
Anticoncentration Regularizers for Stochastic Combinatorial Problems

Shiva Kaul

CMU

Geoff Gordon

CMU

THE MACHINE LEARNING TRAGEDY

————— DRAMATIC STRUCTURE IN IV ACTS —————

I. THE GODS POSE A LEARNING PROBLEM

II. THE PROTAGONIST FINDS A STATISTICALLY
IDEAL WAY TO SOLVE IT

III. HIS WAY IS CURSED BY NP-HARDNESS

IV. HE IS BANISHED TO STATISTICAL INFERIORITY

The sparsity recovery problem

The setting:

$h^* \in \mathbb{R}^D$ is unknown

$\text{supp}(h^*)$ are the positions of its S nonzero entries

S is a constant fraction of D

Given for $1 \leq m \leq M$:

$x_m \sim N(0,1)^D$

$y_m = \langle h^*, x_m \rangle + g_m$ for some $g_m \sim N(0, \sigma^2)$

Asymptotic reliability:

$\mathbb{P}(\text{supp}(h) \neq \text{supp}(h^*)) \rightarrow 0$ as $M, S, D \rightarrow \infty$

with respect to the randomness of x_m and g_m

A direct approach is optimal... [W09]

$M = \Omega(S)$ is *necessary* for asymptotic reliability.

$M = O(S)$ is *sufficient* for asymptotic reliability of:

$$\min_h \frac{1}{M} \sum_m (\langle h, x_m \rangle - y_m)^2 \quad \text{s.t.} \quad \|h\|_0 \leq S \quad (\text{DIRECT}_2)$$

Still NP-hard when the numerical values of the inputs are polynomial in their bit-lengths.

...but (generally) [↑] strongly NP-hard.

Convex relaxation is suboptimal.

[W06]: $M = \Theta(S \log(D - S))$ is necessary for asymptotic reliability of

$$\min_h \frac{1}{M} \sum_m (\langle h, x_m \rangle - y_m)^2 + \lambda \|h\|_1 \quad (\text{LASSO})$$

Let's tweak DIRECT₂

Actually, a linear variant inheriting strong NP-hardness. (It's easier to write down.)

$$\min_h \frac{1}{M} \sum_m |\langle h, x_m \rangle - y_m| \quad \text{s.t.} \quad \|h\|_0 \leq S \quad (\text{DIRECT}_1)$$

Assumption for talk: DIRECT₁ \cong DIRECT₂



Asymptotically reliable at same rate, up to constants.

Theorem: a randomized polynomial-time algorithm is asymptotically reliable given $M = O(S)$.

1. $\text{DIRECT}_1 \cong \text{ROUND}$, where the decision variables are rounded to take polynomially many values..

use LASSO as a starting point to seed the algorithm

DIRECT₁

$$\text{min. } \sum_m \ell_m / M$$

$$\text{s.t. } \forall m \in \{1, \dots, M\}, d \in \{1, \dots, 2D\}, d' \in \{1, \dots, D\}$$

