

Individualized Inference using Bayesian Quantile Directed Acyclic Graphical Models

Ksheera Sagar K. N**

Department of Statistics
Purdue University

August 11, 2022

**collaborative work with Veera B (UMich), Yang Ni (Texas A&M) and Anindya Bhadra (Purdue).

① Quantile DAGs

- Introduction and motivation

- Intuition to the model

- Related works and our contributions

- qDAGx: The model

- Prior formulation

- Synthetic data and simulation settings

- Simulation results

- Application of qDAGx in lung cancer data

② Conclusion and future scope

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

Challenges with GGMs:

- What if the interacting variables are not jointly Gaussian?
- And consequent lack of robustness in model misspecification.

Our solution:

- Circumvent the Gaussian assumption on the likelihood.
- Model association between variables at any given quantile level, $\tau \in (0, 1)$.

Real life motivation:

- Directed acyclic graphs are important to understand protein–protein interaction (PPI) networks.
- Personalized PPI's can help in a better understanding of diseases like cancer, and therefore finding applications in precision medicine.

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

INTUITION TO THE MODEL

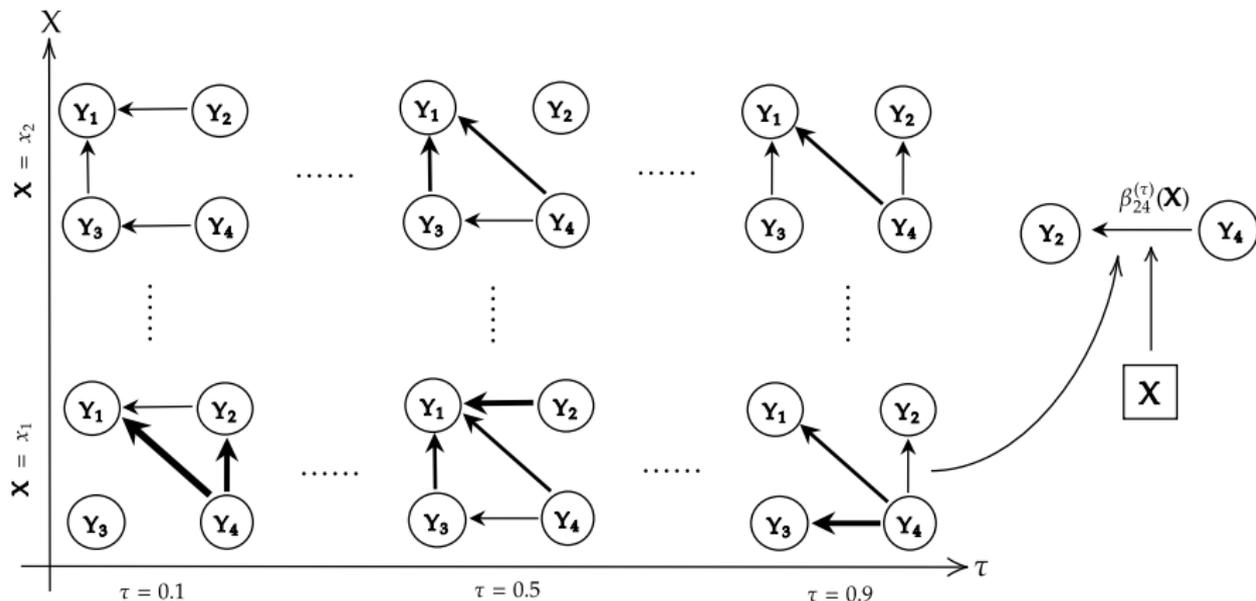


Figure: The directed acyclic graph $QG^{(\tau)}$, for two observations, on four vertices $Y = (Y_1, Y_2, Y_3, Y_4)$ is presented, for univariate X and for given quantile levels $\tau = 0.1, 0.5$ and 0.9 .

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

RELATED WORKS AND OUR CONTRIBUTIONS

Related works inspired by varying coefficient models (Hastie and Tibshirani, 1993):

- DAG inference using node conditional varying coefficient models (Ni et al., 2019). **Drawbacks:** Gaussian likelihood and known ordering assumption.
- Varying coefficient Bayesian quantile regression (Das et al., 2021). **Drawbacks:** lack of support to graphical models.

Related work in quantile-graphs:

- Undirected quantile-graphs constructed from node-wise quantile regression (Guha et al., 2020). **Drawbacks:** not individualized.

RELATED WORKS AND OUR CONTRIBUTIONS (CONT.)

