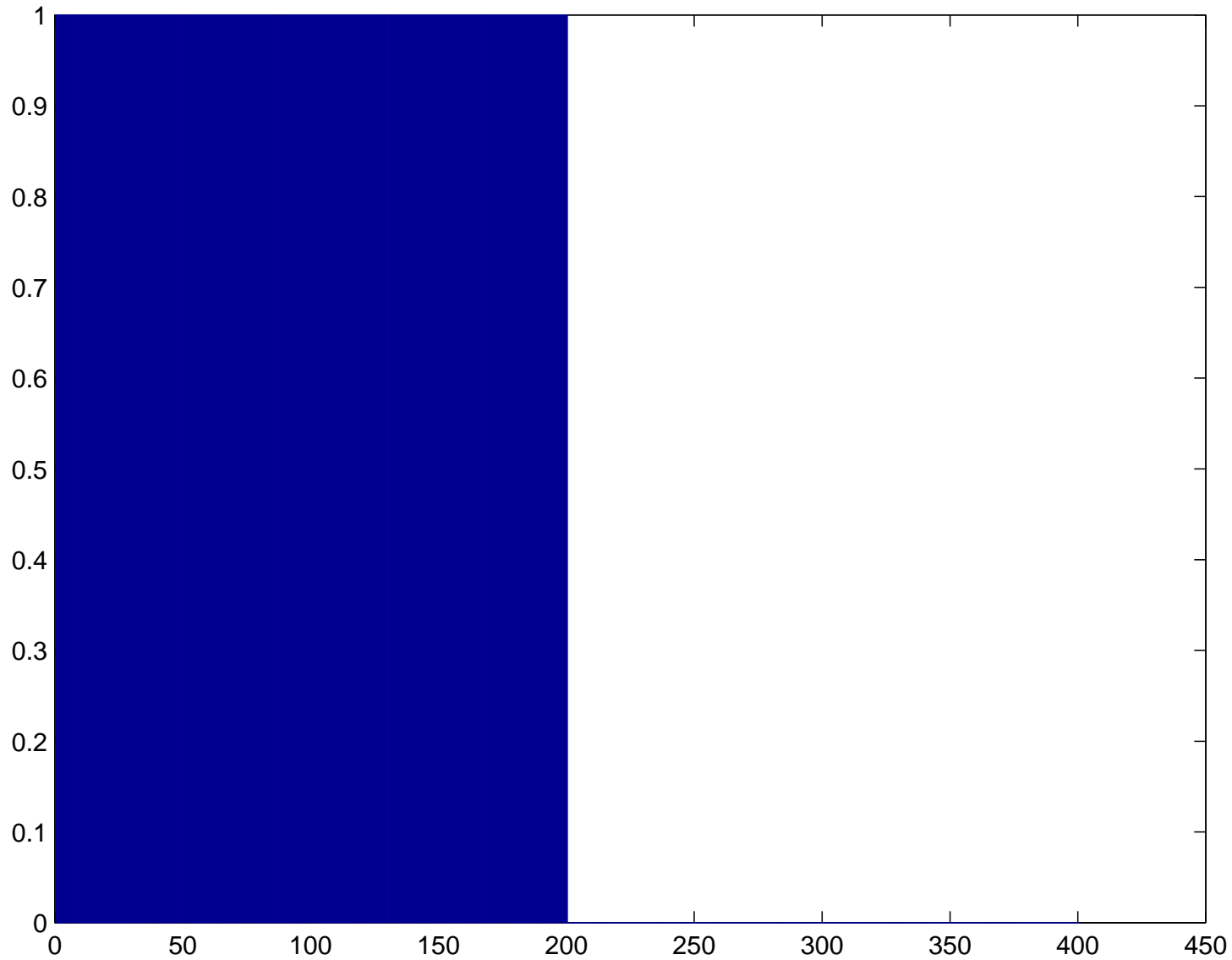


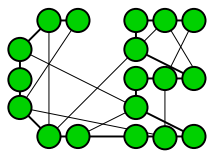
Max-Cut



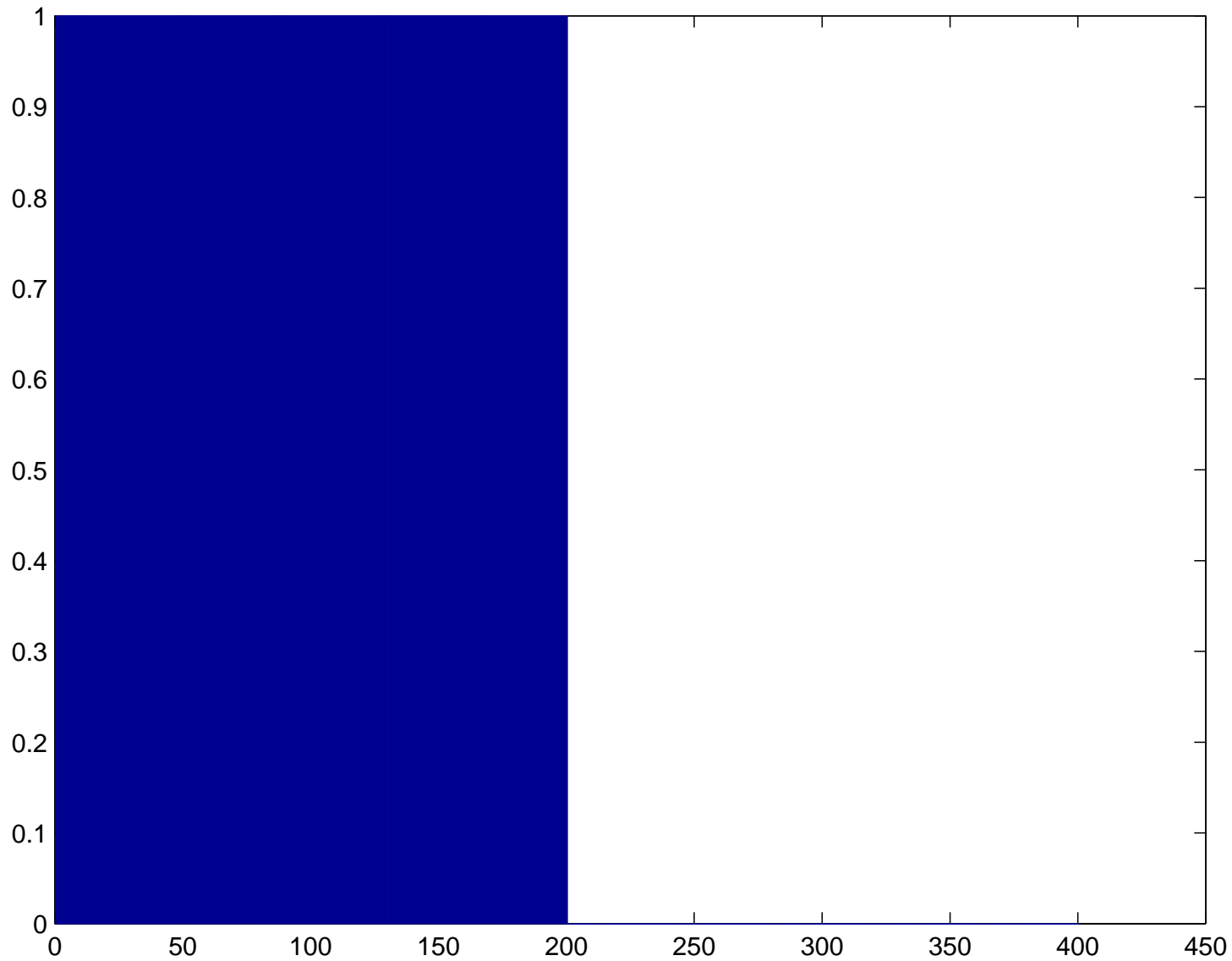