$$\ell_m = p_m + n_m$$

$$\langle h, x_m \rangle - y_m = p_m - n_m$$

$$h_d \leq I_d \cdot ?$$

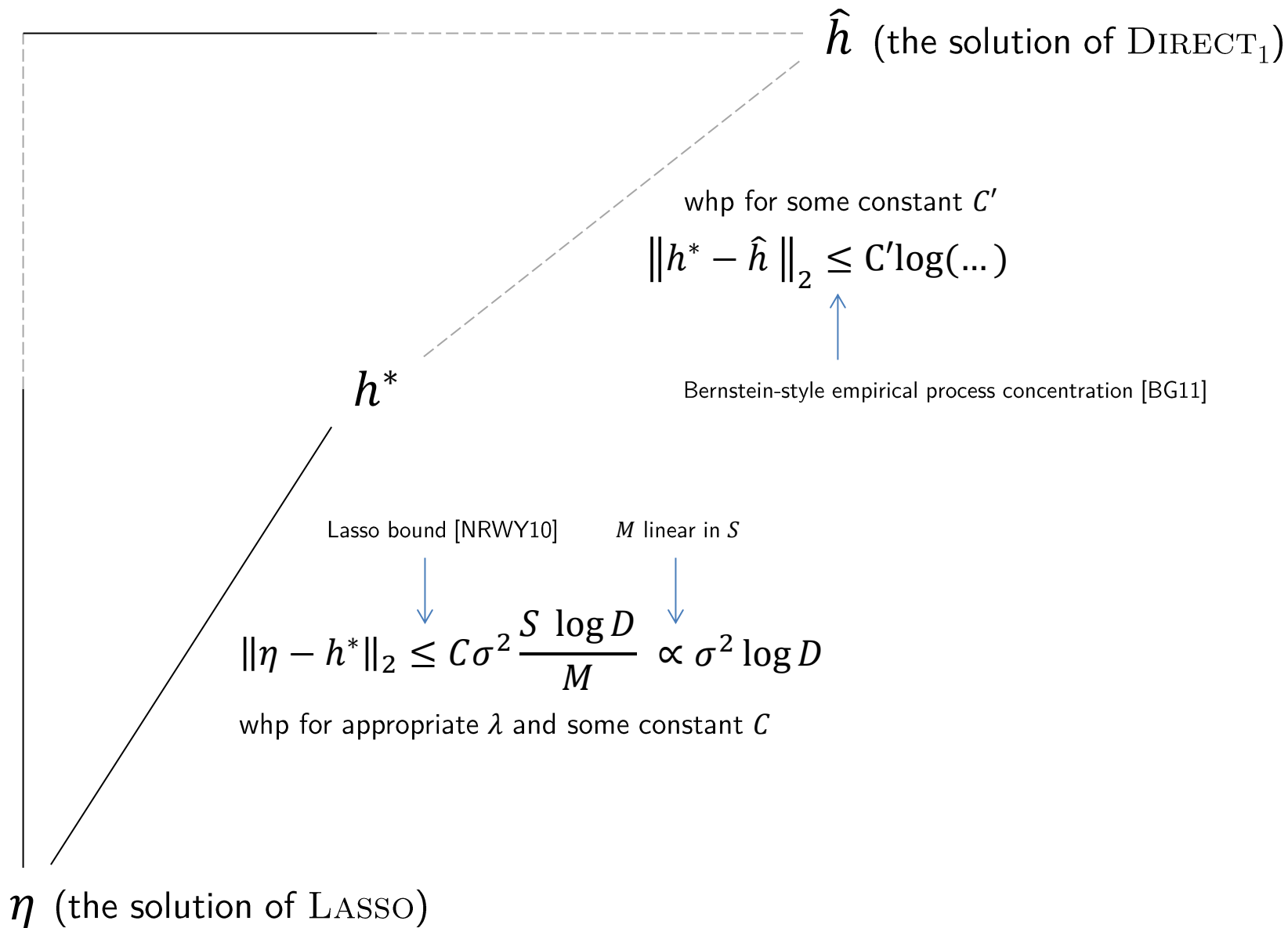
$$I_{d'} + I_{d'+D} \leq 1$$

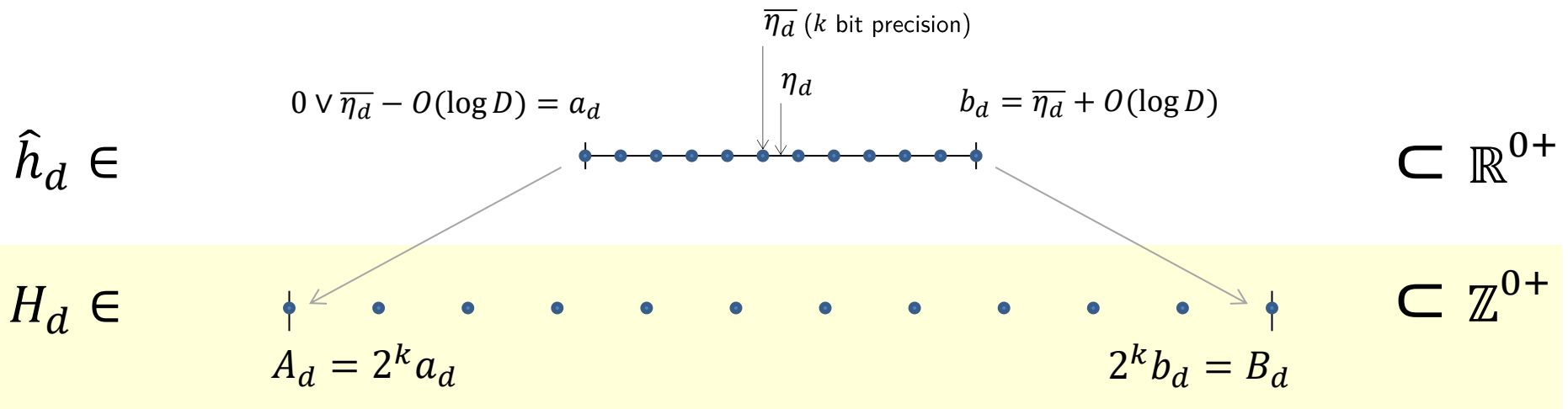
$$\sum_d I_d \leq S$$

$$\ell_m, p_m, n_m, h_d \geq 0$$

$$I_d \in \{0,1\}$$

The real-valued decision variables can take uncountably many values.





