

Intelligent Agents: Web-mining Agents

Probabilistic Graphical Models

Continuous Space

Tanya Braun



UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

Probabilistic Graphical Models (PGMs)

1. Recap: **Propositional** modelling

- Factor model, Bayesian network, Markov network
- Semantics, inference tasks + algorithms + complexity

2. **Probabilistic relational models (PRMs)**

- Parameterised models, Markov logic networks
- Semantics, inference tasks

3. **Lifted inference**

- LVE, LJT, FOKC
- Theoretical analysis

4. **Lifted learning**

- Recap: propositional learning
- From ground to lifted models
- Direct lifted learning

5. **Approximate Inference: Sampling**

- Importance sampling
- MCMC methods

6. **Sequential models & inference**

- Dynamic PRMs
- Semantics, inference tasks + algorithms + complexity, learning

7. **Decision making**

- (Dynamic) Decision PRMs
- Semantics, inference tasks + algorithms, learning

8. **Continuous Space**

- Gaussian distributions and Bayesian networks
- Probabilistic soft logic

Models with Continuous Variables

- Discretisation of continuous variables
 - Discrete model again
 - Own set of problems
 - Hard to find good discretisation
 - High granularity might be necessary
→ large ranges → large factors
 - Lose characteristics of variable
 - Not each value necessarily associated with a probability
 - Nearby values have similar probabilities → hard to capture in a discrete distribution (no notion of closeness between range values)
- Therefore, use models with continuous variables

Outline: 8. Continuous Space

A. *Basics*

- Continuous variables, probability density function, cumulative probability distribution
- Joint distribution, marginal density, conditional density

B. *Gaussian models*

- (Multivariate) Gaussian distribution
- (Parameterised) Gaussian Bayesian networks

C. *Probabilistic Soft Logic (PSL)*

- Modelling, semantics, inference task

Probability Density Function

- Continuous random variable R
 - Range $\mathcal{R}(R) = [0,1]$
- Function $p : \mathbb{R} \rightarrow \mathbb{R}$ is a **probability density function (PDF)** for R if it is a non-negative, integrable function s.t.

$$\int_{\mathcal{R}(R)} p(r) dr = 1$$

- For any a (and b) in event space

$$P(R \leq a) = \int_{-\infty}^a p(r) dr \quad P(a \leq R \leq b) = \int_a^b p(r) dr$$

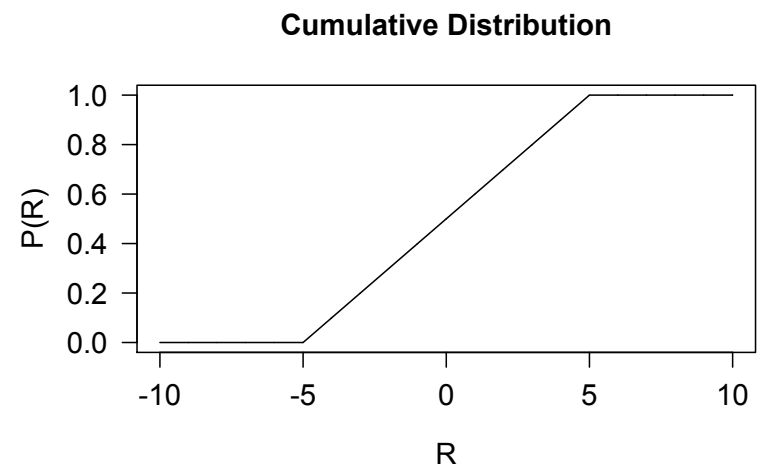
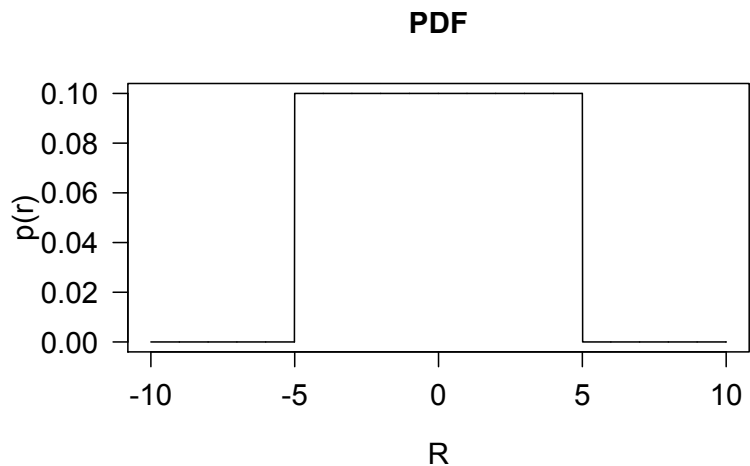
- Function P is a cumulative distribution for R
- Intuitively, value of $p(r)$ at point r is the incremental amount that r adds to the cumulative distribution during integration

PDFs: Uniform Distribution

- Continuous random variable R has a uniform distribution over $[a, b]$, denoted $R \sim \text{Unif}[a, b]$, if it has the PDF

$$p(r) = \begin{cases} \frac{1}{b-a} & b \geq r \geq a \\ 0 & \text{otherwise} \end{cases}$$

- Density can be larger than 1 if $b - a < 1$
 - Can be legal if the total area under the pdf is 1



Joint/Multivariate Distribution

- Let P be a joint distribution over continuous randvars R_1, \dots, R_n
- Function $p(r_1, \dots, r_n)$ is a joint density function of R_1, \dots, R_n if
 - $p(r_1, \dots, r_n) \geq 0$ for all values r_1, \dots, r_n of R_1, \dots, R_n
 - p is an integrable function
 - For any choice a_1, \dots, a_n and b_1, \dots, b_n ,

$$P(a_1 \leq R_1 \leq b_1, \dots, a_n \leq R_n \leq b_n)$$

$$= \int_{a_1}^{b_1} \dots \int_{a_n}^{b_n} p(r_1, \dots, r_n) dr_1 \dots r_n$$

Marginal Density

- Given a joint density, integrate out the non-query randvars
 - E.g., given $p(r, s)$ a joint density for randvars R, S , then

$$p(r) = \int_{-\infty}^{+\infty} p(r, s) ds$$

- Shorthand notations
 - $p_R = p(r)$ marginal density
 - $p_{R,S} = p(r, s)$ joint density

Conditional Density Function

- Discrete case: $P(S|R = r) = \frac{P(S, R=r)}{P(R=r)}$
 - Problem in continuous case: $P(R = r) = 0$
→ $P(S|R = r)$ undefined
- To avoid problem, condition on event $r - \epsilon \leq R \leq r + \epsilon$ and consider limit when $\epsilon \rightarrow 0$
$$P(S|r) = \lim_{\epsilon \rightarrow 0} P(S|r - \epsilon \leq R \leq r + \epsilon)$$
- If a continuous joint density $p(r, s)$ exists, derive form of this expression:

$$p(s|r) = \frac{p(r, s)}{p(r)}$$

- If $p(r) = 0$, conditional density undefined
- Chain rule and Bayes' rule hold as well:

$$p(r, s) = p(r)p(s|r) \qquad p(s|r) = \frac{p(s)p(r|s)}{p(r)}$$

Outline: 8. Continuous Space

A. *Basics*

- Continuous variables, probability density function, cumulative probability distribution
- Joint distribution, marginal density, conditional density

B. ***Gaussian models***

- (Multivariate) Gaussian distribution
- (Parameterised) Gaussian Bayesian networks

C. *Probabilistic Soft Logic (PSL)*

- Modelling, semantics, inference task

Models with Continuous Variables

- Problem: Space of possible parameterisation essentially unbounded
- Special case: (Multivariate) Gaussian distributions
 - Two parameters per variable: mean, variance
 - Strong assumptions, e.g.,
 - Exponential decay away from its mean
 - Linearity of interactions between randvars
 - Assumptions often invalid but still work as a good approximation for many real-world distributions
 - Many generalisations exist which use Gaussians as a foundation
 - Non-linear interactions
 - Mixture of Gaussians

PDFs: Gaussian/Normal Distribution

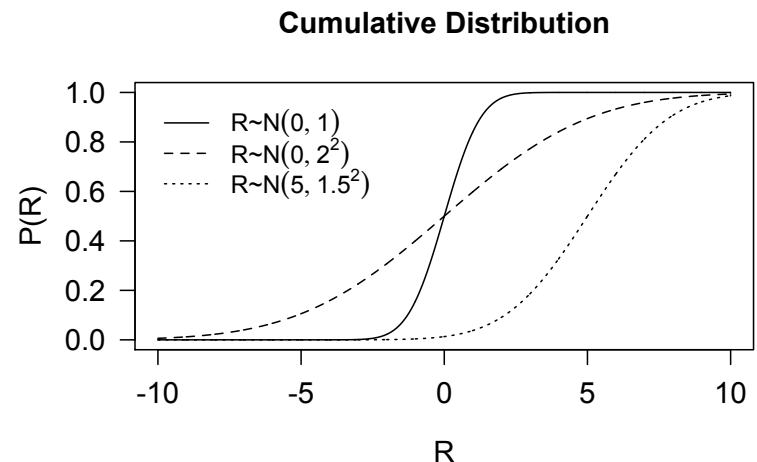
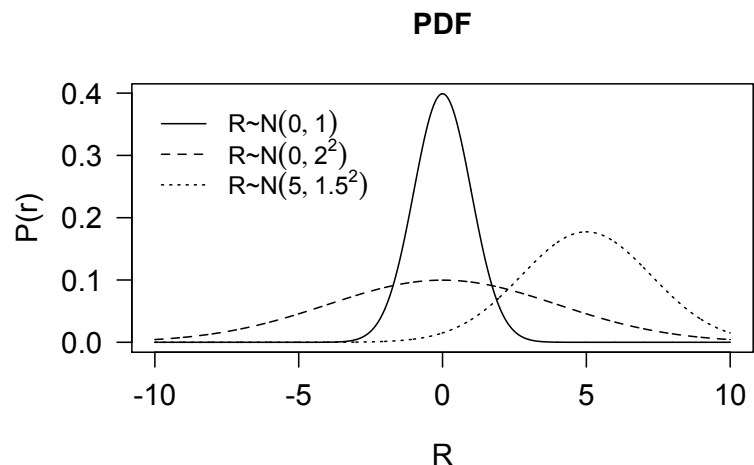
- Continuous random variable R has a Gaussian distribution with mean μ and variance σ^2 , denoted $R \sim \mathcal{N}(\mu, \sigma^2)$, if it has the PDF

$$p(r) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(r-\mu)^2}{2\sigma^2}}$$

Standard Gaussian $R \sim \mathcal{N}(\mu = 0, \sigma^2 = 1)$:

$$p(r) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(r)^2}{2}}$$

- Expected value and variance of R given by μ and σ^2
 - Standard deviation: σ



Multivariate Gaussian

- Univariate Gaussian: two parameters
 - Mean μ and variance σ^2
- Multivariate Gaussian distribution over continuous randvars R_1, \dots, R_n characterised by

- n -dimensional **mean vector** μ
- Symmetric $n \times n$ **covariance matrix** Σ
- I.e., $\mathcal{N}(\mu; \Sigma)$

Standard multivariate Gaussian
 R_1, \dots, R_n with

- $\mu = \mathbf{0}$ (all-zero vector)
- $\Sigma = I$ (identity matrix)

- Density function defined as

$$p(\mathbf{r}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left[-\frac{1}{2} (\mathbf{r} - \mu)^T \Sigma^{-1} (\mathbf{r} - \mu) \right]$$

- $\mathbf{r} = (r_1, \dots, r_n)^T$
- $|\Sigma|$ determinant of Σ
- To induce a well-defined density, Σ must be **positive-definite**
 - For any $\mathbf{r} \in \mathbb{R}^n$ s.t. $\mathbf{r} \neq \mathbf{0} : \mathbf{r}^T \Sigma \mathbf{r} > 0$
 - Guaranteed to be non-singular \rightarrow **non-zero** determinant

