

# Latent Tree Copulas

Sergey Kirshner

skirshne@purdue.edu

Purdue University

West Lafayette, IN, USA



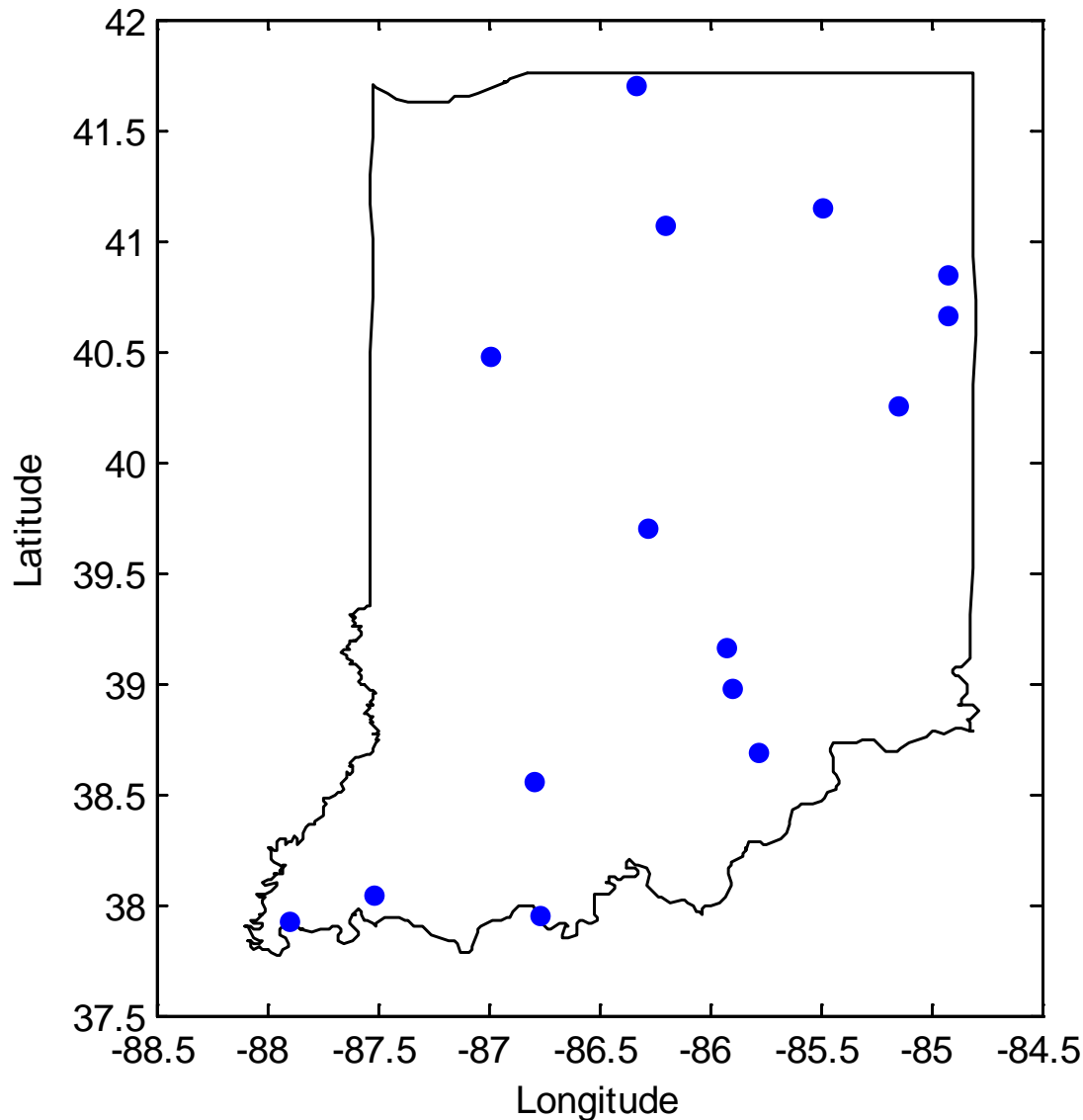
Granada, Spain,  
September 19, 2012



# Coming Attractions

- Want to fit density to model multivariate data?
  - and organize real-valued data into a hierarchy of features?
- New density estimation model based on tree-structured dependence with latent variables
  - Distribution = Univariate Marginals + Copula
  - Hierarchy of variables as a latent tree-copula
  - Parameter estimation and structure learning
    - Efficient inference for Gaussian copulas (100s of variables), several structure learning approaches
    - Variational inference for other copulas (10-30 variables)

# Building a Hierarchy of Rainfall Stations

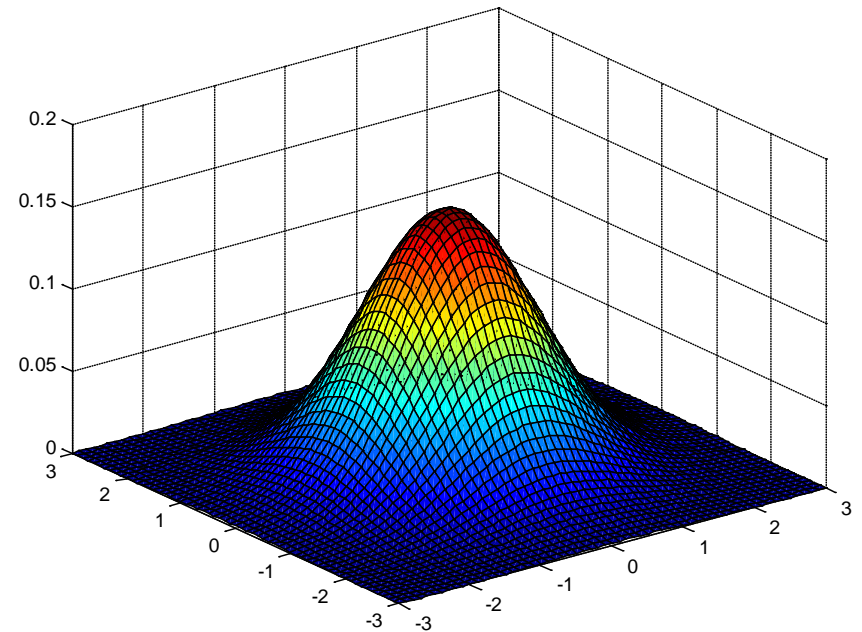


State of Indiana  
(USA)

