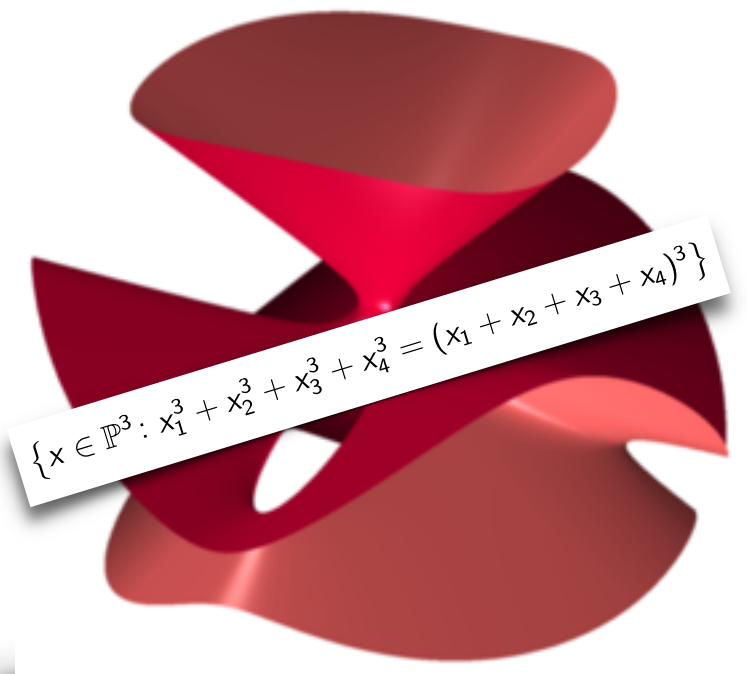


Algebraic Statistics & Quantifier Elimination

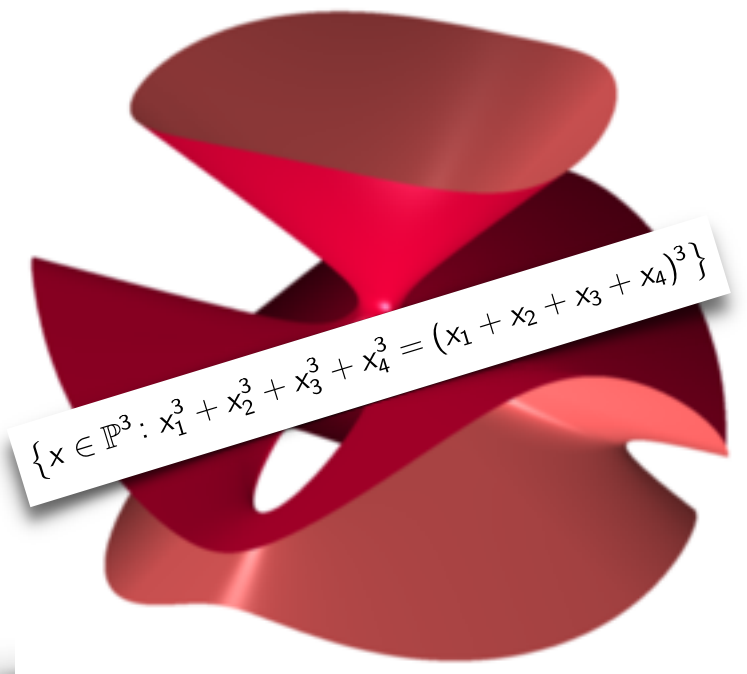
Daniel Suess

What is Algebraic Statistics?



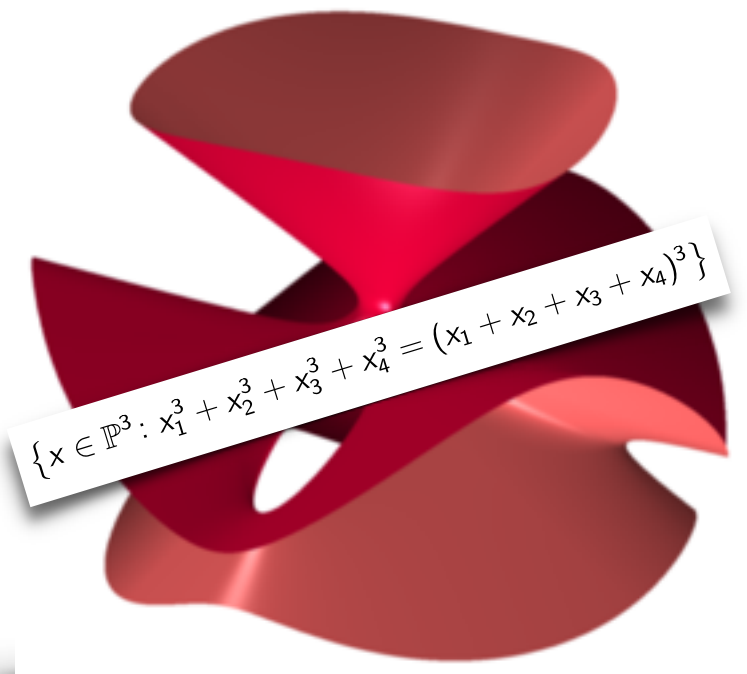
What is Algebraic Statistics?

understand algebraic
structure of statistical models



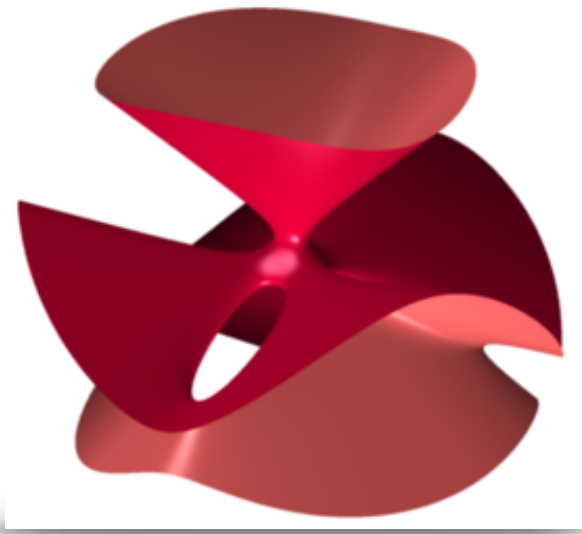
What is Algebraic Statistics?

understand algebraic
structure of statistical models



application of techniques from
computational algebra,
algebraic geometry, ...

