

# Column generation for exact Bayesian network learning: Work in Progress

James Cussens, University of York

CoCoMile 2012

# Encoding graphs with binary 'family' variables

- ▶ Suppose there are  $|V|$  fully-observed discrete variables in some dataset. We want to learn a BN (with  $|V|$  vertices) from this data.

# Encoding graphs with binary 'family' variables

- ▶ Suppose there are  $|V|$  fully-observed discrete variables in some dataset. We want to learn a BN (with  $|V|$  vertices) from this data.
- ▶ Can encode any graph by creating a binary variable  $I(W \rightarrow u)$  for each BN variable  $u \in V$  and each candidate parent set  $W$ .
- ▶ Each  $I(W \rightarrow u)$  has a *local score*  $c'(u, W)$ .

# Encoding graphs with binary 'family' variables

- ▶ Suppose there are  $|V|$  fully-observed discrete variables in some dataset. We want to learn a BN (with  $|V|$  vertices) from this data.
- ▶ Can encode any graph by creating a binary variable  $I(W \rightarrow u)$  for each BN variable  $u \in V$  and each candidate parent set  $W$ .
- ▶ Each  $I(W \rightarrow u)$  has a *local score*  $c'(u, W)$ .
- ▶ **Big problem already:** That could be a lot of  $I(W \rightarrow u)$  variables.

# Encoding graphs with binary 'family' variables

- ▶ Suppose there are  $|V|$  fully-observed discrete variables in some dataset. We want to learn a BN (with  $|V|$  vertices) from this data.
- ▶ Can encode any graph by creating a binary variable  $I(W \rightarrow u)$  for each BN variable  $u \in V$  and each candidate parent set  $W$ .
- ▶ Each  $I(W \rightarrow u)$  has a *local score*  $c'(u, W)$ .
- ▶ **Big problem already:** That could be a lot of  $I(W \rightarrow u)$  variables.
- ▶ **What might save us:** Only  $|V|$  of these  $I(W \rightarrow u)$  variables will be non-zero in any solution.

# BN learning as constrained optimisation

Instantiate the  $I(W \rightarrow u)$  to maximise:

$$\sum_{u,W} c'(u, W) I(W \rightarrow u) \quad (1)$$

subject to the  $I(W \rightarrow u)$  representing a DAG.

Set  $c(u, W) = -c'(u, W)$  and consider

Instantiate the  $I(W \rightarrow u)$  to minimise:

$$\sum_{u,W} c(u, W) I(W \rightarrow u)$$

subject to the  $I(W \rightarrow u)$  representing a DAG.

# An integer linear programming approach

$$\forall u \in V : \sum_W I(W \rightarrow u) = 1 \quad (2)$$

$$\text{Where } C \subseteq V : \sum_{u \in C} \sum_{W: W \cap C = \emptyset} I(W \rightarrow u) \geq 1 \quad (3)$$

Cluster constraints (3) added 'on the fly' as cutting planes.  
Introduced in [1].

# Slack variable representation

$$\sum_{u \in C} \sum_{W: W \cap C = \emptyset} I(W \rightarrow u) \geq 1 \quad (4)$$

$$-w_C + \sum_{u \in C} \sum_{W: W \cap C = \emptyset} I(W \rightarrow u) = 1 \quad (5)$$

where  $w_C \geq 0$ . Let  $n$  be the number of  $I(W \rightarrow u)$  variables then can represent the feasible region using:

$$\mathbf{Ax} = \mathbf{b} \quad (6)$$

where  $\mathbf{A}$  is a  $(|V| + |C|) \times (n + |C|)$  matrix.



# Dictionary representation of an empty graph solution

Let  $\zeta$  be the variable for the objective function and  $\bar{\zeta}$  a constant which is the score of the empty graph.

$$\zeta = \bar{\zeta} + \sum_{u, W: W \neq \emptyset} [c(u, W) - c(u, \emptyset)] I(W \rightarrow u) \quad (7)$$

$$I(\emptyset \rightarrow u) = 1 - \sum_{u, W: W \neq \emptyset} I(W \rightarrow u) \quad (8)$$

$$w_C = |C| - 1 - \sum_{u \in C} \sum_{W: W \cap C \neq \emptyset} I(W \rightarrow u) \quad (9)$$