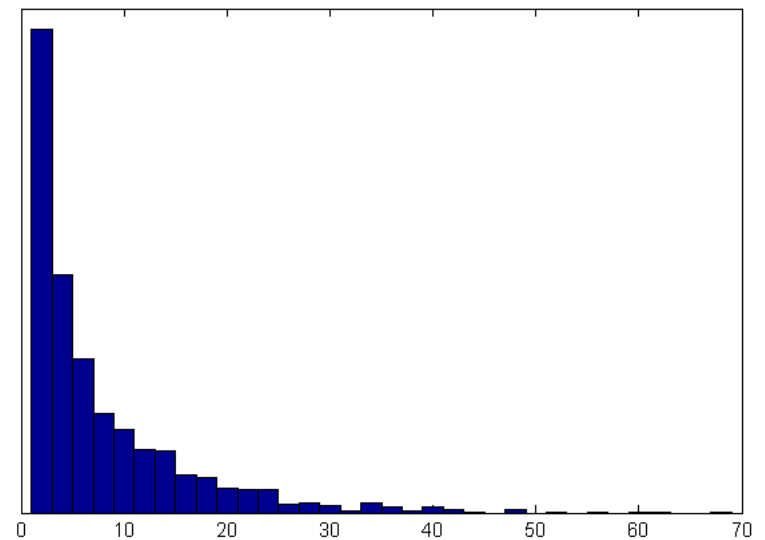
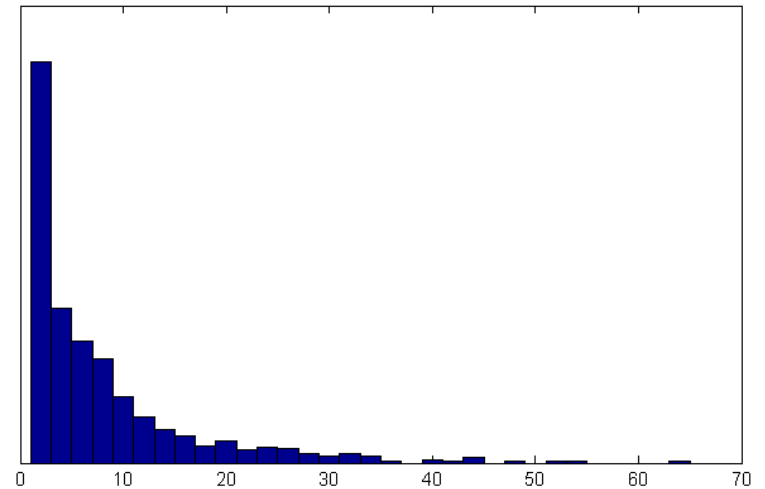
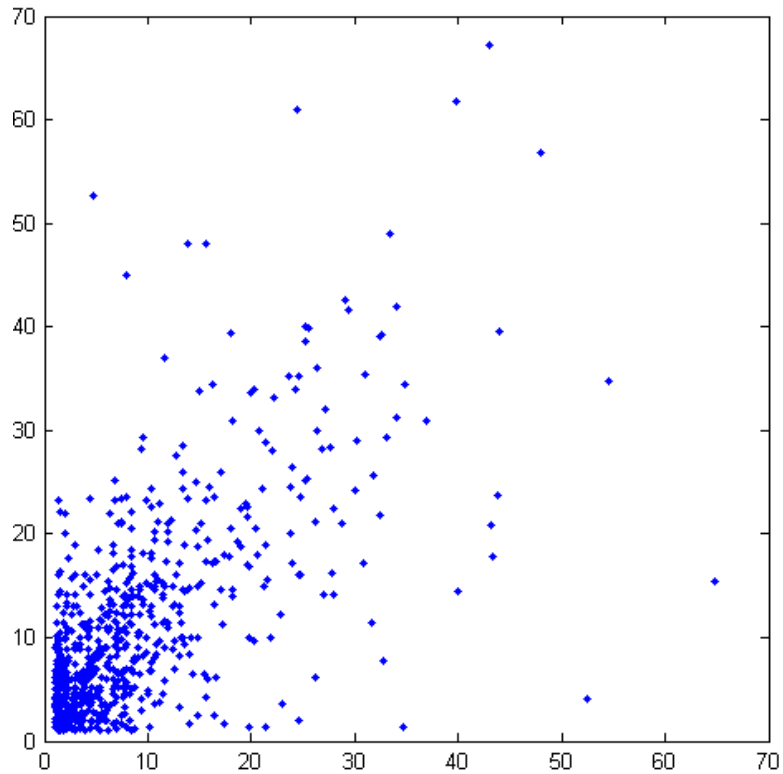
Average monthly  
observations for  
15 rainfall stations  
1951-1996 (47  
years)

# Most Popular Distribution...

- Interpretable
- Closed under taking marginals
- Generalizes to multiple dimensions
- Models pairwise dependence
- Tractable
- 245 pages out of 691 from *Continuous Multivariate Distributions* by Kotz, Balakrishnan, and Johnson



# What If the Data Is NOT Gaussian?



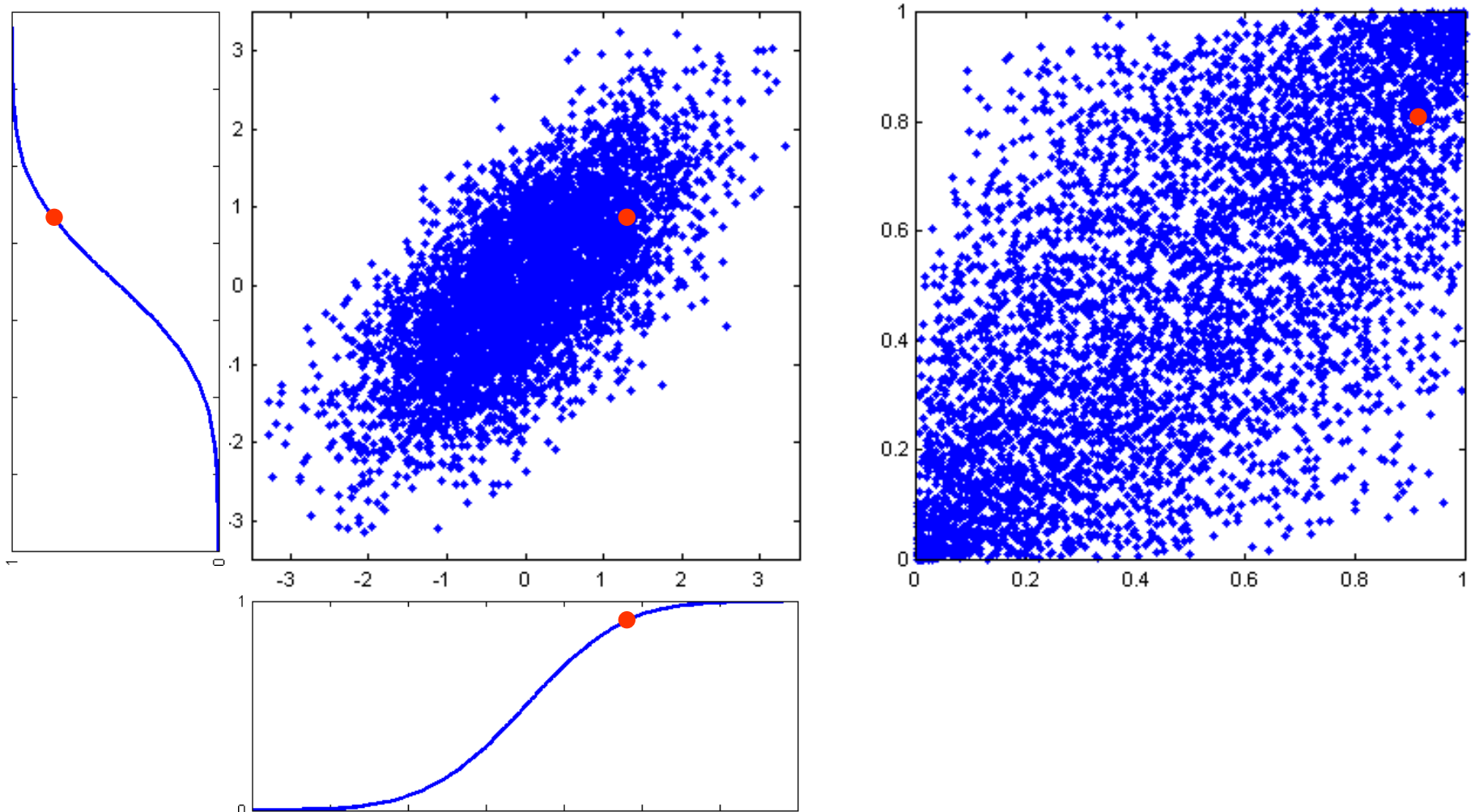
# Separating Univariate Marginals

$$\ln p(\mathbf{x}) = \sum_{i=1}^d \ln p_i(x_i) + \ln \frac{p(\mathbf{x})}{p_1(x_1) \dots p_d(x_d)}$$

univariate marginals,  
independent variables,

multivariate dependence term,  
copula

# Monotonic Transformation of the Variables



# Copula

**Copula**  $C$  is a multivariate distribution (cdf) defined on a unit hypercube with uniform univariate marginals:

$$C : [0, 1]^d \rightarrow [0, 1]$$

$$C_i(a_i) = a_i, \quad i = 1, \dots, d$$

$$F(\mathbf{x}) = C(F_1(x_1), \dots, F_d(x_d)), \quad \mathbf{x} \in \mathbb{R}^d$$