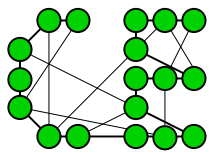
Max-Cut



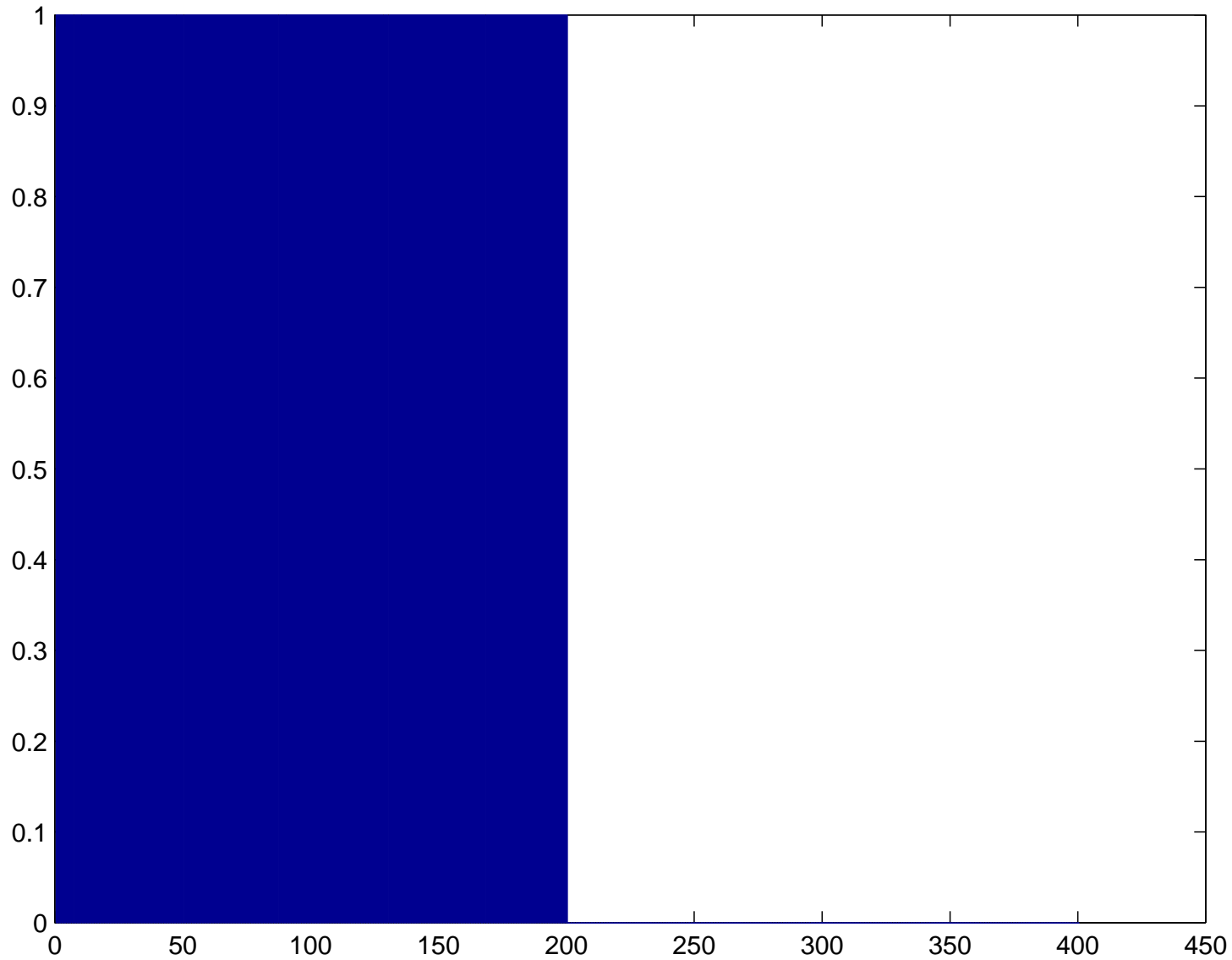


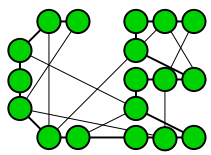