We can choose k so that the number of points is polynomial, yet $\|\hat{h} - H/2^k\|_2$ is exponentially small.

— ROUND (in progress) —

$$\text{min. } \sum_m \ell_m / M$$

$$\text{s.t. } \forall m \in \{1, \dots, M\}, d \in \{1, \dots, 2D\}, d' \in \{1, \dots, D\}$$

$$\ell_m = p_m + n_m$$

$$\langle H/2^k, x_m \rangle - y_m = p_m - n_m$$

$$H_d \leq I_d B_d$$

$$I_{d'} + I_{d'+D} \leq 1$$

$$\sum_d I_d \leq S$$

$$\ell_m, p_m, n_m \geq 0$$

$$H_d \in \{A_d, \dots, B_d\}$$

$$I_d \in \{0, 1\}$$

ROUND (almost)

$$\text{min. } \sum_m L_m / M$$

$$\text{s.t. } \forall m \in \{1, \dots, M\}, d \in \{1, \dots, 2D\}, d' \in \{1, \dots, D\}$$

$$L_m = P_m + N_m$$

$$\langle H, x_m 2^k \rangle - y_m 2^k \approx P_m - N_m$$

$$H_d \leq I_d B_d$$

$$I_{d'} + I_{d'+D} \leq 1$$

$$\sum_d I_d \leq S$$

$$L_m, P_m, N_m \in \{\dots\}$$

$$H_d \in \{A_d, \dots, B_d\}$$

$$I_d \in \{0, 1\}$$

Actually need to
replace with
inequalities and
introduce more
variables. Won't
bother, since it will
simplify anyway.

ROUND (almost)

$$\text{min. } \sum_m L_m / M$$

$$\text{s.t. } \forall m \in \{1, \dots, M\}, d \in \{1, \dots, 2D\}, d' \in \{1, \dots, D\}$$

$$L_m = P_m + N_m$$

$$\langle H, x_m 2^k \rangle - y_m 2^k \approx P_m - N_m$$

$$H_d \leq I_d B_d$$

$$I_{d'} + I_{d'+D} \leq 1$$

$$\sum_d I_d \leq S$$

$$L_m, P_m, N_m \in \{\dots\}$$

$$H_d \in \{A_d, \dots, B_d\}$$

$$I_d \in \{0, 1\}$$

1. $\text{DIRECT}_1 \cong \text{ROUND}$, where the decision variables are rounded to take polynomially many values.
use LASSO as a starting point to seed the algorithm

2. $\text{ROUND} \cong \text{SMOOTHROUND}$, which is smoothed by a random perturbation.
may perturb a random program

SMOOTHROUND (almost)

$$\text{min. } \sum_m L_m / M$$

$$\text{s.t. } \forall m \in \{1, \dots, M\}, d \in \{1, \dots, 2D\}, d' \in \{1, \dots, D\}$$

$$L_m = P_m + N_m$$

$$\langle H, x_m 2^k \rangle - z_m 2^k \approx P_m - N_m$$

$$H_d \leq I_d B_d$$

$$I_{d'} + I_{d'+D} \leq 1$$

$$\sum_d I_d \leq S$$

$$L_m, P_m, N_m \in \{\dots\}$$

$$H_d \in \{A_d, B_d\}$$

$$I_d \in \{0, 1\}$$

$z_m = y_m + \rho_m$ where $\rho_m \sim N(0, r^2)$. Adding it in allows us to 'restart' without drawing new sample.

1. $\text{DIRECT}_1 \cong \text{ROUND}$, where the decision variables are rounded to take polynomially many values.

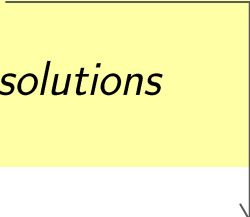
use LASSO as a starting point to seed the algorithm

2. $\text{ROUND} \cong \text{SMOOTHROUND}$, which is smoothed by a random perturbation.

