# Graphical Models

# Stochastic Computation for Large-Scale Graphical Models
—
## Exploratory Analysis of Gene Expression Data

Mike West, Duke University

Adrian Dobra          Beatrix Jones
Carlos Carvalho        Chris Hans

Institute of Statistics and Decision Sciences
&
Computational and Applied Genomics Program

Quanli Wang      Joseph Nevins      Guang Yao

# High-Dimensional Graphical Models

- Exploratory data analysis & visualisation
- Observational data:  Varying contexts
- Uses in prediction


- Sparsity

# Exploring High-Dimensional Observational Data

Exploring & summarising associations:

- structure in p($x_1,...,x_p$)

- visualisation, clues

Predictive models:

- regressions and classification models

- compositional regressions for retrospective data

$$p(y|x) \propto p(x_1|y)p(x_2|y, x_1)p(x_3|y, x_1, x_2)$$

# Graphical Models: p(x)
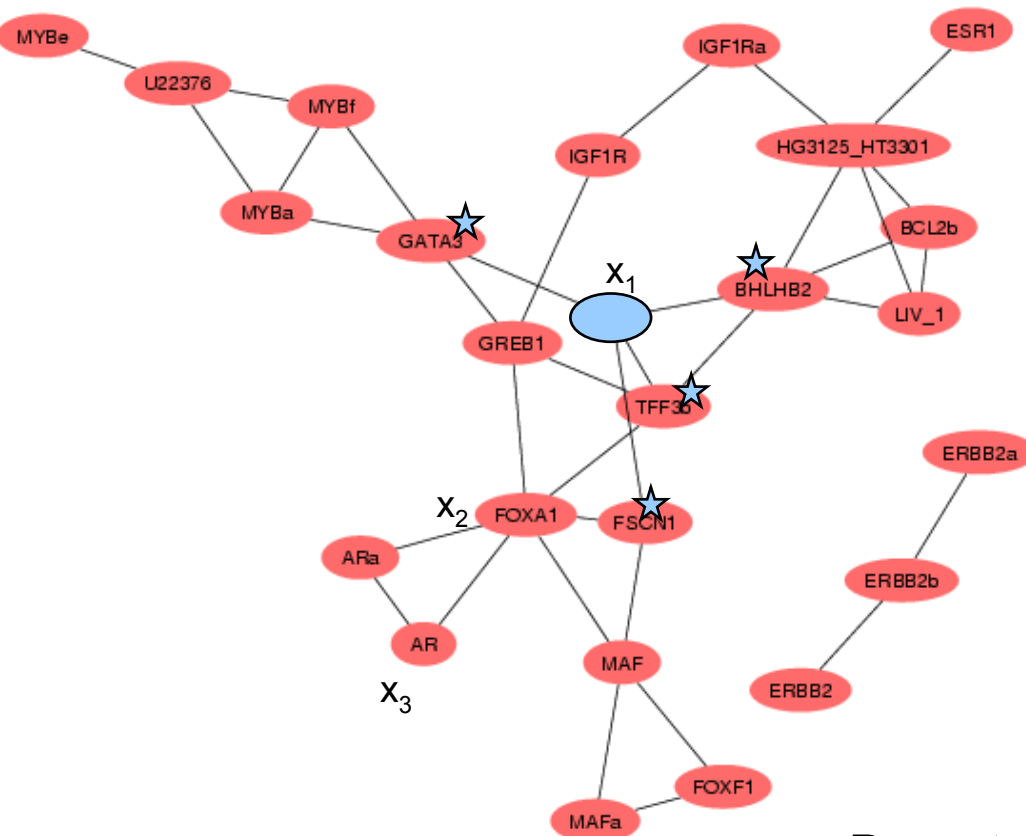
Nodes=p variables
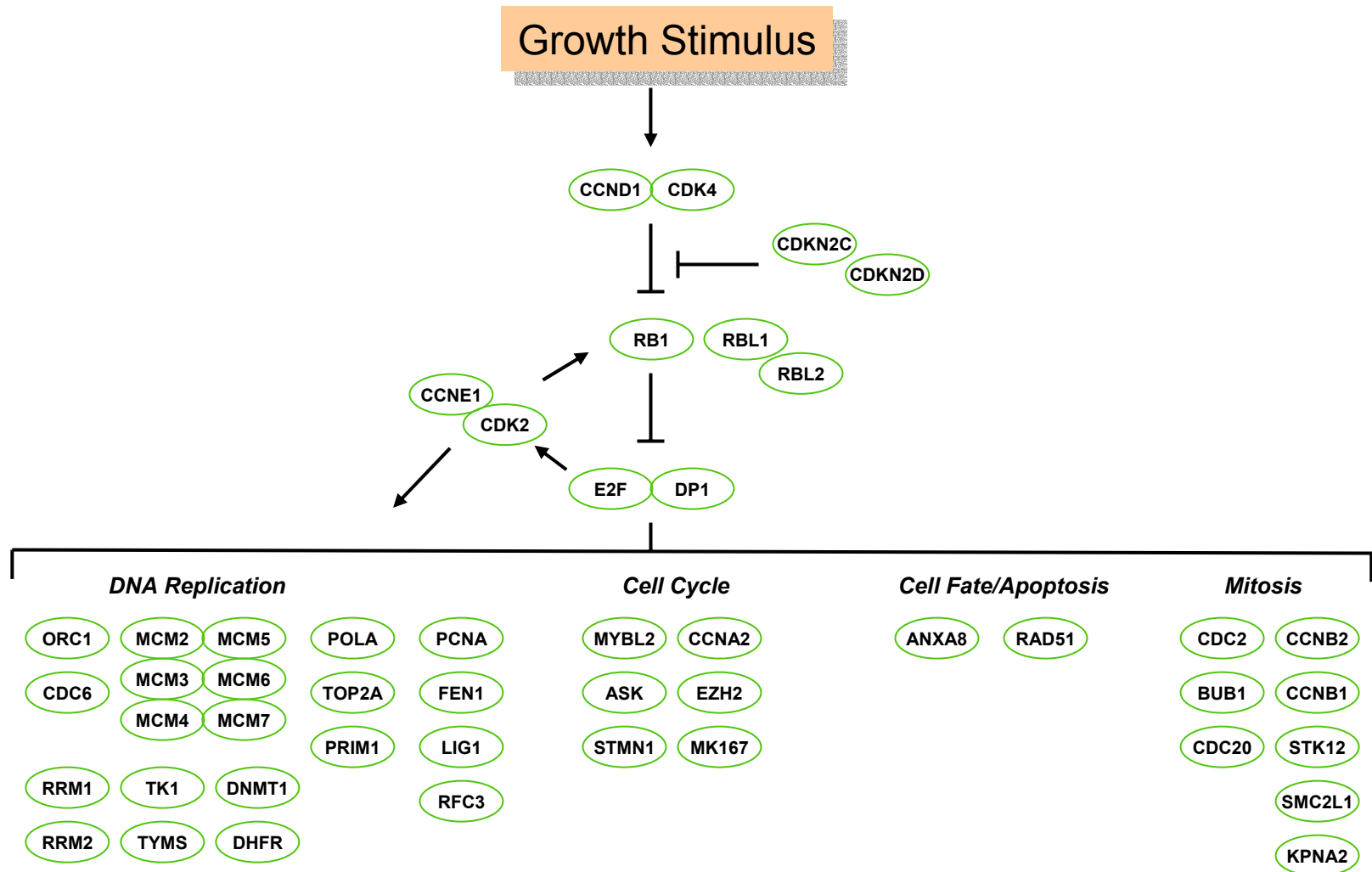
Edges: ne(i)=neighbours