$$C(\mathbf{a}) = F(F_1^{-1}(a_1), \dots, F_d^{-1}(a_d))$$

$$a_i = F_i(x_i), \quad i = 1, \dots, d$$

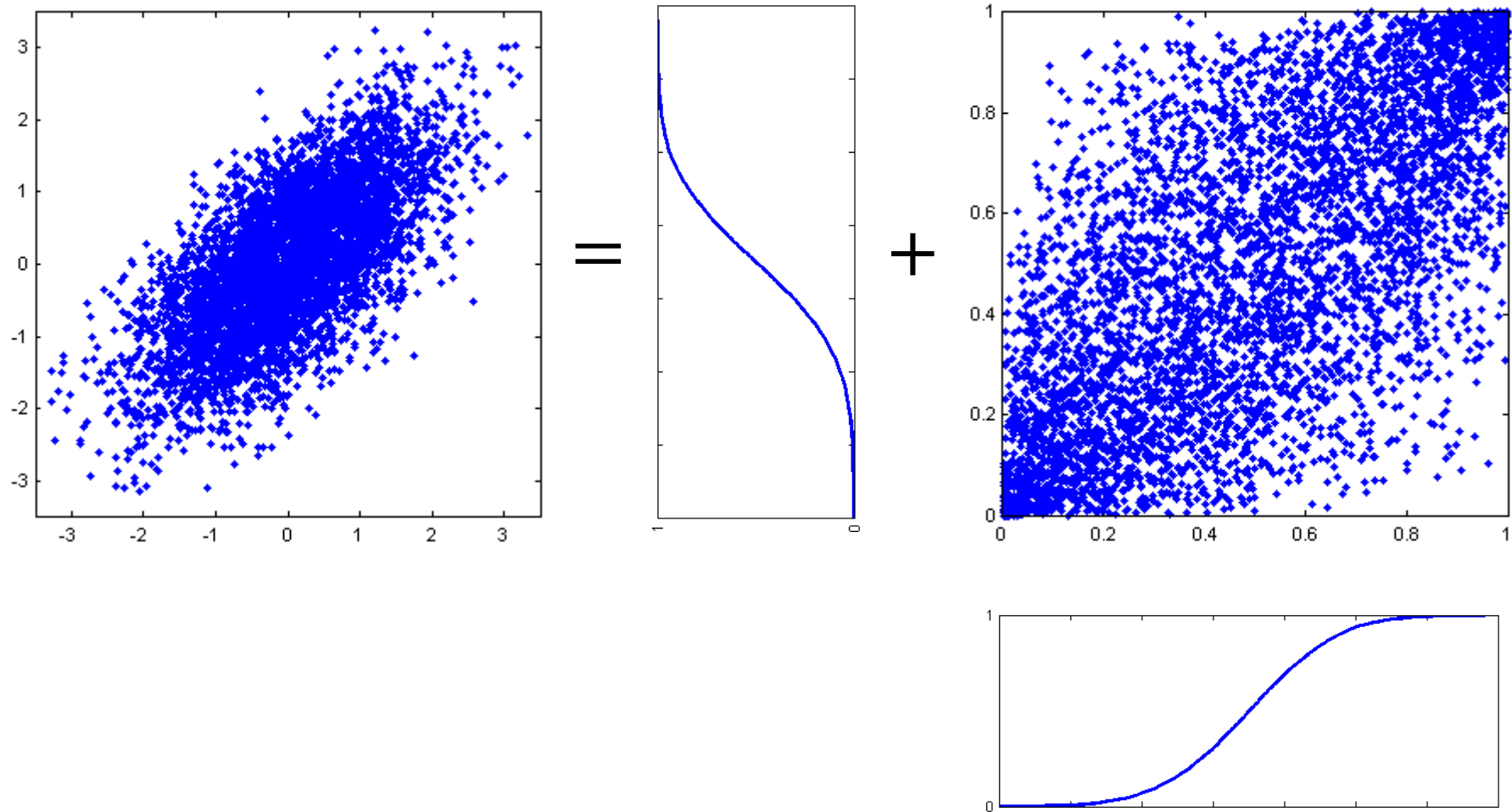
$$x_i = F_i^{-1}(a_i), \quad i = 1, \dots, d$$

$$c(\mathbf{a}) = \frac{\partial^d C(\mathbf{a})}{\partial a_1 \dots \partial a_d} = \frac{p(\mathbf{x})}{\prod_{i=1}^d p_i(x_i)}$$



# Sklar's Theorem

[Sklar 59]



# Example: Multivariate Gaussian Copula

$$F(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_d \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11}^2 & \cdots & \sigma_{11}\sigma_{dd}\rho_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{11}\sigma_{dd}\rho_{d1} & \cdots & \sigma_{dd}^2 \end{pmatrix}$$

$$F(x_1, \dots, x_d) = \Phi_{\mathbf{R}} \left( \frac{x_1 - \mu_1}{\sigma_{11}}, \dots, \frac{x_d - \mu_d}{\sigma_{dd}} \right)$$

$$a_1 = F_1(x_1) = \Phi \left( \frac{x_1 - \mu_1}{\sigma_{11}} \right), \dots, a_d = F_d(x_d) = \Phi \left( \frac{x_d - \mu_d}{\sigma_{dd}} \right)$$

$$C(\mathbf{a}) = F \left( F_1^{-1}(a_1), \dots, F_d^{-1}(a_d) \right) = \Phi_{\mathbf{R}} \left( \Phi^{-1}(a_1), \dots, \Phi^{-1}(a_d) \right)$$

# Separating Univariate Marginals

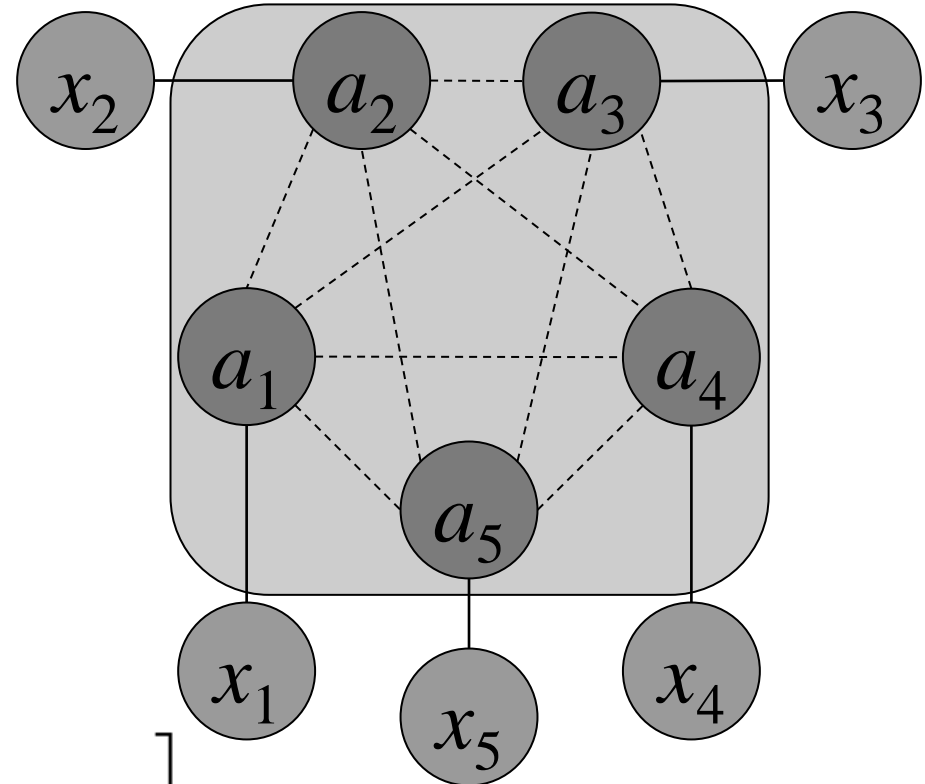
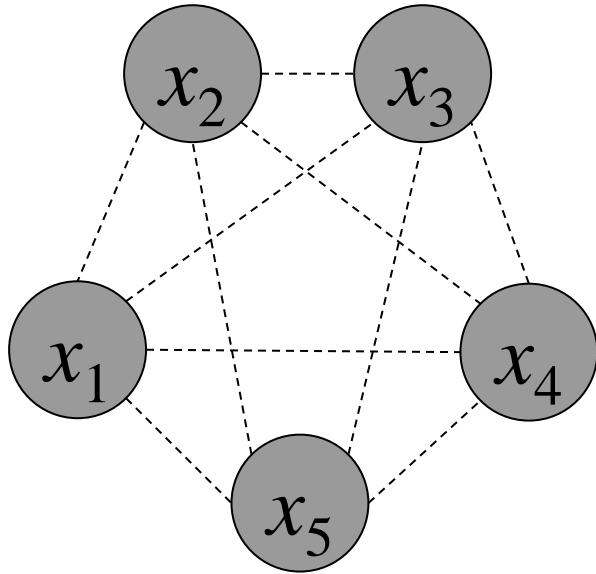
$$\ln p(\mathbf{x}) = \sum_{i=1}^d \ln p_i(x_i) + \ln \frac{p(\mathbf{x})}{p_1(x_1) \dots p_d(x_d)}$$