Max-Cut



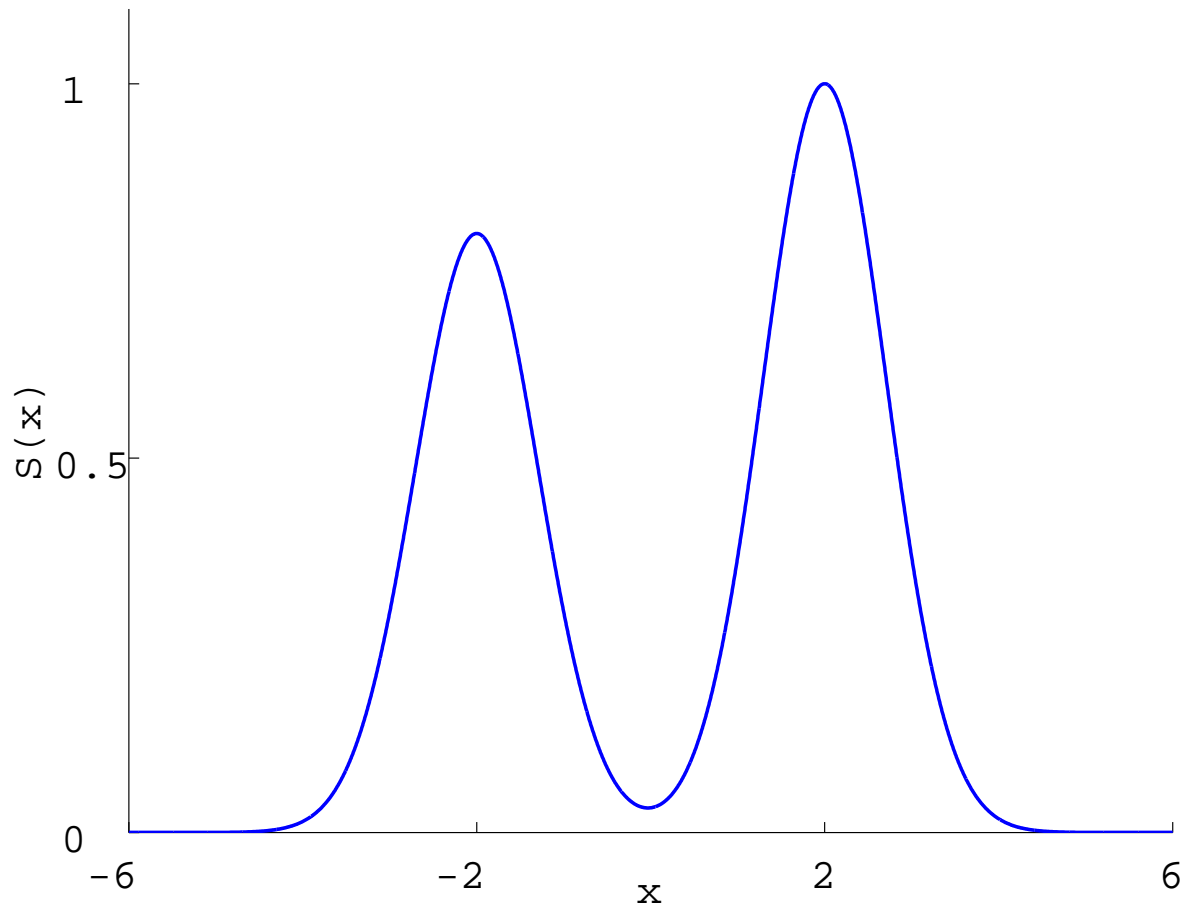


Max-Cut

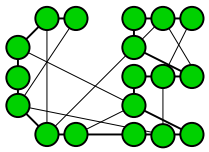




Example: Continuous Optimization