Conditional (in)dependencies

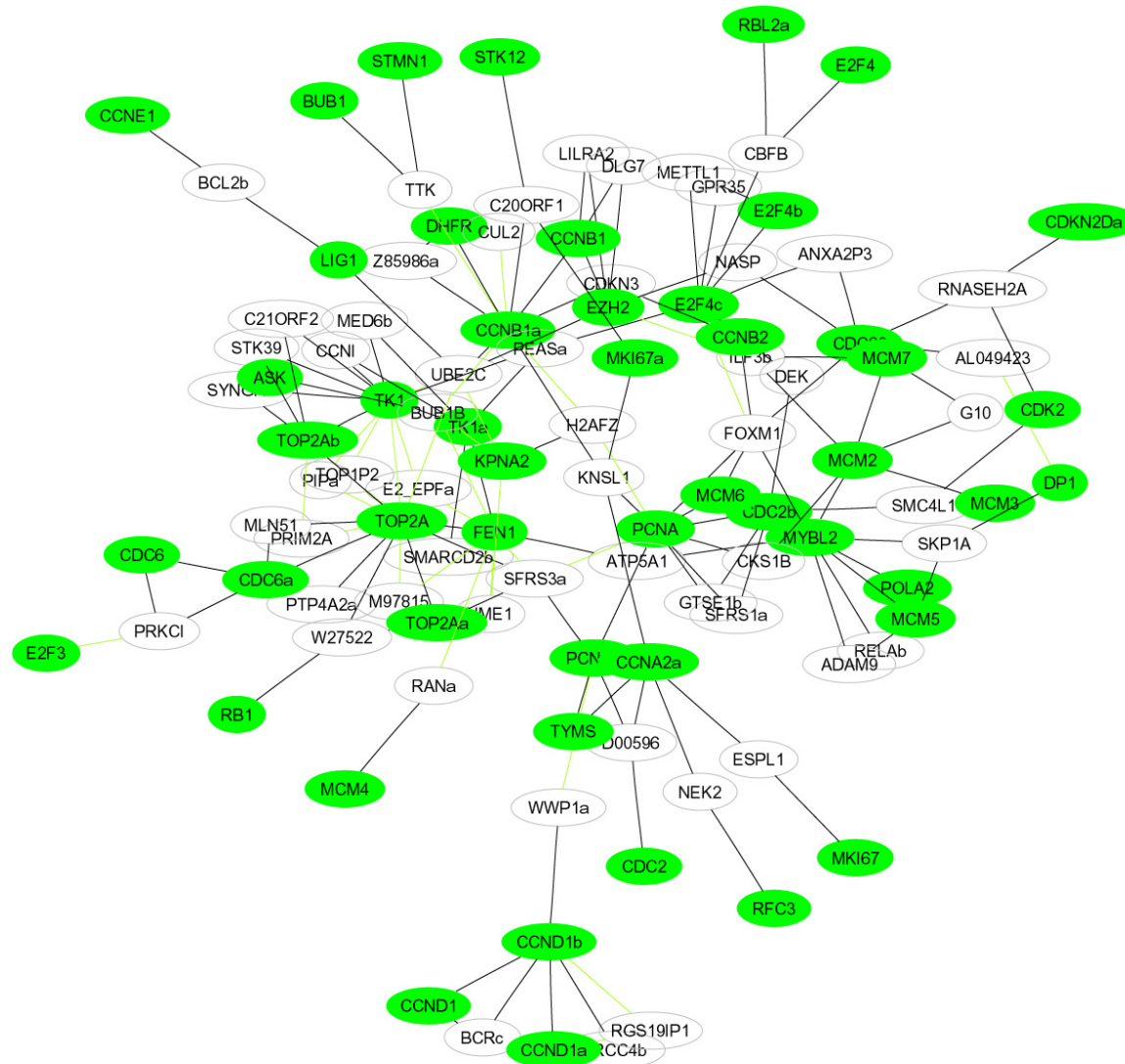$$p(x_i|x_{-i}) = p(x_i|x_{ne(i)})$$

(prediction: $y=x_0$)

