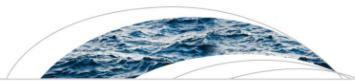# Gaussian conditional independence beyond graphical models

Tobias Boege

Max-Planck Institute for Mathematics in the Sciences, Leipzig

Annual meeting of the IMS
London, 27 June 2022

# A Statistical Graphical Model of the California Reservoir System

**A. Taeb[1] ⓘ, J. T. Reager[2] ⓘ, M. Turmon[2] ⓘ, and V. Chandrasekaran[3]**

[1]Department of Electrical Engineering, California Institute of Technology, Pasadena, CA, USA, [2]Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA, [3]Department of Computing and Mathematical Sciences and Department of Electrical Engineering, California Institute of Technology, Pasadena, CA, USA

**Abstract** The recent California drought has highlighted the potential vulnerability of the state's water management infrastructure to multiyear dry intervals. Due to the high complexity of the network, dynamic storage changes in California reservoirs on a state-wide scale have previously been difficult to model using either traditional statistical or physical approaches. Indeed, although there is a significant line of research on exploring models for single (or a small number of) reservoirs, these approaches are not amenable to a system-wide modeling of the California reservoir network due to the spatial and hydrological heterogeneities of the system. In this work, we develop a state-wide statistical graphical model to characterize the dependencies among a collection of 55 major California reservoirs across the state; this model is defined with respect to a graph in which the nodes refer to reservoirs and the edges specify the

## Gaussian conditional independence

Assume $\xi = (\xi_i : i \in N)$ are jointly Gaussian with covariance matrix $\Sigma \in PD_N$.

### Definition

The polynomial $\Sigma[K] \coloneqq \det \Sigma_{K,K}$ is a *principal minor* of $\Sigma$ and $\Sigma[ij|K] \coloneqq \det \Sigma_{iK,jK}$ is an *almost-principal minor*.

- $\Sigma$ is PD if and only if $\Sigma[K] > 0$ for all $K \subseteq N$.

- $[\xi_i \perp\!\!\!\perp \xi_j \mid \xi_K]$ holds if and only if $\Sigma[ij|K] = 0$.

# Gaussian CI models

## Definition

A *CI constraint* is a CI statement $[\xi_i \perp\!\!\!\perp \xi_j \mid \xi_K]$ or its negation $\neg[\xi_i \perp\!\!\!\perp \xi_j \mid \xi_K]$.
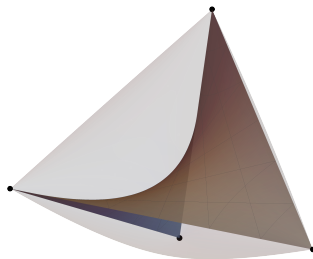The *model* of a set of CI constraints is the set of all PD matrices which satisfy them.



Figure: Model of $\Sigma[12 \mid 3] = a - bc = 0$ in the space of $3 \times 3$ correlation matrices.

## Basic questions

- How hard is it to decide if the model specification is inconsistent?

- How hard is it to *certify* consistency by showing a point in the model?

- What is the geometric structure of the models?

## Basic questions

- How hard is it to decide if the model specification is inconsistent?

- How hard is it to *certify* consistency by showing a point in the model?

- What is the geometric structure of the models?

What is the model of $[X \perp\!\!\!\perp Y] \wedge [X \perp\!\!\!\perp Z \mid Y] \wedge \neg[X \perp\!\!\!\perp Y \mid Z]$ ?