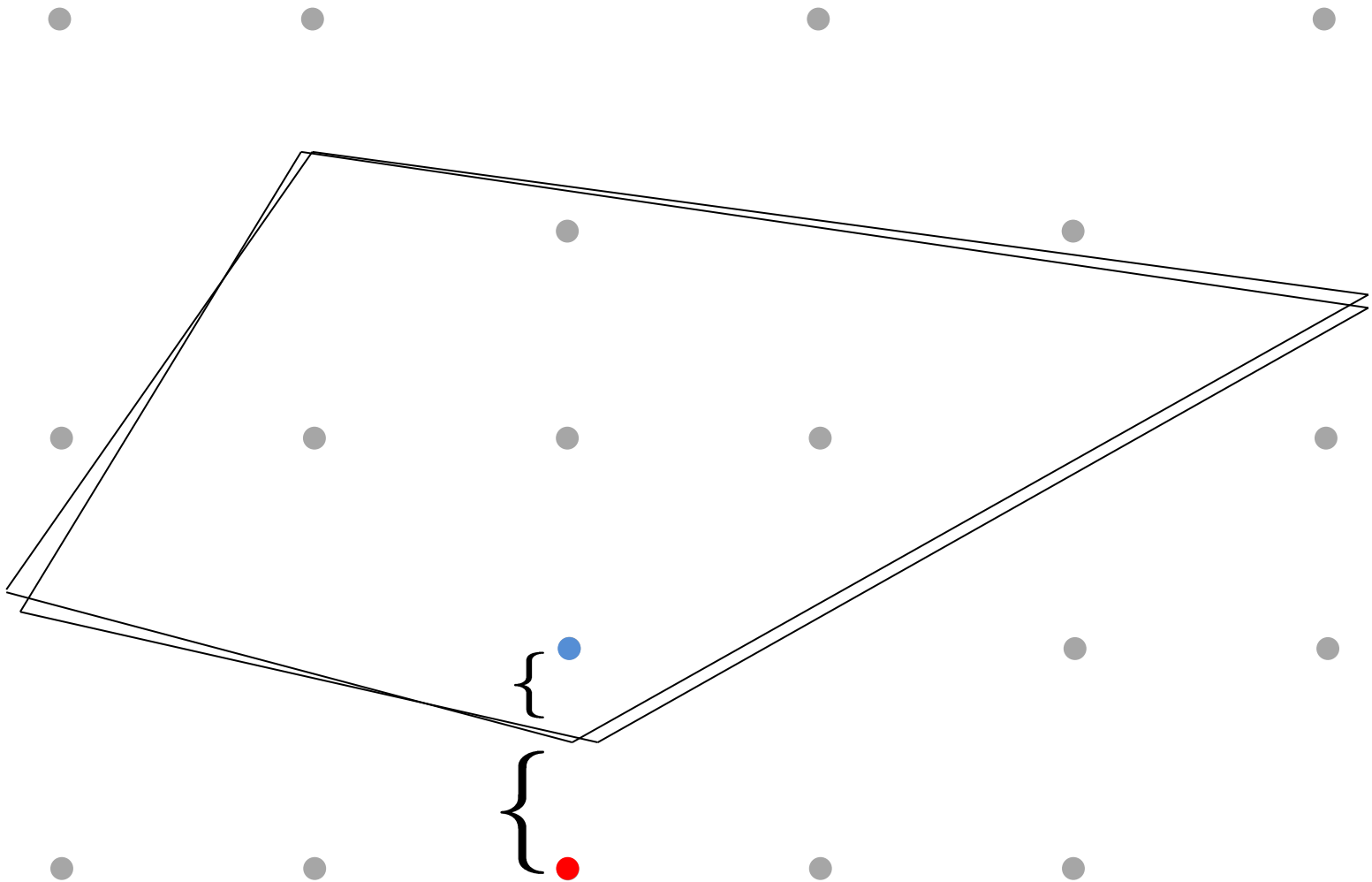
may perturb a random program

3. SMOOTHROUND can be solved in polynomial time if ROUND can be solved in pseudopolynomial time.

perturbed combinatorial problems have few optimal solutions



Polynomial in the numerical values of the inputs.
With this power we can solve some weakly NP-hard problems.



Smoothing leads to poly-size margins; the solution will still be ● even if the inputs (i.e. x_m and z_m) are truncated to logarithmic length i.e. polynomial value.

