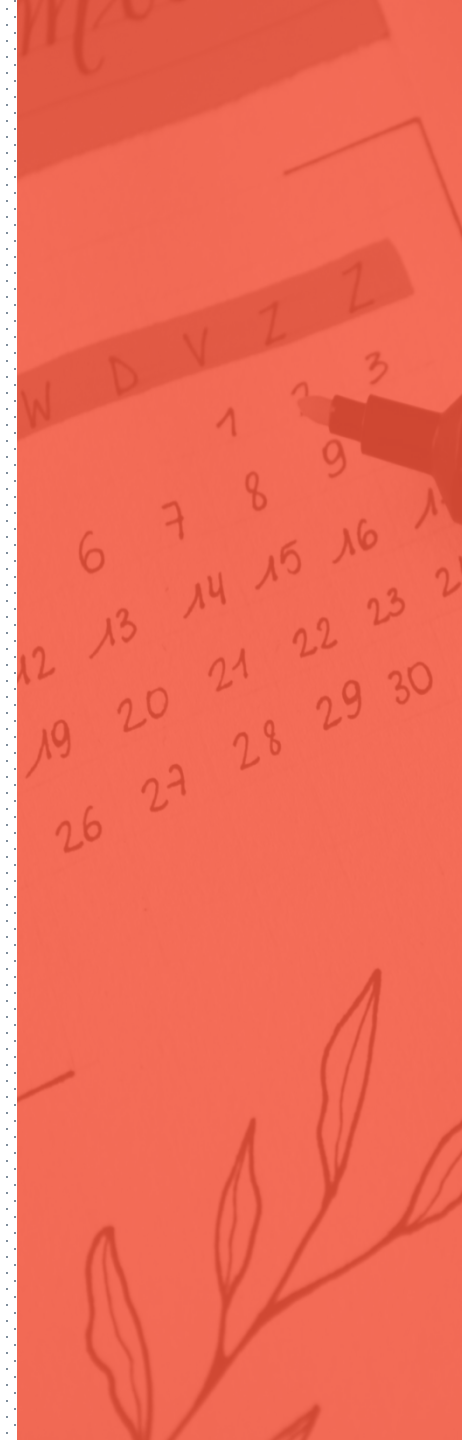


Implementation study of DISCOPOLIS

an algorithm for uniform sampling of metabolic flux distributions via iterative sequences of linear programs

Name: Hongxing NIU

Promotor: Prof. Philippe Bogaerts

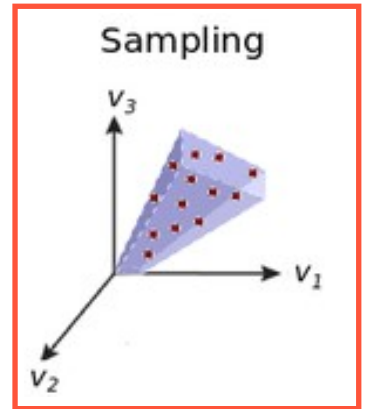
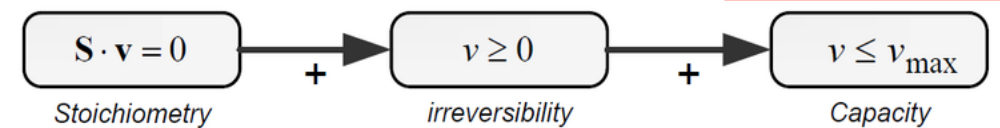
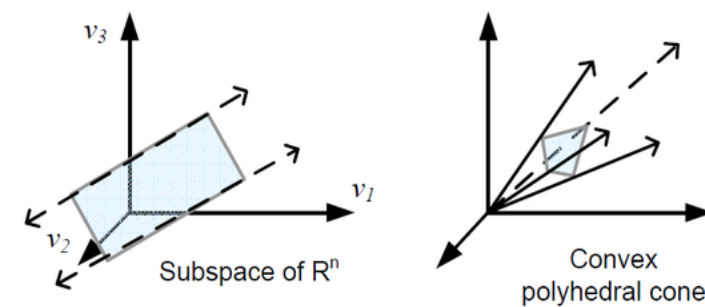
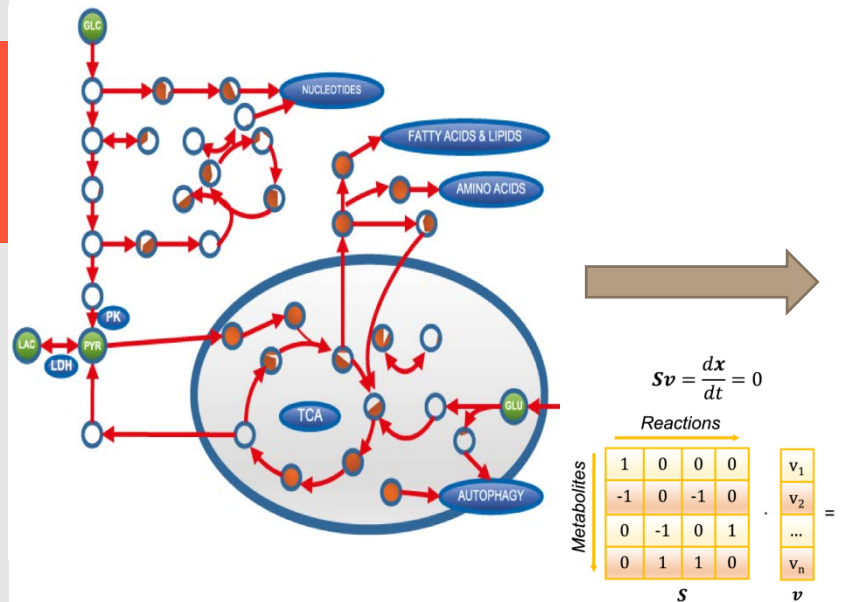


We are interested in metabolic networks, most of which are constrained under-determined systems



Problem

Solving constrained under-determined systems by random sampling



(source: http://2014.igem.org/Team:Valencia_UPV/Modeling/fba)

Approaches to the under-determined system by random sampling

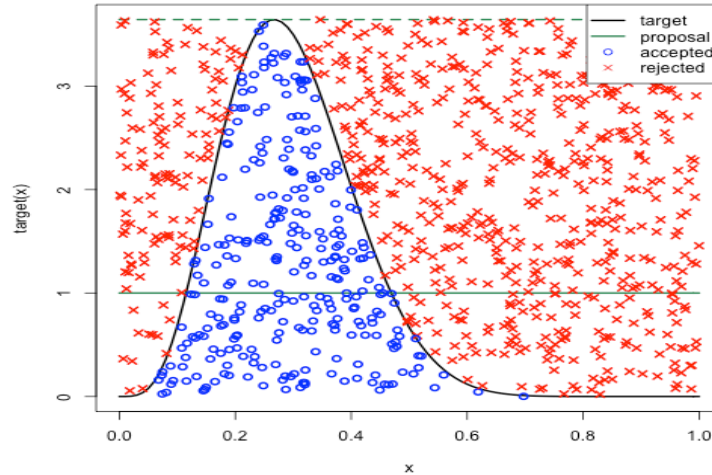


Problem

Solving constrained under-determined systems by random sampling

(1). Acceptance-Rejection

High computation load under high-dimen



(3). DISCOPOLIS

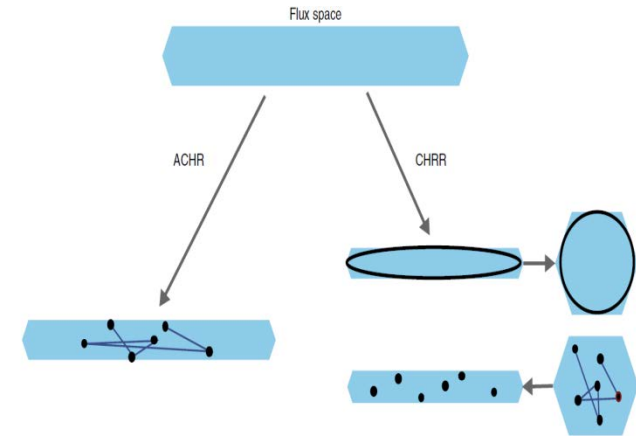
Discrete Sampling of CONvex POLytope via Linear program Iterative Sequences

(Bogaerts and Rومان, 2019)

(2). Hit-and-Run (MCMC, Metropolis-Hastings)

Getting stuck under high-dimen

(Schellendberger, 2011; Haraldsdottir, 2017)



Approaches to the under-determined system by random sampling

(3). DISCOPOLIS : Discrete Sampling of CONvex POLytope via Linear program Iterative Sequences (Bogaerts and Rooman, 2019).