Example

- Joint Standard Gaussian distribution over two randvars R_1, R_2 , i.e.,

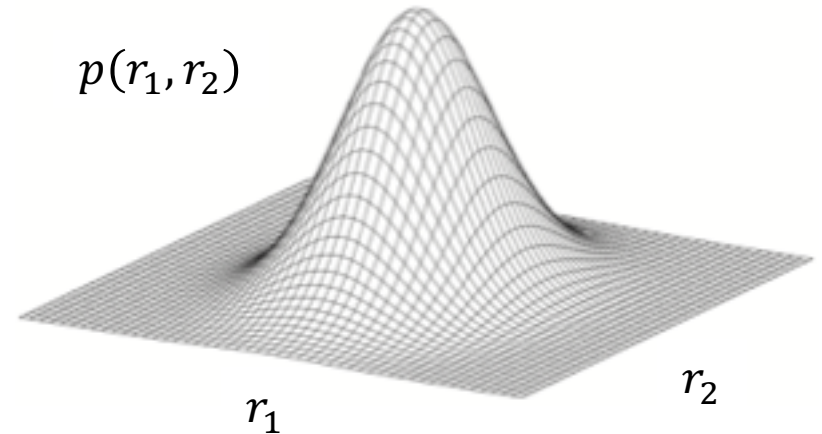
- $\mu = (0 \ 0)^T, \Sigma = I_2$

- Joint Gaussian distribution over three randvars R_1, R_2, R_3

- Mean vector, covariance matrix:

$$\mu = \begin{pmatrix} 1 \\ -3 \\ 4 \end{pmatrix} \quad \Sigma = \begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ -2 & -5 & 8 \end{pmatrix}$$

- Covariances $Cov[R_1; R_3]$ and $Cov[R_2; R_3]$ negative, i.e., R_3 negatively correlated with R_1 (and R_2)
 - When R_1 (R_2) goes up, R_3 goes down



Marginalisation

- Trivial with covariance matrix:
 - Compute pairwise covariances, i.e., generating the distribution in its covariance form
 - Given covariance form Σ : Read off from $\boldsymbol{\mu}, \Sigma$
- Assume a joint Gaussian distribution over $\{\mathbf{R}, \mathbf{T}\}$ where $\mathbf{R} \in \mathbb{R}^n$ and $\mathbf{T} \in \mathbb{R}^m$

- One can decompose mean and covariance:

$$p(\mathbf{r}, \mathbf{t}) = \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_R \\ \boldsymbol{\mu}_T \end{pmatrix}; \begin{bmatrix} \Sigma_{RR} & \Sigma_{RT} \\ \Sigma_{TR} & \Sigma_{TT} \end{bmatrix} \right)$$

- where

- $\boldsymbol{\mu}_R \in \mathbb{R}^n, \boldsymbol{\mu}_T \in \mathbb{R}^m,$
- Σ_{RR} an $n \times n$ matrix, Σ_{RT} an $n \times m$ matrix,
 $\Sigma_{TR} = \Sigma_{RT}^T$ an $m \times n$ matrix, Σ_{TT} a $m \times m$ matrix

- Then, marginal distribution over \mathbf{T} given by Gaussian distribution of $\mathcal{N}(\boldsymbol{\mu}_T; \Sigma_{TT})$

Example

- Given joint Gaussian distribution over three randvars R_1, R_2, R_3

- Mean vector, covariance matrix:

$$\mu = \begin{pmatrix} 1 \\ -3 \\ 4 \end{pmatrix} \quad \Sigma = \begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ -2 & -5 & 8 \end{pmatrix}$$

- $p(R_1, R_2)$ given by Gaussian distribution with

$$\mu = \begin{pmatrix} 1 \\ -3 \end{pmatrix} \quad \Sigma = \begin{pmatrix} 4 & 2 \\ 2 & 5 \end{pmatrix}$$

Dual: Information/Precision Form

- Rewrite $\exp \left[-\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{r} - \boldsymbol{\mu}) \right]$ by setting $\Gamma = \Sigma^{-1}$ and multiplying out:

$$-\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu})^T \Gamma (\mathbf{r} - \boldsymbol{\mu}) = -\frac{1}{2} [\mathbf{r}^T \Gamma \mathbf{r} - 2\mathbf{r}^T \Gamma \boldsymbol{\mu} + \boldsymbol{\mu}^T \Gamma \boldsymbol{\mu}]$$

- $\boldsymbol{\mu}^T \Gamma \boldsymbol{\mu}$ is constant over the different \mathbf{r} , therefore,

$$p(\mathbf{r}) \propto \exp \left(-\frac{1}{2} [\mathbf{r}^T \Gamma \mathbf{r} - 2\mathbf{r}^T \Gamma \boldsymbol{\mu}] \right)$$

$$= \exp \left[-\frac{1}{2} \mathbf{r}^T \Gamma \mathbf{r} + \mathbf{r}^T \Gamma \boldsymbol{\mu} \right]$$

$$= \exp \left[-\frac{1}{2} \mathbf{r}^T \Gamma \mathbf{r} + (\Gamma \boldsymbol{\mu})^T \mathbf{r} \right]$$

- $\Gamma \boldsymbol{\mu}$ called potential vector

$$\begin{aligned} & \mathbf{r}^T \Gamma \boldsymbol{\mu} \\ &= (\mathbf{r}^T \Gamma \boldsymbol{\mu})^{TT} \quad \triangleright A^{TT} = A \\ &= \left((\Gamma \boldsymbol{\mu})^T \mathbf{r}^{TT} \right)^T \quad \triangleright (AB)^T = B^T A^T \\ &= \left((\Gamma \boldsymbol{\mu})^T \mathbf{r} \right)^T \quad \triangleright A^{TT} = A \\ &= (\Gamma \boldsymbol{\mu})^T \mathbf{r} \quad \triangleright k^T = k, k \text{ a scalar} \end{aligned}$$

Dual: Information/Precision Form

- For a decomposition $\{\mathbf{R}, \mathbf{T}\}$ where $\mathbf{R} \in \mathbb{R}^n$ and $\mathbf{T} \in \mathbb{R}^m$:

$$\Gamma = \Sigma^{-1} = \begin{bmatrix} \Sigma_{RR} & \Sigma_{RT} \\ \Sigma_{TR} & \Sigma_{TT} \end{bmatrix}^{-1} = \begin{bmatrix} \Gamma_{RR} & \Gamma_{RT} \\ \Gamma_{TR} & \Gamma_{TT} \end{bmatrix}$$

- Getting to Σ

- $\Sigma_{RR} = (\Gamma_{RR} - \Gamma_{RT}\Gamma_{TT}^{-1}\Gamma_{TR})^{-1}$

- $\Sigma_{TT} = (\Gamma_{TT} - \Gamma_{TR}\Gamma_{RR}^{-1}\Gamma_{RT})^{-1}$

- $\Sigma_{RT} = -\Gamma_{RR}^{-1}\Gamma_{RT}(\Gamma_{TT} - \Gamma_{TR}\Gamma_{RR}^{-1}\Gamma_{RT})^{-1} = \Sigma_{TR}^T$

- $\Sigma_{TR} = -\Gamma_{TT}^{-1}\Gamma_{TR}(\Gamma_{RR} - \Gamma_{RT}\Gamma_{TT}^{-1}\Gamma_{TR})^{-1} = \Sigma_{RT}^T$

- Getting to Γ

- $\Gamma_{RR} = (\Sigma_{RR} - \Sigma_{RT}\Sigma_{TT}^{-1}\Sigma_{TR})^{-1}$

- $\Gamma_{TT} = (\Sigma_{TT} - \Sigma_{TR}\Sigma_{RR}^{-1}\Sigma_{RT})^{-1}$

- $\Gamma_{RT} = -\Sigma_{RR}^{-1}\Sigma_{RT}(\Sigma_{TT} - \Sigma_{TR}\Sigma_{RR}^{-1}\Sigma_{RT})^{-1} = \Gamma_{TR}^T$

- $\Gamma_{TR} = -\Sigma_{TT}^{-1}\Sigma_{TR}(\Sigma_{RR} - \Sigma_{RT}\Sigma_{TT}^{-1}\Sigma_{TR})^{-1} = \Gamma_{RT}^T$

Conditioning

- Conditioning a Gaussian on observations $\mathbf{E} = \mathbf{e}$ easy to perform in the information form by setting \mathbf{E} to \mathbf{e} in one of the following

$$\begin{aligned} p(\mathbf{r}) &\propto \exp \left[-\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu})^T \Gamma (\mathbf{r} - \boldsymbol{\mu}) \right] \\ &\propto \exp \left[-\frac{1}{2} \mathbf{r}^T \Gamma \mathbf{r} + (\Gamma \boldsymbol{\mu})^T \mathbf{r} \right] \end{aligned}$$

- Assuming a decomposition into \mathbf{R} and \mathbf{E} , i.e.,
 $p(\mathbf{r}, \mathbf{e})$

$$\begin{aligned} &= \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_R \\ \boldsymbol{\mu}_E \end{pmatrix}; \begin{bmatrix} \Sigma_{RR} & \Sigma_{RE} \\ \Sigma_{ER} & \Sigma_{EE} \end{bmatrix} \right) \\ &\propto \exp \left[-\frac{1}{2} \begin{pmatrix} \mathbf{r} \\ \mathbf{e} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_R \\ \boldsymbol{\mu}_E \end{pmatrix} \right]^T \begin{bmatrix} \Gamma_{RR} & \Gamma_{RT} \\ \Gamma_{TR} & \Gamma_{TT} \end{bmatrix} \begin{pmatrix} \mathbf{r} \\ \mathbf{e} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_R \\ \boldsymbol{\mu}_E \end{pmatrix} \right] \end{aligned}$$

In the exponential function:

$$p(\mathbf{r}) \propto \exp \left[-\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{r} - \boldsymbol{\mu}) \right]$$

$$= \exp \left[-\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu})^T \boldsymbol{\Gamma} (\mathbf{r} - \boldsymbol{\mu}) \right]$$

$$-\frac{1}{2} \left(\begin{pmatrix} \mathbf{r} \\ \mathbf{e} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_R \\ \boldsymbol{\mu}_E \end{pmatrix} \right)^T \begin{bmatrix} \boldsymbol{\Gamma}_{RR} & \boldsymbol{\Gamma}_{RE} \\ \boldsymbol{\Gamma}_{ER} & \boldsymbol{\Gamma}_{EE} \end{bmatrix} \left(\begin{pmatrix} \mathbf{r} \\ \mathbf{e} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_R \\ \boldsymbol{\mu}_E \end{pmatrix} \right)$$

$$= -\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu}_R)^T \begin{bmatrix} \boldsymbol{\Gamma}_{RR} & \boldsymbol{\Gamma}_{RE} \\ \boldsymbol{\Gamma}_{ER} & \boldsymbol{\Gamma}_{EE} \end{bmatrix} (\mathbf{e} - \boldsymbol{\mu}_E)$$

$$= -\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu}_R)^T \boldsymbol{\Gamma}_{RR} (\mathbf{r} - \boldsymbol{\mu}_R) - \frac{1}{2} 2 (\mathbf{r} - \boldsymbol{\mu}_R)^T \boldsymbol{\Gamma}_{RE} (\mathbf{e} - \boldsymbol{\mu}_E) - \frac{1}{2} (\mathbf{e} - \boldsymbol{\mu}_E)^T \boldsymbol{\Gamma}_{EE} (\mathbf{e} - \boldsymbol{\mu}_E)$$

$$\propto -\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu}_R)^T \boldsymbol{\Gamma}_{RR} (\mathbf{r} - \boldsymbol{\mu}_R) - (\mathbf{r} - \boldsymbol{\mu}_R)^T \boldsymbol{\Gamma}_{RE} (\mathbf{e} - \boldsymbol{\mu}_E)$$

Does not depend on \mathbf{r}

$$= -\frac{1}{2} (\mathbf{r} - \boldsymbol{\mu}_R)^T \boldsymbol{\Gamma}_{RR} (\mathbf{r} - \boldsymbol{\mu}_R) - (\mathbf{r} - \boldsymbol{\mu}_R)^T \boldsymbol{\Gamma}_{RE} (\mathbf{e} - \boldsymbol{\mu}_E) - A + A$$