## Models and inference

Consider two sets of CI statements $\mathcal{P}$ and $\mathcal{Q}$:
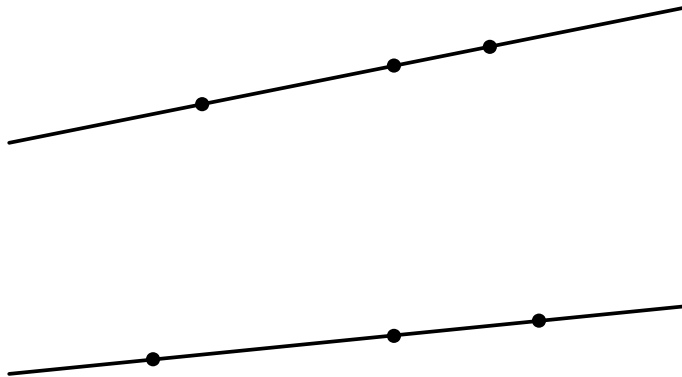
$$\bigwedge \mathcal{P} \Rightarrow \bigvee \mathcal{Q}$$

## Models and inference

Consider two sets of CI statements $\mathcal{P}$ and $\mathcal{Q}$:

$$\bigwedge \mathcal{P} \Rightarrow \bigvee \mathcal{Q}$$
is not valid

$$\iff$$

$$\mathcal{P} \cup \neg \mathcal{Q}$$
has a point

## Models and inference

Consider two sets of CI statements $\mathcal{P}$ and $\mathcal{Q}$:

$$\underset{\text{is not valid}}{\bigwedge \mathcal{P} \Rightarrow \bigvee \mathcal{Q}} \qquad \Longleftrightarrow \qquad \underset{\text{has a point}}{\mathcal{P} \cup \neg \mathcal{Q}}$$

Reasoning about CI statements in normally distributed random variables is the same as reasoning about the vanishing of very special kinds of determinants on very special kinds of varieties inside the positive definite matrices.

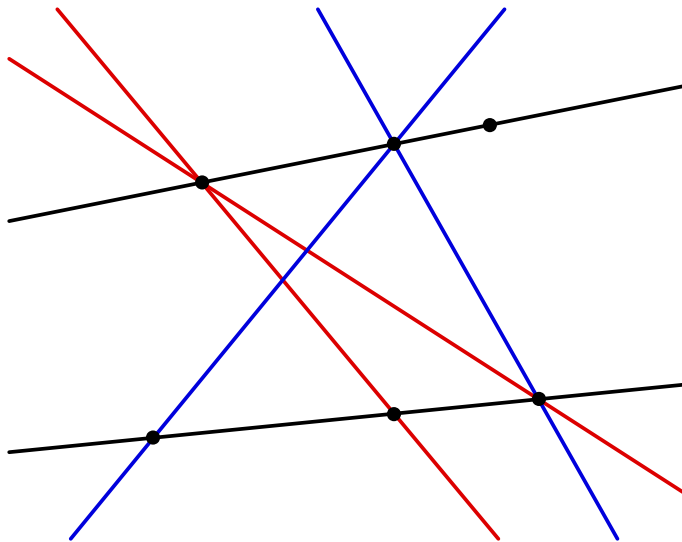# For ancient geometers: conditional independence ≈ collinearity

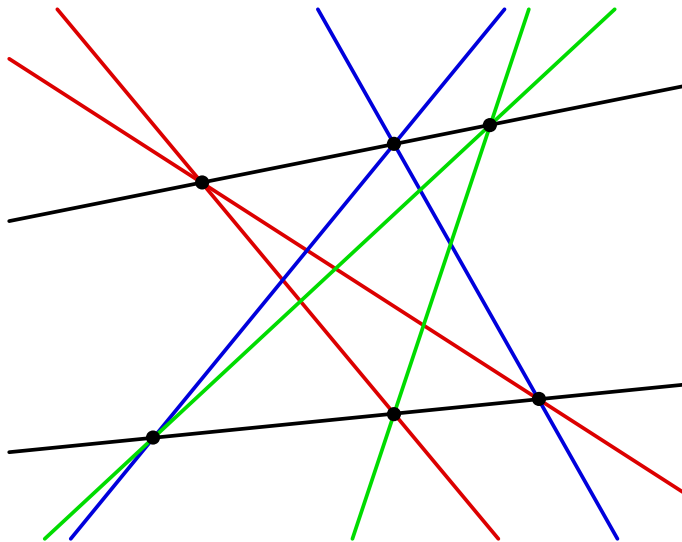# For ancient geometers: conditional independence ≈ collinearity

# For ancient geometers: conditional independence ≈ collinearity

# For ancient geometers: conditional independence ≈ collinearity

## Normal form for proofs and refutations

Let $f_i \in \mathbb{Z}[t_1, \ldots, t_k]$ be integer polynomials in finitely many variables.