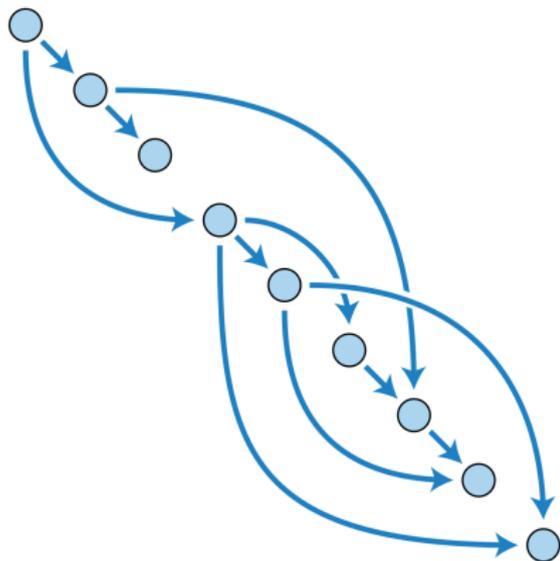


Figure: A sample DAG which is Topologically sorted.

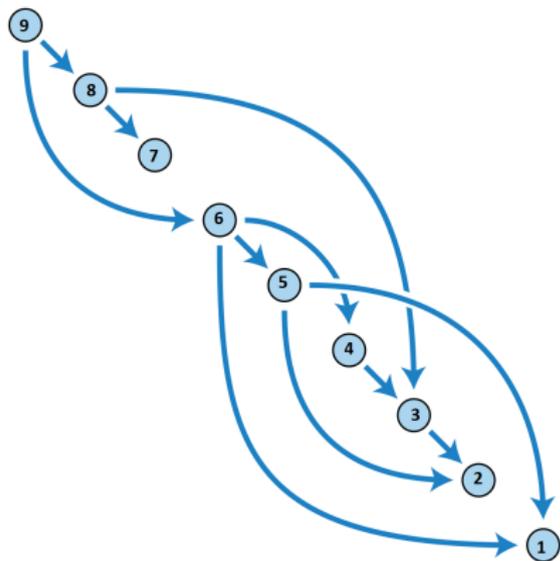


Figure: Labels given to nodes WLOG in the topologically sorted DAG.

Picture credits: https://en.wikipedia.org/wiki/Directed_acyclic_graph

RELATED WORKS AND OUR CONTRIBUTIONS (CONT.)

Our contributions:

- 1 **qDAGx**: Learning *individual-specific* DAG's at any quantile level $\tau \in (0, 1)$, with no assumptions on the data likelihood or on the ordering of nodes.
- 2 Infer for the first time *individual-specific* protein–protein interaction networks in patients with lung adenocarcinoma and lung squamous cell carcinoma.
 - Model the protein-protein association in each patient at a quantile level τ , as a function of external covariates mRNA and methylation.
- 3 Structural identifiability of the quantile-DAGs, properties of prior which aid in sparse quantile DAG discovery and posterior consistency of node conditional fitted densities.

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

qDAGx: THE MODEL

Notations:

- First level covariates: $\mathbf{Y}_1, \dots, \mathbf{Y}_p$ and for $h \in \{1, \dots, p\}$, $\mathbf{Y}_h \in \mathbb{R}^n$.
- Second level covariates: $\mathbf{X}_1, \dots, \mathbf{X}_q$ and for $k \in \{1, \dots, q\}$, $\mathbf{X}_k \in \mathbb{R}^n$.
- Y_{ih} , X_{ik} : first, second level covariate values for i^{th} observation, $i \in \{1, \dots, n\}$.

Model for conditional quantile:

$$Q_{Y_{ih}}(\tau | Y_{ij}, \mathbf{X}_{i\cdot}) = \beta_{h0}^{(\tau)}(\mathbf{X}_{i\cdot}) + \sum_{j \in pa(h)} Y_{ij} \beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})$$

qDAGx: THE MODEL

Notations:

- First level covariates: $\mathbf{Y}_1, \dots, \mathbf{Y}_p$ and for $h \in \{1, \dots, p\}$, $\mathbf{Y}_h \in \mathbb{R}^n$.
- Second level covariates: $\mathbf{X}_1, \dots, \mathbf{X}_q$ and for $k \in \{1, \dots, q\}$, $\mathbf{X}_k \in \mathbb{R}^n$.
- Y_{ih} , X_{ik} : first, second level covariate values for i^{th} observation, $i \in \{1, \dots, n\}$.