$$\ln p(\mathcal{D}) = \sum_{n=1}^N \sum_{i=1}^d \ln p_i(x_i^n) + \sum_{n=1}^N \ln c(F_1(x_1^n), \dots, F_d(x_d^n))$$

1. Fit univariate marginals (parametric or non-parametric)
2. Replace data points with cdf's of the marginals
3. Estimate copula density

Inference for the margins [Joe and Xu 96]; canonical maximum likelihood [Genest et al 95]

# Graphical Model Using a Copula

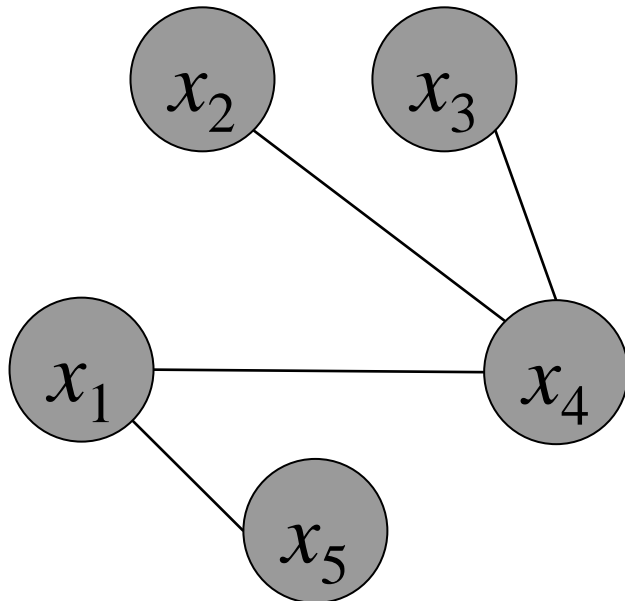


$$\begin{aligned}
 f(x_1, \dots, x_d) &= \left[ \prod_{i=1}^d f_i(x_i) \right] c(F_1(x_1), \dots, F_d(x_d)) \\
 &= \left[ \prod_{i=1}^d \phi_i(x_i, a_i) \right] c(a_1, \dots, a_d)
 \end{aligned}$$

# Graphical Model Approaches to Estimating Copulas

- Vines [[Bedford and Cooke 02](#)]
- Trees [[Kirshner 08](#)]
- Nonparanormal model [[Liu et al 09](#)]
- Copula Bayesian networks [[Elidan 10](#)]

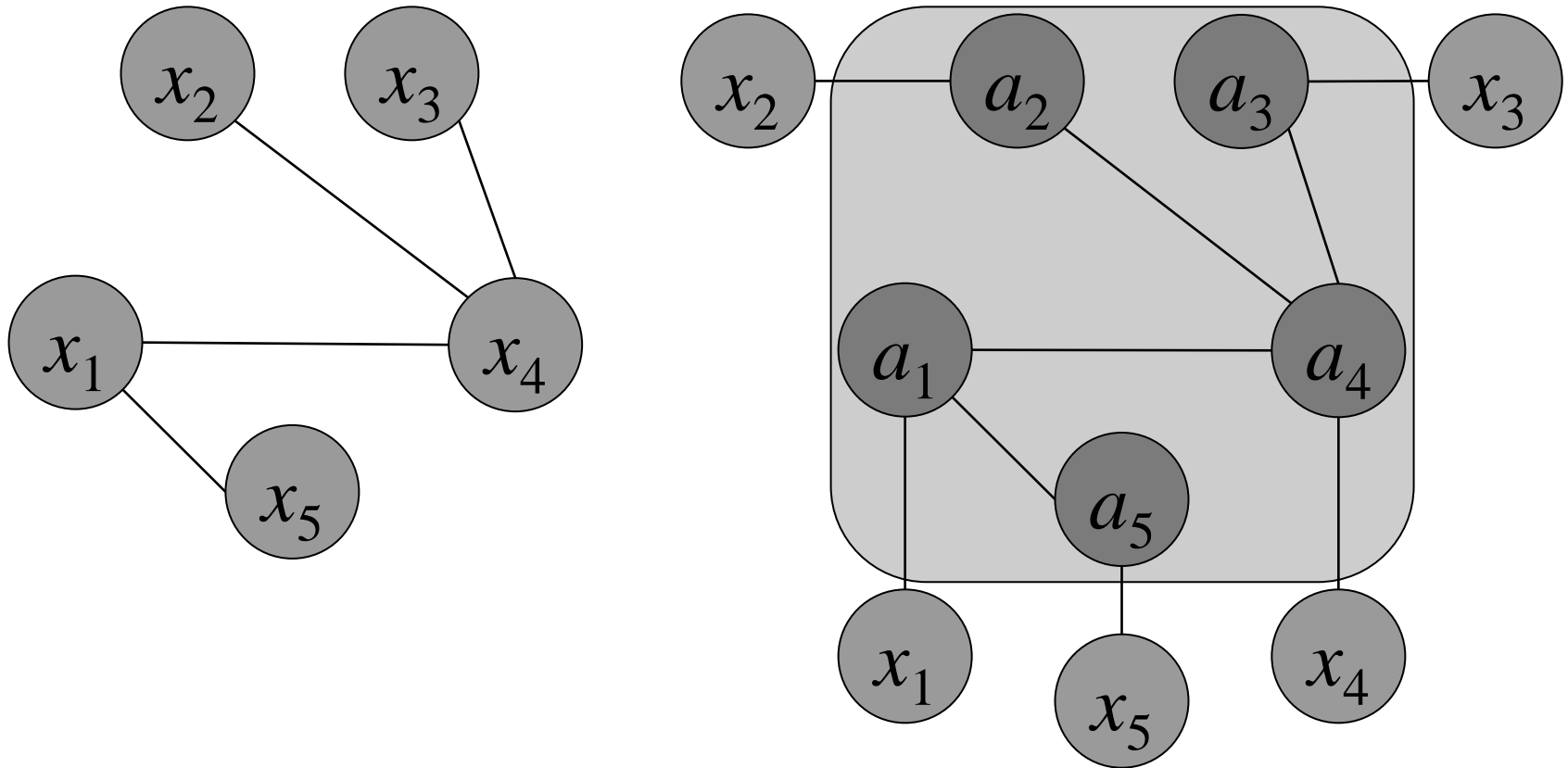
# Tree-Structured Densities



$$p_{\mathcal{E}}(\mathbf{x}) = \left[ \prod_{v \in \mathcal{V}} p_v(x_v) \right] \left[ \prod_{\{u,v\} \in \mathcal{E}} \frac{p_{uv}(x_u, x_v)}{p_u(x_u) p_v(x_v)} \right]$$

# Tree-Structured Copulas

[Kirshner 08]