Algebraic Structure in Statistics

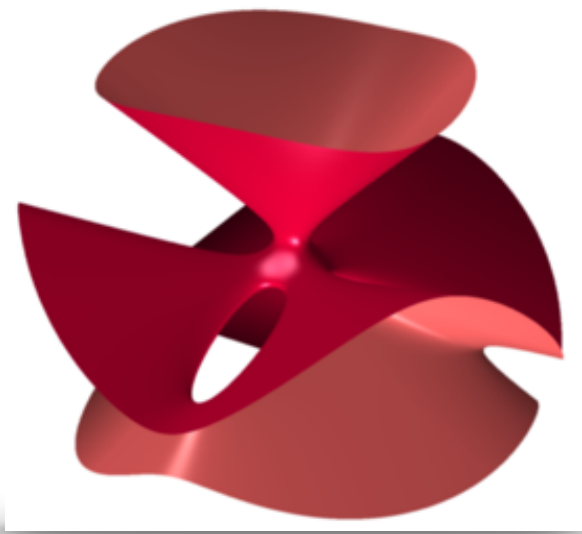


Semi-algebraic set $S \subset \mathbb{R}^n$

$$S = \left\{ x \in \mathbb{R}^n : \begin{aligned} &p_i(x_1, \dots, x_n) = 0, \\ &q_j(x_1, \dots, x_n) \leq 0 \end{aligned} \right\}$$

for finitely many polynomials p_i ,
 q_j with coefficients in \mathbb{Q}

Algebraic Structure in Statistics



Semi-algebraic set $S \subset \mathbb{R}^n$

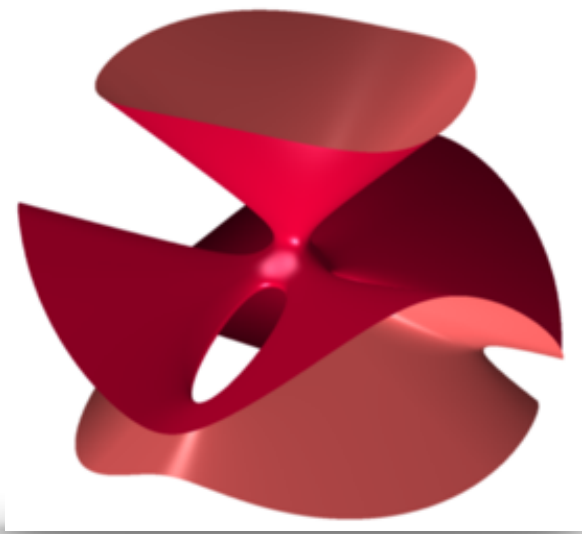
$$S = \left\{ x \in \mathbb{R}^n : \begin{aligned} &p_i(x_1, \dots, x_n) = 0, \\ &q_j(x_1, \dots, x_n) \leq 0 \end{aligned} \right\}$$

for finitely many polynomials p_i ,
 q_j with coefficients in \mathbb{Q}



- probability simplex
$$\Delta_{n-1} = \left\{ x \in \mathbb{R}^n : \sum_i x_i = 1, x_i \geq 0 \right\}$$
- joint prob. of discrete, indep. RVs
- max. likelihood estimates
- ...

Algebraic Structure in Statistics



Semi-algebraic set $S \subset \mathbb{R}^n$

$$S = \left\{ x \in \mathbb{R}^n : \begin{aligned} &p_i(x_1, \dots, x_n) = 0, \\ &q_j(x_1, \dots, x_n) \leq 0 \end{aligned} \right\}$$

for finitely many polynomials p_i , q_j with coefficients in \mathbb{Q}

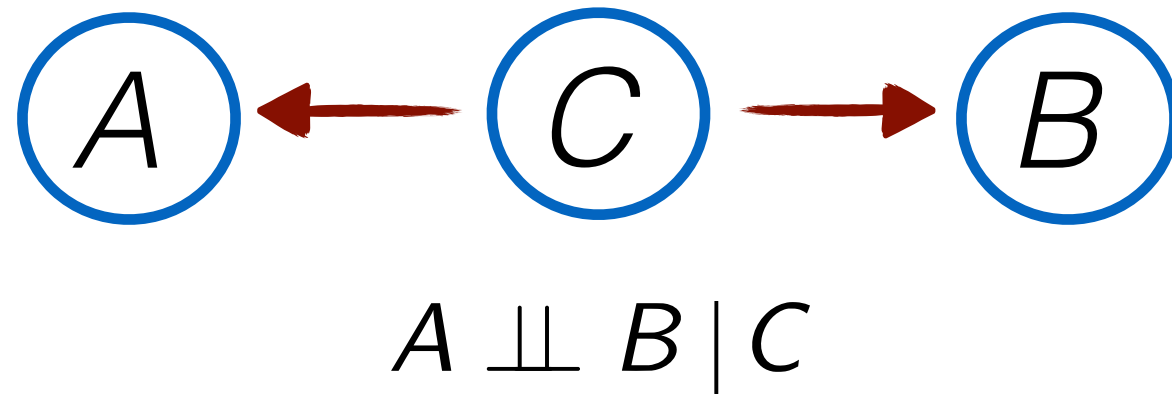


- probability simplex

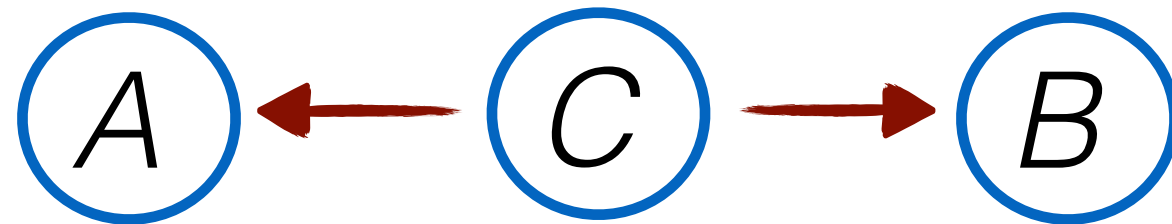
$$\Delta_{n-1} = \left\{ x \in \mathbb{R}^n : \sum_i x_i = 1, x_i \geq 0 \right\}$$

- joint prob. of discrete, indep. RVs
- max. likelihood estimates
- ...