Model for conditional quantile:

$$Q_{Y_{ih}}(\tau | Y_{ij}, \mathbf{X}_{i\cdot}) = \beta_{h0}^{(\tau)}(\mathbf{X}_{i\cdot}) + \sum_{j \in pa(h)} Y_{ij} \beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})$$

$$\beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) = \theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) \cdot \mathbb{1}(|\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})| > t_{hj}) \text{ where,}$$

qDAGx: THE MODEL

Notations:

- First level covariates: $\mathbf{Y}_1, \dots, \mathbf{Y}_p$ and for $h \in \{1, \dots, p\}$, $\mathbf{Y}_h \in \mathbb{R}^n$.
- Second level covariates: $\mathbf{X}_1, \dots, \mathbf{X}_q$ and for $k \in \{1, \dots, q\}$, $\mathbf{X}_k \in \mathbb{R}^n$.
- Y_{ih} , X_{ik} : first, second level covariate values for i^{th} observation, $i \in \{1, \dots, n\}$.

Model for conditional quantile:

$$Q_{Y_{ih}}(\tau | Y_{ij}, \mathbf{X}_{i\cdot}) = \beta_{h0}^{(\tau)}(\mathbf{X}_{i\cdot}) + \sum_{j \in \text{pa}(h)} Y_{ij} \beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})$$

$$\beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) = \theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) \cdot \mathbb{1}(|\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})| > t_{hj}) \text{ where,}$$

$$\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) = \sum_{k=1}^q f_{hjk}^{(\tau)}(X_{ik}).$$

qDAGx: THE MODEL (CONT.)

Union-DAG condition: Let $\mathcal{Q}\mathcal{G}_i^{(\tau)}$ be the adjacency matrix of quantile-DAG of i^{th} observation at quantile level τ . Then,

$$\mathcal{Q}\mathcal{G}_u^{(\tau)} = \bigcup_{i=1}^n \mathcal{Q}\mathcal{G}_i^{(\tau)} \text{ is a DAG, where } \mathcal{Q}\mathcal{G}_i^{(\tau)} = \left((\beta_{hj}(\mathbf{X}_{i\cdot}) \neq 0) \right).$$

loss function \rightarrow negative log-likelihood of ALD (Koenker and Bassett Jr, 1978) \rightarrow Joint likelihood.

Joint likelihood:

$$\begin{aligned} \pi(\mathbf{Y} \mid \mathbf{X}, \boldsymbol{\beta}^{(\tau)}, \tau) &= \prod_{i=1}^n \prod_{h=1}^p \tau(1-\tau) \exp \left(-\psi_{\tau} \left(Y_{ih} - \beta_{h0}^{(\tau)}(\mathbf{X}_{i\cdot}) - \sum_{j \in \text{pa}(h)} Y_{ij} \beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) \right) \right) \\ &\quad \times \mathbb{1} \left(\mathcal{Q}\mathcal{G}_u^{(\tau)} \text{ is a DAG} \right). \end{aligned}$$

Where, $\psi_{\tau}(x) = \tau \mathbb{1}(x \geq 0) - (1 - \tau) \mathbb{1}(x < 0)$.

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

PRIOR FORMULATION

$$\boldsymbol{\theta}_{hj}^{(\tau)}(\mathbf{X}) = \sum_{k=1}^q f_{hjk}^{(\tau)}(\mathbf{X}_k) = \mu_{hj} \mathbf{1}_n + \sum_{k=1}^q \widetilde{\mathbf{X}}_k^* \boldsymbol{\alpha}_{hjk}^* + \sum_{k=1}^q \mathbf{X}_k \alpha_{hjk}^0,$$

$$\text{peNMHS: } \begin{cases} \boldsymbol{\alpha}_{hjk}^* & = \eta_{hjk} \boldsymbol{\xi}_{hjk}, \eta_{hjk} \sim \mathcal{N}(0, T_{hj}^2 L_{hjk}^2), \\ \boldsymbol{\xi}_{hjk} & = \left(\xi_{hjk}^{(1)}, \dots, \xi_{hjk}^{(B_k^*)} \right)^T, \\ \xi_{hjk}^{(l)} & \sim \mathcal{N}(m_{hjk}^{(l)}, 1), \text{ for } l \in \{1, \dots, B_k^*\}, \\ m_{hjk}^{(l)} & \sim 0.5 \cdot \delta_1(m_{hjk}^{(l)}) + 0.5 \cdot \delta_{-1}(m_{hjk}^{(l)}), \\ T_{hj} & \sim \mathcal{C}^+(0, 1), L_{hjk} \sim \mathcal{C}^+(0, 1). \end{cases}$$

$\alpha_{hjk}^0 \sim$ peNMHS prior analogous to $\boldsymbol{\alpha}_{hjk}^*$.

$\mu_{hj} \sim \mathcal{N}(0, \sigma_\mu^2)$ and $t_{hj} \sim \text{Gamma}(\text{shape} = a, \text{rate} = b)$, $1 \leq h, j \leq p$.