Problem

Solving constrained under-determined systems by random sampling

Input : solution polytope defined by A and b ; **number of samples N** ; **number of grid points S** ; minimum and maximum values of the fluxes v_i^{MIN} and v_i^{MAX} ($i \in [1, n]$) obtained with Flux Variability Analysis
Output : N samples $v(k) \in \mathbb{R}^n$ ($k \in [1, N]$) with their weights $w(k)$

```

1   $A_{eq} = \emptyset$ ;  $b_{eq} = \emptyset$ ; /* initialize empty matrices for equality constraints
2   $L_i = (v_i^{MAX} - v_i^{MIN}) / (S - 1)$ ; /* compute for each flux  $v_i$  the interval
   between 2 grid points
3  for  $k = 1$  to  $N$  do
4     $w(k) = 1$ ; /* initialize weight of the  $k^{th}$  sample
5     $I = [1, n]$ ; /* set of all indexes  $i$  of all the fluxes  $v_i \in v$ 
6    Randomly select an index  $i$  in  $I$ ;
7    Remove index  $i$  from set  $I$ ;
8    Generate one index  $g$  from a uniform distribution on  $[1, S]$ ;
9     $v_i = v_i^{MIN} + (v_i^{MAX} - v_i^{MIN}) * (g - 1) / (S - 1)$ ; /* discrete uniform
   sampling of  $v_i$  corresponding to the  $g^{th}$  grid point
10   while  $I \neq \emptyset$  do
11     Augment  $A_{eq}$  and  $b_{eq}$  to account for last fixed  $v_i$ ;
12     Randomly select an index  $i$  in  $I$ ;
13     Remove index  $i$  from set  $I$ ;
14      $v_i^{MINnew} = \min_v v_i$  computed with LP subject to  $A * v \leq b$ 
   and  $A_{eq} * v = b_{eq}$ ;
15      $v_i^{MAXnew} = \max_v v_i$  computed with LP subject to  $A * v \leq b$ 
   and  $A_{eq} * v = b_{eq}$ ;
16      $S^{new} = 1 + \text{floor}((v_i^{MAXnew} - v_i^{MINnew}) / L_i)$ ; /* number of grid
   points remaining in the new constrained solution interval
17     if  $S^{new} > 1$  then
18       Generate one index  $g$  from a uniform distribution on  $[1, S^{new}]$ ;
19        $v_i = v_i^{MINnew} + (v_i^{MAXnew} - v_i^{MINnew}) * (g - 1) / (S^{new} - 1)$ ; /* discrete
   uniform sampling of  $v_i$  corresponding to the  $g^{th}$  grid point
20     else
21        $v_i = (v_i^{MAXnew} + v_i^{MINnew}) / 2$ ; /* use of the center of the new
   solution interval in case of only 1 remaining grid point
22     end
23      $w(k) = w(k) * S^{new} / S$ ; /* update weight of the  $k^{th}$  sample
24   end
25 end

```

→ N iterations ($nSample$)

→ Discretized with S grid points ($nGrid$)

Iteratively, discrete sampling of convex polytope over which the objective function is optimized:

$$v_{Min,Max}(i) = Min,Max[v(i)] \quad \forall i \in [1, n]$$



Problem

Solving constrained under-determined systems by random sampling



Objectives & Methods

Investigating the DISCOPOLIS algorithm and optimizing the parameter settings:
 1) monitoring convergence;
 2) getting a “fit” solution distribution by appropriate discretization.

1) Monitoring convergence of averages ($nSample$)

----Univariate Gelman and Rubin diagnostic by potential scale reduction factor (PSRF)

$$PSRF = \hat{R} = \sqrt{\frac{\widehat{var}(X)}{W}},$$

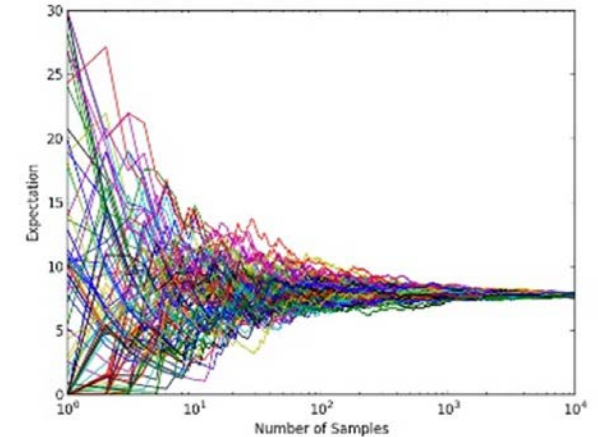
with

$$\widehat{var}(X) = \frac{N-1}{N}W + \frac{1}{N}B$$

----Multivariate extension (MPSRF)

$$MPSRF = \hat{R}^n = \sqrt{\frac{N-1}{N} + \frac{\lambda_{\max}(W^{-1}B)}{N}},$$

$\lambda_{\max}()$ is the largest eigenvalue



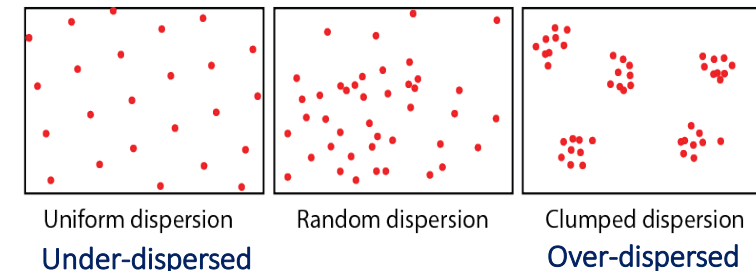
2) Choosing an appropriate number of grid points ($nGrid$)

----Generalized variance (the determinant of covariance matrix S of fluxes)

----Total sample variance ($trace(S)$)

----P (99.9% weight)

(Bogaerts and Rooman, 2019 ,
 the percentage of samples whose sum
 of weights =99.9% total sum of weights)





Problem

Solving constrained under-determined systems by random sampling



Objectives & Methods

Investigating the DISCOPOLIS algorithm and optimizing the parameter settings:
1)convergence;
2)discretized sampling