Use $-A$ to get expression into the form $(\mathbf{r} - \boldsymbol{\mu})^T \boldsymbol{\Gamma} (\mathbf{r} - \boldsymbol{\mu})$ by factoring out $\boldsymbol{\Gamma}_{RR}$

$$A = \frac{1}{2} (\mathbf{e} - \boldsymbol{\mu}_E)^T \boldsymbol{\Gamma}_{ER} \boldsymbol{\Gamma}_{RR}^{-1} \boldsymbol{\Gamma}_{RR} \boldsymbol{\Gamma}_{RR}^{-1} \boldsymbol{\Gamma}_{RE} (\mathbf{e} - \boldsymbol{\mu}_E)$$

$$\exp \left[-\frac{1}{2} \left(\left(\mathbf{r} - \boldsymbol{\mu}_R + \boldsymbol{\Gamma}_{RR}^{-1} \boldsymbol{\Gamma}_{ER} (\mathbf{e} - \boldsymbol{\mu}_E) \right)^T \boldsymbol{\Gamma}_{RR} \left(\mathbf{r} - \boldsymbol{\mu}_R + \boldsymbol{\Gamma}_{RR}^{-1} \boldsymbol{\Gamma}_{RE} (\mathbf{e} - \boldsymbol{\mu}_E) \right) \right) \right] \exp[A]$$

$$\propto \exp \left[-\frac{1}{2} \left(\left(\mathbf{r} - \boldsymbol{\mu}_R + \boldsymbol{\Gamma}_{RR}^{-1} \boldsymbol{\Gamma}_{ER} (\mathbf{e} - \boldsymbol{\mu}_E) \right)^T \boldsymbol{\Gamma}_{RR} \left(\mathbf{r} - \boldsymbol{\mu}_R + \boldsymbol{\Gamma}_{RR}^{-1} \boldsymbol{\Gamma}_{ER} (\mathbf{e} - \boldsymbol{\mu}_E) \right) \right) \right]$$

$$\boldsymbol{\mu}^* = \boldsymbol{\mu}_R - \boldsymbol{\Gamma}_{RR}^{-1} \boldsymbol{\Gamma}_{ER} (\mathbf{e} - \boldsymbol{\mu}_E) \quad \boldsymbol{\Sigma}^* = \boldsymbol{\Gamma}_{RR}$$

Conditioning

- Conditioning a Gaussian on observations $\mathbf{E} = \mathbf{e}$ with remaining randvars \mathbf{R}
- Result:

$$\mathbf{R}|\mathbf{E} = \mathbf{e} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

- Information form:
 - $\boldsymbol{\mu}^* = \boldsymbol{\mu}_R - \boldsymbol{\Gamma}_R^{-1}\boldsymbol{\Gamma}_{ER}(\mathbf{e} - \boldsymbol{\mu}_E)$
 - $\boldsymbol{\Sigma}^* = \boldsymbol{\Gamma}_{RR}$
- Covariance form:
 - $\boldsymbol{\mu}^* = \boldsymbol{\mu}_R + \boldsymbol{\Sigma}_{RE}\boldsymbol{\Sigma}_{EE}^{-1}(\mathbf{e} - \boldsymbol{\mu}_E)$
 - $\boldsymbol{\Sigma}^* = \boldsymbol{\Sigma}_{RR} - \boldsymbol{\Sigma}_{RE}\boldsymbol{\Sigma}_{EE}^{-1}\boldsymbol{\Sigma}_{ER}$
- Mean moved from $\boldsymbol{\mu}_R$ according to correlation and variance on observations $\boldsymbol{\Sigma}_{RE}\boldsymbol{\Sigma}_{EE}^{-1}(\mathbf{e} - \boldsymbol{\mu}_E)$
- Covariance does not depend on observations \mathbf{e}

Query Answering: Summary

- For marginalisation, read off parameters in covariance form
 - Marginal query for \mathbf{T} : $\mathcal{N}(\boldsymbol{\mu}_T; \Sigma_{TT})$
- For conditioning, one needs to invert the covariance matrix to obtain the information form
 - Conditioning on $\mathbf{E} = \mathbf{e}$:
$$\mathbf{R}|\mathbf{E} = \mathbf{e} \sim \mathcal{N}(\boldsymbol{\mu}^*, \Sigma^*)$$
 - In covariance form
 - $\boldsymbol{\mu}^* = \boldsymbol{\mu}_R + \Sigma_{RE}\Sigma_{EE}^{-1}(\mathbf{e} - \boldsymbol{\mu}_E)$
 - $\Sigma^* = \Sigma_{RR} - \Sigma_{RE}\Sigma_{EE}^{-1}\Sigma_{ER}$
 - Matrix inversion can be very costly!

Linear Gaussian Model

- Let S be a continuous random variable with continuous parents R_1, \dots, R_k
- S has a **linear Gaussian model** if there are parameters β_0, \dots, β_k and σ^2 such that
$$p(S|r_1, \dots, r_k) = \mathcal{N}(\beta_0 + \beta_1 r_1 + \dots + \beta_k r_k; \sigma^2)$$
$$= \mathcal{N}(\beta_0 + \boldsymbol{\beta}^T \mathbf{r}; \sigma^2) \longleftarrow \text{(vector notation)}$$
- $p(S|r_1, \dots, r_k)$ a **conditional probability distribution (CPD)**
- Interpretations
 - β_0 is an initial mean μ_0 that is moved according to the influences by the parents
 - S is a linear function of R_1, \dots, R_k with the addition of Gaussian noise: $S = \beta_0 + \beta_1 r_1 + \dots + \beta_k r_k + \epsilon$
 - ϵ a Gaussian random variable with mean 0 and variance σ^2 , representing the noise in the system
- Does not allow σ^2 to depend on parent values
 - But can be a useful approximation

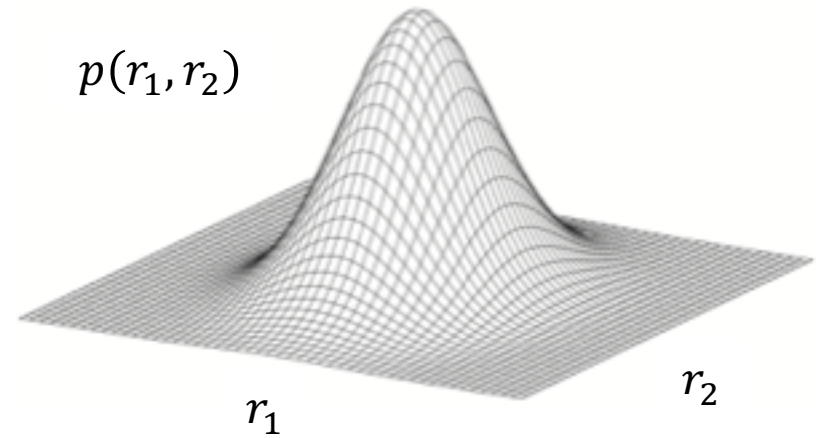
Independencies in Gaussians

- Let randvars R_1, \dots, R_n have a joint distribution $\mathcal{N}(\boldsymbol{\mu}; \Sigma)$
- Then, R_i, R_j independent iff $\Sigma_{ij} = 0$
 - Joint distribution needs to be Gaussian for this equivalence to hold
 - If the distribution is not Gaussian, $\Sigma_{ij} = 0$ might be the case and there still might be a dependence between R_i, R_j
- Conditional independence can be read of in the inverse of the covariance matrix, Σ^{-1}
 - Given a Gaussian distribution $p(r_1, \dots, r_n) = \mathcal{N}(\boldsymbol{\mu}; \Sigma)$
 - Then, $\Sigma_{ij}^{-1} = 0$ iff $p \models (R_i \perp R_j | \{R_1, \dots, R_n\} \setminus \{R_i, R_j\})$

Example

- Joint Standard Gaussian distribution over two randvars R_1, R_2 , i.e.,

- $\mu = (0 \ 0)^T, \Sigma = \mathbf{I}_2$
- R_1, R_2 independent as $\Sigma_{ij} = \Sigma_{ji} = 0$



- Gaussian for R_1, R_2, R_3 from before

- Covariance and inverse covariance matrix:

$$\Sigma = \begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ -2 & -5 & 8 \end{pmatrix} \quad \Sigma^{-1} = \begin{pmatrix} 0.3125 & -0.125 & 0 \\ -0.125 & 0.5833 & 0.3333 \\ 0 & 0.3333 & 0.3333 \end{pmatrix}$$

- R_1, R_3 conditionally independent given R_2

- $\Sigma_{13}^{-1} = 0$ iff

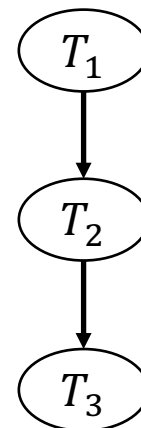
$$p \models (R_1 \perp R_3 | \{R_1, R_2, R_3\} \setminus \{R_1, R_3\}) = (R_1 \perp R_3 | R_2)$$

Gaussian Bayesian Network (GBN)

- Factorisation of a joint distribution into factors also possible with linear Gaussians as local CPDs
- A BN is a directed acyclic graph G whose nodes are discrete randvars $\{R_1, \dots, R_n\}$ and whose full joint P_G factorises according to the local CPTs, i.e.,

$$P_G = \prod_i P(R_i | \text{parents}(R_i))$$

- **Gaussian BN** is a BN where
 - R_i are continuous randvars
 - All CPDs are linear Gaussians
 - E.g., $T_1 \rightarrow T_2 \rightarrow T_3$ (also depicted right)
 - $p(T_1) = \mathcal{N}(1; 4)$
 - $p(T_2|T_1) = \mathcal{N}(-3.5 + 0.5 \cdot T_1; 4)$
 - $p(T_3|T_2) = \mathcal{N}(1 + (-1) \cdot T_2; 3)$



Connection to Multivariate Gaussian

- Linear GBN an alternative representation to multivariate Gaussian distribution
 - A linear Gaussian BN always defines a joint multivariate Gaussian distribution
- Let S be a linear Gaussian of its parents R_1, \dots, R_k
 - $\mathcal{N}(\beta_0 + \boldsymbol{\beta}^T \mathbf{r}; \sigma^2) = \mathcal{N}(\beta_0 + \beta_1 r_1 + \dots + \beta_k r_k; \sigma^2)$
 - R_1, \dots, R_k jointly Gaussian with $\mathcal{N}(\boldsymbol{\mu}; \boldsymbol{\Sigma})$
- Distribution of S is a Gaussian $p(S) = \mathcal{N}(\mu_S; \sigma_S^2)$ with
$$\begin{aligned}\mu_S &= \beta_0 + \boldsymbol{\beta}^T \mathbf{r} \\ \sigma_S^2 &= \sigma^2 + \boldsymbol{\beta}^T \boldsymbol{\Sigma} \boldsymbol{\beta}\end{aligned}$$
- Joint distribution over $\{R_1, \dots, R_k, S\}$ is a Gaussian with
$$\text{Cov}[R_i; S] = \sum_{j=1}^k \beta_j \Sigma_{ij}$$

General Procedure for Conversion

- Let (R_1, \dots, R_n) be the randvars of a GBN
 - Each R_i is a Gaussian $\mathcal{N}(\beta_0 + \boldsymbol{\beta}^T \mathbf{r}; \sigma^2)$ conditional on its parents $parents(R_i)$
 - (R_1, \dots, R_n) follows a topological ordering θ s.t.
 $\forall R_j \in \{R_1, \dots, R_n\} : \forall R_i \in parents(R_j) : R_i \prec_{\theta} R_j$
 - Build a matrix $B^{n \times n}$ that has a non-zero entry β_{ij} if there exists a parent-child relation $R_i \rightarrow R_j$ with β_{ij} being the factor for R_i in the $\boldsymbol{\beta}$ of R_j

$$B = \begin{pmatrix} 0 & \beta_{12} & \dots & \beta_{1n} \\ 0 & 0 & & \beta_{2n} \\ & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}$$

- i chooses the row(s), j chooses the column(s)
- B is upper-triangular because no loops allowed in BNs
 - Including self-loops $\rightarrow \beta_{ii} = 0$ as well

General Procedure for Conversion

- Joint distribution $p(r_1, \dots, r_n)$ given by $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$

- Means

$$\boldsymbol{\mu} = (\mu_1, \beta_{0,2} + \boldsymbol{\beta}_2^T \mathbf{r}, \dots, \beta_{0,n} + \boldsymbol{\beta}_n^T \mathbf{r})^T$$

- Covariance (*recursive rules*): $j \in \{2, \dots, n\}, i = 1 \dots j - 1$

$$\Sigma_{11} \leftarrow \sigma_1^2$$

$$\Sigma_{ij} \leftarrow \Sigma_{ii} B_{ij}$$

$$\Sigma_{ji} \leftarrow \Sigma_{ij}^T$$

$$\Sigma_{jj} \leftarrow \sigma_j^2 + \Sigma_{ji} B_{ij}$$

- First index chooses the row(s), second index chooses the column(s)