Conditional Independence: Explicit & Implicit Form



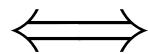
Conditional Independence: Explicit & Implicit Form



$$A \perp\!\!\!\perp B \mid C$$

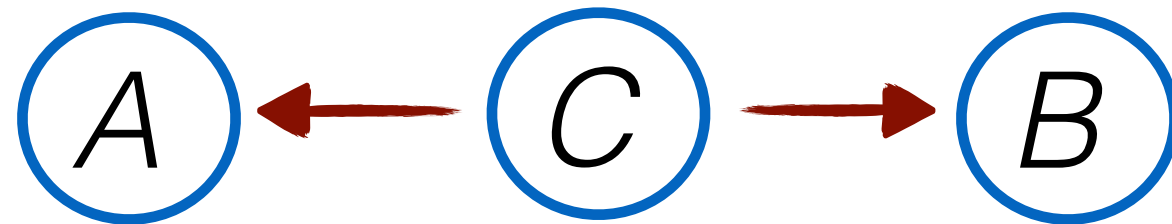


$$P(A, B, C) = P(A|C) P(B|C) P(C)$$



$$p_{abc} = q_{ac} r_{bc} s_c$$

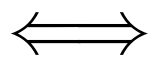
Conditional Independence: Explicit & Implicit Form



$$A \perp\!\!\!\perp B \mid C$$



$$P(A, B, C) = P(A|C) P(B|C) P(C)$$



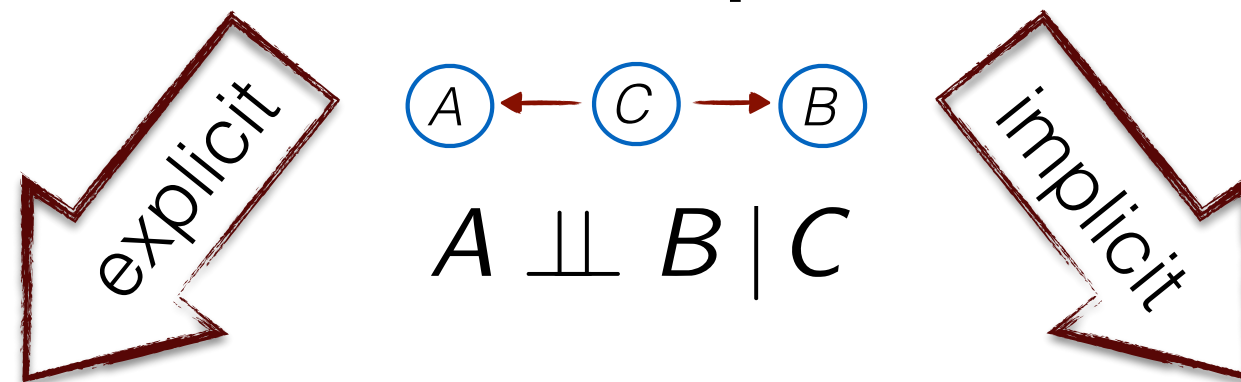
$$p_{abc} = q_{ac} r_{bc} s_c$$



$$p_{abc} p_{a'b'c} - p_{ab'c} p_{a'bc} = 0$$

Additional constraint for the semi-algebraic set of probabilities consistent with DAG.

Conditional Independence: Explicit & Implicit Form



$$P(A, B, C) = P(A|C) P(B|C) P(C)$$

\iff

$$p_{abc} = q_{ac} r_{bc} s_c$$

$$p_{abc} p_{a'b'c} - p_{ab'c} p_{a'bc} = 0$$

Additional constraint for the semi-algebraic set of probabilities consistent with DAG.

$$p_{abc} = q_{ac} r_{bc} s_c$$

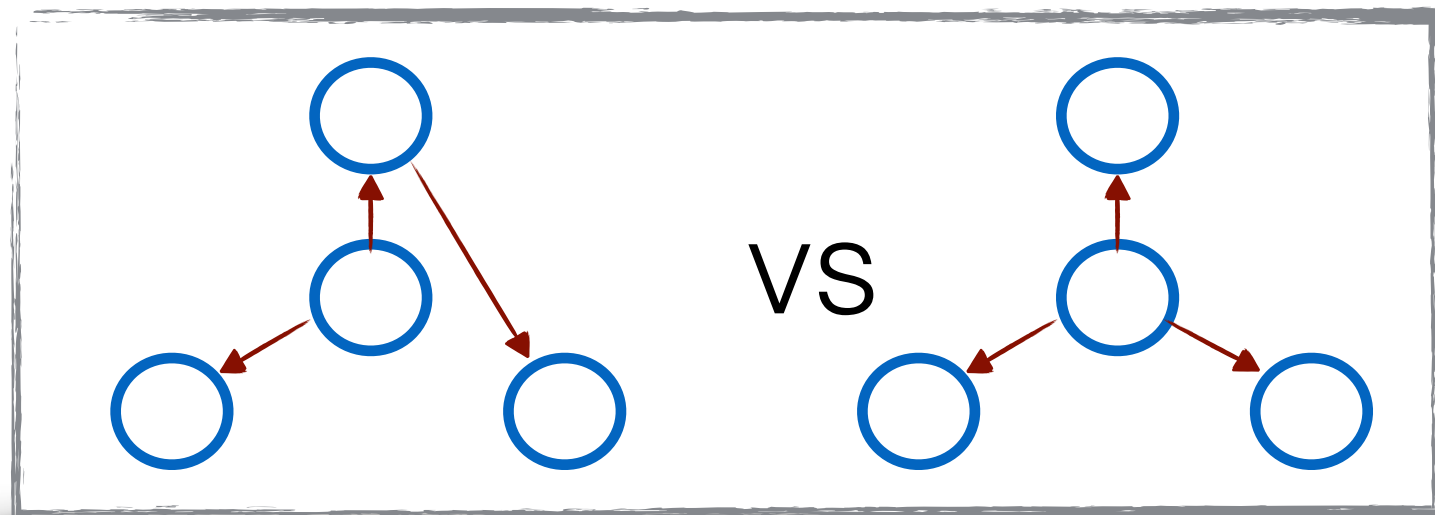
\iff for all c , $(p_{abc})_{ab}$ is rank 1