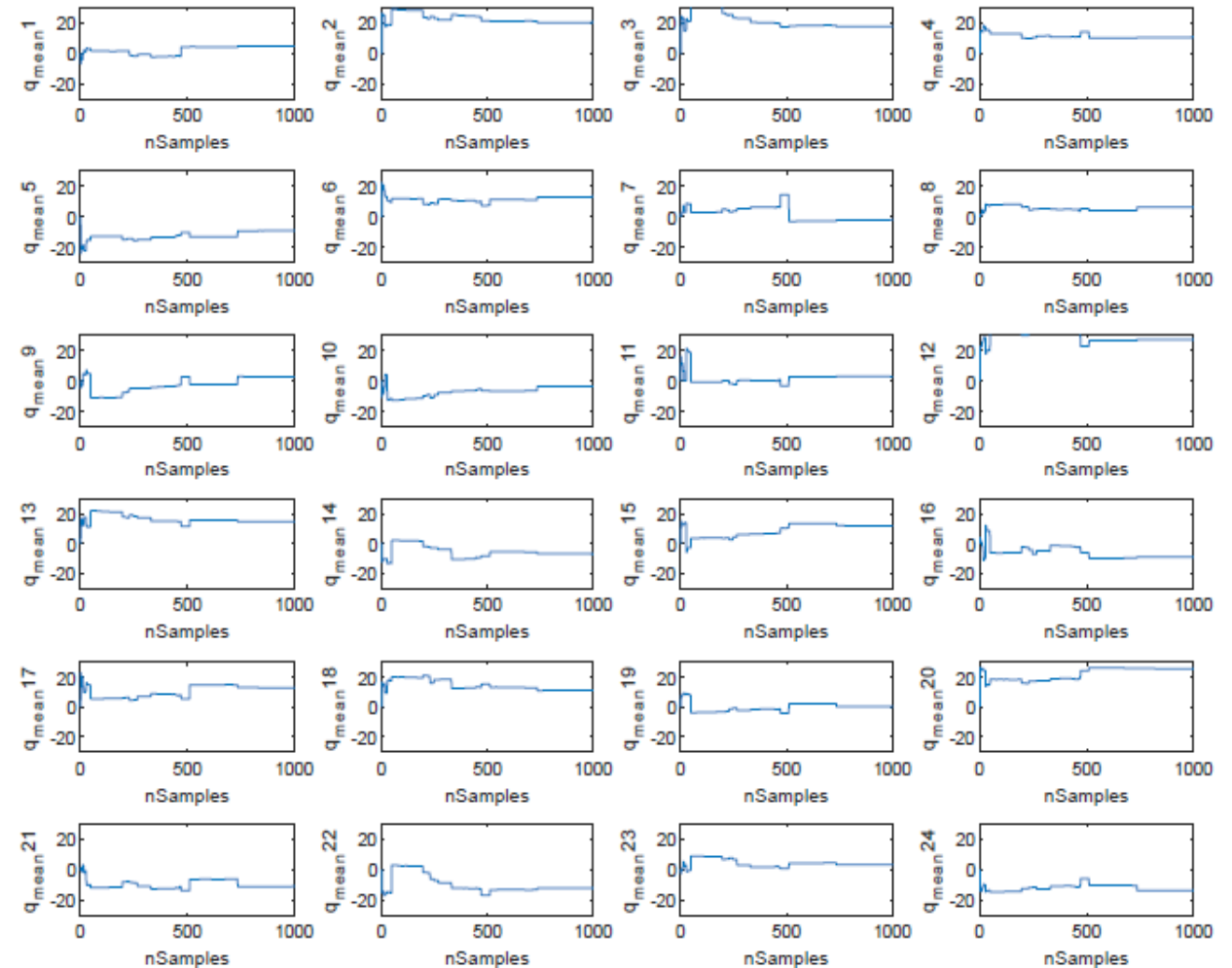


Case Studies

Toy example (extra slides);

Core metabolic network of *Escherichia coli*.

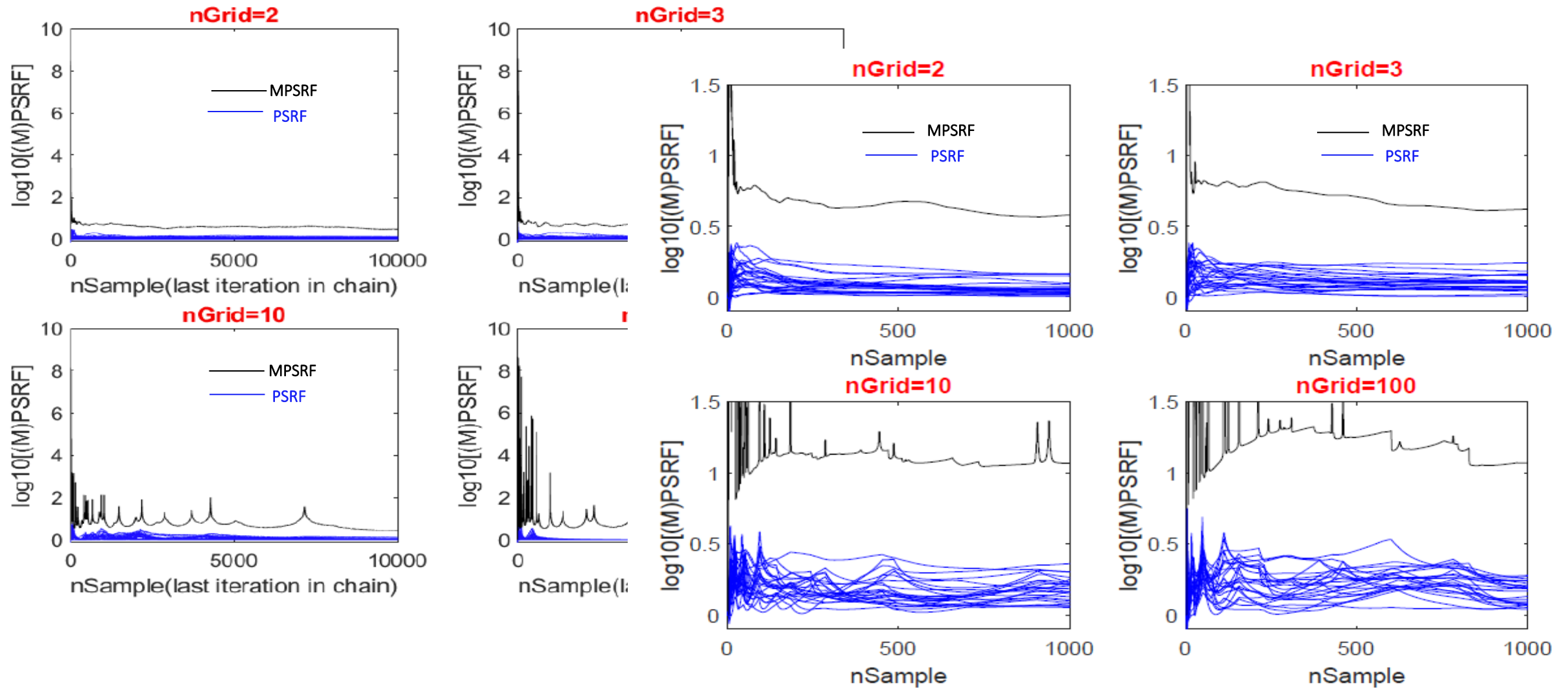
Core metabolic network of *E.coli*:
after elimination of equality constraints,
 $A'q \leq b'$, where $A' \in \mathcal{R}^{172 \times 24}$



The flux means over *nSample* by the algorithm with *nGrid*=10

Convergence monitoring of means with changing *nGrid* w.r.t *nSample*.