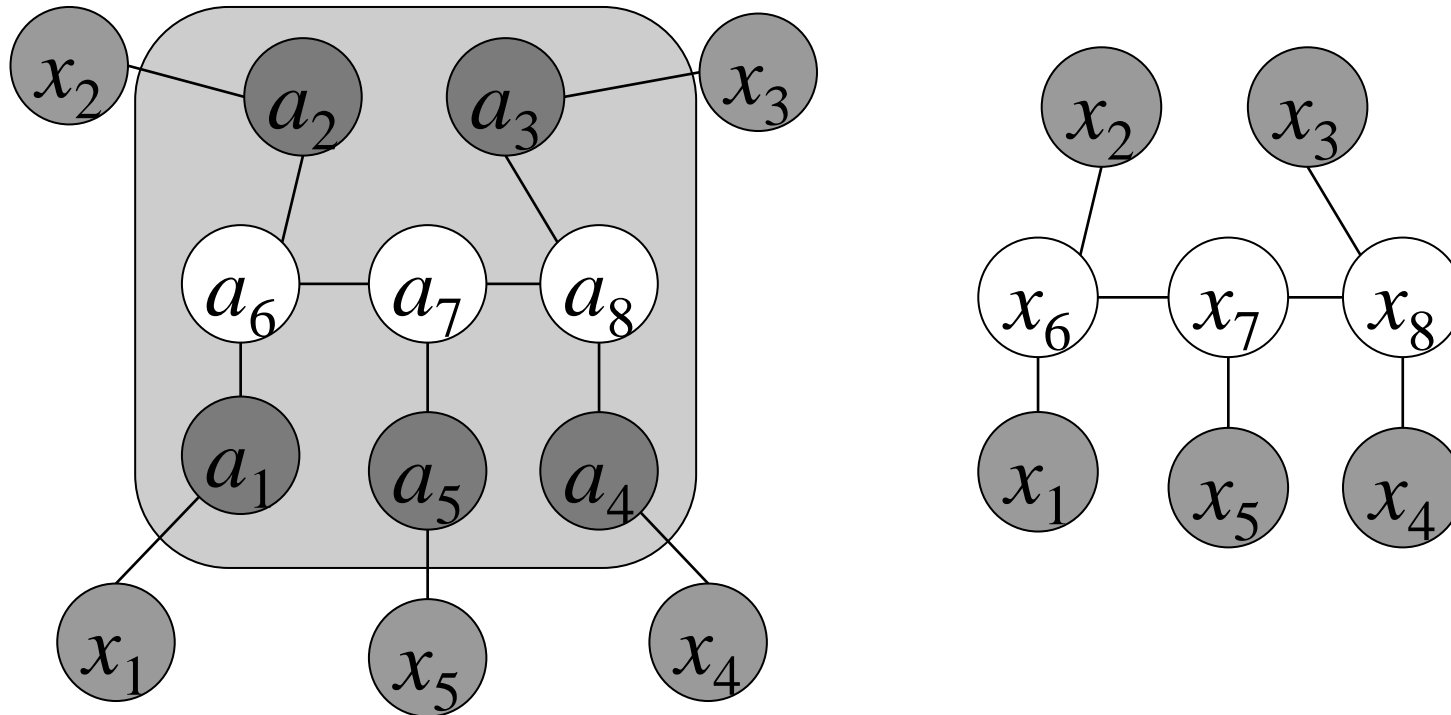
$$c_{\mathcal{E}}(\mathbf{a}) = \prod_{\{u,v\} \in \mathcal{E}} c_{uv}(a_u, a_v)$$

# Using Tree-Structured Copulas

- Tree-structured copulas are convenient, but are restrictive
  - True distribution may require much larger cliques to decompose
- Can approximate other dependencies using latent variables
  - Mixtures [\[Kirshner 08\]](#): discrete latent variables
  - [Latent tree copulas](#): continuous random variables embedded in copula trees



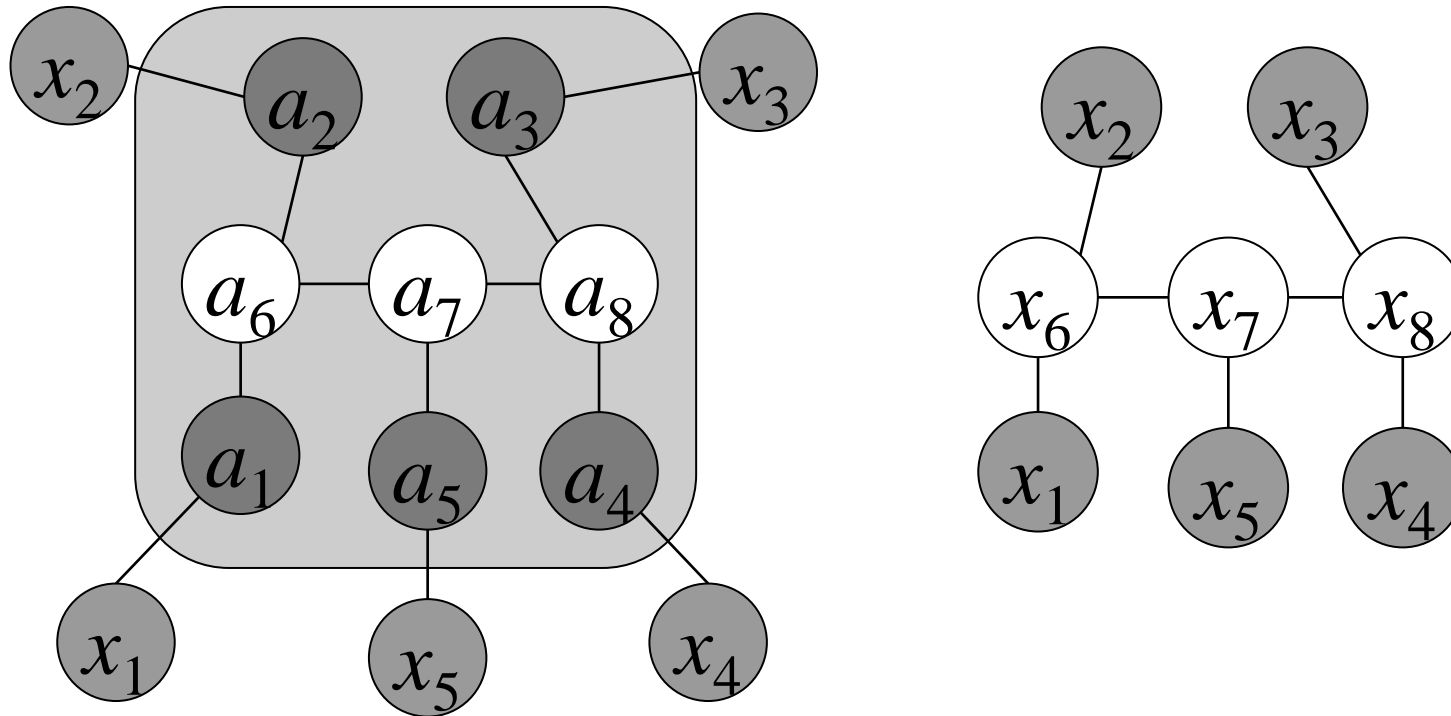
# Latent Tree Copulas



- Defined as a tuple of variables, tree structure, and bivariate copulas

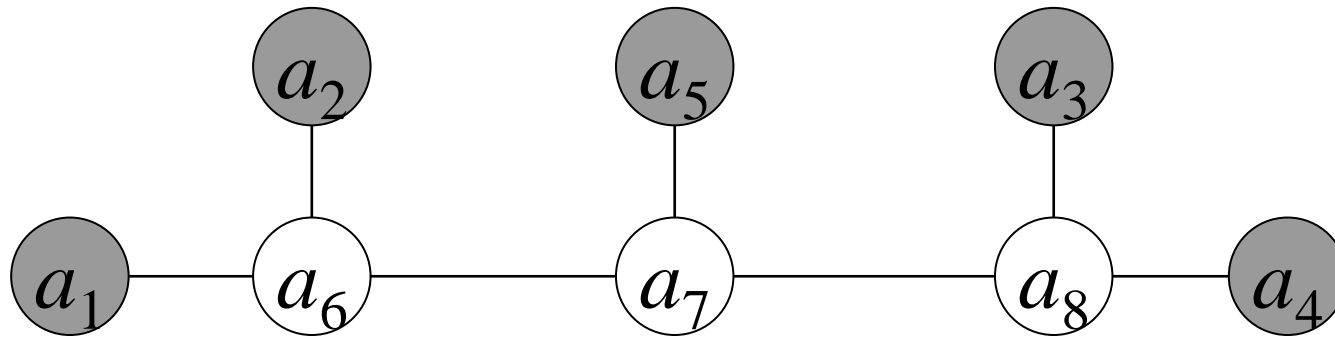
$$c_{LT}(\mathbf{a}) = \iint \prod_{\{u,v\} \in \mathcal{E}} c_{uv}(a_u, a_v) da_{d+1} \dots da_t$$