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

SYNTHETIC DATA AND SIMULATION SETTINGS

Problem dimensions: $n \in \{100, 250\}$, $p \in \{25, 50, 100\}$, $q \in \{2, 5\}$.

Data generating mechanism:

- $\mathbf{X}_1, \dots, \mathbf{X}_q$ from a multivariate normal $\mathcal{N}(0, \mathbf{I}_q)$.
- WLOG, true order is $\mathbf{Y}_1, \dots, \mathbf{Y}_p$ and for each \mathbf{Y}_h choose $\max\left\{1, \lfloor \frac{p-h}{5} \rfloor\right\}$ number of parents from $\{\mathbf{Y}_{h+1}, \dots, \mathbf{Y}_p\}$.
- Compute values of corresponding $\theta_{hj}(\cdot)$'s where $j \in pa(h)$, which are functions of τ , \mathbf{X}
- When $q = 2$, all thresholds = 0.5 and when $q = 5$, all thresholds = 1.
- Compute $\beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) = \theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}) \cdot \mathbb{1}(|\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})| > t_{hj})$
- Generate n random samples for \mathbf{Y}_h :

$$Q_{Y_{ih}}(\tau | Y_{ij}, \mathbf{X}_{i\cdot}) = \beta_{h0}^{(\tau)}(\mathbf{X}_{i\cdot}) + \sum_{j \in pa(h)} Y_{ij} \beta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot}).$$

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

SIMULATION RESULTS

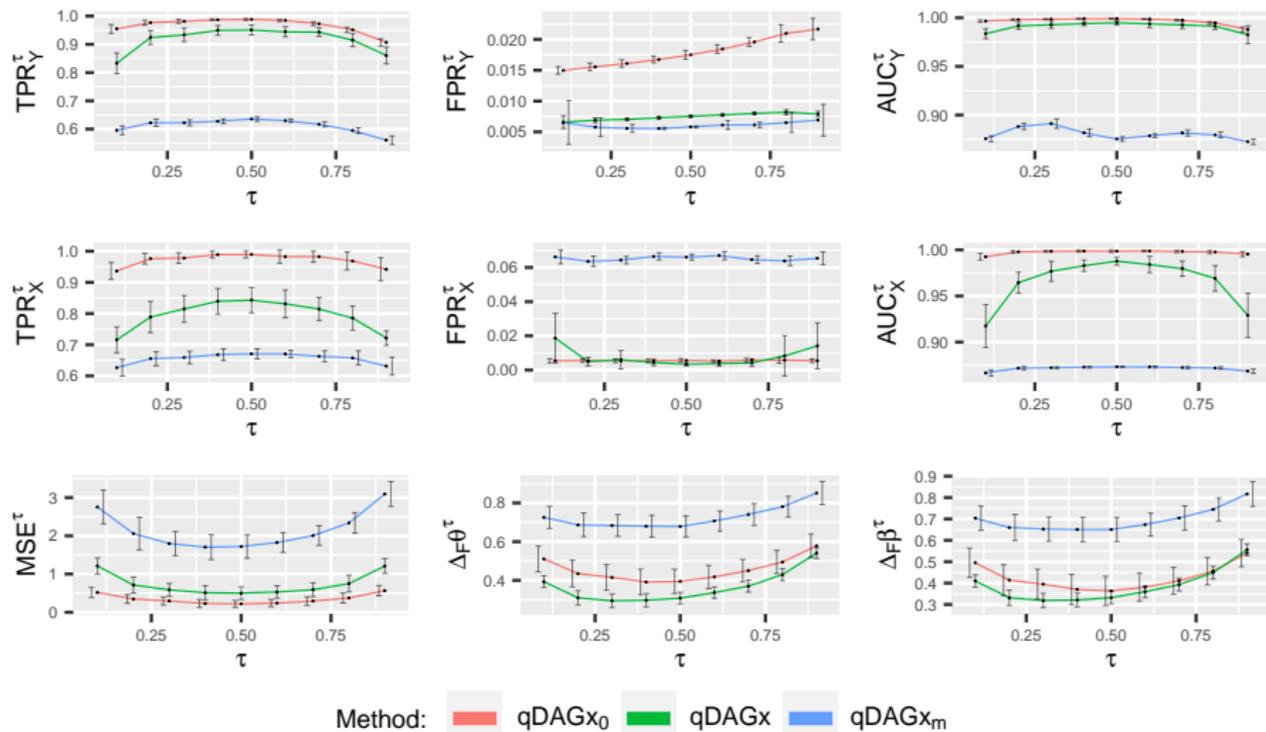


Figure: $p = 25, q = 5, n = 250$. Kendall's' T for the misspecified sequence is 0.5

① Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

② Conclusion and future scope

Overview:

- Protein expressions of 67 proteins, analogous to $\mathbf{Y}_1, \dots, \mathbf{Y}_{p=67}$.
- $\mathbf{X}_1, \mathbf{X}_2$: mRNA expression and methylation.
- $n = 306$ patients: Lung adenocarcinoma (LUAD).
- $n = 278$ patients: Lung squamous cell carcinoma (LUSC).
- Estimate quantile-DAGs at $\tau \in \{0.1, \dots, 0.9\}$.
- Aggregate DAGs at each quantile level *for visualization purposes* and show edges present in $\geq n/2$ patients and node size \propto in-degree.

qDAGx IN LUNG CANCER DATA (CONT.)

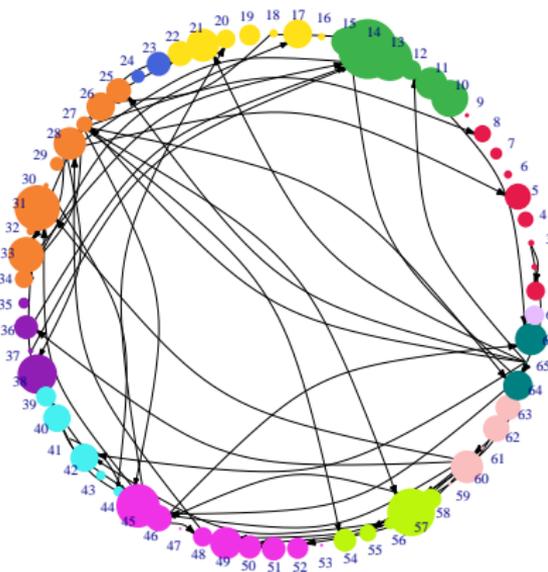


Figure: Graph of $E_{\text{LUAD}}^{(0.1)}$.

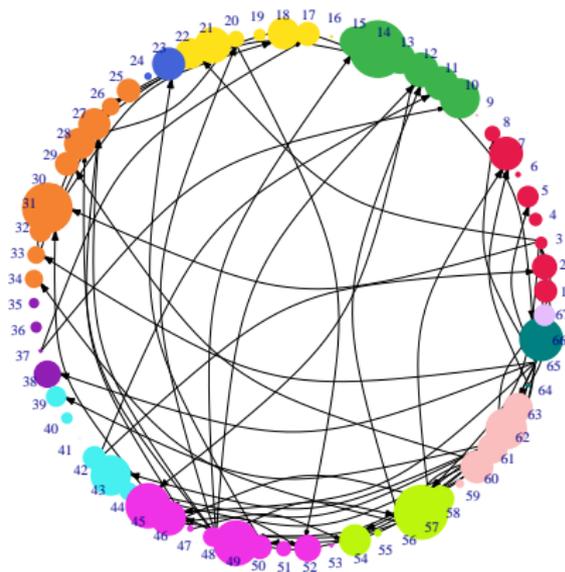


Figure: Graph of $E_{\text{LUSC}}^{(0.1)}$.

qDAGx IN LUNG CANCER DATA (CONT.)

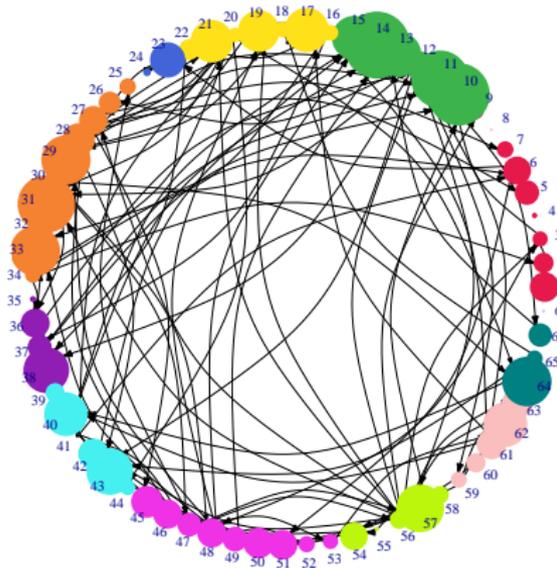


Figure: Graph of $E_{LUAD}^{(0.5)}$.

