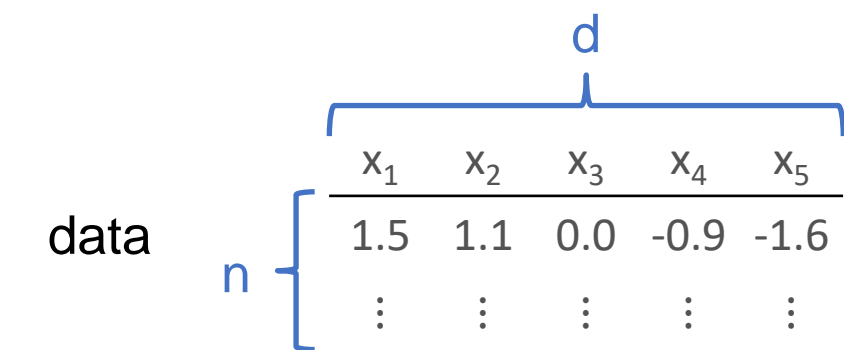


# DAGs with No Fears: A Closer Look at Continuous Optimization for Learning Bayesian Networks

Presented by: **Yue Yu**, Department of Mathematics, Lehigh University

Joint work with: **Dennis Wei, Tian Gao**, IBM research

## Bayesian Networks via Continuous Optimization



Given  $n$  samples of  $X \in \mathbb{R}^d$ , learn Bayesian network parametrized by  $W$

Score-based learning:

$$\begin{aligned} & \text{minimize} && \text{score } F(W) \\ & \text{subject to} && G(W) \in \text{DAGs} \end{aligned}$$

combinatorial

NOTEARS [1]:

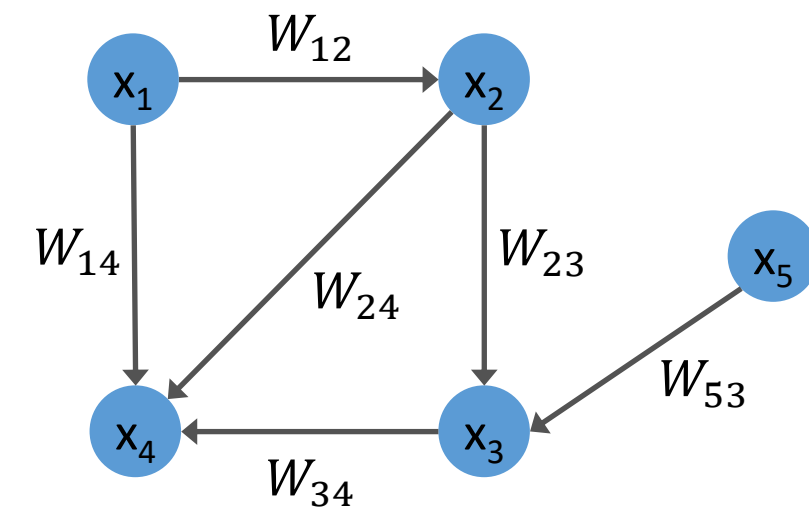
$$\begin{aligned} & \text{minimize} && \text{score } F(W) \\ & \text{subject to} && h(A) = 0 \end{aligned}$$

continuous, differentiable

adjacency matrix

## 1. Better Understanding of NOTEARS

Assumption: Each edge corresponds to one parameter  $W_{ij}$   
e.g. generalized linear SEM



$$\begin{aligned} & \text{minimize} && F(W) \\ & \text{subject to} && h(W \circ W) = 0 \end{aligned}$$

quadratic adjacency matrix  $A = W \circ W$

– Acyclic solutions cannot satisfy KKT conditions

– Augmented Lagrangian algorithm cannot converge to acyclic solution even with high penalty

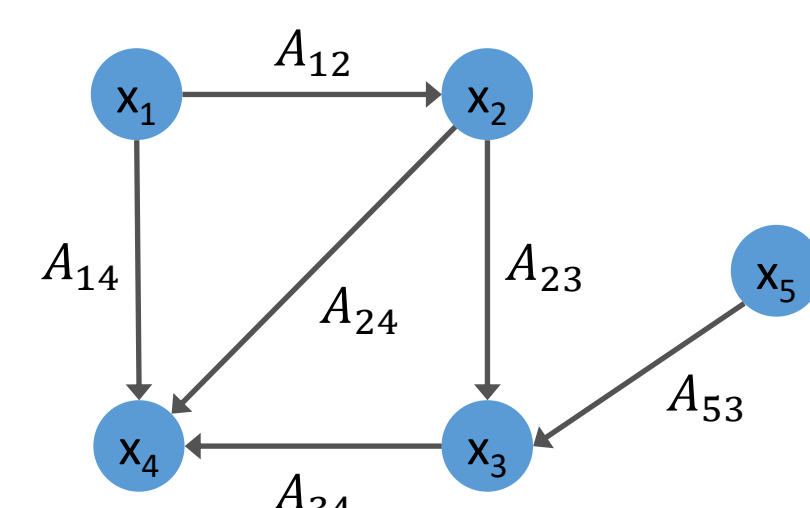
## Our Contributions

Theoretical understanding of continuous optimization framework, leading to significant algorithmic improvements