1. $\text{DIRECT}_1 \cong \text{ROUND}$, where the decision variables are rounded to take polynomially many values.
use LASSO as a starting point to seed the algorithm
2. $\text{ROUND} \cong \text{SMOOTHROUND}$, which is smoothed by a random perturbation.
may perturb a random program
3. SMOOTHROUND can be solved in polynomial time if ROUND can be solved in pseudopolynomial time.
perturbed combinatorial problems have few optimal solutions
4. ROUND can be solved in pseudopolynomial time.
input can't encode complex dependencies

ROUND

$$\text{min. } \sum_m L_m / M$$

$$\text{s.t. } \forall m \in \{1, \dots, M\}, d \in \{1, \dots, 2D\}, d' \in \{1, \dots, D\}$$

$$L_m = P_m + N_m$$

$$\langle H, \bar{x}_m 2^k \rangle - \bar{y}_m 2^k = P_m - N_m$$

$$H_d + \phi_d = I_d B_d$$

$$I_{d'} + I_{d'+D} + \psi_{d'} = 1$$

$$\sum_d I_d + \Psi = S$$

$$L_m, P_m, N_m \in \{\dots\}$$

$$H_d \in \{A_d, \dots, B_d\}$$

$$I_d \in \{0, 1\}$$

$$\phi_d, \psi_{d'}, \Psi \in \mathbb{Z}^{0+}$$

Since $M = \Omega(S) = \Omega(D)$ we may assume there are as many rows as columns.

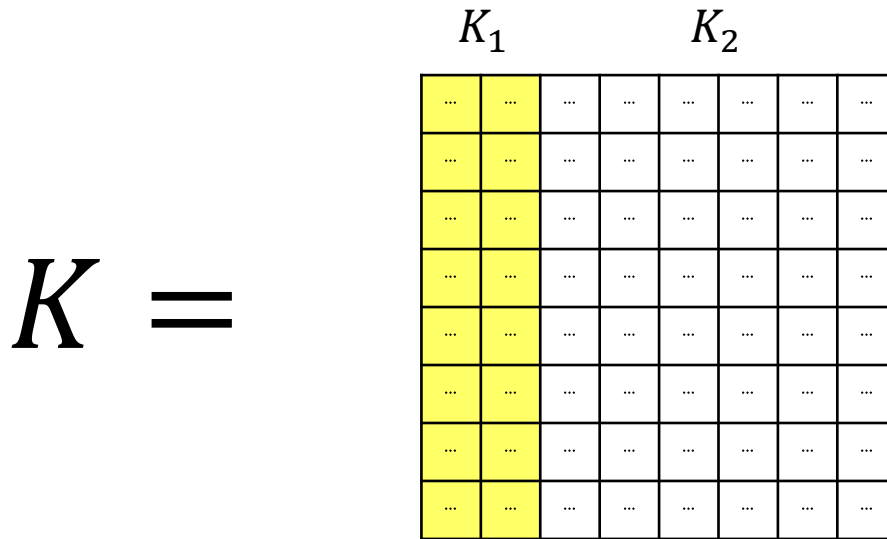
$m \in \{1, \dots, M\}$

$m \in \{1, \dots, M\}$			$d' \in \{1, \dots, D\}$						
L_m	N_m	P_m	$H_{d'}$	$H_{d'+D}$	$I_{d'}$	$I_{d'+D}$	$\phi_{d'}$	$\psi_{d'}$	Ψ
1	-1	-1							
			1	1	-1		1		
					1	1		1	
	1	-1	$\bar{x}_{m,d'}$	$\bar{x}_{m,d'+D}$					
...

0
0
0
\bar{y}_m
...

It mostly encodes a matrix of iid $N(0,1)$ random variables.

The inability to encode complex dependencies is captured by constant branchwidth.



A branch decomposition is a binary tree on the columns

Cutting an edge partitions the columns into K_1 and K_2

$$\text{branchwidth} = \min_{\text{decompositions}} \max_{\text{cuts}} (\text{rank}(K_1) + \text{rank}(K_2) - \text{rank}(K) + 1)$$

[CG06]: An integer linear program in equational form can be solved in pseudopolynomial time if its decision variables take polynomially many values and its constraint matrix has constant branchwidth.

Done.

Conclusions

Don't let worst-case hardness scare you away from average-case problems.

In order to obtain better statistical guarantees, you can exploit:

- the huge amount of work on relaxations (don't just toss it out!),
- the instable, random nature of the optimization program,
- simple structure of the input.