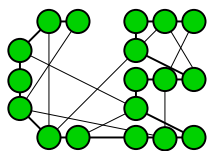


$$S(x) = e^{-(x-2)^2} + 0.8e^{-(x+2)^2}$$

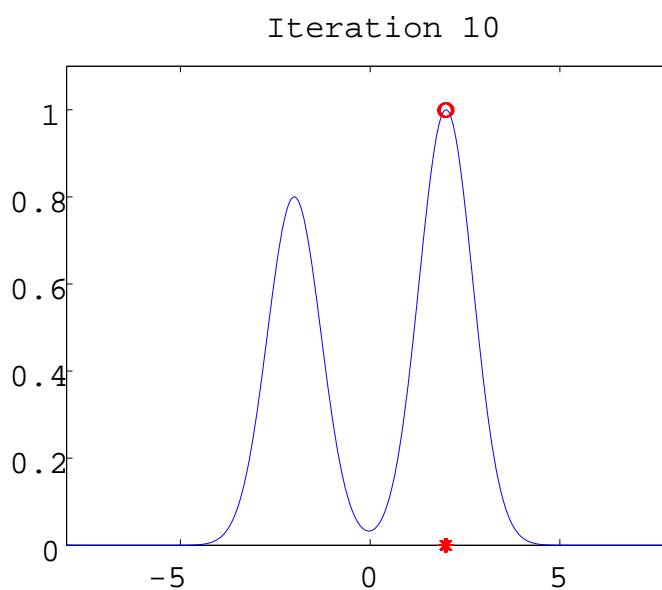