Theorem (Tarski's transfer principle)

*If a polynomial system $\{f_i \bowtie_i 0\}$, where $\bowtie_i \in \{=, \neq, <, \leq, \geq, >\}$, has a solution over $\mathbb{R}$, then it has a solution in a finite real extension of $\mathbb{Q}$.*

## Normal form for proofs and refutations

Let $f_i \in \mathbb{Z}[t_1, \ldots, t_k]$ be integer polynomials in finitely many variables.

Theorem (Tarski's transfer principle)

*If a polynomial system $\{f_i \bowtie_i 0\}$, where $\bowtie_i \in \{=, \neq, <, \leq, \geq, >\}$, has a solution over $\mathbb{R}$, then it has a solution in a finite real extension of $\mathbb{Q}$.*

$\rightarrow$ If $\bigwedge \mathcal{P} \Rightarrow \bigvee \mathcal{Q}$ is false, there exists a counterexample matrix $\Sigma$ with algebraic entries.

$[12|\,] \wedge [12|3] \Rightarrow [13|\,]$ is false and a counterexample is

$$\begin{pmatrix} 1 & 0 & {}^1\!/2 \\ 0 & 1 & 0 \\ {}^1\!/2 & 0 & 1 \end{pmatrix}.$$

## Normal form for proofs and refutations

Let $f_i, g_j, h_k \in \mathbb{Z}[t_1, \ldots, t_k]$ be integer polynomials in finitely many variables.

Theorem (Positivstellensatz)

*A polynomial system $\{f_i = 0, g_j \geq 0, h_k \neq 0\}$ is infeasible if and only if there exist $f \in \text{ideal}(f_i)$, $g \in \text{cone}(g_j)$ and $h \in \text{monoid}(h_k)$ such that $g + h^2 = f$.*

## Normal form for proofs and refutations

Let $f_i, g_j, h_k \in \mathbb{Z}[t_1, \ldots, t_k]$ be integer polynomials in finitely many variables.

Theorem (Positivstellensatz)

*A polynomial system $\{f_i = 0, g_j \geq 0, h_k \neq 0\}$ is infeasible if and only if there exist $f \in \text{ideal}(f_i)$, $g \in \text{cone}(g_j)$ and $h \in \text{monoid}(h_k)$ such that $g + h^2 = f$.*

$\rightarrow$ If $\bigwedge \mathcal{P} \Rightarrow \bigvee \mathcal{Q}$ is true, there exists an algebraic proof for it with integer coefficients.

$[12|\,] \wedge [12|3] \Rightarrow [13|\,] \vee [23|\,]$ is true and a proof is the final polynomial

$$\Sigma[13|\,] \cdot \Sigma[23|\,] = \Sigma[3] \cdot \Sigma[12|\,] - \Sigma[12|3].$$

## A 5 × 5 final polynomial

The following inference rule is valid for all positive definite 5 × 5 matrices:

$[12\,|\,] \wedge [14\,|\,5] \wedge [23\,|\,5] \wedge [35\,|\,1] \wedge [45\,|\,2] \wedge [15\,|\,23] \wedge [34\,|\,12] \wedge [24\,|\,135] \; \Rightarrow \; [25\,|\,] \vee [34\,|\,].$