$\iff (p_{abc})_{ab}$ has determinantal rank 1

\iff determinant of all 2×2 minors = 0

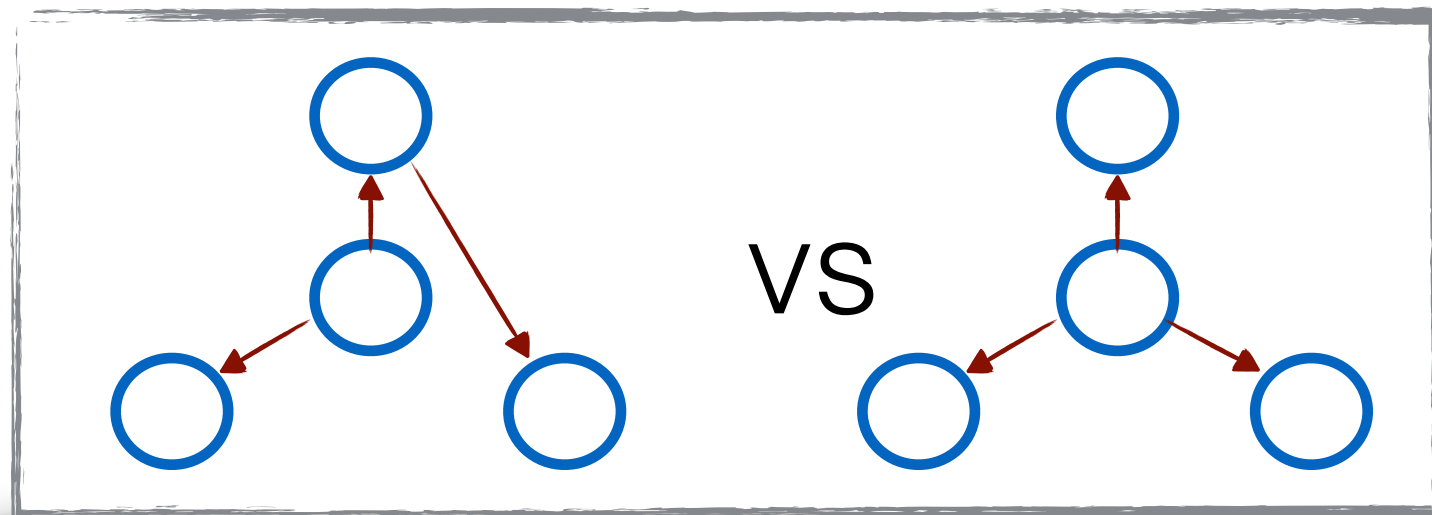
$\iff \forall a, a', b, b', c: p_{abc} p_{a'b'c} - p_{ab'c} p_{a'bc} = 0$

Possible Applications



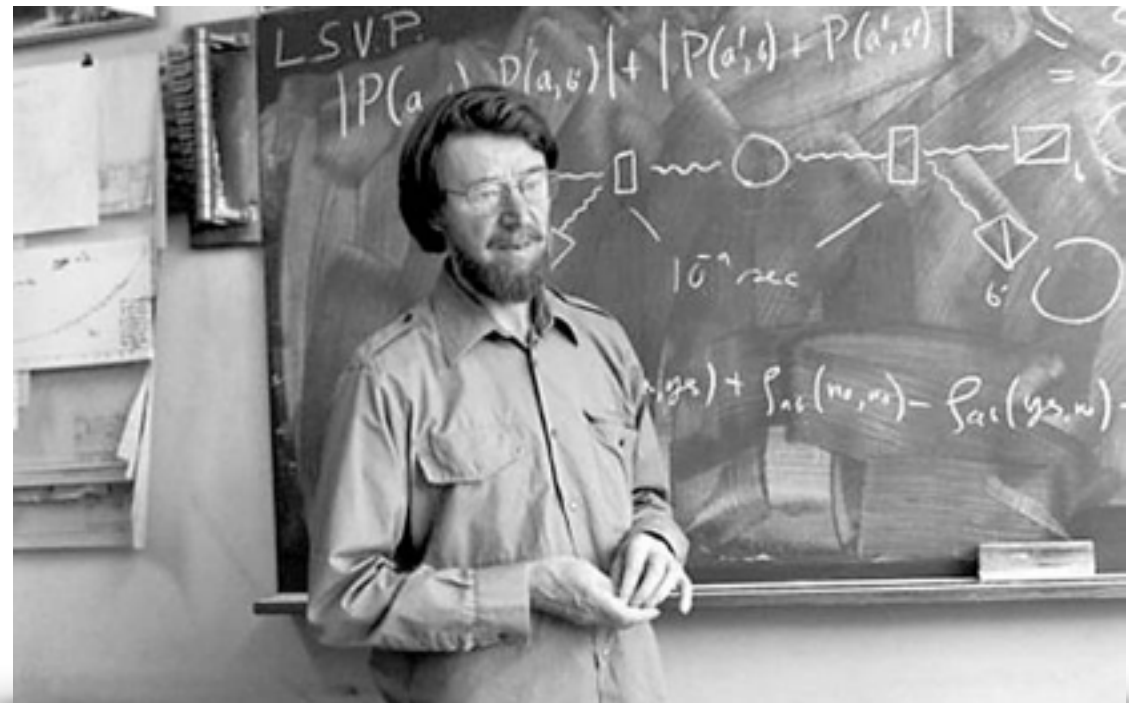
distinguishing Bayesian
networks from
observational data