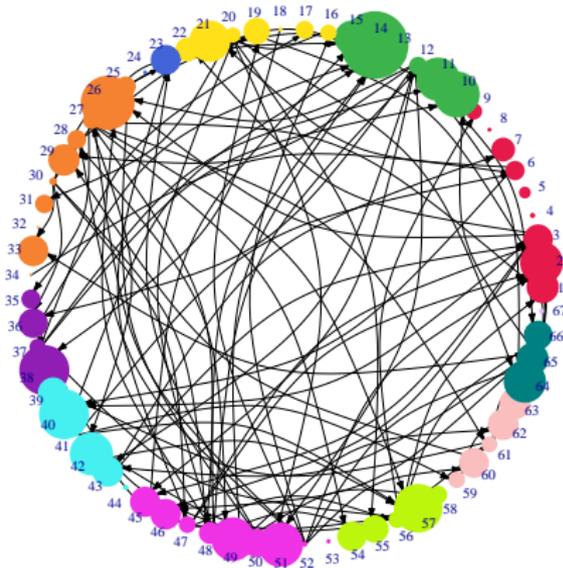


Figure: Graph of $E_{LUSC}^{(0.5)}$.

qDAGx IN LUNG CANCER DATA (CONT.)

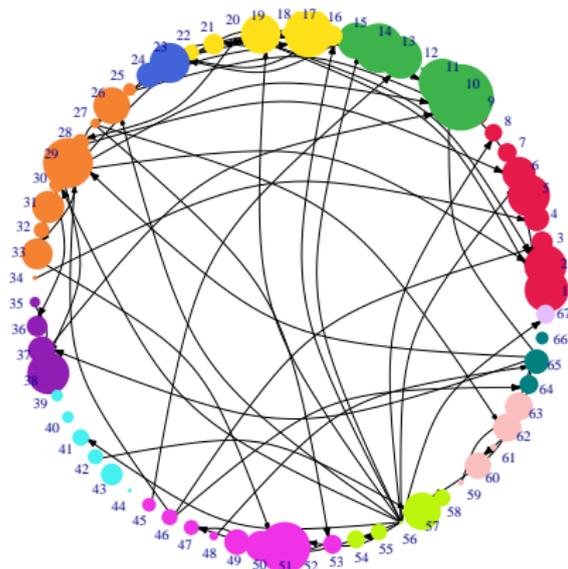


Figure: Graph of $E_{LUAD}^{(0.9)}$.

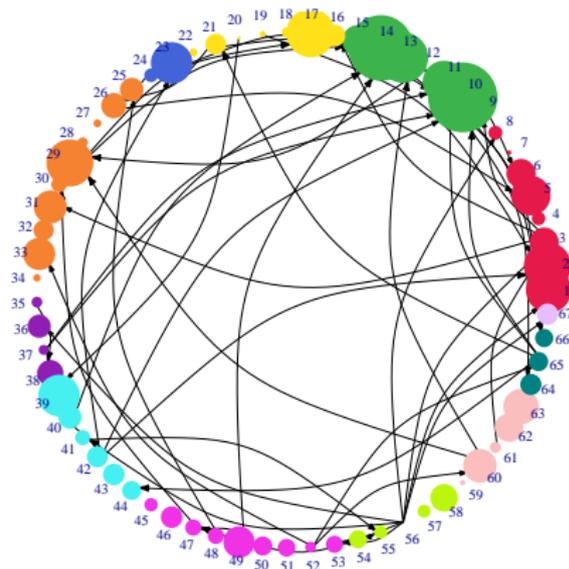


Figure: Graph of $E_{LUSC}^{(0.9)}$.

qDAGx IN LUNG CANCER DATA (CONT.)

Table: Directed edges in quantile-DAG estimates which are present in at least 50% of patients and across five out of nine quantile levels, $\tau \in \{0.1, \dots, 0.9\}$. Common edges in LUAD and LUSC are highlighted.

Lung adenocarcinoma (LUAD)			Lung squamous cell carcinoma (LUSC)		
BAK1←BID	BAD←ATK1S1	BID←ERBB3	BAK1←BID	AKT1, AKT2, AKT3←AKT1S1	CAV1←PGR
CAV1←COL6A1	EGFR←ERBB2	GAPDH←CDH2	CAV1←COL6A1	EGFR←ERBB2	CCNB1←COL6A1
JUN←ERBB3	MAPK1, MAPK3←MAP2K1	MYH11←COL6A1	MTOR←PGR	MAPK1, MAPK3←MAP2K1	MYH11←COL6A1
PCNA←CHEK1	RPS6KB1←PGR		MYH11←FOXM1	RPS6KB1←PGR	RAD51←PGR

Table: Percentage of edges, mean (sd), influenced by second level covariates in quantile-DAG estimates of all patients, across the quantiles $\tau \in \{0.1, \dots, 0.9\}$.

	only mRNA	only methylation	both
LUAD	13.7 (0.78)	28 (0.86)	58.3 (1.55)
LUSC	13.6 (0.66)	28.1 (0.61)	58.3 (0.85)

qDAGx IN LUNG CANCER DATA (CONT.)