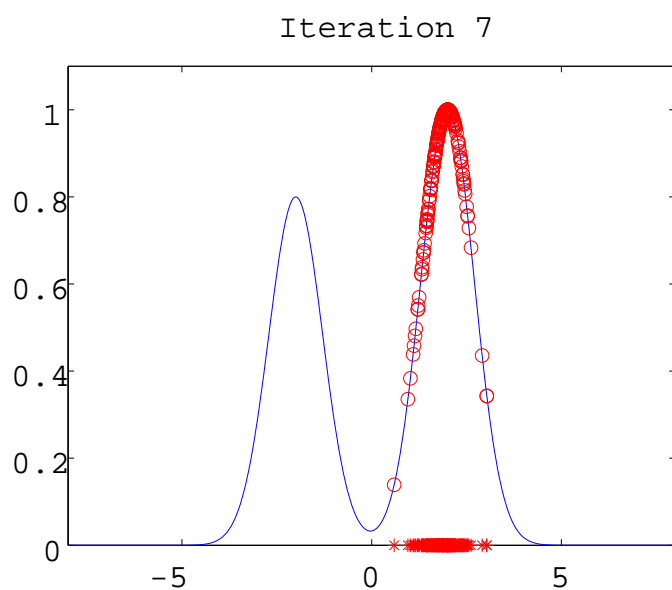
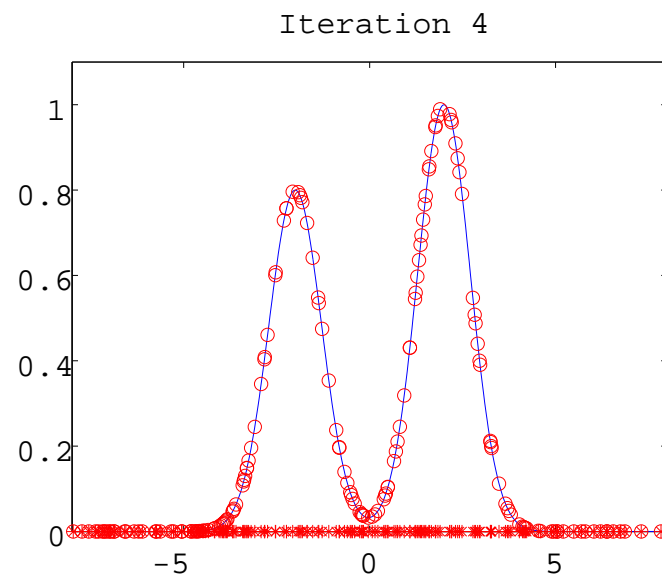
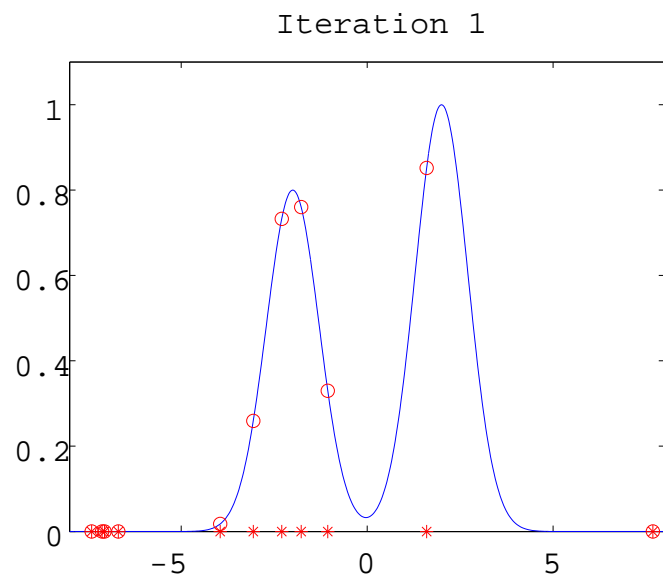