Possible Applications

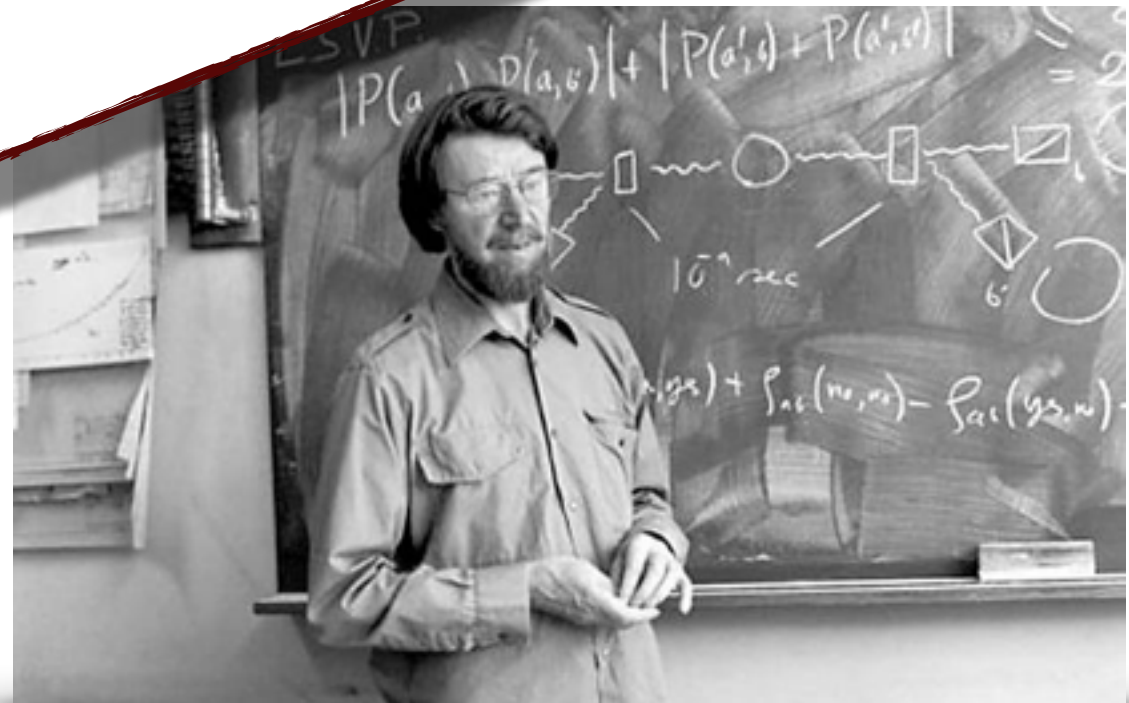
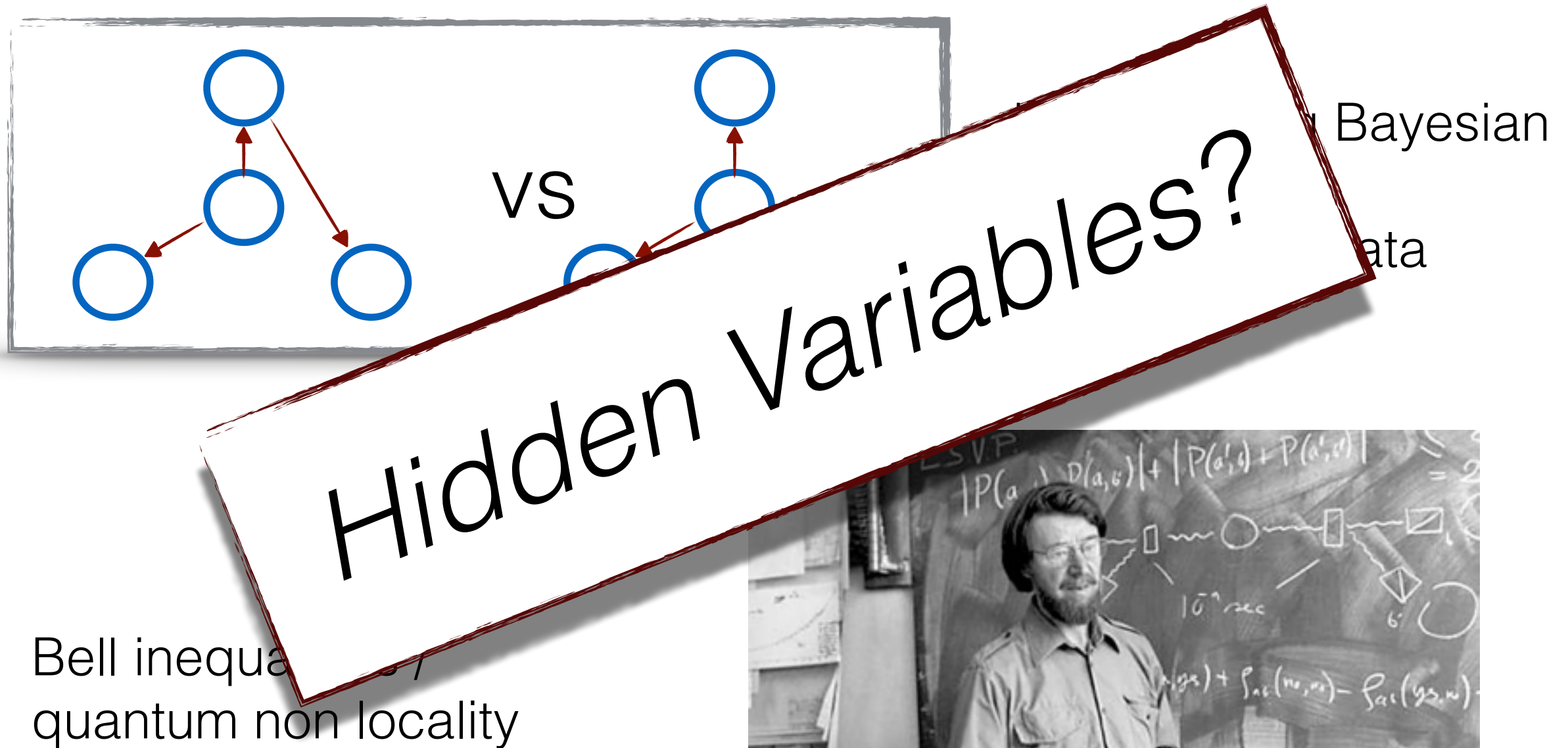


distinguishing Bayesian networks from observational data

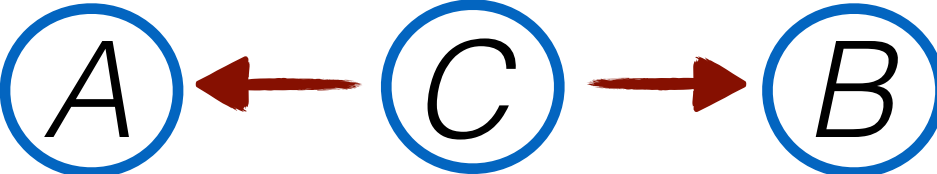
Bell inequalities /
quantum non locality