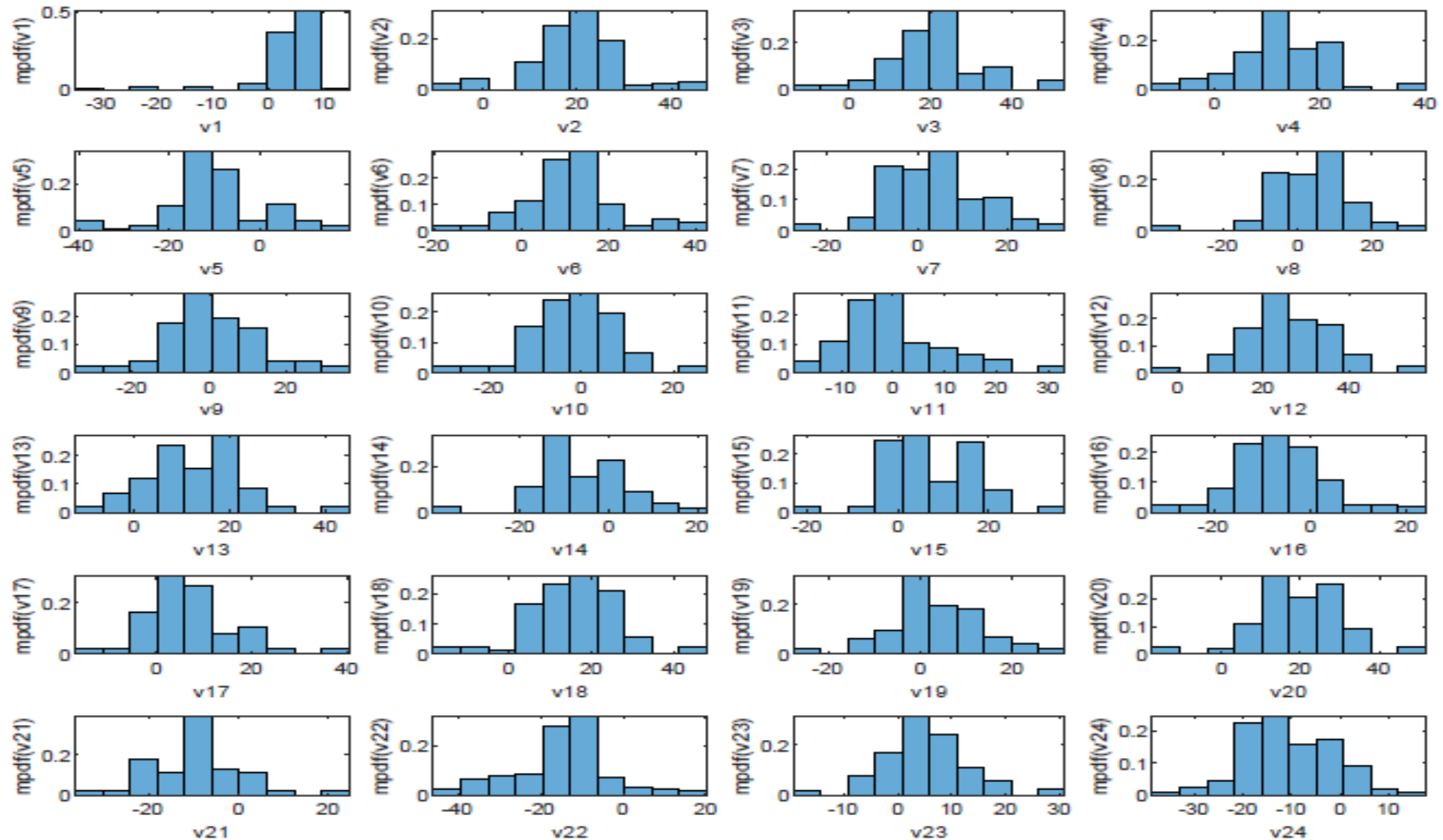
10 random seeds



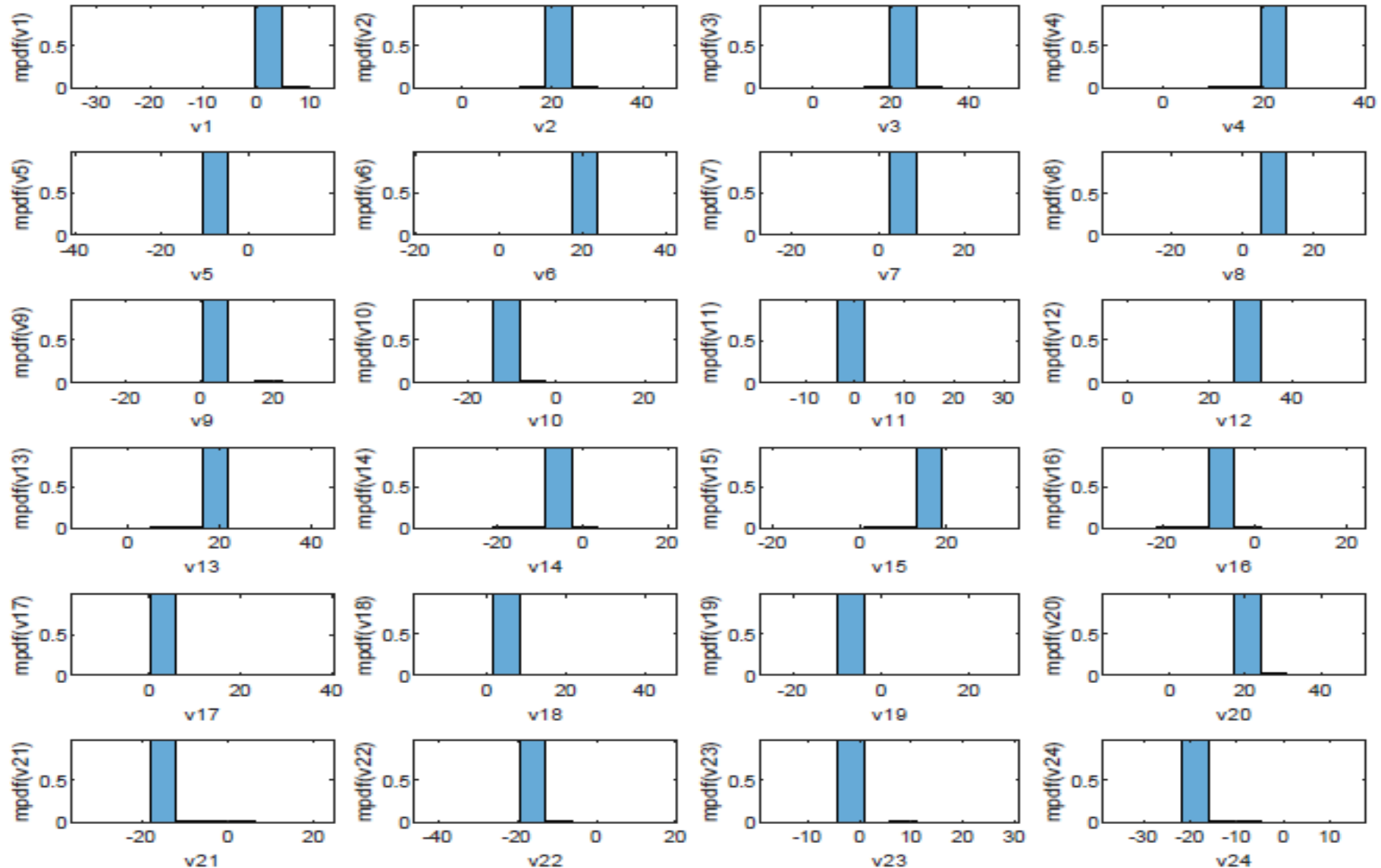
MPSRF >> *PSRF* in high-dimen, thus *MPSRF* is a conservative termination criterion.

Marginal distributions of the fluxes by the DISCOPOLIS algorithm, $nSample=10,000$.

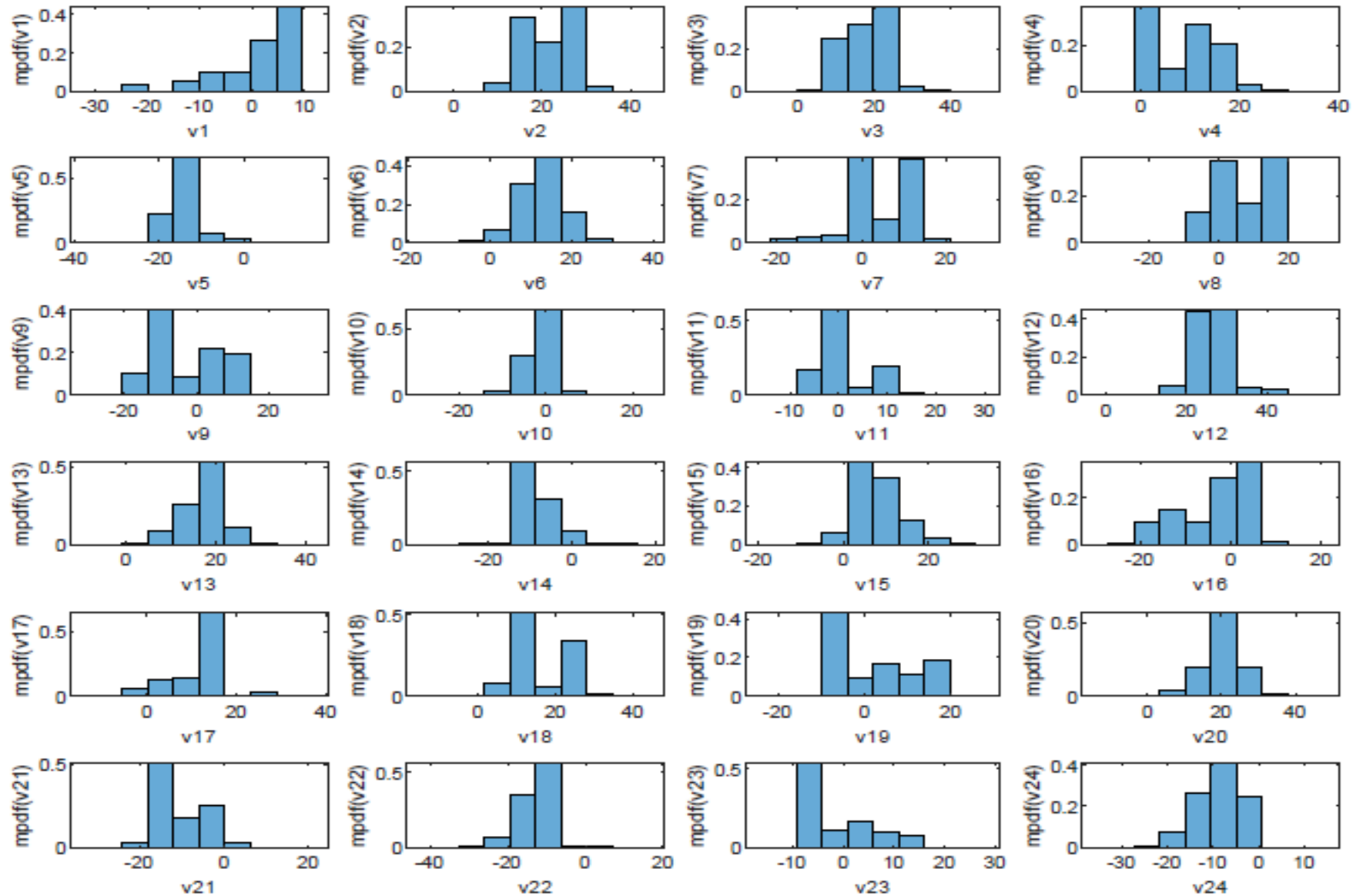
(a) $nGrid=2$, the solutions could move to the tails