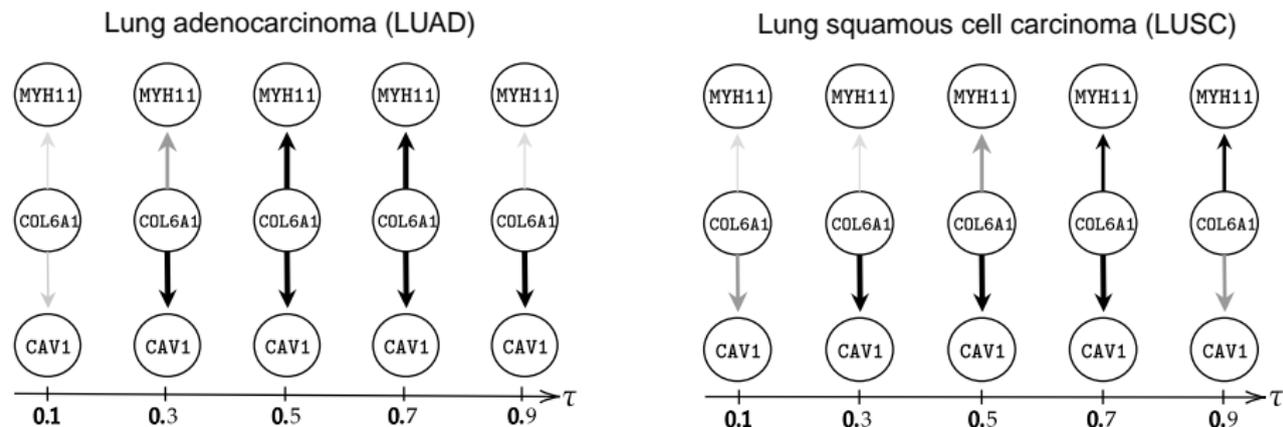


Figure: Prevalence of $CAV1 \leftarrow COL6A1$ and $MYH11 \leftarrow COL6A1$ in LUAD and LUSC. Boldness of the edge is proportional to the number of patients in whom the edge was inferred at the specific quantile level τ .

1 Quantile DAGs

Introduction and motivation

Intuition to the model

Related works and our contributions

qDAGx: The model

Prior formulation

Synthetic data and simulation settings

Simulation results

Application of qDAGx in lung cancer data

2 Conclusion and future scope

CONCLUSION AND FUTURE SCOPE

Contributions:

- **qDAGx**: quantile-DAG learning framework with neither assumptions on likelihood nor on ordering of the nodes.
- Individualized inference via a varying coefficient framework.
- Demonstration of qDAGx in patients with LUAD and LUSC → usefulness in precision medicine.

Future work:

- Theoretical guarantees of estimating quantile-DAG structure similar to Cao et al. (2019); DAG estimation consistency in GGMs.
- Incorporate conditions for preserving increasing nature of quantile estimates (Ali et al., 2016; Yang and Tokdar, 2017).
- Mixture quantile-DAG modeling, possibly by modeling threshold parameter $t_{h,j}$ as function of categorical variables like type of cancer, gender etc. And/or relaxing the union-DAG condition.

- Ali, A., Kolter, J. Z., and Tibshirani, R. J. (2016). The multiple quantile graphical model. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Cao, X., Khare, K., and Ghosh, M. (2019). Posterior graph selection and estimation consistency for high-dimensional bayesian dag models. *The Annals of Statistics*, 47(1):319–348.
- Das, P., Christine, B. P., Ni, Y., Reuben, A., Zhang, J., Zhang, J., Do, K.-A., and Baladandayuthapani, V. (2021). Bayesian hierarchical quantile regression for precision immuno-oncology. *TBD*, 00(0):123–123.
- Guha, N., Baladandyuthapani, V., and Mallick, B. (2020). Quantile graphical models: a bayesian approach. *Journal of machine learning research*.
- Hastie, T. and Tibshirani, R. (1993). Varying-coefficient models. *Journal of the Royal Statistical Society: Series B (Methodological)*, 55(4):757–779.
- Koenker, R. and Bassett Jr, G. (1978). Regression quantiles. *Econometrica: journal of the Econometric Society*, pages 33–50.
- Ni, Y., Stingo, F. C., and Baladandayuthapani, V. (2019). Bayesian graphical regression. *Journal of the American Statistical Association*, 114(525):184–197.
- Yang, Y. and Tokdar, S. T. (2017). Joint estimation of quantile planes over arbitrary predictor spaces. *Journal of the American Statistical Association*, 112(519):1107–1120.