Possible Applications



Quantifier Elimination: Conditional Independence

p_{abc} compatible with 

$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

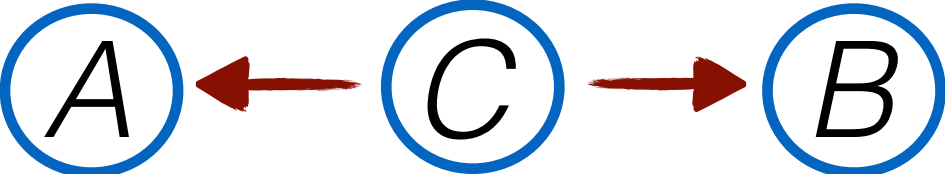
Quantifier Elimination: Conditional Independence

p_{abc} compatible with $A \leftarrow C \rightarrow B$

$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

$$P(A, B, C) = P(A|C) P(B|C) P(C)$$

Quantifier Elimination: Conditional Independence

p_{abc} compatible with 

$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

$$= \pi \left(\{ (p_{abc}, q_{ab}, r_{bc}, s_c) : \dots \} \right)$$

$$\text{with } \pi : \mathbb{R}^{n_1+n_2} \rightarrow \mathbb{R}^{n_1}, (p_{abc}, q_{ab}, r_{bc}, s_c) \mapsto (p_{abc})$$

Quantifier Elimination: Conditional Independence

p_{abc} compatible with 

$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

$$= \pi \left(\{ (p_{abc}, q_{ab}, r_{bc}, s_c) : \dots \} \right)$$

$$\text{with } \pi : \mathbb{R}^{n_1+n_2} \rightarrow \mathbb{R}^{n_1}, (p_{abc}, q_{ab}, r_{bc}, s_c) \mapsto (p_{abc})$$

Theorem (Tarski-Seidenberg) The image of a semi-algebraic set under a projection map π is a semi-algebraic set.

Quantifier Elimination: Conditional Independence

p_{abc} compatible with 

$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

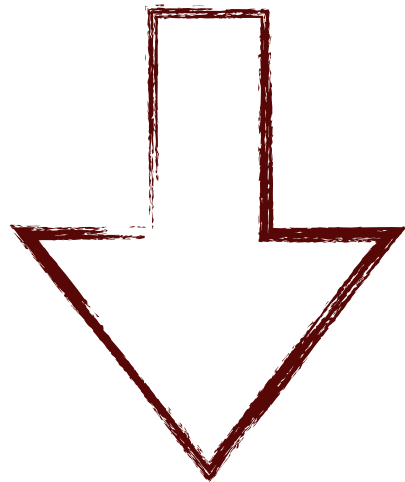
$$= \pi \left(\{ (p_{abc}, q_{ab}, r_{bc}, s_c) : \dots \} \right)$$

$$\text{with } \pi : \mathbb{R}^{n_1+n_2} \rightarrow \mathbb{R}^{n_1}, (p_{abc}, q_{ab}, r_{bc}, s_c) \mapsto (p_{abc})$$

$$\begin{aligned} & \{(p, q) \in \mathbb{R}^2 : (\exists x \in \mathbb{R}) x^2 + px + q = 0\} \\ &= \{(p, q) \in \mathbb{R}^2 : p^2 \geq 4q\} \end{aligned}$$

Quantifier Elimination: Algorithms

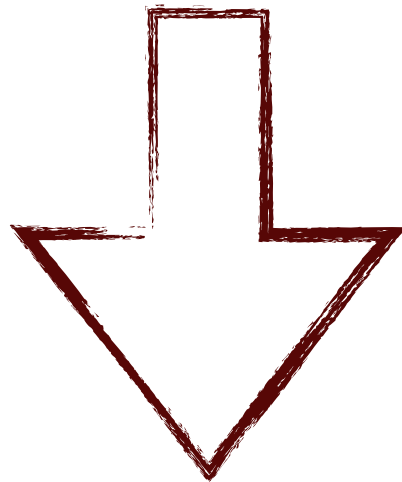
Input: formula ψ



Output: equivalent,
quantifier-free formula ψ'

Quantifier Elimination: Algorithms

Input: formula ψ

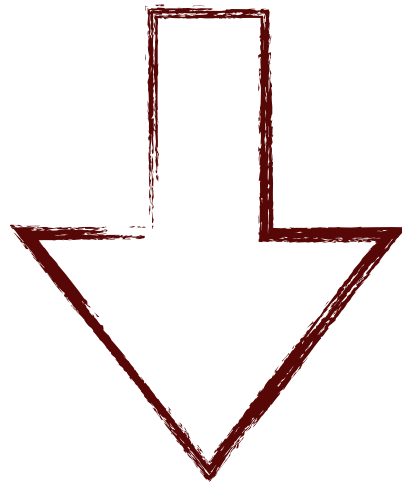


Output: equivalent,
quantifier-free formula ψ'

- real valued variables
- rational constants
- operations (+, -, ×)
- binary relations (=, ≠, <, ≤)
- logical connectives
($\wedge, \vee, \neg, \Rightarrow, \Leftrightarrow$)
- quantifiers (\forall, \exists)

Quantifier Elimination: Algorithms

Input: formula ψ

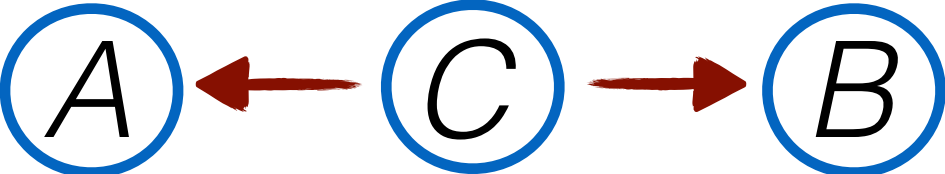


Output: equivalent,
quantifier-free formula ψ'

- real valued variables
- rational constants
- operations (+, -, ×)
- binary relations (=, ≠, <, ≤)
- logical connectives ($\wedge, \vee, \neg, \Rightarrow, \Leftrightarrow$)
- quantifiers (\forall, \exists)

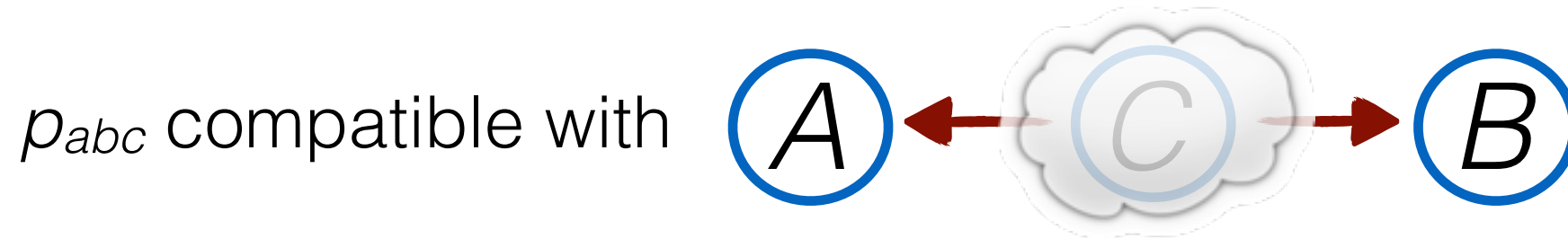
- Tarski's algorithm
- Cylinder Algebraic decomposition