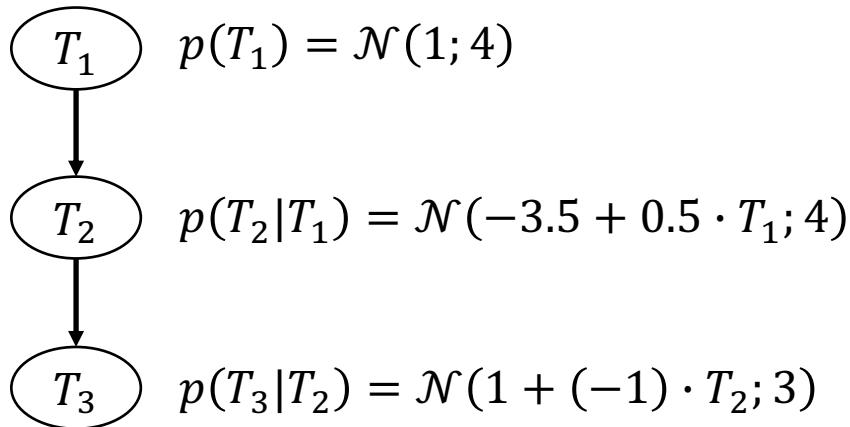
- given B

filling Σ layer-wise:

$$B = \begin{pmatrix} 0 & \beta_{12} & \dots & \beta_{1n} \\ 0 & 0 & & \beta_{2n} \\ & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \dots & \Sigma_{1n} \\ \Sigma_{21} & \Sigma_{22} & & \Sigma_{2n} \\ & \vdots & \ddots & \vdots \\ \Sigma_{n1} & \Sigma_{n2} & \dots & \Sigma_{nn} \end{pmatrix}$$

GBN: Conversion Example

- GBN



- Goal: Joint distribution

$$p(t_1, t_2, t_3) = \mathcal{N}(\boldsymbol{\mu}; \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix}$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

- Matrix B

- $B_{12}: T_1 \rightarrow T_2, \beta_1 = 0.5$
- $B_{23}: T_2 \rightarrow T_3, \beta_1 = -1$
- Rest: zeroes
- Result:

$$\bullet B = \begin{pmatrix} 0 & 0.5 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{pmatrix}$$

- Means

- $\mu_1 = 1$
- $\mu_2 = -3.5 + 0.5 \cdot \mu_1 = -3$
- $\mu_3 = 1 + (-1) \cdot \mu_2 = 4$
- Result:

$$\bullet \boldsymbol{\mu} = \begin{pmatrix} 1 \\ -3 \\ 4 \end{pmatrix}$$

GBN: Conversion Example

- Filling Σ :

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

- Need B and the recursive rules

$$B = \begin{pmatrix} 0 & 0.5 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{pmatrix}$$

$$\Sigma_{11} \leftarrow \sigma_1^2$$

$$\Sigma_{ij} \leftarrow \Sigma_{ii} B_{ij}$$

$$\Sigma_{ji} \leftarrow \Sigma_{ij}^T$$

$$\Sigma_{jj} \leftarrow \sigma_j^2 + \Sigma_{ji} B_{ij}$$

- First index: row(s)
- Second index: column(s)

- $\Sigma_{11} = \sigma_1^2 = 4$

- $j = 2, i = 1$

- Σ_{12}

$$= \Sigma_{11} B_{12}$$

$$= 4 \cdot 0.5$$

$$= 2$$

- Σ_{21}

$$= \Sigma_{12}^T$$

$$= 2^T$$

$$= 2$$

- Σ_{22}

$$= \sigma_2^2 + \Sigma_{21} B_{12}$$

$$= 4 + 2 \cdot 0.5$$

$$= 5$$

$$\begin{pmatrix} 4 & \Sigma_{12} & \Sigma_{13} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

$$\begin{pmatrix} 4 & 2 & \Sigma_{13} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

$$\begin{pmatrix} 4 & 2 & \Sigma_{13} \\ 2 & \Sigma_{22} & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

$$\begin{pmatrix} 4 & 2 & \Sigma_{13} \\ 2 & 5 & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

GBN: Conversion Example

- Remaining goal:

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

- Need B and the recursive rules

$$B = \begin{pmatrix} 0 & 0.5 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{pmatrix}$$

$$\Sigma_{11} \leftarrow \sigma_1^2$$

$$\Sigma_{ij} \leftarrow \Sigma_{ii} B_{ij}$$

$$\Sigma_{ji} \leftarrow \Sigma_{ij}^T$$

$$\Sigma_{jj} \leftarrow \sigma_j^2 + \Sigma_{ji} B_{ij}$$

- First index: row(s)
- Second index: column(s)

- $j = 3, i = 12$

$$\begin{aligned} & \bullet \Sigma_{(12)3} \\ &= \Sigma_{(12)(12)} B_{(12)3} \\ &= \begin{pmatrix} 4 & 2 \\ 2 & 5 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \end{pmatrix} \\ &= \begin{pmatrix} -2 \\ -5 \end{pmatrix} \end{aligned}$$

$$\begin{pmatrix} 4 & 2 & \Sigma_{13} \\ 2 & 5 & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

$$\begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{pmatrix}$$

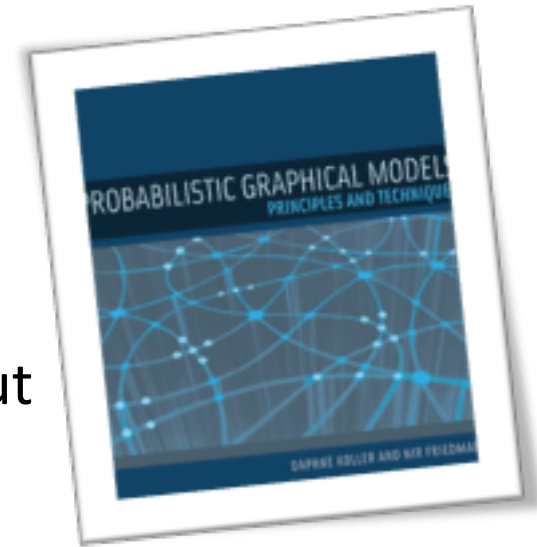
$$\begin{aligned} & \bullet \Sigma_{3(12)} \\ &= \Sigma_{(12)3}^T \\ &= \begin{pmatrix} -2 \\ -5 \end{pmatrix}^T \\ &= (-2 \quad -5) \end{aligned}$$

$$\begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ -2 & -5 & \Sigma_{33} \end{pmatrix}$$

$$\begin{aligned} & \bullet \Sigma_{33} \\ &= \sigma_3^2 + \Sigma_{3(12)} B_{(12)3} \\ &= 3 + (-2 \quad -5) \begin{pmatrix} 0 \\ -1 \end{pmatrix} \begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & -5 \\ -2 & -5 & 8 \end{pmatrix} \\ &= 3 + 5 = 8 \end{aligned}$$

Inference in GBNs

- Inference in linear Gaussians with Variable Elimination
 - Representation through linear Gaussian CPDs instead of CPTs/factors
 - Modified operations for multiply/sum-out
- Message passing formulation
 - Approximate belief propagation
- Sampling in the continuous space
 - Rejection sampling, importance sampling, MCMC methods for GBNs
- Actually using the full joint
 - Marginalisation, conditioning as sketched in Basics



See Ch. 14 of PGM book for further information

Lifting the Full Joint

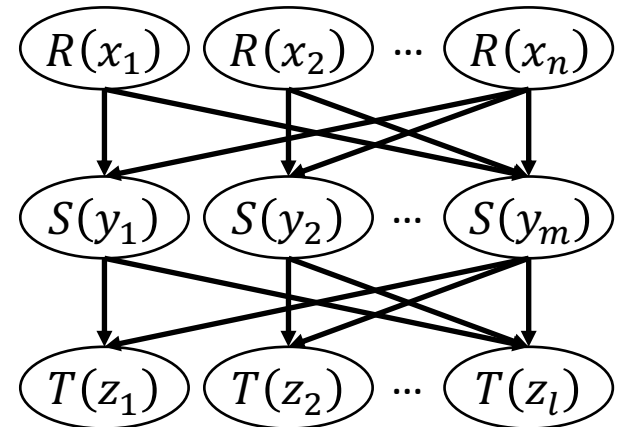
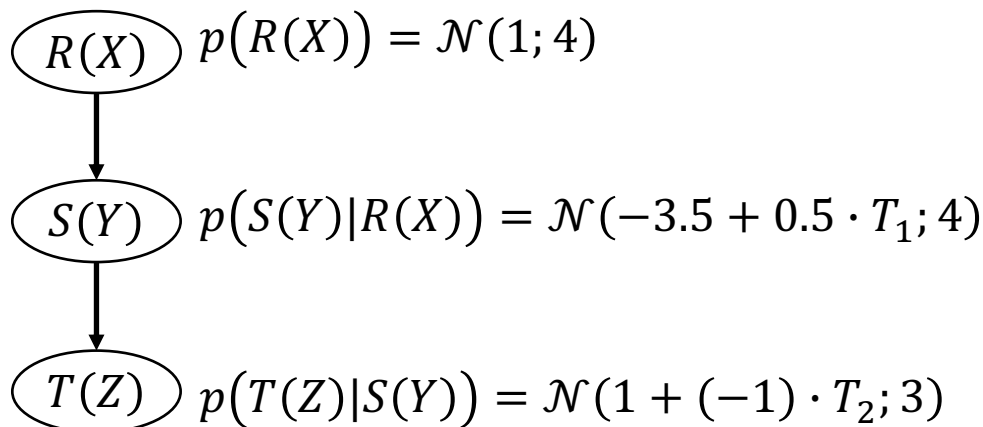
- Lifting conversion approach by Shachter and Kenley for parameterised GBNs

$$\begin{aligned} \Sigma_{11} &\leftarrow \sigma_1^2 \\ \Sigma_{ij} &\leftarrow \Sigma_{ii} B_{ij} \\ \Sigma_{ji} &\leftarrow \Sigma_{ij}^T \\ \Sigma_{jj} &\leftarrow \sigma_j^2 + \Sigma_{ji} B_{ij} \end{aligned}$$

- GBN with PRVs A_1, \dots, A_m as nodes
 - PDF for each A_i applies to each $R \in gr(A_i)$
 - $m \ll n, n = |\cup_i gr(A_i)|$
 - Semantics: grounding and forming full joint $p(\cup_i gr(A_i))$
 - Simple case for GBNs (*general case under review*):

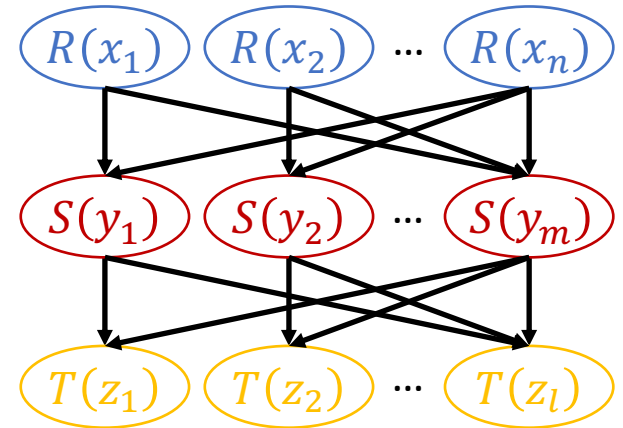
For all parent-child relations $R(\mathbf{X}) \rightarrow S(\mathbf{Y})$,
it holds that $\mathbf{X} \cap \mathbf{Y} = \emptyset$

- Each child instance has the same parent instances as its siblings



Lifting the Full Joint: Simple Case

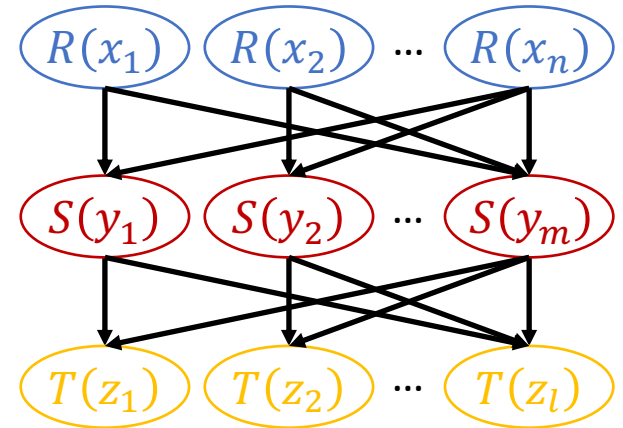
- With PRVs, matrix B and covariance matrix have liftable blocks for each PRV
 - Given the case of no overlaps in logvars: B



$$\begin{array}{c}
 R(X) \\
 S(Y) \\
 T(Z)
 \end{array}
 \begin{pmatrix}
 & R(X) & & S(Y) & & T(Z) \\
 0 & \dots & 0 & \beta_{r_1 s_1} & \dots & \beta_{r_1 s_m} & 0 & \dots & 0 \\
 \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 0 & \dots & 0 & \beta_{r_n s_1} & \dots & \beta_{r_n s_m} & 0 & \dots & 0 \\
 0 & \dots & 0 & 0 & \dots & 0 & \beta_{s_1 t_1} & \dots & \beta_{s_1 t_l} \\
 \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 0 & \dots & 0 & 0 & \dots & 0 & \beta_{s_m t_1} & \dots & \beta_{s_m t_l} \\
 0 & \dots & 0 & 0 & \dots & 0 & 0 & \dots & 0 \\
 \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 0 & \dots & 0 & 0 & \dots & 0 & 0 & \dots & 0
 \end{pmatrix}$$