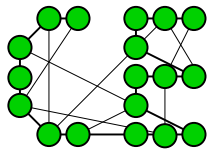
Matlab Program

```
S = inline('exp(-(x-2).^2) + 0.8*exp(-(x+2).^2)');
mu = -10; sigma = 10; rho = 0.1; N = 100; eps = 1E-3;
t=0; % iteration counter
while sigma > eps
    t = t+1;
    x = mu + sigma*randn(N,1);
    SX = S(x); % Compute the performance.
    sortSX = sortrows([x SX],2);
    mu = mean(sortSX((1-rho)*N:N,1));
    sigma = std(sortSX((1-rho)*N:N,1));
    fprintf('%g %6.9f %6.9f %6.9f \n', t, S(mu),mu, sigma)
```



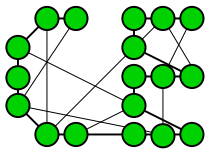
Numerical Result





Cross-Entropy: Some Theory

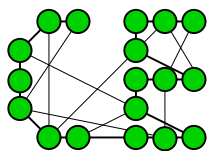
Estimate $\ell := \mathbb{P}_{\mathbf{u}}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}_{\mathbf{u}} I_{\{S(\mathbf{X}) \geq \gamma\}}$ via



Cross-Entropy: Some Theory

Estimate $\ell := \mathbb{P}_{\mathbf{u}}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}_{\mathbf{u}} I_{\{S(\mathbf{X}) \geq \gamma\}}$ via

$$\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} \frac{f(\mathbf{X}_i; \mathbf{u})}{g(\mathbf{X}_i)} .$$



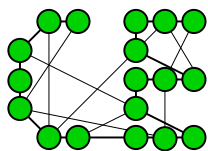
Cross-Entropy: Some Theory

Estimate $\ell := \mathbb{P}_{\mathbf{u}}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}_{\mathbf{u}} I_{\{S(\mathbf{X}) \geq \gamma\}}$ via

$$\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \gamma\}} \frac{f(\mathbf{X}_i; \mathbf{u})}{g(\mathbf{X}_i)} .$$

The best density (zero variance estimator!) is

$$g^*(\mathbf{x}) := \frac{I_{\{S(\mathbf{x}) \geq \gamma\}} f(\mathbf{x}; \mathbf{u})}{\ell} .$$



Cross-Entropy: Some Theory

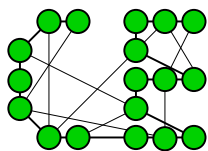
Estimate $\ell := \mathbb{P}_{\mathbf{u}}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}_{\mathbf{u}} I_{\{S(\mathbf{X}) \geq \gamma\}}$ via

$$\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} \frac{f(\mathbf{X}_i; \mathbf{u})}{g(\mathbf{X}_i)} .$$

The best density (zero variance estimator!) is

$$g^*(\mathbf{x}) := \frac{I_{\{S(\mathbf{x}) \geq \gamma\}} f(\mathbf{x}; \mathbf{u})}{\ell} .$$

Problem: g^* depends on the unknown ℓ .



Cross-Entropy: Some Theory

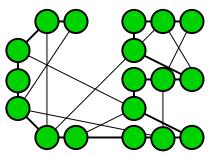
Estimate $\ell := \mathbb{P}_{\mathbf{u}}(S(\mathbf{X}) \geq \gamma) = \mathbb{E}_{\mathbf{u}} I_{\{S(\mathbf{X}) \geq \gamma\}}$ via

$$\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} \frac{f(\mathbf{X}_i; \mathbf{u})}{g(\mathbf{X}_i)} .$$

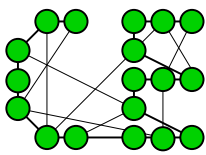
The best density (zero variance estimator!) is

$$g^*(\mathbf{x}) := \frac{I_{\{S(\mathbf{x}) \geq \gamma\}} f(\mathbf{x}; \mathbf{u})}{\ell} .$$

Problem: g^* depends on the unknown ℓ .



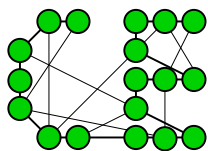
Idea: choose $g = f(\cdot; \boldsymbol{v})$ such that the “distance” between the densities g^* and $f(\cdot; \boldsymbol{v})$ is minimal.



Idea: choose $g = f(\cdot; \mathbf{v})$ such that the “distance” between the densities g^* and $f(\cdot; \mathbf{v})$ is minimal.