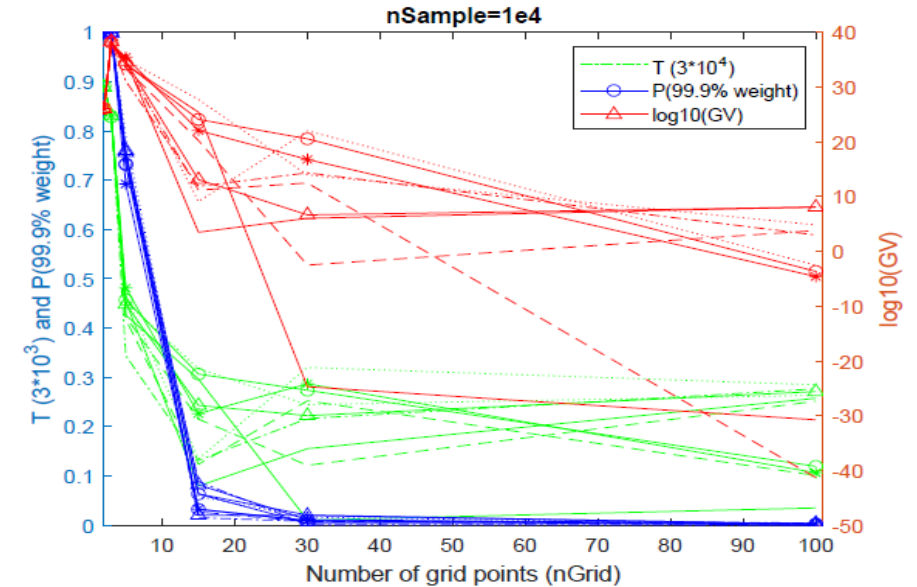
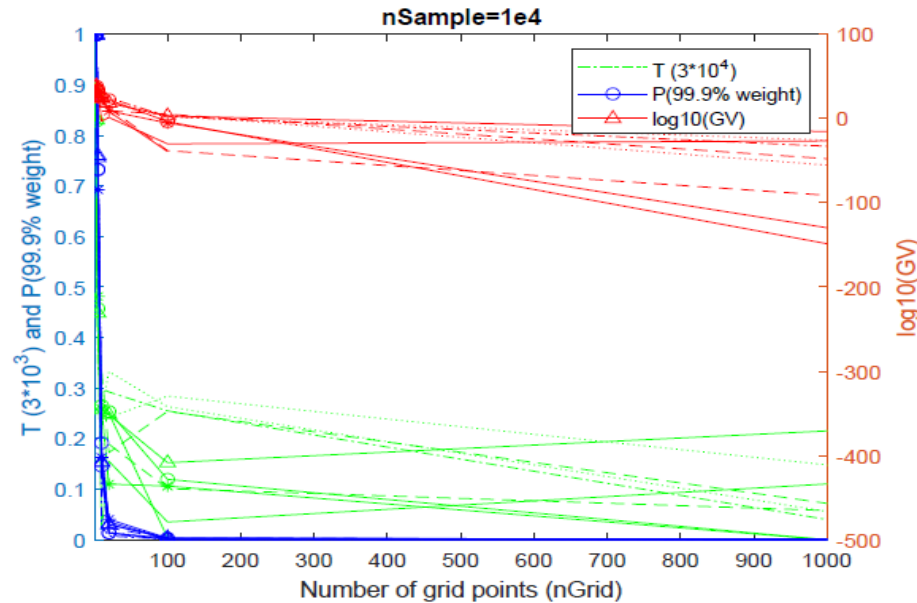
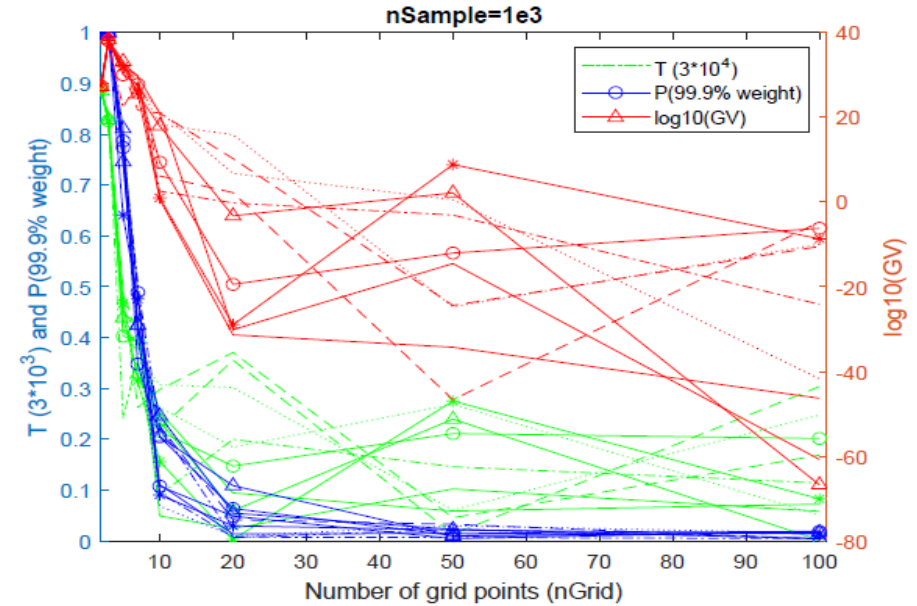
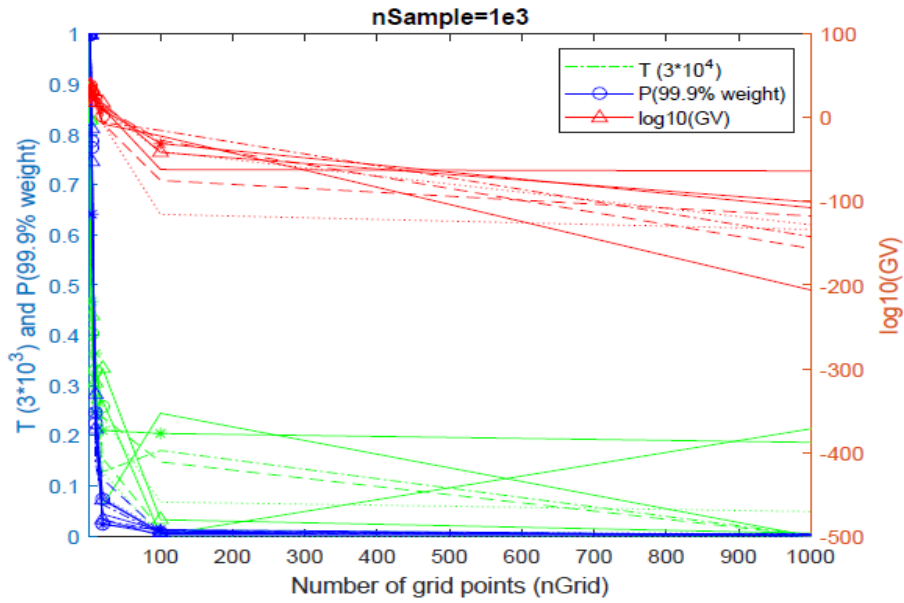
Marginal distributions of the fluxes by the DISCOPOLIS algorithm, $nSample=10,000$.
(b) $nGrid=1000$, under-dispersion, more regular than random