Quantifier Elimination: Hidden Variables

p_{abc} compatible with 

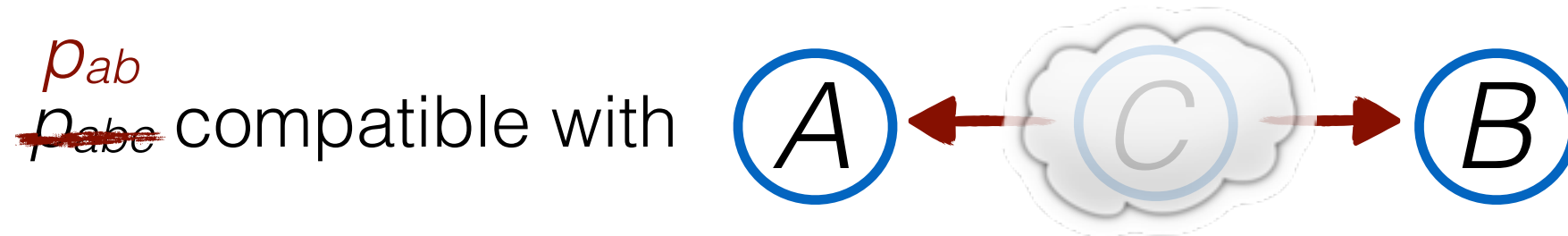
$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

Quantifier Elimination: Hidden Variables



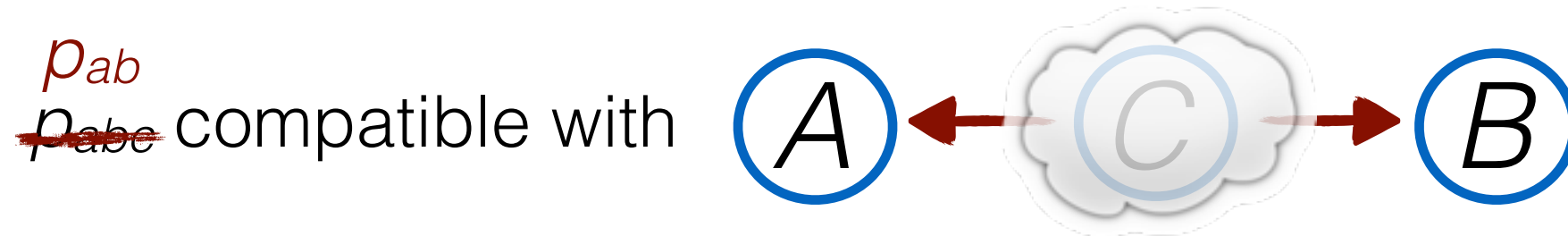
$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

Quantifier Elimination: Hidden Variables



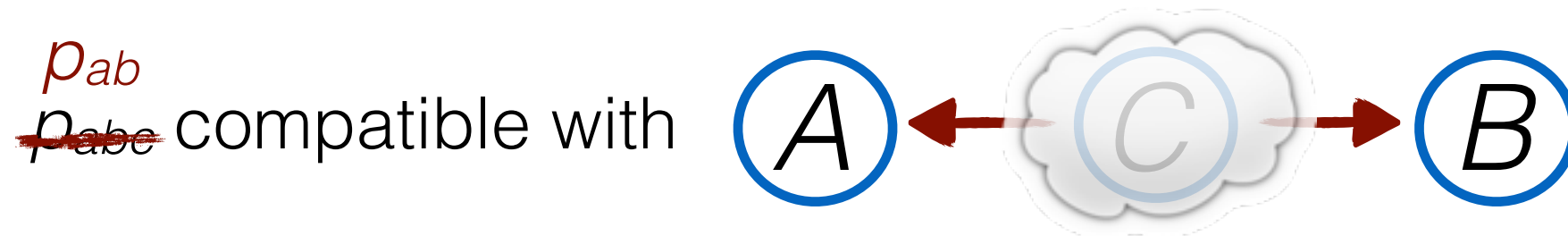
$$\left\{ \begin{array}{l} (p_{abc}) : \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ (\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

Quantifier Elimination: Hidden Variables



$$\left\{ \begin{array}{l} \textcolor{red}{(p_{ab})} \\ \textcolor{red}{(\cancel{p_{abc}})}: \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\wedge_c \sum_a q_{ac} = 1) \wedge (\wedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\wedge_{ac} q_{ac} \geq 0) \wedge (\wedge_{bc} r_{bc} \geq 0) \wedge (\wedge_c s_c \geq 0) \wedge \\ (\wedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c) \end{array} \right\}$$

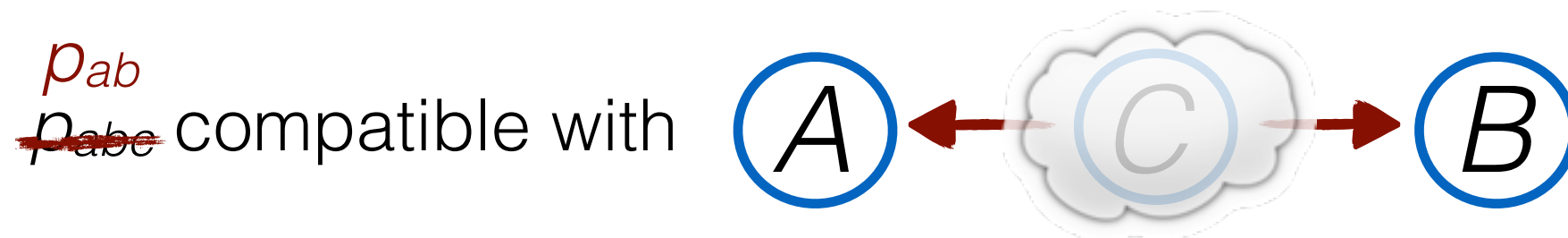
Quantifier Elimination: Hidden Variables



$$\left\{ \begin{array}{l} \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ \cancel{(\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c)} \end{array} \right\}$$