The **Kullback-Leibler** or **cross-entropy** distance is defined as:

$$\begin{aligned}\mathcal{D}(g, h) &= \mathbb{E}_g \log \frac{g(\mathbf{X})}{h(\mathbf{X})} \\ &= \int g(\mathbf{x}) \log g(\mathbf{x}) d\mathbf{x} - \int g(\mathbf{x}) \log h(\mathbf{x}) d\mathbf{x} .\end{aligned}$$

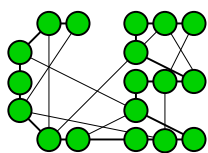


Idea: choose $g = f(\cdot; \mathbf{v})$ such that the “distance” between the densities g^* and $f(\cdot; \mathbf{v})$ is minimal.

The **Kullback-Leibler** or **cross-entropy** distance is defined as:

$$\begin{aligned}\mathcal{D}(g, h) &= \mathbb{E}_g \log \frac{g(\mathbf{X})}{h(\mathbf{X})} \\ &= \int g(\mathbf{x}) \log g(\mathbf{x}) d\mathbf{x} - \int g(\mathbf{x}) \log h(\mathbf{x}) d\mathbf{x} .\end{aligned}$$

Determine the optimal \mathbf{v}^* from $\min_{\mathbf{v}} \mathcal{D}(g^*, f(\cdot; \mathbf{v}))$.



This is equivalent to solving

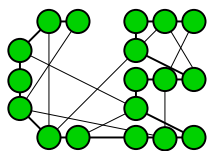
$$\max_{\boldsymbol{v}} \mathbb{E}_{\boldsymbol{u}} I_{\{S(\boldsymbol{X}) \geq \gamma\}} \log f(\boldsymbol{X}; \boldsymbol{v}) .$$

Using again IS, we can rewrite this as

$$\max_{\boldsymbol{v}} \mathbb{E}_{\boldsymbol{w}} I_{\{S(\boldsymbol{X}) \geq \gamma\}} W(\boldsymbol{X}; \boldsymbol{u}, \boldsymbol{w}) \log f(\boldsymbol{X}; \boldsymbol{v}),$$

for *any* reference parameter \boldsymbol{w} , where

$$W(\boldsymbol{x}; \boldsymbol{u}, \boldsymbol{w}) = \frac{f(\boldsymbol{x}; \boldsymbol{u})}{f(\boldsymbol{x}; \boldsymbol{w})}$$



We may *estimate* the optimal solution \mathbf{v}^* by solving the following stochastic counterpart:

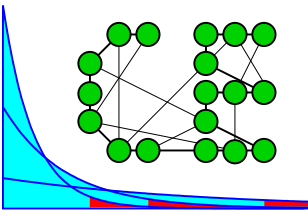
$$\max_{\mathbf{v}} \frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} W(\mathbf{X}_i; \mathbf{u}, \mathbf{w}) \log f(\mathbf{X}_i; \mathbf{v}) ,$$

where $\mathbf{X}_1, \dots, \mathbf{X}_N$ is a random sample from $f(\cdot; \mathbf{w})$.

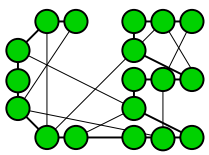
Alternatively, solve:

$$\frac{1}{N} \sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} W(\mathbf{X}_i; \mathbf{u}, \mathbf{w}) \nabla \log f(\mathbf{X}_i; \mathbf{v}) = \mathbf{0},$$

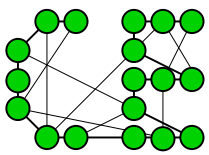
where the gradient is with respect to \mathbf{v} .



- The solution to the CE program can often be calculated **analytically**.

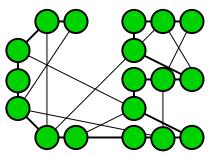


- The solution to the CE program can often be calculated **analytically**.
- Note that the CE program is useful only when under w the event $\{S(\mathbf{X}) \geq \gamma\}$ is **not too rare**, say $\geq 10^{-5}$.



- The solution to the CE program can often be calculated **analytically**.
- Note that the CE program is useful only when under w the event $\{S(\mathbf{X}) \geq \gamma\}$ is **not too rare**, say $\geq 10^{-5}$.