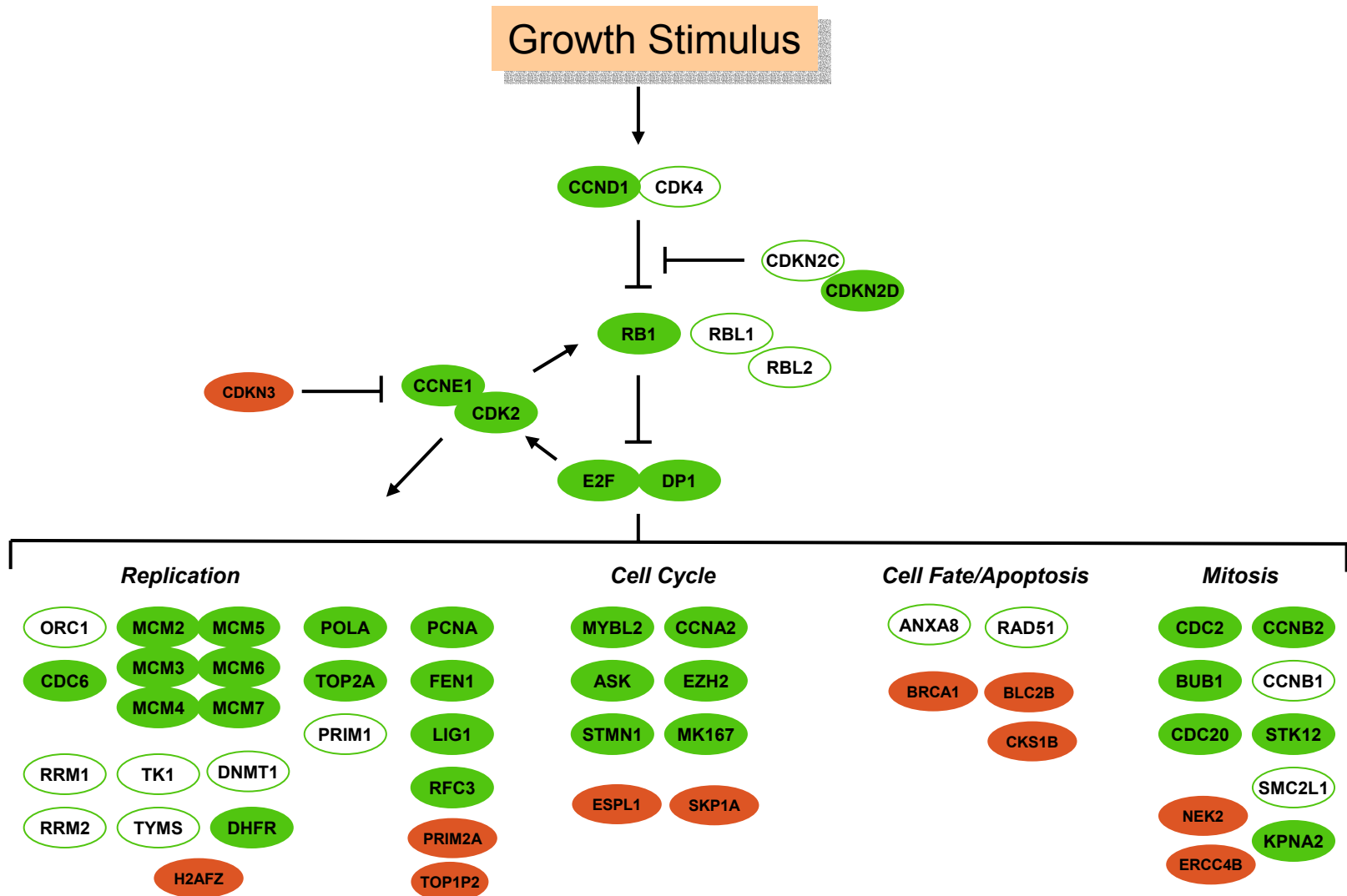Breast cancer - ER subgraph:  p=12558

# The Rb-E2F Pathway

# A Subgraph of an Rb-E2F Association Graph
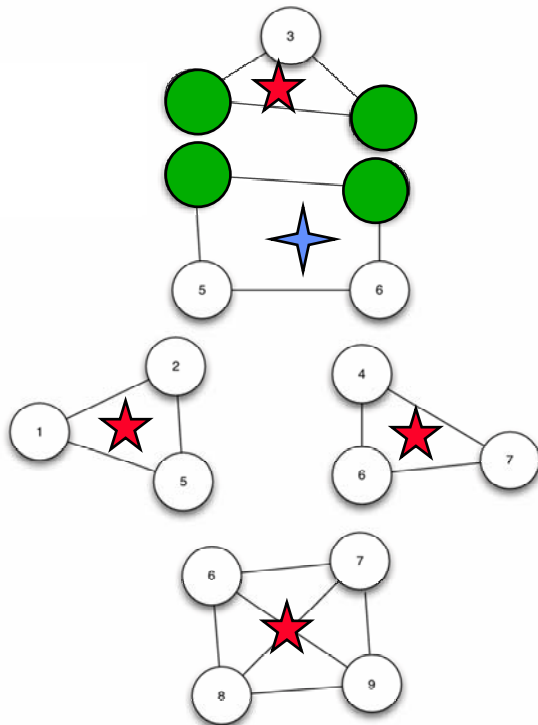## - Pathway