Marginal distributions of the fluxes by the DISCOPOLIS algorithm, $nSample=10,000$.
(c) $nGrid=10$



GV , T , and P (99.9% weight) are calculated and plotted against $nGrid$. 10 random seeds are taken under each scenario.



Summary

			
Problem	Objectives & Methods	Case Studies	Conclusions
Solving constrained under-determined systems by random sampling	Investigating the DISCOPOLIS algorithm and optimizing the parameter settings: 1)convergence; 2)discretized sampling	1)Toy example; 2)Core metabolic network of <i>Escherichia coli</i>	The configuration of DISCOPOLIS algorithm, <i>i.e.</i> , $nSamples + nGrid$ were tuned to improve its performance

1) Gelman and Rubin diagnostic: the trend of $PSRF/MPSRF$ from the \bar{v}^T tells if and when the results converge;

2) Two measures (GV and T) quantify the statistical dispersion of solution cloud: in the two cases the optimal $nGrid=10\sim 20$ with a neither over-dispersed nor under-dispersed flux distribution.



THANK
YOU

Extra Slides

Structure Overview



Problem

Solving constrained under-determined systems by random sampling



Objectives & Methods

Investigating the DISCOPOLIS algorithm and optimizing the parameter settings:
1)convergence;
2)discretized sampling



Case Studies

Toy example (extra slides);

Core metabolic network of *Escherichia coli*.



Conclusions

The configuration of DISCOPOLIS algorithm, *i.e.*, $nSamples + nGrid$ were tuned to improve its performance

Structure Overview



Problem

Solving constrained under-determined systems by random sampling



Objectives & Methods

Investigating the DISCOPOLIS algorithm and optimizing the parameter settings:
1)convergence;
2)discretized sampling



Case Studies

Toy example (extra slides);

Core metabolic network of *Escherichia coli*.



Conclusions

The configuration of DISCOPOLIS algorithm, *i.e.*, $nSamples + nGrid$ were tuned to improve its performance



Problem

Solving constrained under-determined systems by random sampling



Objectives & Methods

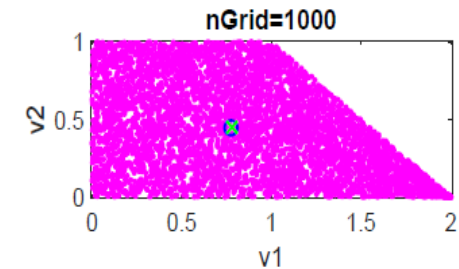
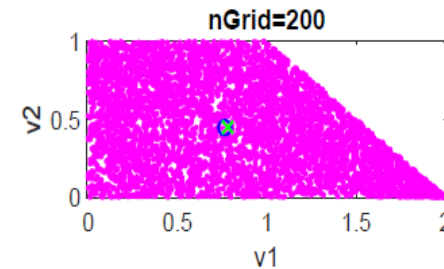
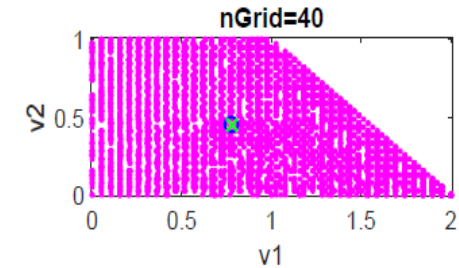
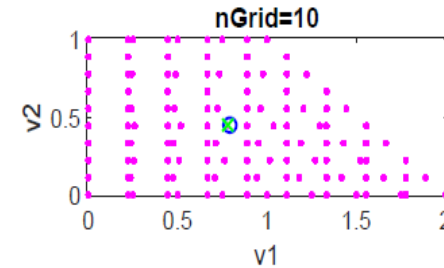
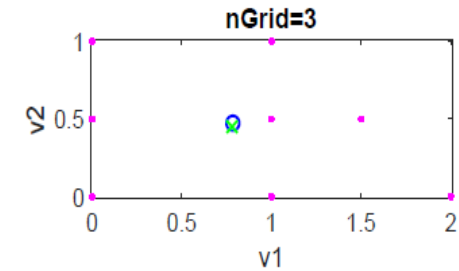
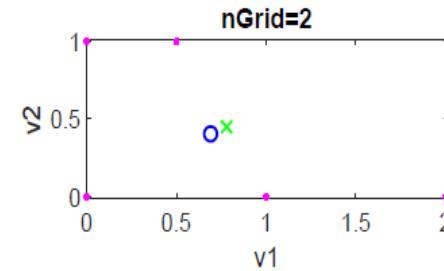
Investigating the DISCOPOLIS algorithm and optimizing the parameter settings:
1)convergence;
2)discretized sampling