## A $5 \times 5$ final polynomial

The following inference rule is valid for all positive definite $5 \times 5$ matrices:

$$[12\,|\,] \wedge [14\,|\,5] \wedge [23\,|\,5] \wedge [35\,|\,1] \wedge [45\,|\,2] \wedge [15\,|\,23] \wedge [34\,|\,12] \wedge [24\,|\,135] \;\Rightarrow\; [25\,|\,] \vee [34\,|\,].$$

$$[25\,|\,][34\,|\,] \cdot [1][2][3][15] =$$

$$\Big( cd^2egr + bd^2fgr - ad^2grh - 2cd^2e^2i - 2bd^2efi - 2pdfgri + 2ad^2ehi + 2pdefi^2 - 2pdqhi^2 + 2pcqi^3 +$$

$$2pdqrij - 2pbqi^2j - pcegrt + pbfgrt + pagrht + 2pce^2it - 2pcqrit + 2pbqhit - 2paehit \Big) \cdot [12\,|\,] +$$

$$\Big( pdqer + pbqgr - 2pbqei \Big) \cdot [14\,|\,5] - \Big( pcdqr + p^2fgr - 2pbcqi + 2pb^2qj - 2p^2qrj \Big) \cdot [23\,|\,5] +$$

$$\Big( cdqgr - 2cdqei + 2pqghi - 2pqfi^2 - pqgrj + 2pqeij - 2pe^2ft + 2pqfrt \Big) \cdot [35\,|\,1] +$$

$$\Big( pd^2er - 2pbdei + p^2gri + 2pb^2et - 2p^2ert \Big) \cdot [45\,|\,2] - \Big( 2pdfi - 2pbft \Big) \cdot [15\,|\,23] -$$

$$\Big( d^2gr - 2d^2ei - pgrt + 2peit \Big) \cdot [34\,|\,12] - 2pqi \cdot [24\,|\,135].$$

## A 5 × 5 final polynomial

```
R = QQ[p,a,b,c,d, q,e,f,g, r,h,i, s,j, t];
X = genericSymmetricMatrix(R,p,5);
I = ideal(
  det X_{0}^{1}, det X_{0,3}^{2,3}, det X_{0,4}^{3,4},
  det X_{1,4}^{2,4}, det X_{2,0}^{4,0}, det X_{3,1}^{4,1},
  det X_{0,1,2}^{4,1,2}, det X_{2,0,1}^{3,0,1},
  det X_{1,0,2,4}^{3,0,2,4}
);
U = g*h*p*q*r*(p*t-d^2); -- [25|][34|]·[1][2][3][15] ∈ monoid(𝒱)
U % I --> 0, meaning monoid(𝒱) ∩ ideal(𝒱) ≠ ∅ in ℚ[X]
-- Get a proof that U is in I:
G = gens I; -- the equations generating ideal(𝒱)
H = U // G; -- linear combinators for U from G
U == G*H --> true
```

## Consistency checking is hard

The complexity class $\exists \mathbb{R}$ contains all decision problems which can be reduced in polynomial time to the feasibility of a semialgebraic set:

- polynomial optimization
- computational geometry
- algebraic statistics …

## Consistency checking is hard

The complexity class $\exists\mathbb{R}$ contains all decision problems which can be reduced in polynomial time to the feasibility of a semialgebraic set:

- polynomial optimization
- computational geometry
- algebraic statistics …

### Theorem

*The problem of deciding whether a general CI model is non-empty is complete for $\exists\mathbb{R}$.*

(Graphical models are always consistent.)

# Consistency certification is hard

Šimeček's Question

*Does every non-empty Gaussian CI model contain a rational point?*

## Consistency certification is hard

### Šimeček's Question

*Does every non-empty Gaussian CI model contain a rational point?*

### Theorem

*For every finite real extension $\mathbb{K}$ of $\mathbb{Q}$ there exists a CI model $\mathcal{M}$ such that $\mathcal{M} \cap \mathrm{PD}_N(\mathbb{K}) \neq \varnothing$ but $\mathcal{M} \cap \mathrm{PD}_N(\mathbb{L}) = \varnothing$ for all proper subfields $\mathbb{L} \subsetneq \mathbb{K}$.*

(Graphical models always have rational points.)

# Model topology can be bad

An oriented CI model is specified by sign constraints on partial correlations.

### Theorem

*For every primary basic semialgebraic set $Z$ there exists an oriented CI model $\mathcal{M}$ which is homotopy-equivalent to $Z$.*

(Graphical models are always contractible.)