Lifting the Full Joint: Simple Case

- With PRVs, matrix B and covariance matrix have liftable blocks for each PRV
 - Given the case of no overlaps in logvars: B



$$R(X) \begin{pmatrix}
 & R(X) & & S(Y) & & T(Z) \\
 0 & \dots & 0 & \beta_{r_1 s_1} & \dots & \beta_{r_1 s_m} & 0 & \dots & 0 \\
 \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 0 & \dots & 0 & \beta_{r_n s_1} & \dots & \beta_{r_n s_m} & \beta_{s_1 t_1} & \dots & \beta_{s_1 t_l} \\
 & & & & & & \vdots & & \vdots \\
 & & & & & & \beta_{s_m t_1} & \dots & \beta_{s_m t_l} \\
 & & & & & & 0 & \dots & 0 \\
 & & & & & & \vdots & \ddots & \vdots \\
 & & & & & & 0 & \dots & 0
 \end{pmatrix}$$

Lifted B'

$$\begin{pmatrix}
 0 & \beta_{rs} & 0 \\
 0 & 0 & \beta_{st} \\
 0 & 0 & 0
 \end{pmatrix}$$

Lifted B'

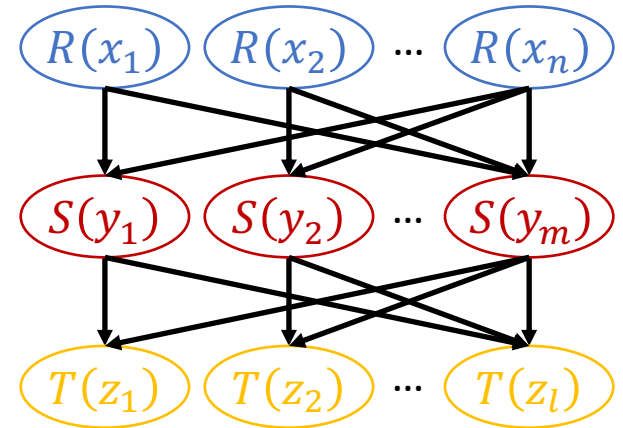
$$\begin{pmatrix}
 0 & 0.5 & 0 \\
 0 & 0 & -1 \\
 0 & 0 & 0
 \end{pmatrix}$$

Each $R(x_i)$ has the same influence on each $S(y_j)$.
 Given $P(s(Y)|r(X)) = \mathcal{N}(\beta_0 + \beta_1 r(X); \sigma^2)$,
 $\beta_{r_i s_j} = \beta_1$
 for all $i \in \{1, \dots, n\}, j \in \{1, \dots, m\}$.
 The same holds for $S(y_j)$ and $T(z_k)$ (as well as
 $R(x_i)$ and $T(z_k)$, which has $\beta_{r_i t_k} = 0$).



Lifting the Full Joint: Simple Case

- With PRVs, matrix B and covariance matrix have liftable blocks for each PRV
 - Given the case of no overlaps in logvars: Σ



$$\begin{array}{c}
 R(X) \\
 S(Y) \\
 T(Z)
 \end{array}
 \begin{pmatrix}
 \begin{array}{ccc}
 R(X) & & \\
 & S(Y) & \\
 & & T(Z)
 \end{array} \\
 \left(\begin{array}{cccccc}
 \Sigma_{r_1 r_1} & \cdots & \Sigma_{r_1 r_n} & \Sigma_{r_1 s_1} & \cdots & \Sigma_{r_1 s_m} \\
 \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 \Sigma_{r_n r_1} & \cdots & \Sigma_{r_n r_n} & \Sigma_{r_n s_1} & \cdots & \Sigma_{r_n s_m} \\
 \Sigma_{s_1 r_1} & \cdots & \Sigma_{s_1 r_n} & \Sigma_{s_1 s_1} & \cdots & \Sigma_{s_1 s_m} \\
 \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 \Sigma_{s_m r_1} & \cdots & \Sigma_{s_m r_n} & \Sigma_{s_m s_1} & \cdots & \Sigma_{s_m s_m} \\
 \Sigma_{t_1 r_1} & \cdots & \Sigma_{t_1 r_n} & \Sigma_{t_1 s_1} & \cdots & \Sigma_{t_1 s_m} \\
 \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 \Sigma_{t_l r_1} & \cdots & \Sigma_{t_l r_n} & \Sigma_{t_l s_1} & \cdots & \Sigma_{t_l s_m}
 \end{array} \right)
 \end{pmatrix}$$

Lifting the Full Joint: Simple Case

$$\begin{pmatrix} \Sigma_{r_1 r_1} & \dots & \Sigma_{r_1 r_n} & \Sigma_{r_1 s_1} & \dots & \Sigma_{r_1 s_m} & \Sigma_{r_1 t_1} & \dots & \Sigma_{r_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{r_n r_1} & \dots & \Sigma_{r_n r_n} & \Sigma_{r_n s_1} & \dots & \Sigma_{r_n s_m} & \Sigma_{r_n t_1} & \dots & \Sigma_{r_n t_l} \\ \Sigma_{s_1 r_1} & \dots & \Sigma_{s_1 r_n} & \Sigma_{s_1 s_1} & \dots & \Sigma_{s_1 s_m} & \Sigma_{s_1 t_1} & \dots & \Sigma_{s_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{s_m r_1} & \dots & \Sigma_{s_m r_n} & \Sigma_{s_m s_1} & \dots & \Sigma_{s_m s_m} & \Sigma_{s_m t_1} & \dots & \Sigma_{s_m t_l} \\ \Sigma_{t_1 r_1} & \dots & \Sigma_{t_1 r_n} & \Sigma_{t_1 s_1} & \dots & \Sigma_{t_1 s_m} & \Sigma_{t_1 t_1} & \dots & \Sigma_{t_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{t_l r_1} & \dots & \Sigma_{t_l r_n} & \Sigma_{t_l s_1} & \dots & \Sigma_{t_l s_m} & \Sigma_{t_l t_1} & \dots & \Sigma_{t_l t_l} \end{pmatrix}$$

$$\Sigma_{r_1 r_1} = \sigma_{R(X)}^2$$

$$\Sigma_{r_1 r_2} = \sigma_{R(X)}^2 B_{r_1 r_2} = \sigma_{R(X)}^2 B'_{11} = \sigma_{R(X)}^2 \cdot 0 = 0$$

$$\Sigma_{r_2 r_1} = 0$$

$$\Sigma_{r_2 r_2} = \sigma_{R(X)}^2 + \Sigma_{r_2 r_1} B_{r_1 r_2} = \sigma_{R(X)}^2 + 0 = \sigma_{R(X)}^2$$

...

$R(X)$

$$R(X) \begin{pmatrix} \sigma_{R(X)}^2 & \dots & 0 & 4 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{R(X)}^2 & 0 & \dots & 4 \end{pmatrix} \rightarrow \begin{array}{l} \text{on-diagonal: } \sigma_{R(X)}^2 \\ \text{off-diagonal: } 0 \end{array}$$

$$\begin{array}{l} \Sigma_{11} \leftarrow \sigma_1^2 \\ \Sigma_{ij} \leftarrow \Sigma_{ii} B_{ij} \\ \Sigma_{ji} \leftarrow \Sigma_{ij}^T \\ \Sigma_{jj} \leftarrow \sigma_j^2 + \Sigma_{ji} B_{ij} \end{array}$$

Lifted B'

$$\begin{pmatrix} 0 & \beta_{rs} & 0 \\ 0 & 0 & \beta_{st} \\ 0 & 0 & 0 \end{pmatrix}$$

Lifting the Full Joint: Simple Case

$$\begin{pmatrix} \sigma_{R(X)}^2 & \dots & 0 & \Sigma_{r_1 s_1} & \dots & \Sigma_{r_1 s_m} & \Sigma_{r_1 t_1} & \dots & \Sigma_{r_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{R(X)}^2 & \Sigma_{r_n s_1} & \dots & \Sigma_{r_n s_m} & \Sigma_{r_n t_1} & \dots & \Sigma_{r_n t_l} \\ \Sigma_{s_1 r_1} & \dots & \Sigma_{s_1 r_n} & \Sigma_{s_1 s_1} & \dots & \Sigma_{s_1 s_m} & \Sigma_{s_1 t_1} & \dots & \Sigma_{s_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{s_m r_1} & \dots & \Sigma_{s_m r_n} & \Sigma_{s_m s} & \dots & \Sigma_{s_m s_m} & \Sigma_{s_m t_1} & \dots & \Sigma_{s_m t_n} \\ \Sigma_{t_1 r_1} & \dots & \Sigma_{t_1 r_n} & \Sigma_{t_1 s_1} & \dots & \Sigma_{t_1 s_m} & \Sigma_{t_1 t_1} & \dots & \Sigma_{t_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{t_l r_1} & \dots & \Sigma_{t_l r_n} & \Sigma_{t_l s_1} & \dots & \Sigma_{t_l s_m} & \Sigma_{t_l t_1} & \dots & \Sigma_{t_l t_l} \end{pmatrix}$$

$$\begin{aligned} & \Sigma_{(r_1 \dots r_n) s_1} \\ &= \Sigma_{(r_1 \dots r_n)(r_1 \dots r_n)} B_{(r_1 \dots r_n) s_1} \\ &= \begin{pmatrix} \sigma_{R(X)}^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{R(X)}^2 \end{pmatrix} \begin{pmatrix} \beta_{rs} \\ \vdots \\ \beta_{rs} \end{pmatrix} = \begin{pmatrix} \sigma_{R(X)}^2 \beta_{rs} \\ \vdots \\ \sigma_{R(X)}^2 \beta_{rs} \end{pmatrix} = \begin{pmatrix} 4 \cdot 0.5 \\ \vdots \\ 4 \cdot 0.5 \end{pmatrix} = \begin{pmatrix} 2 \\ \vdots \\ 2 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} \Sigma_{11} &\leftarrow \sigma_1^2 \\ \Sigma_{ij} &\leftarrow \Sigma_{ii} B_{ij} \\ \Sigma_{ji} &\leftarrow \Sigma_{ij}^T \\ \Sigma_{jj} &\leftarrow \sigma_j^2 + \Sigma_{ji} B_{ij} \end{aligned}$$

$$\begin{aligned} & \Sigma_{s_1 s_1} \\ &= \sigma_{S(Y)}^2 + \Sigma_{s_1(r_1 \dots r_n)} B_{(r_1 \dots r_n) s_1} \\ &= \sigma_{S(Y)}^2 + (\sigma_{R(X)}^2 \beta_{rs} \quad \dots \quad \sigma_{R(X)}^2 \beta_{rs}) \begin{pmatrix} \beta_{rs} \\ \vdots \\ \beta_{rs} \end{pmatrix} = \sigma_{S(Y)}^2 + n \sigma_{R(X)}^2 \beta_{rs}^2 \\ &= 4 + n \cdot 4 \cdot 0.5^2 = 4 + n \end{aligned}$$

$$\begin{aligned} & \text{Lifted } B' \\ & \begin{pmatrix} 0 & \beta_{rs} & 0 \\ 0 & 0 & \beta_{st} \\ 0 & 0 & 0 \end{pmatrix} \end{aligned}$$