Case Studies

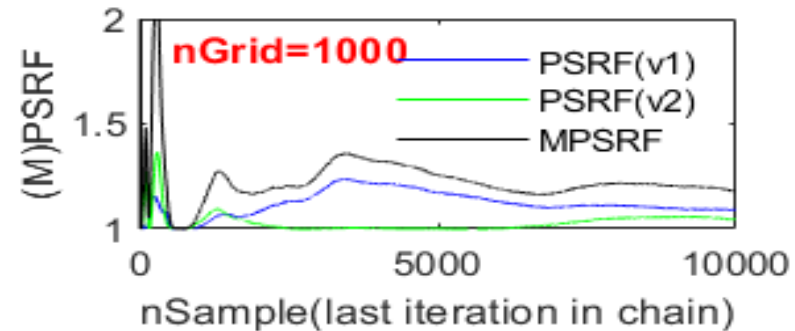
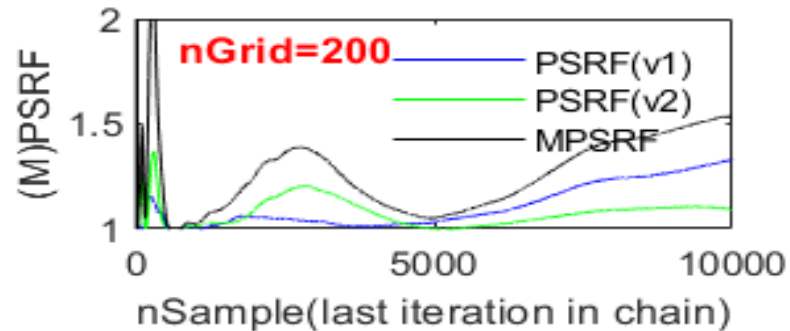
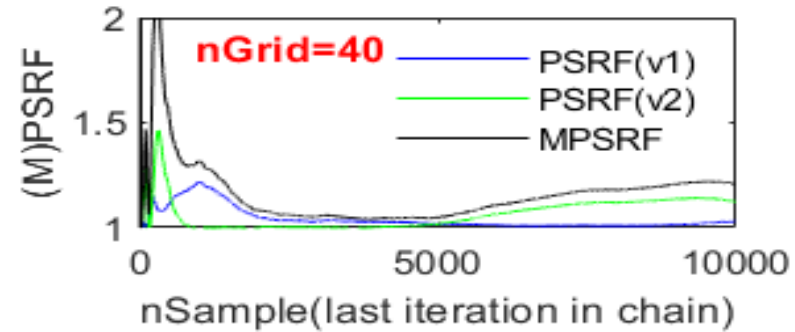
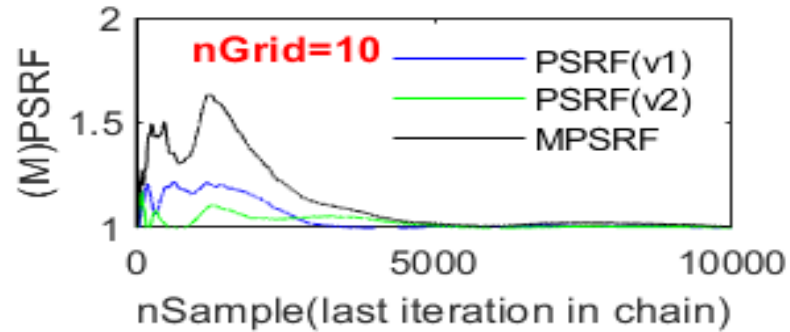
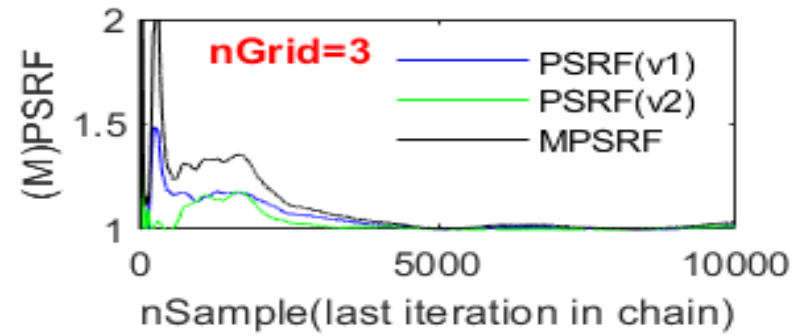
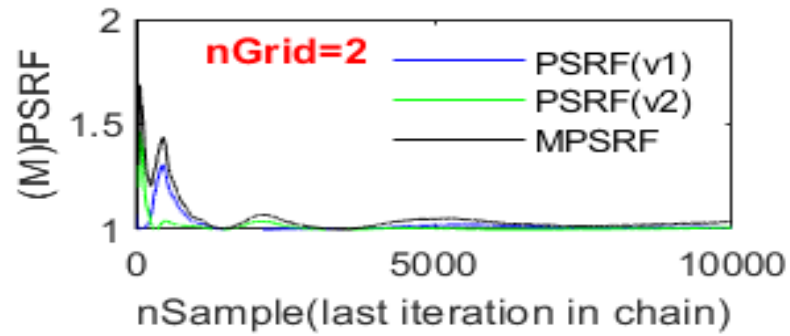
Toy example;
Core metabolic network of *Escherichia coli*

1) Toy example: $A = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}, b = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 2 \end{bmatrix}$

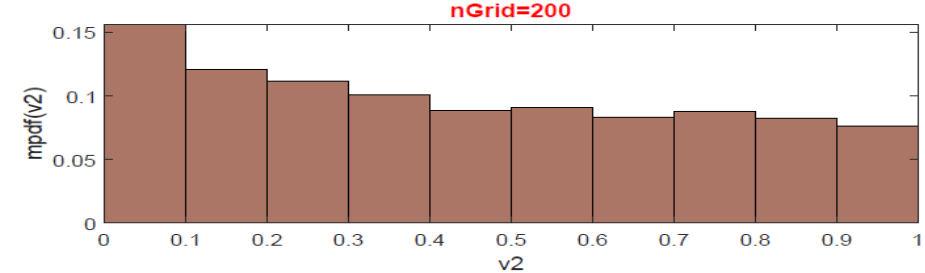
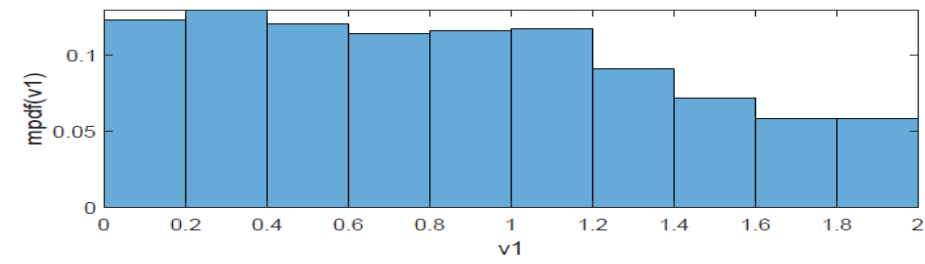
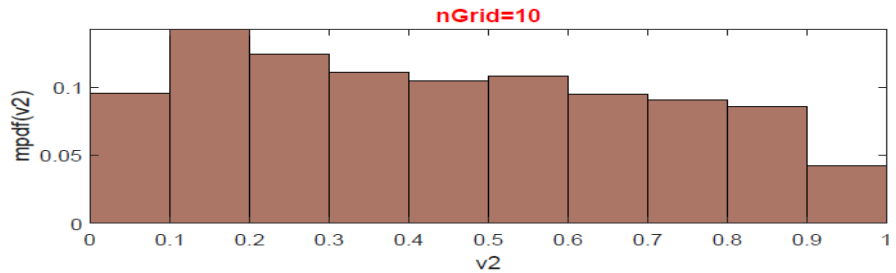
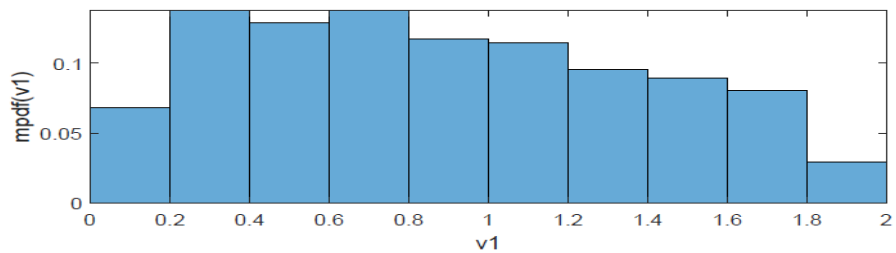
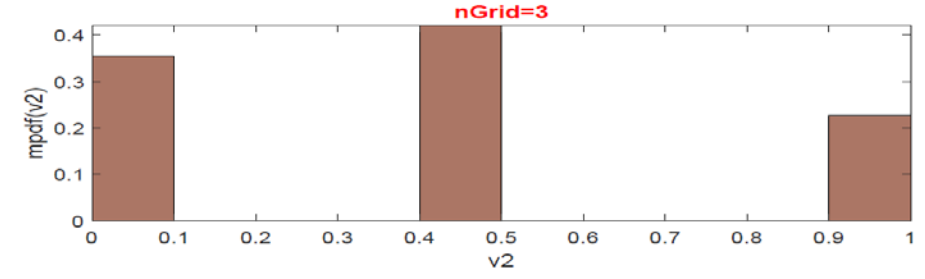
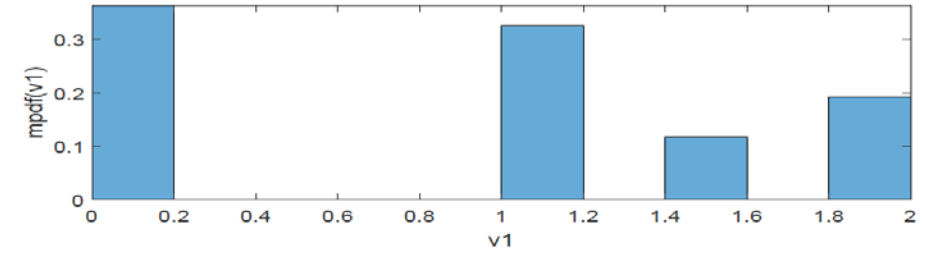
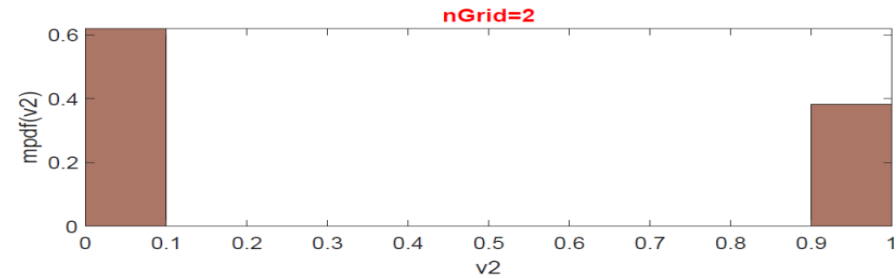
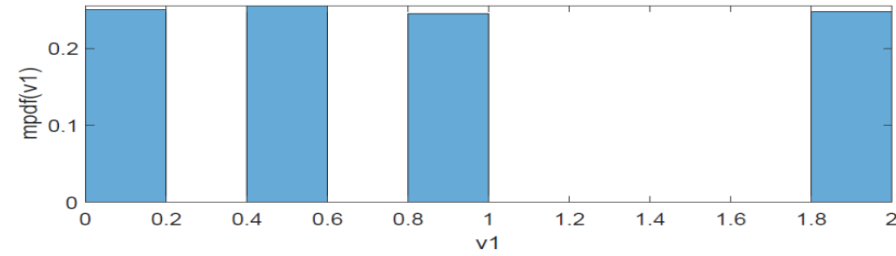