1. Better understanding of NOTEARS
2. Understanding of KKT optimality conditions for reformulation
3. KKT-search post-processing improves all tested algorithms

## General Class of Acyclicity Constraints

Adjacency matrix  $A$ :  
 $A_{ij} > 0 \Leftrightarrow \text{edge}$



NOTEARS [1]:  $h(A) = \text{tr}(e^A) - d$

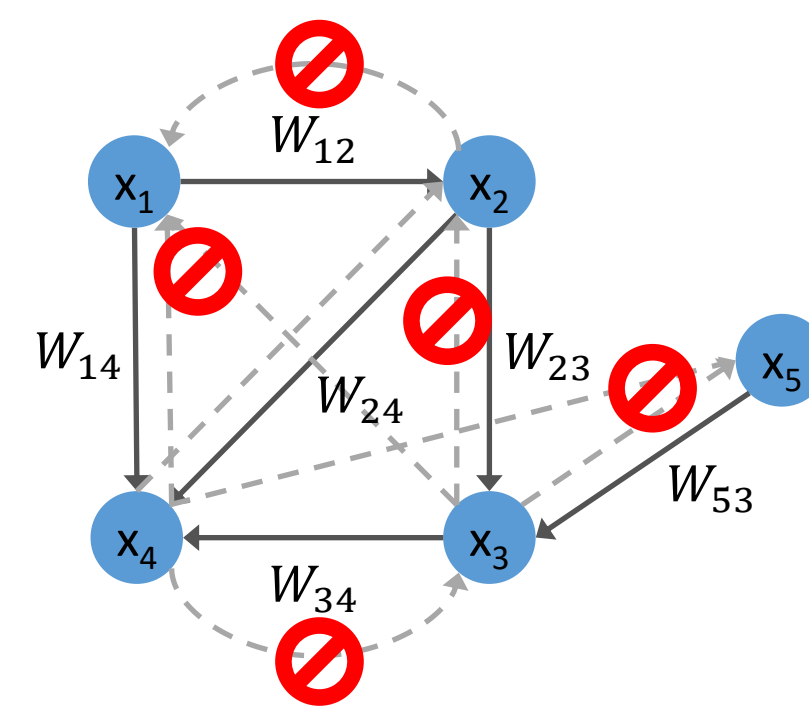
DAG-GNN [2]:  $h(A) = \text{tr}((I + A/d)^d) - d$

Our generalization:

$$h(A) = \text{tr} \left( \sum_{p=1}^d c_p A^p \right) \quad c_p > 0$$

$\nabla h(A)$  has a directed walk interpretation

## 2. KKT Conditions for Reformulation



$$\begin{aligned} & \text{minimize} && F(W) \\ & \text{subject to} && h(|W|) = 0 \end{aligned}$$

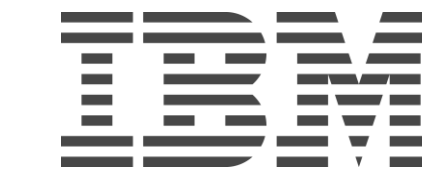
absolute adjacency matrix  $A = |W|$

– KKT conditions are indeed necessary: Local minima must satisfy them

– If score  $F(W)$  is convex, KKT conditions **sufficient** for local minimality, despite non-convexity of constraint

– Understanding of KKT conditions in terms of **edge absence constraints**: Collectively sufficient, individually necessary in preventing cycles

$$\begin{aligned} & \text{minimize} && F(W) \\ & \text{subject to} && W_{ij} = 0, \quad (i, j) \in \mathcal{Z} \end{aligned}$$