Lifting the Full Joint: Simple Case

$$\begin{aligned}
 & \Sigma_{(r_1 \dots r_n s_1) s_2} \\
 &= \Sigma_{(r_1 \dots r_n s_1)(r_1 \dots r_n s_1)} B_{(r_1 \dots r_n s_1) s_2} \\
 &= \begin{pmatrix} \sigma_{R(X)}^2 & \dots & 0 & \sigma_{R(X)}^2 \beta_{rs} \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & \sigma_{R(X)}^2 & \sigma_{R(X)}^2 \beta_{rs} \\ \sigma_{R(X)}^2 \beta_{rs} & \dots & \sigma_{R(X)}^2 \beta_{rs} & \sigma_{S(Y)}^2 + n \sigma_{R(X)}^2 \beta_{rs}^2 \end{pmatrix} \begin{pmatrix} \beta_{rs} \\ \vdots \\ \beta_{rs} \\ 0 \end{pmatrix} \\
 &= \begin{pmatrix} \sigma_{R(X)}^2 \beta_{rs} \\ \vdots \\ \sigma_{R(X)}^2 \beta_{rs} \\ n \sigma_{R(X)}^2 \beta_{rs}^2 \end{pmatrix} = \begin{pmatrix} 4 \cdot 0.5 \\ \vdots \\ 4 \cdot 0.5 \\ n \cdot 4 \cdot 0.5^2 \end{pmatrix} = \begin{pmatrix} 2 \\ \vdots \\ 2 \\ n \end{pmatrix}
 \end{aligned}$$

$$\begin{pmatrix} \Sigma_{r_1 s_1} & \dots & \Sigma_{r_1 s_m} & \Sigma_{r_1 t_1} & \dots & \Sigma_{r_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{r_n s_1} & \dots & \Sigma_{r_n s_m} & \Sigma_{r_n t_1} & \dots & \Sigma_{r_n t_l} \\ \Sigma_{s_1 s_1} & \dots & \Sigma_{s_1 s_m} & \Sigma_{s_1 t_1} & \dots & \Sigma_{s_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{s_m s} & \dots & \Sigma_{s_m s_m} & \Sigma_{s_m t_1} & \dots & \Sigma_{s_m t_n} \\ \Sigma_{t_1 s_1} & \dots & \Sigma_{t_1 s_m} & \Sigma_{t_1 t_1} & \dots & \Sigma_{t_1 t_l} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_{t_l s_1} & \dots & \Sigma_{t_l s_m} & \Sigma_{t_l t_1} & \dots & \Sigma_{t_l t_l} \end{pmatrix}$$

$$\begin{aligned}
 \Sigma_{11} &\leftarrow \sigma_1^2 \\
 \Sigma_{ij} &\leftarrow \Sigma_{ii} B_{ij} \\
 \Sigma_{ji} &\leftarrow \Sigma_{ij}^T \\
 \Sigma_{jj} &\leftarrow \sigma_j^2 + \Sigma_{ji} B_{ij}
 \end{aligned}$$

$$\begin{aligned}
 & \Sigma_{s_2 s_2} \\
 &= \sigma_{S(Y)}^2 + \Sigma_{s_2(r_1 \dots r_n s_1)} B_{(r_1 \dots r_n s_1) s_2}
 \end{aligned}$$

$$= \sigma_{S(Y)}^2 + \begin{pmatrix} \sigma_{R(X)}^2 \beta_{rs} & \dots & \sigma_{R(X)}^2 \beta_{rs} & n \sigma_{R(X)}^2 \beta_{rs}^2 \end{pmatrix} \begin{pmatrix} \beta_{rs} \\ \vdots \\ \beta_{rs} \\ 0 \end{pmatrix}$$

$$= \sigma_{S(Y)}^2 + n \sigma_{R(X)}^2 \beta_{rs}^2 = 4 + n$$

Lifting the Full Joint: Simple Case

$$\begin{array}{ccc} \sigma_{R(X)}^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{R(X)}^2 \end{array} \rightarrow \begin{array}{l} \text{on-diagonal: } \sigma_{R(X)}^2 \\ \text{off-diagonal: } 0 \end{array}$$

$$\begin{array}{ccc} \sigma_{R(X)}^2 \beta_{rs} & \dots & \sigma_{R(X)}^2 \beta_{rs} \\ \vdots & \ddots & \vdots \\ \sigma_{R(X)}^2 \beta_{rs} & \dots & \sigma_{R(X)}^2 \beta_{rs} \end{array} \rightarrow \text{all: } \sigma_{R(X)}^2 \beta_{rs}$$

$$\begin{array}{ccc} \sigma_{S(Y)}^2 + n\sigma_{R(X)}^2 \beta_{rs}^2 & \dots & n\sigma_{R(X)}^2 \beta_{rs}^2 \\ \vdots & \ddots & \vdots \\ n\sigma_{R(X)}^2 \beta_{rs}^2 & \dots & \sigma_{S(Y)}^2 + n\sigma_{R(X)}^2 \beta_{rs}^2 \end{array} \rightarrow \begin{array}{l} \text{on-diagonal: } \sigma_{S(Y)}^2 + n\sigma_{R(X)}^2 \beta_{rs}^2 \\ \text{off-diagonal: } n\sigma_{R(X)}^2 \beta_{rs}^2 \end{array}$$

$$\begin{array}{l} R(X) \\ S(Y) \\ T(Z) \end{array} \left(\begin{array}{ccccccccc} & \begin{array}{c} R(X) \\ S(Y) \\ T(Z) \end{array} & & & & & & & \\ & \Sigma_{r_1 r_1} & \dots & \Sigma_{r_1 r_n} & \Sigma_{r_1 s_1} & \dots & \Sigma_{r_1 s_m} & \Sigma_{r_1 t_1} & \dots & \Sigma_{r_1 t_l} \\ & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ & \Sigma_{r_n r_1} & \dots & \Sigma_{r_n r_n} & \Sigma_{r_n s_1} & \dots & \Sigma_{r_n s_m} & \Sigma_{r_n t_1} & \dots & \Sigma_{r_n t_l} \\ & \Sigma_{s_1 r_1} & \dots & \Sigma_{s_1 r_n} & \Sigma_{s_1 s_1} & \dots & \Sigma_{s_1 s_m} & \Sigma_{s_1 t_1} & \dots & \Sigma_{s_1 t_l} \\ & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ & \Sigma_{s_m r_1} & \dots & \Sigma_{s_m r_n} & \Sigma_{s_m s_1} & \dots & \Sigma_{s_m s_m} & \Sigma_{s_m t_1} & \dots & \Sigma_{s_m t_l} \\ & \Sigma_{t_1 r_1} & \dots & \Sigma_{t_1 r_n} & \Sigma_{t_1 s_1} & \dots & \Sigma_{t_1 s_m} & \Sigma_{t_1 t_1} & \dots & \Sigma_{t_1 t_l} \\ & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ & \Sigma_{t_l r_1} & \dots & \Sigma_{t_l r_n} & \Sigma_{t_l s_1} & \dots & \Sigma_{t_l s_m} & \Sigma_{t_l t_1} & \dots & \Sigma_{t_l t_l} \end{array} \right)$$

Lifting the Full Joint: Simple Case

$T(Z)$

$$\begin{array}{ccc}
 m\sigma_{R(X)}^2\beta_{rs}\beta_{st} & \dots & m\sigma_{R(X)}^2\beta_{rs}\beta_{st} \\
 \vdots & \ddots & \vdots \\
 m\sigma_{R(X)}^2\beta_{rs}\beta_{rs} & \dots & m\sigma_{R(X)}^2\beta_{rs}\beta_{rs}
 \end{array}
 \rightarrow \text{all: } m\sigma_{R(X)}^2\beta_{rs}\beta_{st} = m \cdot 4 \cdot 0.5 \cdot (-1) = -2m$$

$T(Z)$

$$\begin{array}{ccc}
 (\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2)\beta_{st} & \dots & (\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2)\beta_{st} \\
 \vdots & \ddots & \vdots \\
 (\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2)\beta_{st} & \dots & (\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2)\beta_{st}
 \end{array}
 \rightarrow \text{all: } (\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2)\beta_{st} = -4 - mn$$

$T(Z)$

$$\begin{array}{ccc}
 \sigma_{T(Z)}^2 + m\beta_{st}^2(\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2) & \dots & m\beta_{st}^2(\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2) \\
 \vdots & \ddots & \vdots \\
 m\beta_{st}^2(\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2) & \dots & \sigma_{T(Z)}^2 + m\beta_{st}^2(\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2)
 \end{array}$$

→ on-diagonal: $\sigma_{T(Z)}^2 + m\beta_{st}^2(\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2) = 3 + 4m + m^2n$

→ off-diagonal: $m\beta_{st}^2(\sigma_{S(Y)}^2 + mn\sigma_{R(X)}^2\beta_{rs}^2) = 4m + m^2n$

Lifted Joint

- Only two structures required for covariance matrix

- A matrix

$$\begin{array}{c}
 R(X) \\
 S(Y) \\
 T(Z)
 \end{array}
 \begin{array}{c}
 R(X) \\
 S(Y) \\
 T(Z)
 \end{array}
 \begin{array}{c}
 S(Y) \\
 T(Z) \\
 T(Z)
 \end{array}
 \begin{array}{c}
 0 \\
 \sigma_{R(X)}^2 \beta_{rs} \\
 m \sigma_{R(X)}^2 \beta_{rs} \beta_{st}
 \end{array}
 \begin{array}{c}
 \sigma_{R(X)}^2 \beta_{rs} \\
 n \sigma_{R(X)}^2 \beta_{rs}^2 \\
 (\sigma_{S(Y)}^2 + mn \sigma_{R(X)}^2 \beta_{rs}^2) \beta_{st}
 \end{array}
 \begin{array}{c}
 m \sigma_{R(X)}^2 \beta_{rs} \beta_{st} \\
 (\sigma_{S(Y)}^2 + mn \sigma_{R(X)}^2 \beta_{rs}^2) \beta_{st} \\
 m \beta_{st}^2 (\sigma_{S(Y)}^2 + mn \sigma_{R(X)}^2 \beta_{rs}^2)
 \end{array}
 \Bigg)$$

$$= \begin{pmatrix}
 4 & 2 & -2m \\
 2 & n & -4 - mn \\
 -2m & -4 - mn & 4m + m^2 n
 \end{pmatrix}$$

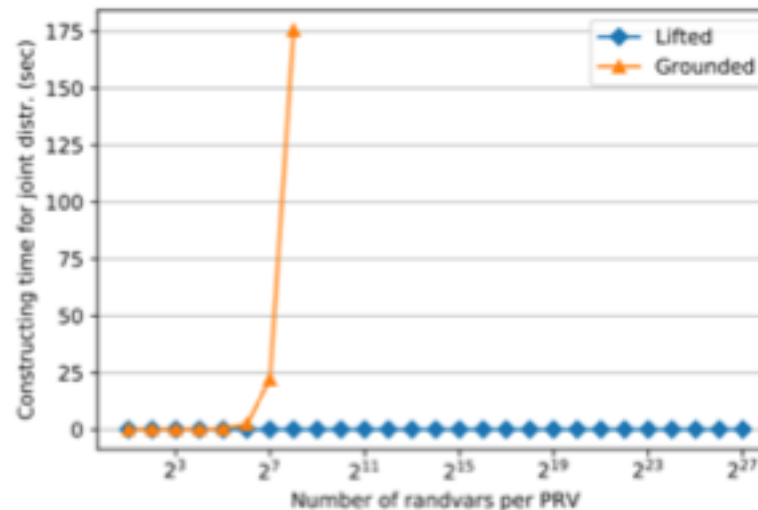
- A vector for on-diagonal covariance entries

- Individual variances
 - Have to be stored anyway

$$\begin{array}{c}
 R(X) \\
 S(Y) \\
 T(Z)
 \end{array}
 \begin{pmatrix}
 \sigma_{R(X)}^2 \\
 \sigma_{S(Y)}^2 \\
 \sigma_{T(Z)}^2
 \end{pmatrix}
 = \begin{pmatrix}
 4 \\
 4 \\
 3
 \end{pmatrix}$$

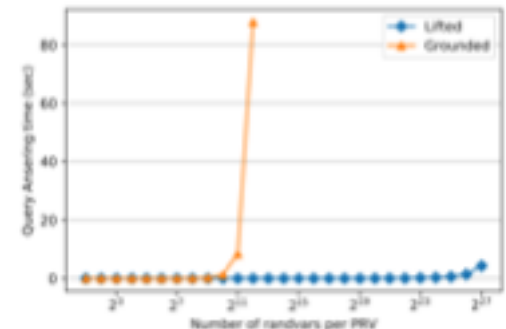
Lifted Joint

- Only two structures required for covariance matrix
- Depend only on the number of PRVs, not the domain sizes!