Exploration -

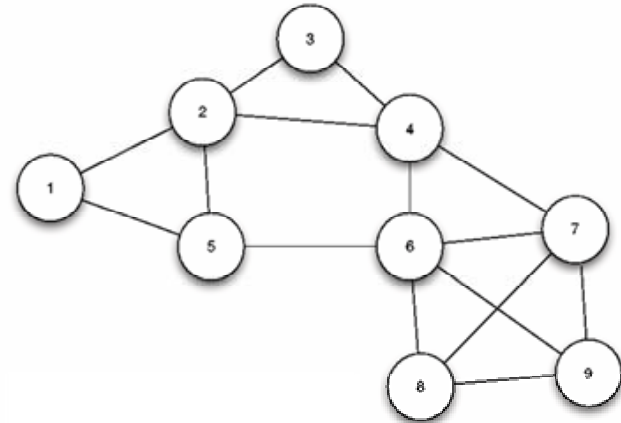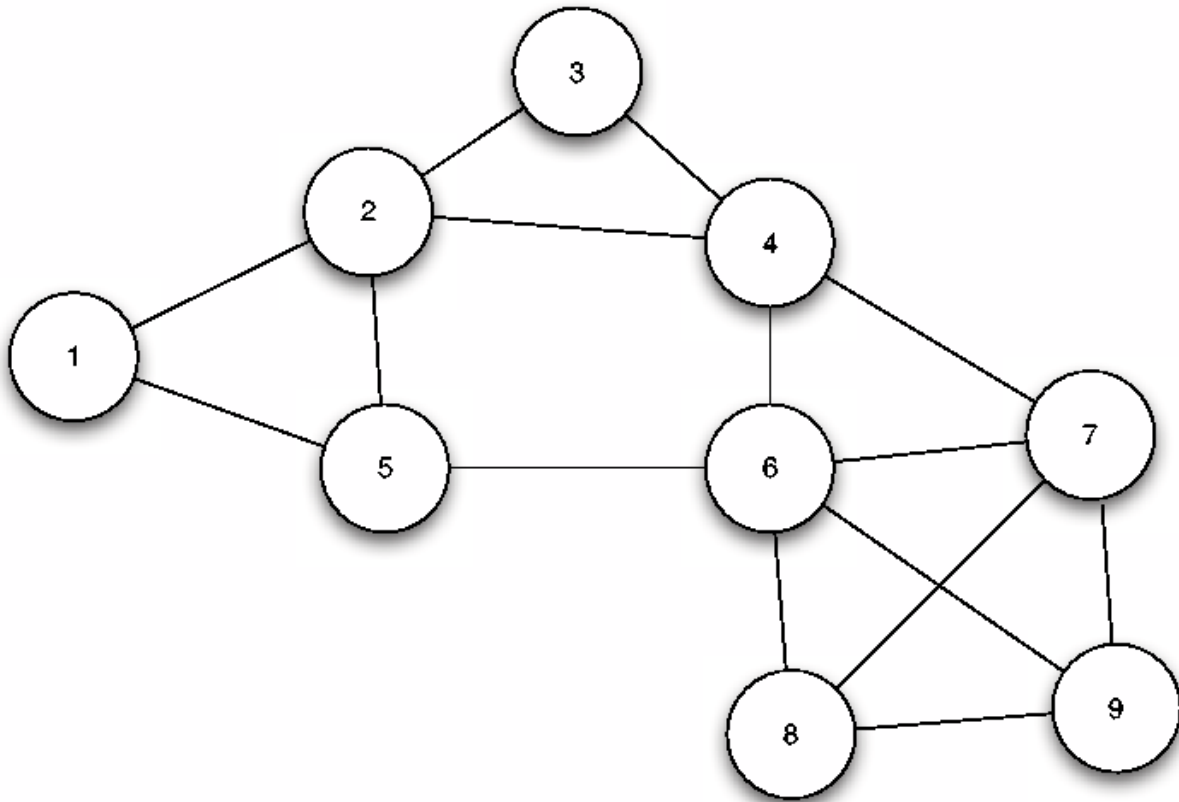# Gene Set Enrichment by Neighbour & Path Exploration