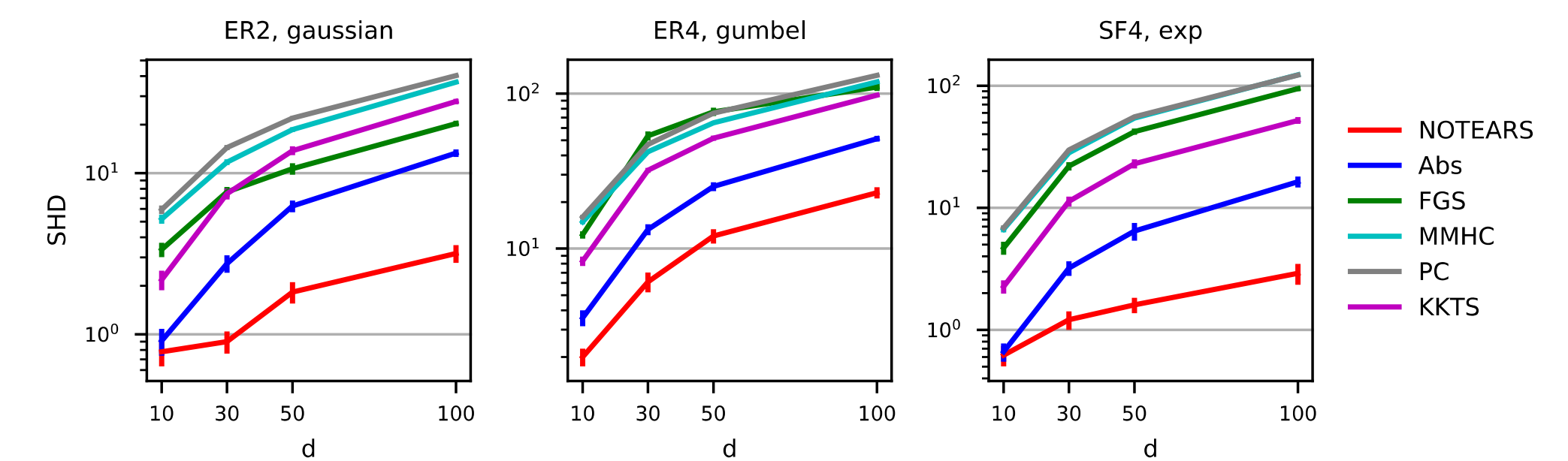


## Proposed Algorithms

1. Augmented Lagrangian with  $A = |W|$  ('Abs')
2. KKT-search to satisfy KKT conditions

## 3. KKT-search Improves Existing Algorithms

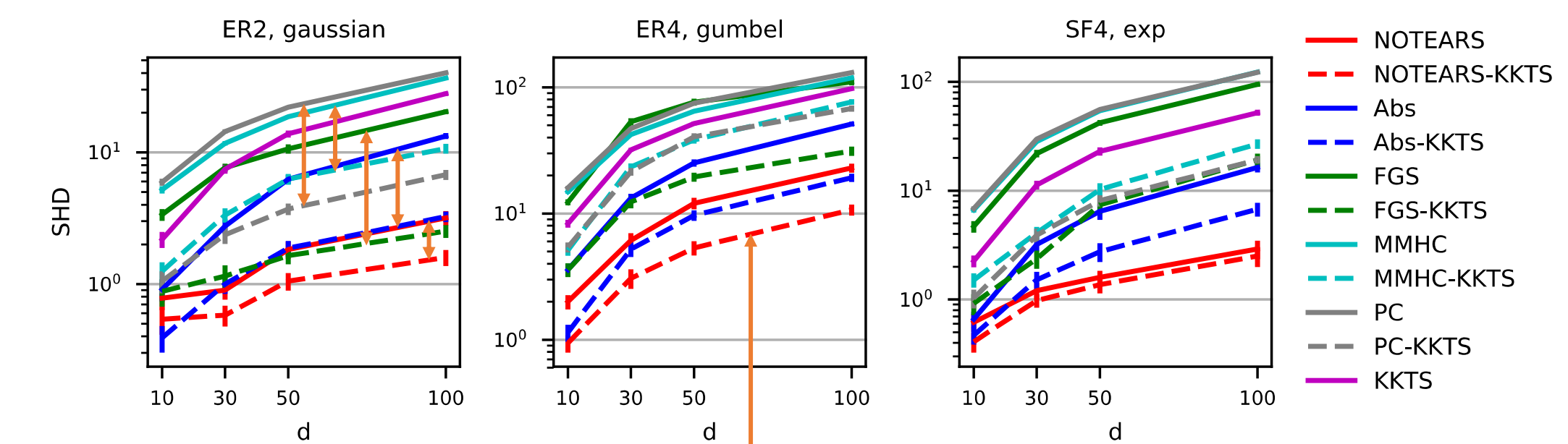
Base algorithms: NOTEARS still best, Abs second



Structural Hamming distance (SHD) with respect to true graph for different graph types and  $n = 1000$  samples

With KKT-search:

Consistent improvement across base algorithms and dimension  $d$



>2X reductions  
state-of-the-art accuracy

## References

- [1] Xun Zheng, Bryon Aragam, Pradeep Ravikumar, Eric P. Xing (2018). "DAGs with NO TEARS: Continuous optimization for structure learning." In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [2] Yue Yu, Jie Chen, Tian Gao, Mo Yu (2019). "DAG-GNN: DAG structure learning with graph neural networks." In *International Conference on Machine Learning (ICML)*.
- [3] Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, Simon Lacoste-Julien (2020). "Gradient-based neural DAG learning". In *International Conference on Learning Representations (ICLR)*.
- [4] Xun Zheng, Chen Dan, Bryon Aragam, Pradeep Ravikumar, Eric P. Xing (2020). "Learning sparse nonparametric DAGs." In *International Conference on Artificial Intelligence and Statistics (AISTATS)*.