$$\bigwedge_{ab} p_{ab} = \sum_c q_{ac} \times r_{bc} \times s_c$$

Quantifier Elimination: Hidden Variables



$$\left\{ \begin{array}{l} \textcolor{red}{(p_{ab})} \\ \textcolor{red}{(p_{abc})}: \quad \exists(q_{ac})\exists(r_{bc})\exists(s_c) \\ (\bigwedge_c \sum_a q_{ac} = 1) \wedge (\bigwedge_c \sum_b r_{bc} = 1) \wedge (\sum_c s_c = 1) \wedge \\ (\bigwedge_{ac} q_{ac} \geq 0) \wedge (\bigwedge_{bc} r_{bc} \geq 0) \wedge (\bigwedge_c s_c \geq 0) \wedge \\ \textcolor{red}{(\bigwedge_{abc} p_{abc} = q_{ac} \times r_{bc} \times s_c)} \end{array} \right\}$$

$$\bigwedge_{ab} \textcolor{red}{p_{ab}} = \sum_c \textcolor{red}{q_{ac}} \times \textcolor{red}{r_{bc}} \times \textcolor{red}{s_c}$$

- no additional constraints on p_{ab} for $N_A = N_B = N_C = 2$
- $N_A = N_B = 3$ did not finish

Quantifier Elimination: Other Applications

Membership Problem

$$(A \perp C) \wedge (A \perp B \mid C) \Rightarrow (A \perp B) ?$$

Quantifier Elimination: Other Applications

Membership Problem

$$\frac{(A \perp C) \wedge (A \perp B \mid C)}{\psi(p)} \Rightarrow \frac{(A \perp B)}{\psi'(p)} ?$$

Quantifier Elimination: Other Applications

Membership Problem

$$\frac{(A \perp C) \wedge (A \perp B \mid C)}{\psi(p)} \Rightarrow \frac{(A \perp B)}{\psi'(p)} ?$$

$$\forall p: \psi(p) \implies \psi'(p)$$

Quantifier Elimination: Other Applications

Membership Problem

$$\frac{(A \perp C) \wedge (A \perp B \mid C)}{\psi(p)} \Rightarrow \frac{(A \perp B)}{\psi'(p)} ?$$

$$\forall p: \psi(p) \implies \psi'(p)$$

Identifiability Problem

Can we uniquely identify the parameter θ in a parametric model from given observation?

Quantifier Elimination: Other Applications

Membership Problem

$$\frac{(A \perp C) \wedge (A \perp B \mid C)}{\psi(p)} \Rightarrow \frac{(A \perp B)}{\psi'(p)} ?$$

$$\forall p: \psi(p) \implies \psi'(p)$$

Identifiability Problem

Can we uniquely identify the parameter θ in a parametric model from given observation?

$$\forall \theta, \theta': (O(\theta) = O(\theta')) \implies (\theta = \theta')$$

Conclusion

Conclusion

- probabilities compatible with directed Bayesian network constitute semi-algebraic set:

Conclusion

- probabilities compatible with directed Bayesian network constitute semi-algebraic set:
- no hidden variables \rightarrow quadratic equations

Conclusion

- probabilities compatible with directed Bayesian network constitute semi-algebraic set:
 - no hidden variables \rightarrow quadratic equations
 - hidden variables \rightarrow possibly QE

Conclusion

- probabilities compatible with directed Bayesian network constitute semi-algebraic set:
 - no hidden variables \rightarrow quadratic equations
 - hidden variables \rightarrow possibly QE
- CAD: best known algorithm for QE, but too slow beyond toy models

Conclusion

- probabilities compatible with directed Bayesian network constitute semi-algebraic set:
 - no hidden variables \rightarrow quadratic equations
 - hidden variables \rightarrow possibly QE
- CAD: best known algorithm for QE, but too slow beyond toy models
- also for continuous RV (exponential families)

Conclusion

- probabilities compatible with directed Bayesian network constitute semi-algebraic set:
 - no hidden variables \rightarrow quadratic equations
 - hidden variables \rightarrow possibly QE
- CAD: best known algorithm for QE, but too slow beyond toy models
- also for continuous RV (exponential families)
- <https://github.com/dseuss/algstat.git>

Theorem (Tarski) Every Tarski sentence (formula without free variables) is equivalent to a quantifier-free formula.

Theorem (Tarski-Seidenberg) Let $f: X \rightarrow Y$ be a semi-algebraic map. Then, the image $f(X) \subseteq Y$ is a semi-algebraic set.

ALGORITHM ComputeCAD(F,j)

Input: $F \subset Q[x_1, \dots, x_j]$.

Output: (K_j, α_j) where K_j is an F -sign-invariant CAD of R^j and α_j is a set of algebraic sample points, one per cell in K_j .

Recurse: If $j > 1$, then do

$\Phi(F) := \Phi_1(F) \cup \Phi_2(F) \cup \Phi_3(F)$

$(K_{j-1}, \alpha_{j-1}) := \text{ComputeCAD}(\Phi(F), j-1),$

else find the roots r_1, \dots, r_m of all polynomials in F and do

$K_1 := \{[-\infty, r_1), [r_1, r_1], (r_1, r_2), \dots, (r_m, +\infty]\}$

$\alpha_1 := \{r_1 - 1, r_1, (r_1 + r_2)/2, \dots, r_m, r_m + 1\}$

Return (K_1, α_1)

Lift: For every cell $C_i \in K_{j-1}$ do

1. Compute the product of all polynomials in F that do not vanish at the sample point α_i of C_i and call the resulting polynomial $\pi(\alpha_i, x)$

2. Find the roots r_1, \dots, r_m of $\pi(\alpha_i, x)$

3. Set $K_{j,i} := \{\{C_i \times [-\infty, r_1)\}, \{C_i \times [r_1, r_1]\}, \{C_i \times (r_1, r_2)\}, \dots, \{C_i \times (r_m, +\infty]\}\}$

Comment: $K_{j,i}$ are the cylinders over C_i .

4. Set $\alpha_{j,i} := \{(\alpha_i, r_1 - 1), (\alpha_i, r_1), (\alpha_i, (r_1 + r_2)/2), \dots, (\alpha_i, r_m + 1)\}$

Comment: $\alpha_{j,i}$ are the algebraic sample points for the cylinders over C_i .

$K_j := \bigcup_i K_{j,i} ; \quad \alpha_j := \bigcup_i \alpha_{j,i}$

Return (K_j, α_j)