Functional formulation of $\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})$ in simulations:

- 1 For $q^* = 0$, $\{\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})\} = (1 + \tau^2)\mathbf{1}_n$, where $\mathbf{1}_n$ is the unit vector of dimension n .
- 2 For $q^* = 1$, $\{\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})\} = \mathbf{X}_{k_1}^2 + \log((1 + \tau^2)\mathbf{1}_n)$, where k_1 is randomly chosen from $\{1, \dots, q\}$.
- 3 For $q^* = 2$, $\{\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})\} = \mathbf{X}_{k_1}^2 + \log((1 + \tau^2)\mathbf{1}_n) + \exp(\mathbf{X}_{k_2})$, where k_1, k_2 are distinct and randomly chosen from $\{1, \dots, q\}$.
- 4 For $q^* = 3$, $\{\theta_{hj}^{(\tau)}(\mathbf{X}_{i\cdot})\} = \mathbf{X}_{k_1}^2 + \log((1 + \tau^2)\mathbf{1}_n) + \exp(\mathbf{X}_{k_2}) + \log|\mathbf{X}_{k_3}|$, where k_1, k_2, k_3 are distinct and randomly chosen from $\{1, \dots, q\}$.

SUPPLEMENTARY MATERIAL (CONT.)

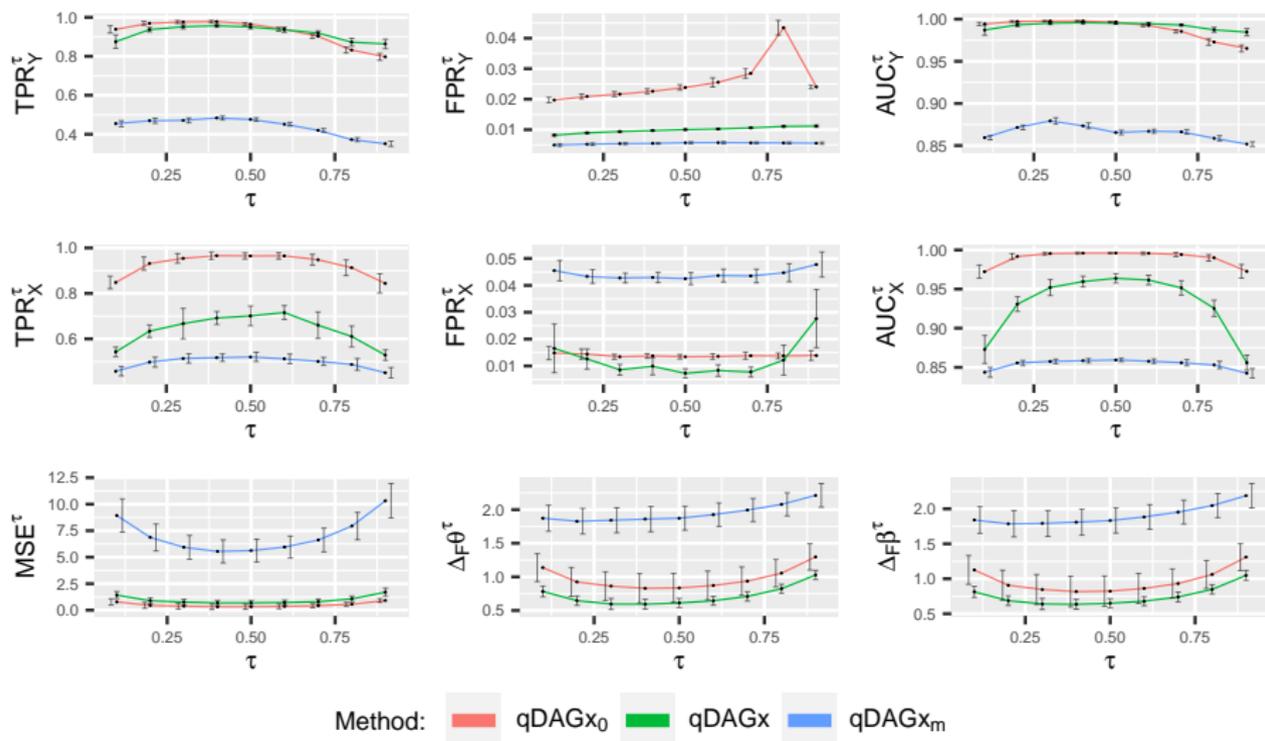


Figure: $p = 50, q = 2, n = 250$. Kendall's' τ for the misspecified sequence is 0.25