# Latent Tree Copulas



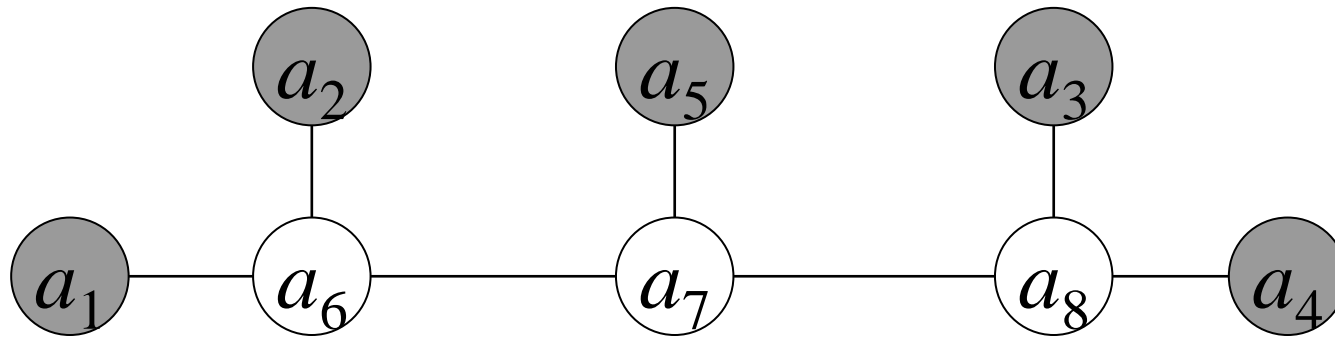
- Defined as a tuple of variables, tree structure, and bivariate copulas
- “Siblings” of latent tree models (LTMs) for categorical variables [e.g., Zhang 02, 04]

# Inference



- **Good news:** posterior distribution is also tree-structured
  - Fairly easy to carry out inference for LTMs
- **Bad news:** Latent variables are continuous: infinite number of possible values
  - Need to estimate the joint posterior densities

# Inference



- Easy for Gaussian copulas
  - Apply inverse standard normal CDF; use belief propagation on jointly Gaussian distribution
- Difficult for non-Gaussian copulas
  - May have no exact form for the posterior!

# Inference for non-Gaussian Case

- Variational approach:
  - Approximate the posterior distribution using a tree-structured distribution over piece-wise uniform variables
  - Essentially, approximate using the tree over categorical variables

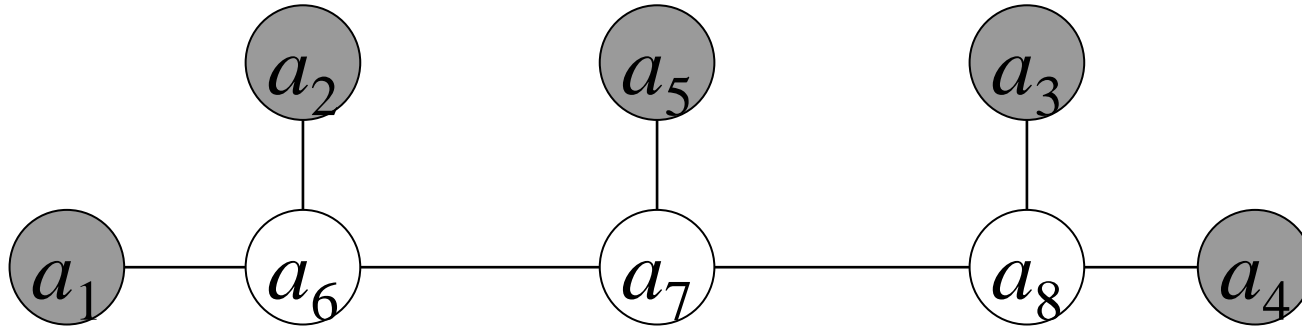
$$q(\mathbf{a}_H) = \prod_{u \in H} q_u(a_u) \left[ \prod_{\{u,v\} \in \mathcal{E}_H} \frac{q_{uv}(a_u, a_v)}{q_u(a_u) q_v(a_v)} \right]$$

$$q_{uv}(a_u, a_v) = p_{uv}(i, j) \geq 0 \text{ for } a_u \in \mathbb{I}_i, a_v \in \mathbb{I}_j,$$

$$q_u(a_u) = p_u(i) \text{ for } a_u \in \mathbb{I}_i, \quad \text{where } \mathbb{I}_i = \left( \frac{i-1}{K}, \frac{i}{K} \right]$$

$$\operatorname{argmin}_{q \in \mathcal{Q}} D \left( q^n(\mathbf{a}_H^n) \parallel c_{LT}(\mathbf{a}_H^n | \mathbf{a}_O^n, \boldsymbol{\theta}') \right)$$

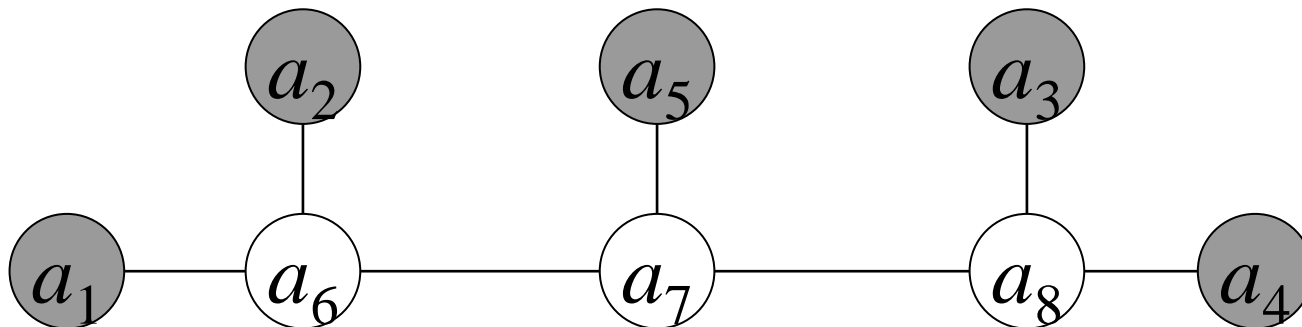
# Parameter Estimation with Known Structure



- (Variational) EM
  - E-step: minimize KL divergence
  - M-step: maximize the expected complete-data log-likelihood

$$l(\theta) = \sum_{i=1}^n \int_{\mathbb{I}^{t-d}} q^n(\mathbf{a}_H^n) \ln \frac{c_{LT}(\mathbf{a}_O^n, \mathbf{a}_H^n | \theta')}{q^n(\mathbf{a}_H^n)} d\mathbf{a}_H^n + \sum_{i=1}^n D(q^n(\mathbf{a}_H^n) \| c_{LT}(\mathbf{a}_H^n | \mathbf{a}_O^n, \theta'))$$