Lifted Query Answering

- Marginal queries
 - Read off values in (lifted) covariance representation
- Conditional queries $\mathbf{R}|\mathbf{E} = \mathbf{e} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$
 - $\boldsymbol{\mu}^* = \boldsymbol{\mu}_R + \boldsymbol{\Sigma}_{RE}\boldsymbol{\Sigma}_{EE}^{-1}(\mathbf{e} - \boldsymbol{\mu}_E)$
 - $\boldsymbol{\Sigma}^* = \boldsymbol{\Sigma}_{RR} - \boldsymbol{\Sigma}_{RE}\boldsymbol{\Sigma}_{EE}^{-1}\boldsymbol{\Sigma}_{ER}$
 - Matrix multiplication, inversion required
 - Possible to compute them in a lifted manner due to block structure
 - Proof in paper by Hartwig and Möller (2020)
 - Evidence is ground
 - Probably no symmetries in observations with real numbers as range values
 - unlikely to get identical observations
 - Fig.: 50% of ground instances get random values assigned as evidence

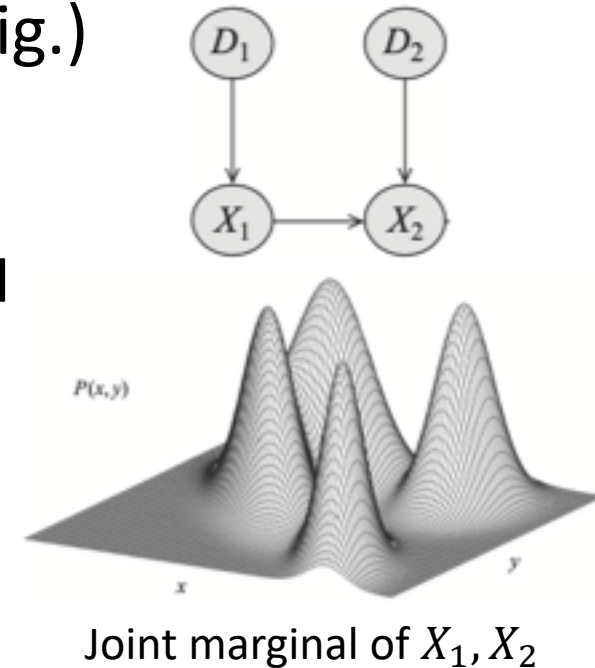


Interim Summary

- Linear Gaussian models
 - Linear dependency between child and parent randvars
 - Full joint given by vector of means and covariance matrix
 - Information form as inverse of covariance form
 - Query answering
 - Marginal using covariance matrix
 - Conditional using information form
- Gaussian BNs
 - Explicitly encode independencies in network structure
 - Conditional linear Gaussian
 - GBN = multivariate Gaussian distribution
 - Lifting for PRVs without an overlap in logvars between parent and child

Hybrid Models

- Models that contain discrete (D_i in fig.) and continuous randvars (X_i in fig.)
- Some general results
 - Even representing the correct marginal distribution in a hybrid network can require space that is exponential in the size of the network
 - Query answering problem is NP-hard even if the GBN is a polytree where all discrete randvars are Boolean-valued and where every continuous randvar has at most one discrete ancestor
 - There are not even approximate algorithms to solve the problem in polynomial time with a useful error bound without further restrictions



Outline: 8. Continuous Space

A. *Basics*

- Continuous variables, probability density function, cumulative probability distribution
- Joint distribution, marginal density, conditional density

B. *Gaussian models*

- (Multivariate) Gaussian distribution
- (Parameterised) Gaussian Bayesian networks

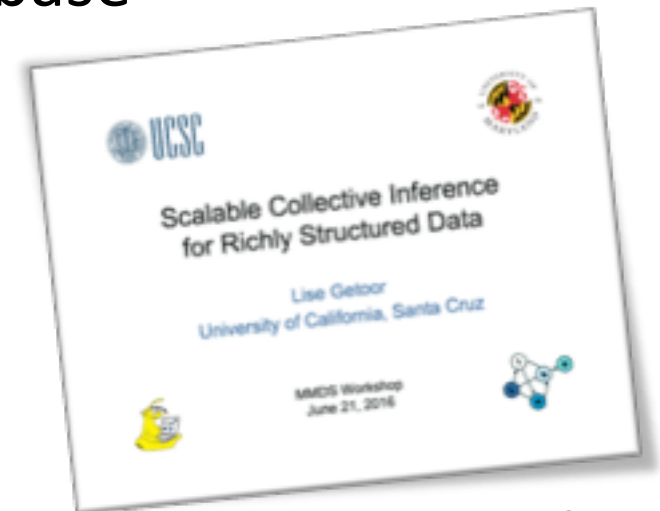
C. ***Probabilistic Soft Logic (PSL)***

- Modelling, semantics, inference task

Probabilistic Soft Logic (PSL)

- Logic-based approach
- Probabilistic programming language
 - Predicate = relationship or property
 - Atom = continuous randvar
 - Rule = dependency or constraint
 - Set = define aggregates
- PSL program = rules + input database
- Implementation:
<https://psl.linqs.org>

Based on slides by Lise Getoor, “Scalable Collective Inference for Richly Structured Data”, MMDS Workshop 2016.



Syntax & Semantics

- Let \mathbf{R} be a set of weighted logical rules, each R_j has the form

$$w_j : \bigwedge_{i \in I_j^-} x_i \Rightarrow \bigvee_{i \in I_j^+} x_i$$

- $w_j \geq 0$
- Sets I_j^-, I_j^+ index conjuncted/disjuncted literals
- Equivalent clausal form:

$$\left(\bigvee_{i \in I_j^+} x_i \right) \vee \left(\bigvee_{i \in I_j^-} \neg x_i \right)$$

- Probability distribution (compare: MLNs)

$$P(\mathbf{x}) \propto \exp \left(\sum_{R_j \in \mathbf{R}} w_j \left(\bigvee_{i \in I_j^+} x_i \right) \vee \left(\bigvee_{i \in I_j^-} \neg x_i \right) \right)$$

MPE Inference

- MPE: Find the most probable assignment to the unobserved randvars

- I.e., given a model ground over an input database,

$$\operatorname{argmax}_x \sum_{R_j \in \mathcal{R}} w_j \left(\bigvee_{i \in I_j^+} x_i \right) \vee \left(\bigvee_{i \in I_j^-} \neg x_i \right)$$

- Combinatorial, NP-hard
- Approximation:
View as optimising rounding probabilities

Expected Score

- Expected score of a clause is the weight times the probability that **at least one literal is true**:

$$w_j \left(1 - \prod_{i \in I_j^+} (1 - p_i) \prod_{i \in I_j^-} p_i \right)$$

At least one literal true
→ or-semantics → trick:
Instead of computing
 $P(A \vee B)$
 $= P(A) + P(B) - P(A \wedge B)$
compute
 $P(\neg \neg(A \vee B))$
 $= 1 - P(\neg A \wedge \neg B)$

- Then, expected total score is

$$\widehat{W} = \sum_{R_j \in \mathbf{R}} w_j \left(1 - \prod_{i \in I_j^+} (1 - p_i) \prod_{i \in I_j^-} p_i \right)$$

- But, $\operatorname{argmax}_p \widehat{W}$ highly non-convex due to product

Approximate Inference

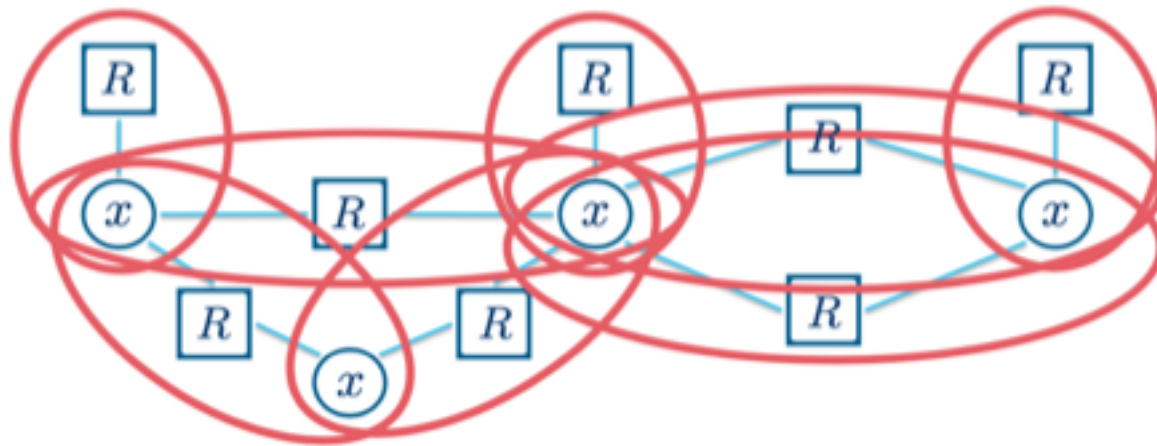
- Instead: Optimise a **linear program** that bounds expected score

$$\sum_{R_j \in \mathbf{R}} w_j \left(1 - \prod_{i \in I_j^+} (1 - p_i) \prod_{i \in I_j^-} p_i \right) \geq \left(1 - \frac{1}{e} \right) \sum_{R_j \in \mathbf{R}} w_j \min \left\{ \sum_{i \in I_j^+} p_i + \sum_{i \in I_j^-} (1 - p_i), 1 \right\}$$

- Can give $\left(1 - \frac{1}{e} \right)$ -optimal *discrete* solution

Scalable Approximate Inference

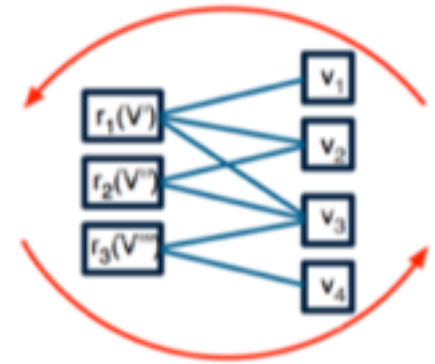
- Linear programming algorithms do not scale well to big probabilistic models



- Instead of solving the problem as one big optimisation, **decompose** the problem based on its graphical structure
 - Compare: cliques/clusters

Consensus Optimisation

- Decompose problem and solve sub-problems independently (in parallel), then merge results
 - Sub-problems are ground rules
 - Auxiliary variables enforce consensus across sub-problems
- Framework:
Alternating direction method of multipliers (ADMM) (Boyd, 2011)
 - Guaranteed to converge for convex problems
 - Inference with ADMM fast, scalable, straightforward to implement (Bach et al, 2017)



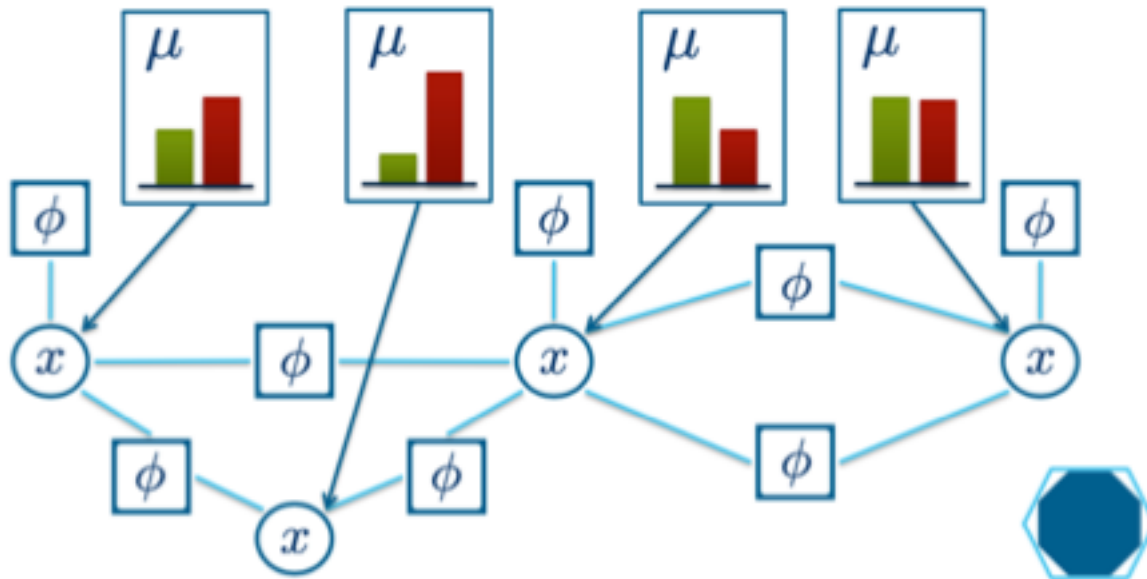
S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed Optimization and Statistical Learning Via the Alternating Direction Method of Multipliers. *Now Publishers*, 2011.

Stephan H. Bach. Matthias Broecheler, Bert Huang, and Lise Getoor: Hinge-Loss Markov Random Fields and Probabilistic Soft Logic. In: *Journal of Machine Learning Research* 18, 2017.