# Basic and non-basic variables

- ▶ **LHS:** The  $I(\emptyset \rightarrow u)$  variables and the slack variables are (currently) *basic* and may be positive.
- ▶ **RHS:** The  $I(W \rightarrow u)$  variables for  $W \neq \emptyset$  are (currently) non-basic and have value zero.

$$\zeta = \bar{\zeta} + \sum_{u, W: W \neq \emptyset} [c(u, W) - c(u, \emptyset)] I(W \rightarrow u) \quad (10)$$

$$I(\emptyset \rightarrow u) = 1 - \sum_{u, W: W \neq \emptyset} I(W \rightarrow u) \quad (11)$$

$$w_C = |C| - 1 - \sum_{u \in C} \sum_{W: W \cap C \neq \emptyset} I(W \rightarrow u) \quad (12)$$

Choose a non-basic variable with *negative reduced cost* to bring into the basis.

# Column generation

- ▶ Column generation = variable generation
- ▶ *It is not necessary to explicitly represent non-basic variables.*
- ▶ Only create them when they are to enter the basis.

# Computing reduced costs

For any basis:

- ▶  $\mathbf{x} = (\mathbf{x}_B, \mathbf{x}_D)$
- ▶  $\mathbf{c} = (\mathbf{c}_B, \mathbf{c}_D)$
- ▶  $\mathbf{B}$  is the (square) submatrix of the original  $\mathbf{A}$  matrix whose columns correspond to  $\mathbf{x}_B$ .  $\mathbf{D}$  is the (non-square) matrix formed from the remaining (non-basic) columns.
- ▶ Compute dual values:  $\lambda^T = \mathbf{c}_B^T \mathbf{B}^{-1}$ .
- ▶ Compute reduced costs:  $\mathbf{r}_D = \mathbf{c}_D - \lambda^T \mathbf{D}$ .

# What we need

- ▶ To compute the reduced cost of a potential new variable we need:
- ▶ its objective coefficient value and
- ▶ its coefficient for each original linear constraint (= row of  $\mathbf{A}$ ).
- ▶ (Note that a row is added to  $\mathbf{A}$  each time a cutting plane is added.)

# An ILP for variable creation

- ▶ A new variable  $I(W \rightarrow u)$  is determined by a choice of the child  $u$  and also the parents  $W$ .
- ▶ Let  $I_{\text{ch}}(u)$  indicate that  $u$  is chosen as the child and let  $I_{\text{pa}}(u)$  represent that  $u$  is chosen as a parent.

$$\sum_{u \in V} I_{\text{ch}}(u) = 1 \quad (13)$$

$$\forall u \in V : I_{\text{ch}}(u) + I_{\text{pa}}(u) \leq 1 \quad (14)$$

Create a variable  $x_C$  for each  $C \in \mathcal{C}$  indicating whether the new variable is involved in the constraint for  $C$ . We have:

$$x_C \geq \sum_{u \in C} I_{\text{ch}}(u) - \sum_{u \in C} I_{\text{pa}}(u) \quad (15)$$

$$\sum_{u \in C} I_{\text{ch}}(u) \geq x_C \quad (16)$$

$$1 - \sum_{u \in C} I_{\text{pa}}(u) \geq x_C \quad (17)$$

# Computing the reduced cost of the new variable

- ▶ Let  $\lambda_C$  be the dual value corresponding to the constraint for cluster  $C$ .
- ▶ Let  $\lambda_u$  be the dual value for the convexity constraint (2) for variable  $u$ .

The reduced cost for a new variable  $I(W \rightarrow u)$  is then:

$$c(u, W) - \sum_u \lambda_u I_{\text{ch}}(u) - \sum_{C \in \mathcal{C}} \lambda_C x_C \quad (18)$$

Hmmm, how to get  $c(u, W)$ ?



# Proposed strategy

- ▶ View  $c(u, W)$  as a real-valued variable.

# Proposed strategy

- ▶ View  $c(u, W)$  as a real-valued variable.
- ▶ Or rather a variable which is a lower bound on this value.
- ▶ This will lead to over-optimistic generation of new variables.
- ▶ Once a new variable is proposed probably worth the effort to consult the data to compute  $c(u, W)$  exactly.



Tommi Jaakkola, David Sontag, Amir Globerson, and Marina Meila.

Learning Bayesian network structure using LP relaxations.

In *Proceedings of 13th International Conference on Artificial Intelligence and Statistics (AISTATS 2010)*, volume 9, pages 358–365, 2010.

Journal of Machine Learning Research Workshop and Conference Proceedings.