# Graph Decompositions:
# The Key to Dealing with Dimension

Dimension 'reduces' via graph decompositions:
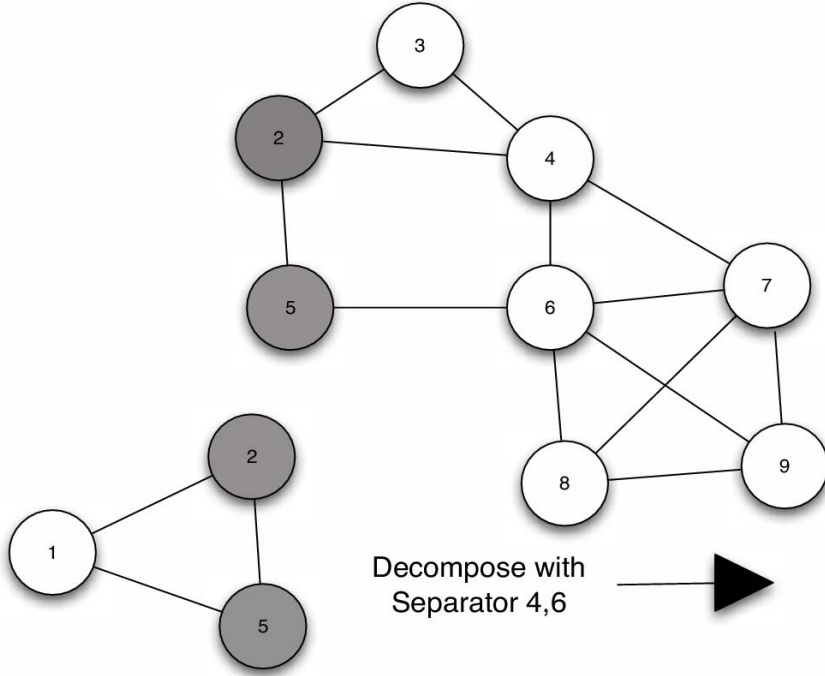
Intersecting or disconnected subgraphs



S :  separator ...
  - complete subgraph that "separates" PCs

PC :  prime component ...
  - either a complete subgraph,
  - or cannot be separated