Local Consistency Relaxation

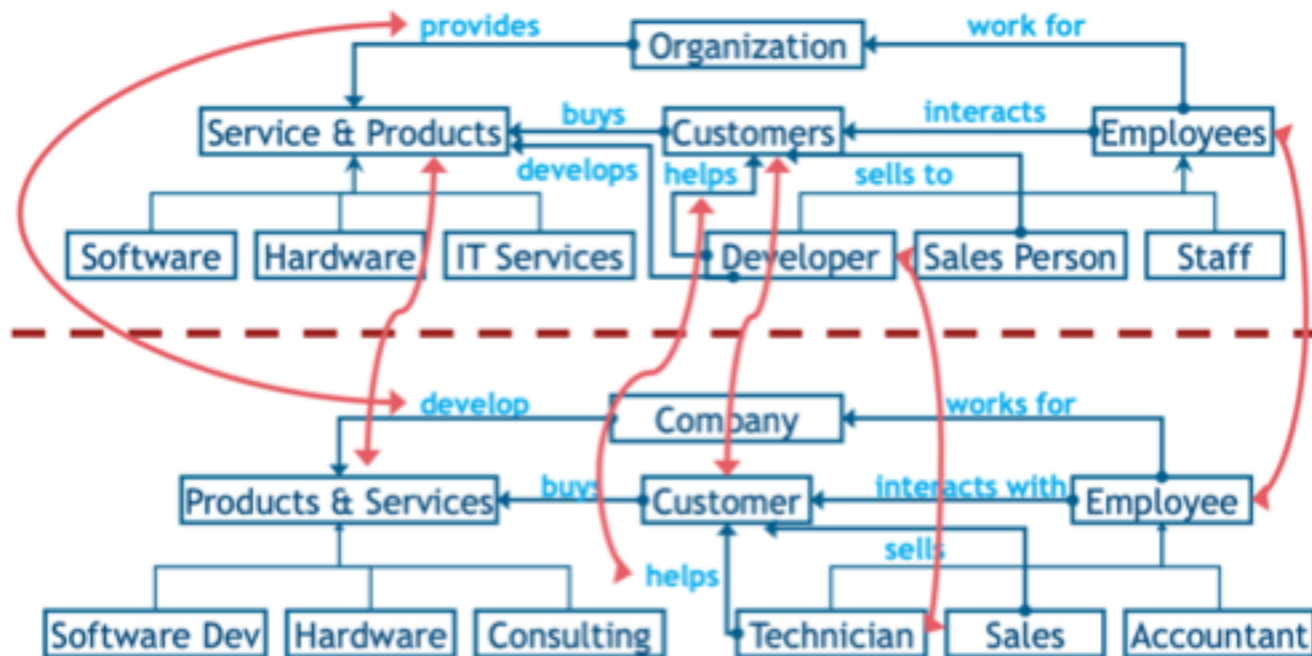
- Relax search over consistent marginals to simpler set

$$\operatorname{argmax}_{\mu \in [0,1]^n} \sum_{R_j \in \mathcal{R}} w_j \min \left\{ \sum_{i \in I_j^+} \mu_i + \sum_{i \in I_j^-} (1 - \mu_i), 1 \right\}$$



Continuous Variables & Similarity

- Continuous values interpreted as similarities
 - E.g., multiple ontologies → alignment
 - Match/Don't match → similar to what extent?



⇒ Soft logic

Soft Logic

- Logical operators defined for continuous values in the $[0,1]$ interval
 - Interpret as similarities or degree of truth
- Łukasiewicz logic
 - $p \wedge q = \max\{p + q - 1, 0\}$
 - $p \vee q = \min\{p + q, 1\}$
 - $\neg p = 1 - p$
- PSL: Use Łukasiewicz logic to interpret rules
 - Hinge-loss MNs (or Markov random fields as called in the publications by the PSL team) formalise this

Hinge-loss MNs

- Relaxed, logic-based MNs can reason about both discrete and continuous graph data scalably and accurately

- General objective

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{y} \in [0,1]^n} P(\mathbf{y}) \\ &= \operatorname{argmax}_{\mathbf{y} \in [0,1]^n} \sum_{j=1}^m w_j \min \left\{ \sum_{i \in I_j^+} y_i + \sum_{i \in I_j^-} (1 - y_i), 1 \right\} \\ &= \operatorname{argmin}_{\mathbf{y} \in [0,1]^n} \sum_{j=1}^m w_j \max \left\{ 1 - \sum_{i \in I_j^+} y_i - \sum_{i \in I_j^-} (1 - y_i), 0 \right\} \end{aligned}$$

- Notion of distance to satisfaction

Distance to Satisfaction

$$\operatorname{argmin}_{\mathbf{y} \in [0,1]^n} \sum_{j=1}^m w_j \max \left\{ 1 - \sum_{i \in I_j^+} y_i - \sum_{i \in I_j^-} (1 - y_i), 0 \right\}$$

- Maximum value of any unweighted term is 1
 - Term is *satisfied*
- Unsatisfied term \rightarrow distance to satisfaction
 - How far it is from achieving its maximum value
 - Each unweighted objective term measures how far the linear constraint is away from being satisfied:

$$1 - \sum_{i \in I_j^+} y_i - \sum_{i \in I_j^-} (1 - y_i) \leq 0$$

Relaxed Linear Constraints

- Instead of requiring logical clauses, each term can be defined using any function $\ell_j(\mathbf{y})$ linear in \mathbf{y}

$$\operatorname{argmin}_{\mathbf{y} \in [0,1]^n} \sum_{j=1}^m w_j \max\{\ell_j(\mathbf{y}), 0\}$$

- Each term represents the distance to satisfaction of a linear constraint $\ell_j(\mathbf{y}) \leq 0$
 - Can use logical clauses or something else based on domain knowledge
 - Also called **hinge losses**
 - Sometimes $\max\{\ell_j(\mathbf{y}), 0\}$ gets squared to better trade off conflicting objective terms
- Weight indicates how important it is to satisfy a constraint relative to others by scaling the distance to satisfaction

Hinge-loss MNs

- Let $\mathbf{y} = (y_1, \dots, y_n)$ be a vector of n randvars and $\mathbf{x} = (x_1, \dots, x_{n'})$ be a vector of n' randvars with joint range $\mathbf{D} = [0,1]^{n+n'}$

- Let $\boldsymbol{\phi} = (\phi_1, \dots, \phi_m)$ be a vector of m continuous potentials of the form

$$\phi_j(\mathbf{y}, \mathbf{x}) = (\max\{\ell_j(\mathbf{y}, \mathbf{x}), 0\})^{p_j}$$

- $\ell_j(\mathbf{y}, \mathbf{x})$ linear function of \mathbf{y}, \mathbf{x}
- $p_j \in \{1,2\}$
- For $(\mathbf{y}, \mathbf{x}) \in \mathbf{D}$ and given a vector of m weights $\mathbf{w} = (w_1, \dots, w_m)$, **constrained hinge-loss energy function** $f_{\mathbf{w}}$ is defined as

$$f_{\mathbf{w}}(\mathbf{y}, \mathbf{x}) = \sum_{j=1}^m w_j \phi_j(\mathbf{y}, \mathbf{x})$$

Hinge-loss MNs

- Let $c = (c_1, \dots, c_r)$ be a vector of linear constraint functions which further restrict the domain \mathbf{D} to \mathbf{D}'
- Hinge-loss MN over randvars \mathbf{y} and conditioned on randvars \mathbf{x} is a PDF defined as follows

- if $(\mathbf{y}, \mathbf{x}) \notin \mathbf{D}'$, then $P(\mathbf{y}|\mathbf{x}) = 0$
- if $(\mathbf{y}, \mathbf{x}) \in \mathbf{D}'$, then

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{w}, \mathbf{x})} \exp(-f_{\mathbf{w}}(\mathbf{y}, \mathbf{x}))$$

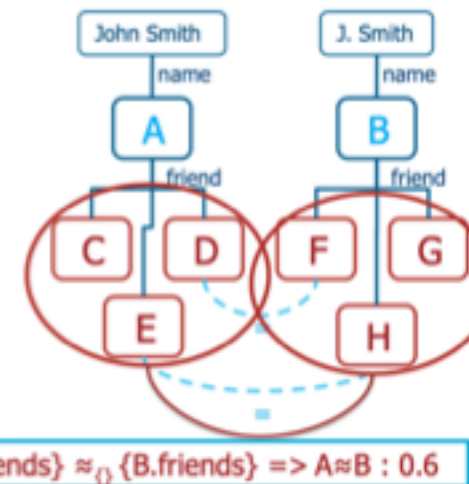
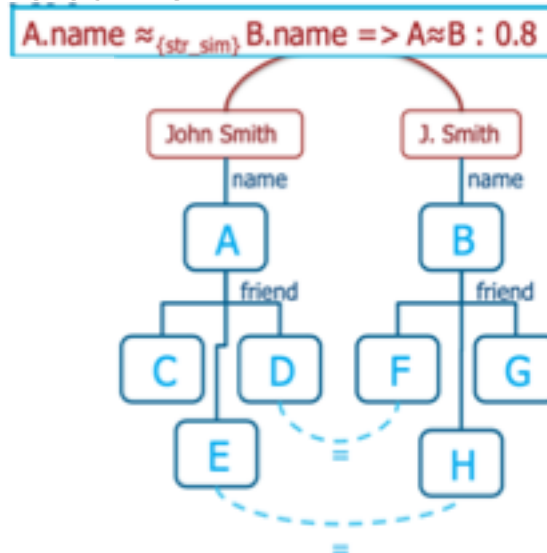
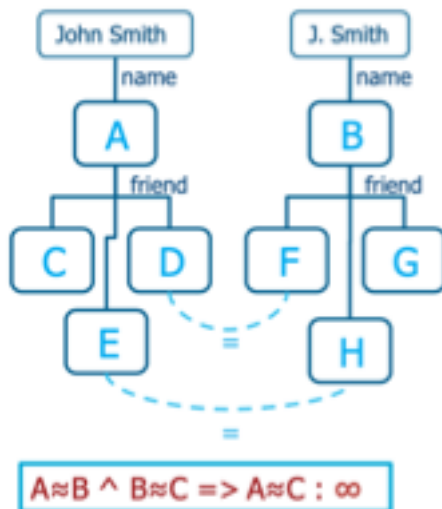
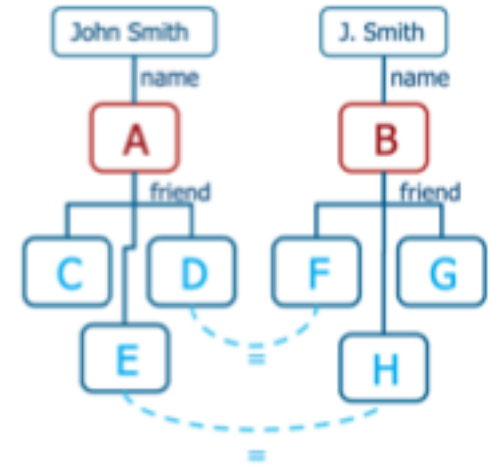
- where

$$Z(\mathbf{w}, \mathbf{x}) = \int_{\mathbf{y} | (\mathbf{y}, \mathbf{x}) \in \mathbf{D}'} \exp(-f_{\mathbf{w}}(\mathbf{y}, \mathbf{x})) d\mathbf{y}$$

- Define hinge-loss MNs using PSL

Application: E.g., Entity Resolution

- Goal: Identify references that denote the same person
- Use model to express dependencies
 - “If $A=B$ and $B=C$, then A and C must also denote the same person”
 - “If two people have similar names, they are probably the same”
 - “If two people have similar friends, they are probably the same”



Interim Summary

- PSL
 - Logic programming language
 - Approximations
 - Linear program that bounds MPE solution from below
 - Decomposition of PGM to optimise set of subproblems (consensus optimisation)
 - Local consistency relaxation
 - Soft logic: Łukasiewicz logic
 - Interpret continuous values as similarities/degree of truth
- Hinge-loss MNs
 - Notion of distance to satisfaction
 - Define using PSL

Outline: 8. Continuous Space

A. *Basics*

- Continuous variables, probability density function, cumulative probability distribution
- Joint distribution, marginal density, conditional density

B. *Gaussian models*

- (Multivariate) Gaussian distribution
- (Parameterised) Gaussian Bayesian networks

C. *Probabilistic Soft Logic (PSL)*

- Modelling, semantics, inference task

The End