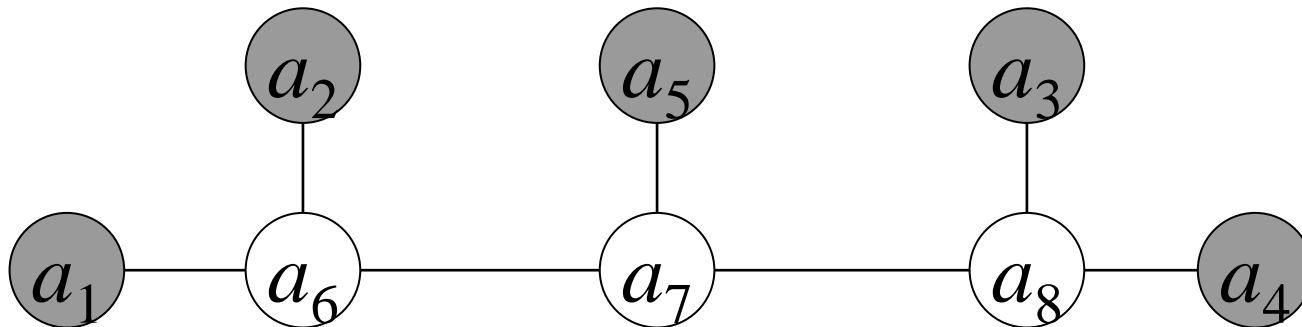
# Parameter Estimation with Known Structure



- Gaussian copula case: EM  $q^n(\mathbf{a}_H^n) = c(\mathbf{a}_H^n | \mathbf{a}_O^n | \boldsymbol{\theta})$ 
  - E-step: closed form inference,  $O(Nt)$  per iteration
  - M-step: maximize the expected complete-data log-likelihood

$$l(\boldsymbol{\theta}) = \sum_{i=1}^n \int_{\mathbb{I}^{t-d}} q^n(\mathbf{a}_H^n) \ln \frac{c_{LT}(\mathbf{a}_O^n, \mathbf{a}_H^n | \boldsymbol{\theta}')}{q^n(\mathbf{a}_H^n)} d\mathbf{a}_H^n + \sum_{i=1}^n D(q^n(\mathbf{a}_H^n) \| c_{LT}(\mathbf{a}_H^n | \mathbf{a}_O^n, \boldsymbol{\theta}'))$$

# Parameter Estimation with Known Structure



- Non-gaussian copula case: variational EM
  - E-step: approximate inference,  $O(sN/E/k^2) + |E|$  bivariate integrals per iteration
  - M-step: approximate maximization, need to update  $|E|$  bivariate copula parameters

$$l(\theta) = \sum_{i=1}^n \int_{\mathbb{I}^{t-d}} q^n(\mathbf{a}_H^n) \ln \frac{c_{LT}(\mathbf{a}_O^n, \mathbf{a}_H^n | \theta')}{q^n(\mathbf{a}_H^n)} d\mathbf{a}_H^n$$

$$+ \sum_{i=1}^n D(q^n(\mathbf{a}_H^n) \| c_{LT}(\mathbf{a}_H^n | \mathbf{a}_O^n, \theta'))$$



# Unknown Structure

- Gaussian LTCs: same as for tree-structured Gaussians
  - Size of possible trees can be limited
  - e.g., can use information distances [Choi et al 2011]
- Non-Gaussian LTCs: need to restrict the space of possible models
  - Very large space of structures/copula families
  - Fix the bivariate copula family
  - Consider only **binary** latent tree copulas
    - Observed nodes = leafs
    - Motivation: Any Gaussian LTC is equivalent to some binary LTC

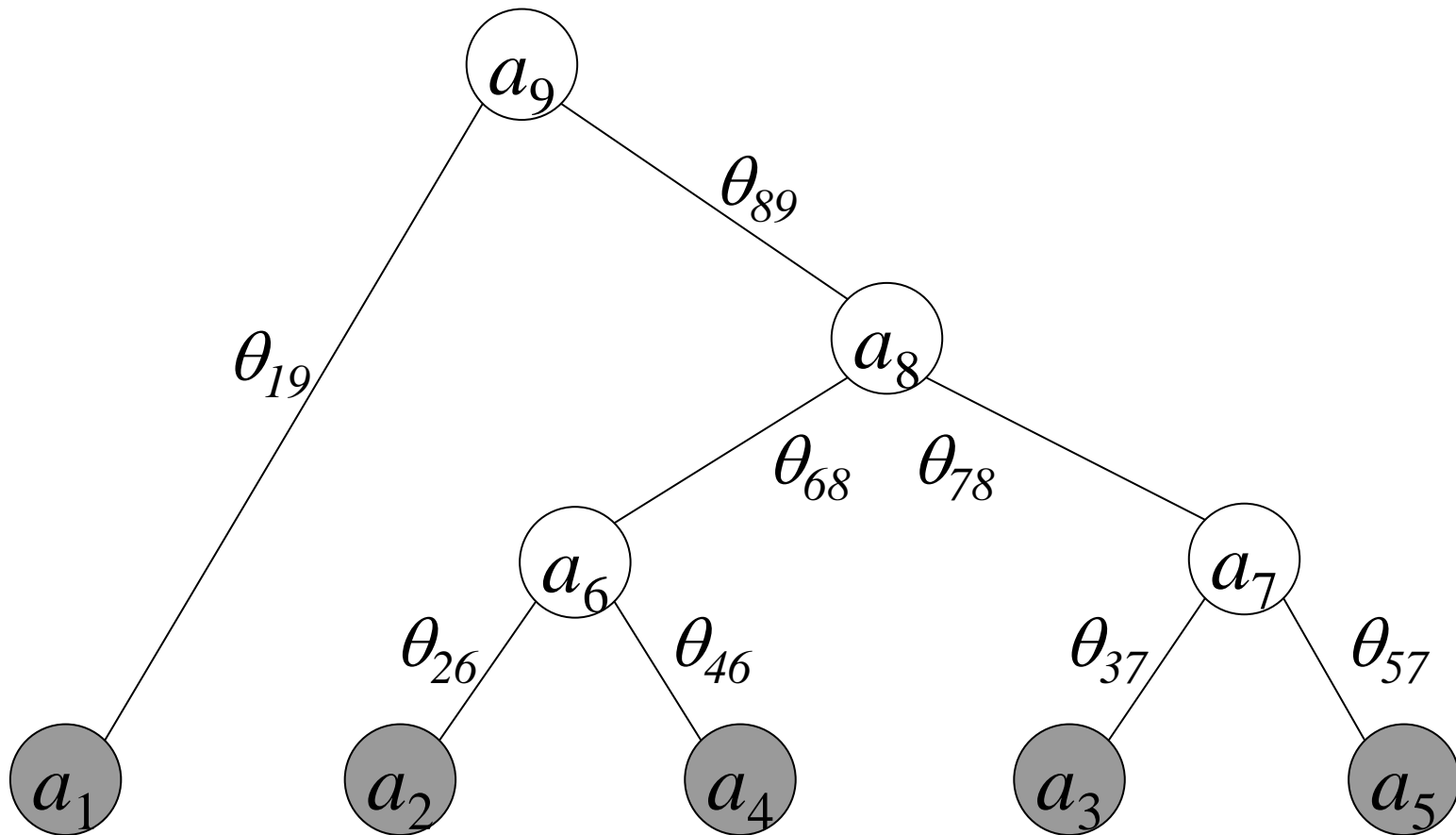
# Bottom-up Binary LTC Learning

[Similar to [Harmeling and Williams 11](#)]