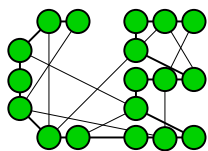
Question: how to choose w so that this is indeed the case?



- The solution to the CE program can often be calculated **analytically**.
- Note that the CE program is useful only when under w the event $\{S(\mathbf{X}) \geq \gamma\}$ is **not too rare**, say $\geq 10^{-5}$.

Question: how to choose w so that this is indeed the case?

Answer: use a multi-level approach.

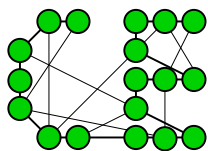


- The solution to the CE program can often be calculated **analytically**.
- Note that the CE program is useful only when under w the event $\{S(\mathbf{X}) \geq \gamma\}$ is **not too rare**, say $\geq 10^{-5}$.

Question: how to choose w so that this is indeed the case?

Answer: use a multi-level approach.

Introduce a sequence of reference parameters $\{\mathbf{v}_t, t \geq 0\}$ and a sequence of levels $\{\gamma_t, t \geq 1\}$, and iterate in both γ_t and \mathbf{v}_t .



Toy example 1 (continued)

Recall

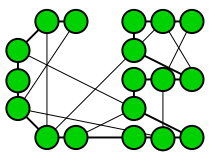
$$f(\mathbf{x}; \mathbf{v}) = \exp \left(- \sum_{j=1}^5 \frac{x_j}{v_j} \right) \prod_{j=1}^5 \frac{1}{v_j}.$$

The optimal \mathbf{v} follows from the system of equations

$$\sum_{i=1}^N I_{\{S(\mathbf{X}_i) \geq \gamma\}} W(\mathbf{X}_i; \mathbf{u}, \mathbf{w}) \nabla \log f(\mathbf{X}_i; \mathbf{v}) = \mathbf{0}.$$

Since

$$\frac{\partial}{\partial v_j} \log f(\mathbf{x}; \mathbf{v}) = \frac{x_j}{v_j^2} - \frac{1}{v_j},$$



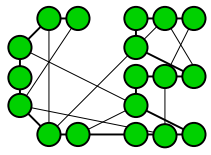
we have for the j th equation

$$\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} W(\mathbf{X}_i; \mathbf{u}, \mathbf{w}) \left(\frac{X_{ij}}{v_j^2} - \frac{1}{v_j} \right) = 0 ,$$

whence,

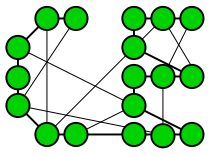
$$v_j = \frac{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} W(\mathbf{X}_i; \mathbf{u}, \mathbf{w}) X_{ij}}{\sum_{i=1}^N I_{\{S(\mathbf{x}_i) \geq \gamma\}} W(\mathbf{X}_i; \mathbf{u}, \mathbf{w})} ,$$

which leads to the updating formula in step 3 of the Algorithm.



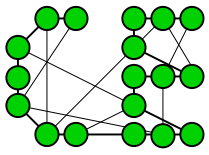
Further research

- Multi-extremal constrained continuous optimization.



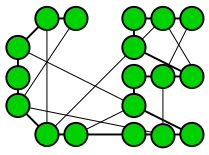
Further research

- Multi-extremal constrained continuous optimization.
- Noisy optimization.



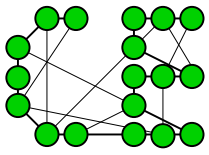
Further research

- Multi-extremal constrained continuous optimization.
- Noisy optimization.
- Simulation with heavy tail distributions.



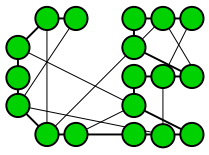
Further research

- Multi-extremal constrained continuous optimization.
- Noisy optimization.
- Simulation with heavy tail distributions.
- Incorporating MaxEnt (MinxEnt), e.g. MCE



Further research

- Multi-extremal constrained continuous optimization.
- Noisy optimization.
- Simulation with heavy tail distributions.
- Incorporating MaxEnt (MinxEnt), e.g. MCE
- Multi-actor games



Further research

- Multi-extremal constrained continuous optimization.
- Noisy optimization.
- Simulation with heavy tail distributions.
- Incorporating MaxEnt (MinxEnt), e.g. MCE
- Multi-actor games
- Convergence of CE algorithm.