Solutions (dots) by DISCOPOLIS, $nSample=5,000$, \bigcirc =the means;
 $\vec{v}^T = [0.78 \ 0.45]$, by the rejection algorithm shown by \times .

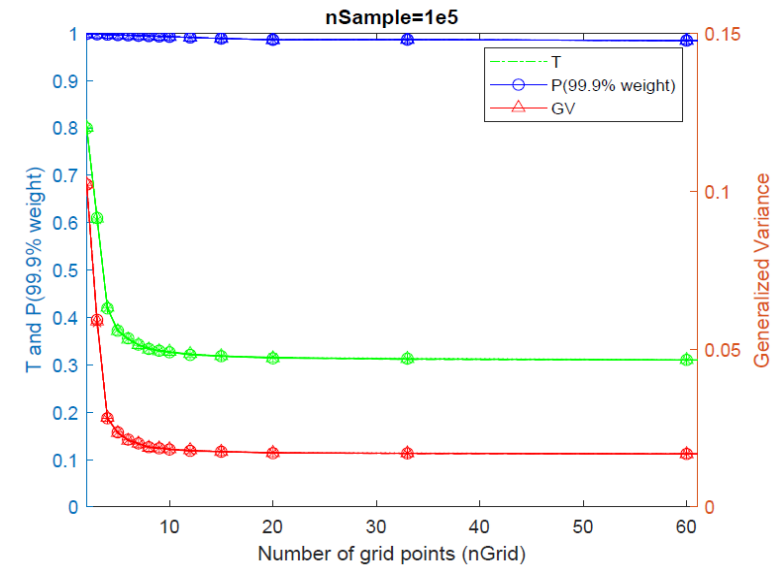
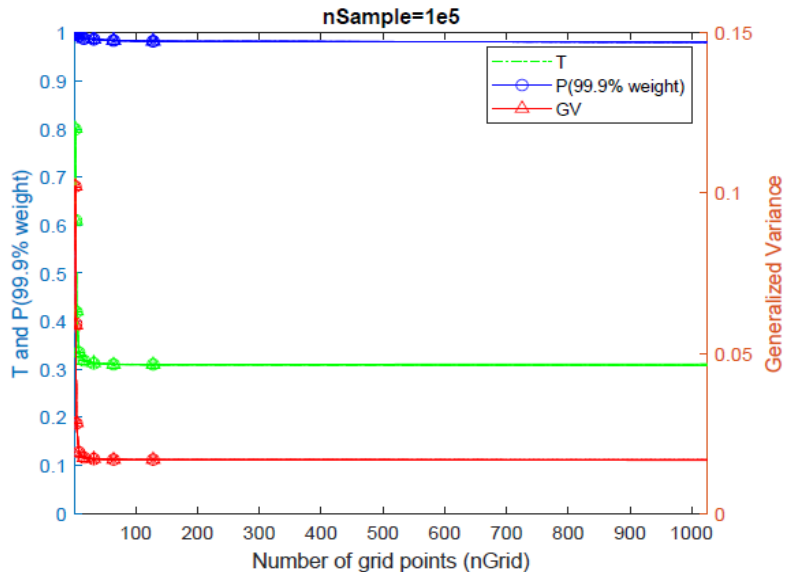
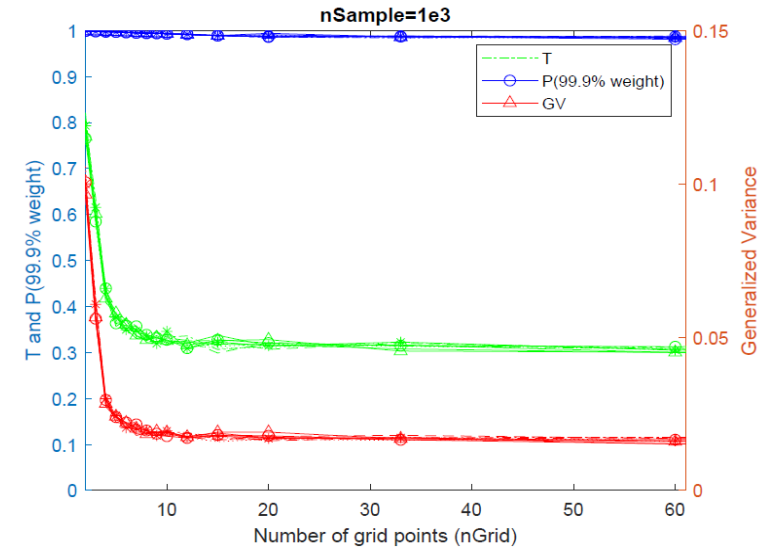
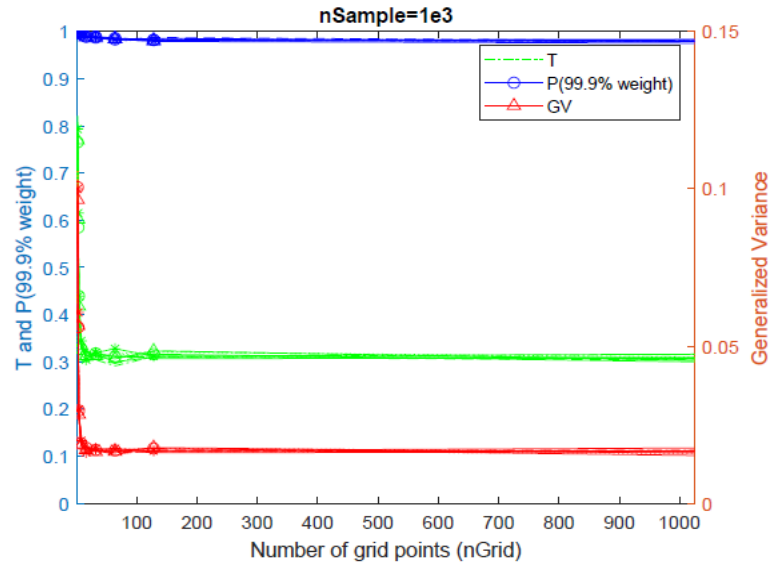
The sequences of PSRF and MPSRF indicate that they approach 1~1.1 over iterations ($nSample=5,000$) under scenarios of $nGrid=2, 3$ or 10 , showing fast convergence.



Marginal distributions of the fluxes of the toy example by the DISCOPOLIS algorithm with changing *nGrid*. *nSample*=5,000.



“Generalized variance” (GV) and **“total sample variance” (T)** are calculated and plotted against *nGrid* with changing *nSample*. **P (99.9% weight)** is also shown. 10 random seeds are taken under each scenario.



Flux mean values of the flux distribution obtained with optimal setting of the DISCOPOLIS algorithm