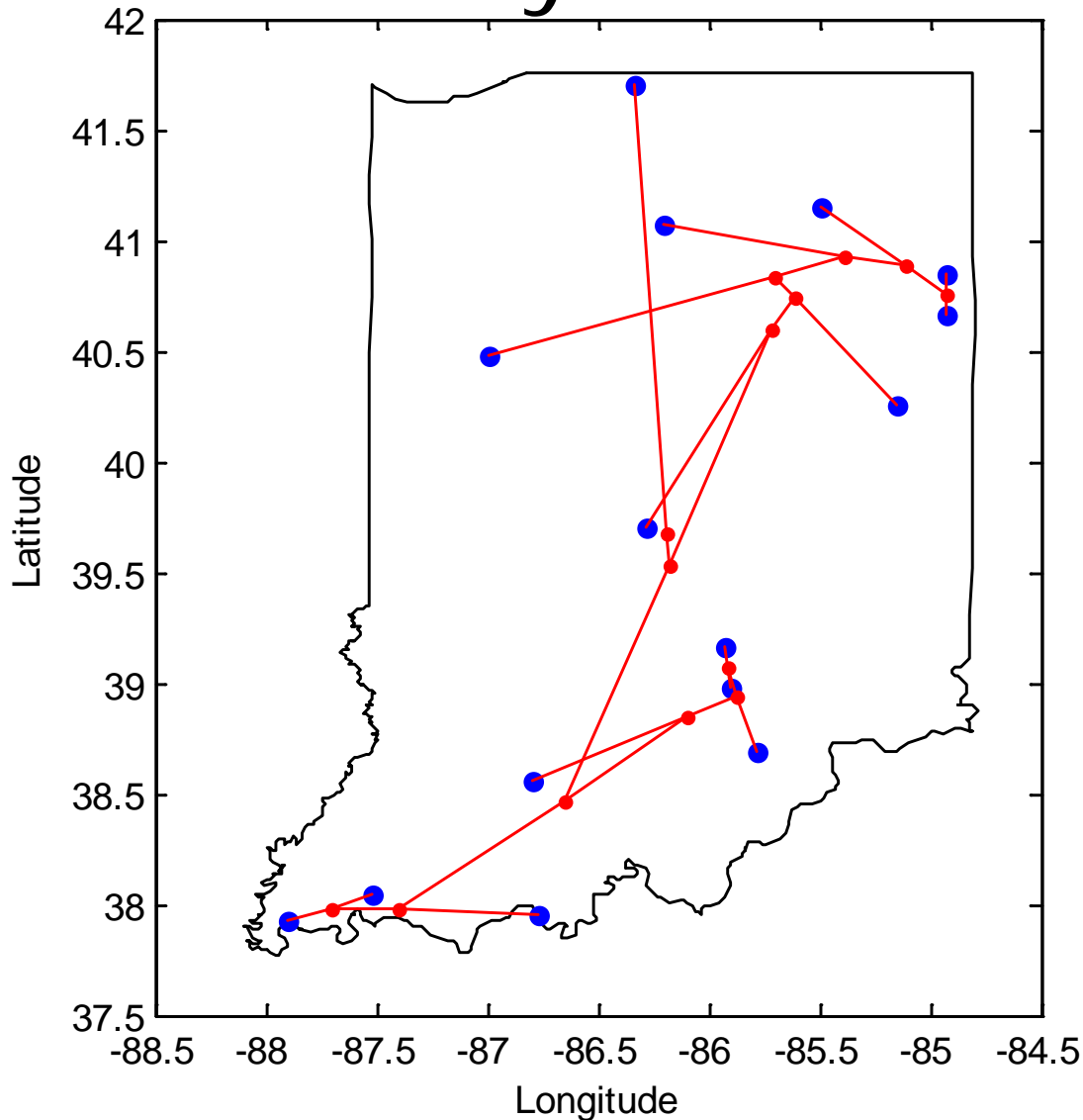
- Initialize the subtrees to consist of individual variables (variable = root of a subtree)
- Iterate until all variables are in one tree
  - Estimate mutual informations (MI) between the root nodes
  - Pick the pair of roots with the largest MI
  - Merge the subtrees by creating a new latent root node
  - Re-estimate parameters (EM)

# Bottom-up Binary LTC Learning

[Similar to [Harmeling and Williams 11](#)]



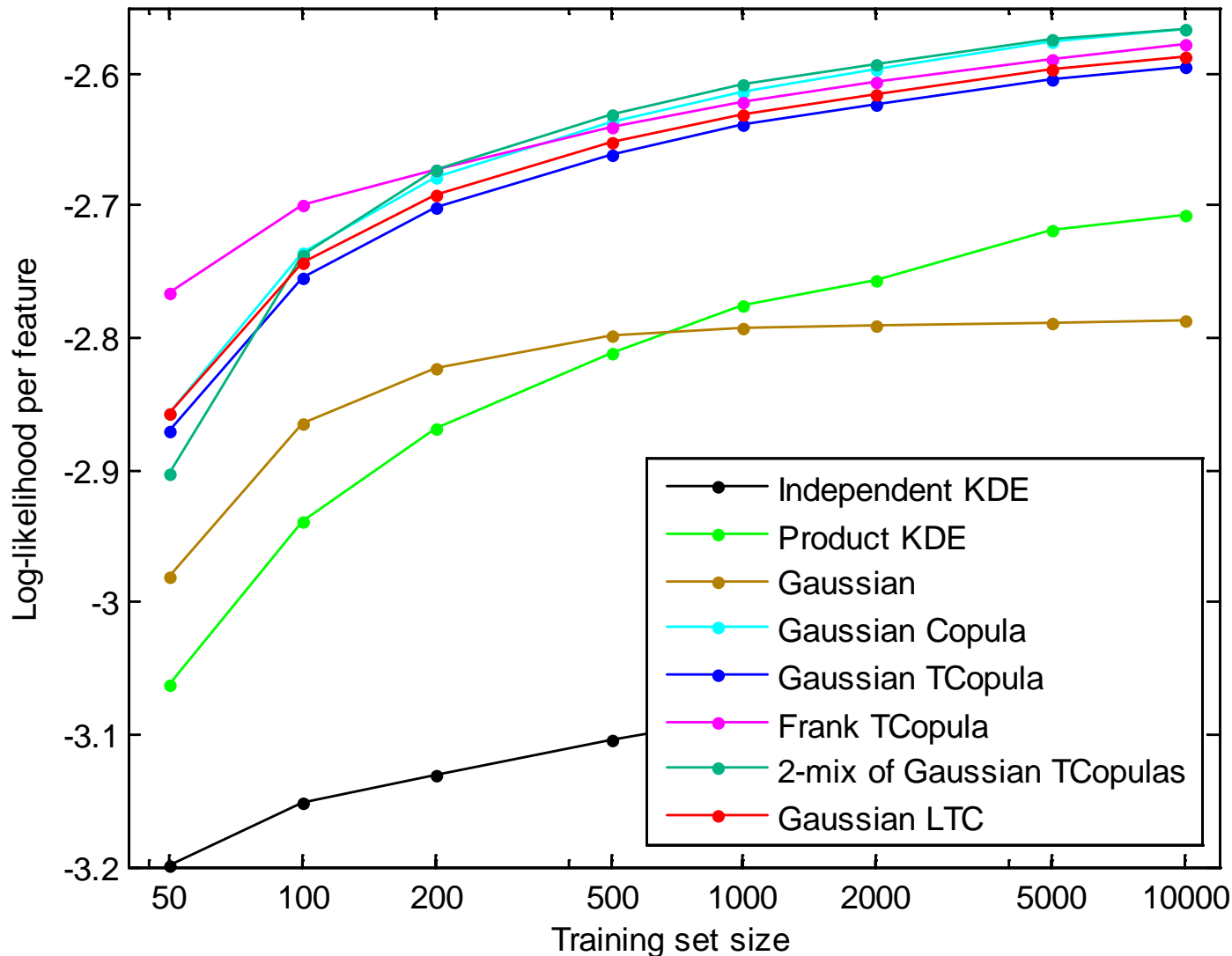
# Illustration for Building of Hierarchy of Rainfall Stations



State of Indiana  
(USA)

Average monthly  
observations for  
15 rainfall stations  
1951-1996 (47  
years)

# Experiments: Log-Likelihood on Test Data



UCI ML  
Repository  
MAGIC data set

12000 10-  
dimensional  
vectors

2000 examples in  
test sets

Average over 10  
partitions

# Summary

- Multivariate distribution = univariate marginals + copula
- New model: tree-structured multivariate distribution with marginally uniform latent variables (latent tree copula, LTC)
  - Sufficient to employ only bivariate copula families!
- Closed form inference for Gaussian copulas (efficient)
- Variational inference for non-Gaussian copulas (slow)
- Parameter estimation using the EM algorithm
- Bottom-up structure learning for bivariate LTCs
- Can be used for parsimonious multivariate density estimation or to structure variables into hierarchies

# Thank you!

Software:

<http://www.stat.purdue.edu/~skirshne/LTC/index.html>

Support:



US National Science Foundation Award AGS-1025430

## Questions?

See me at the poster tonight for more details