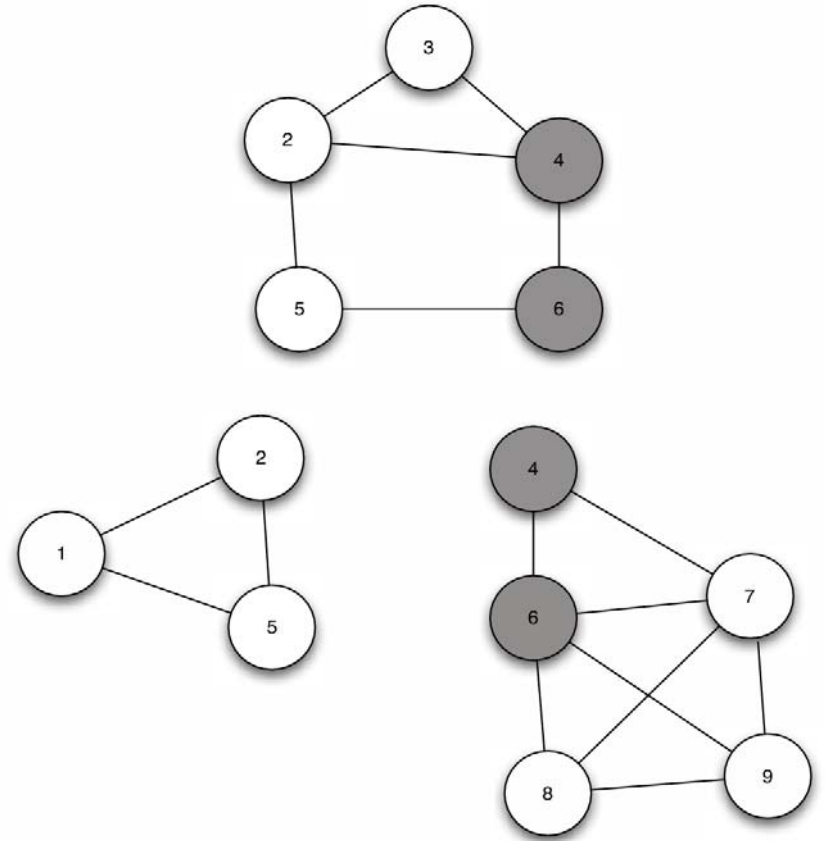Decompose with
Separator 2,5
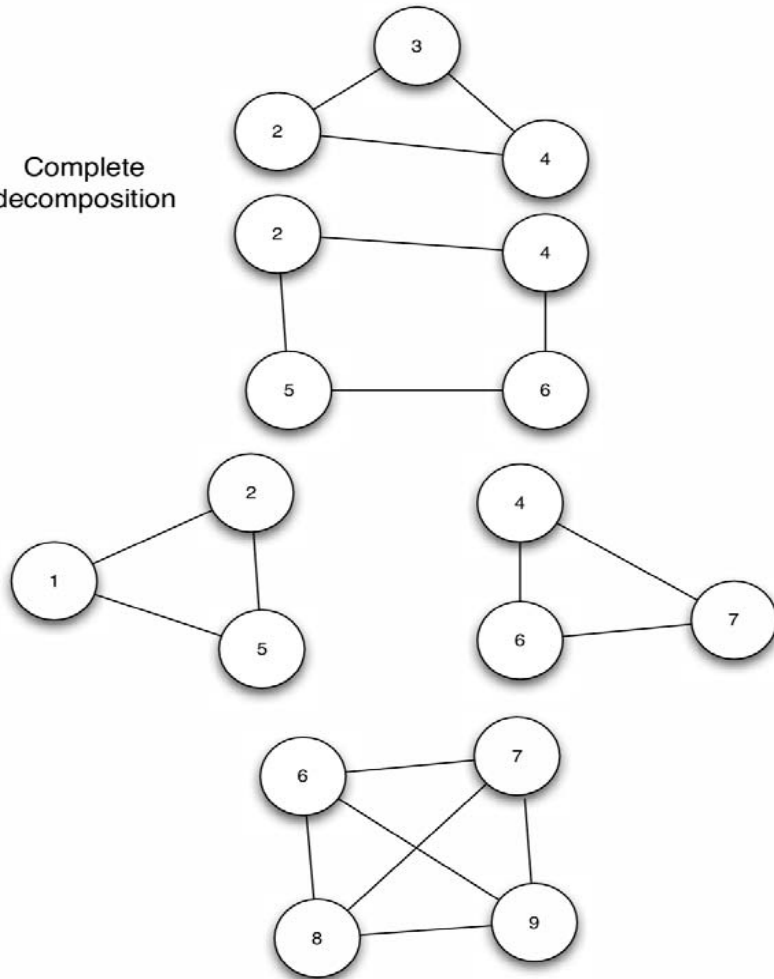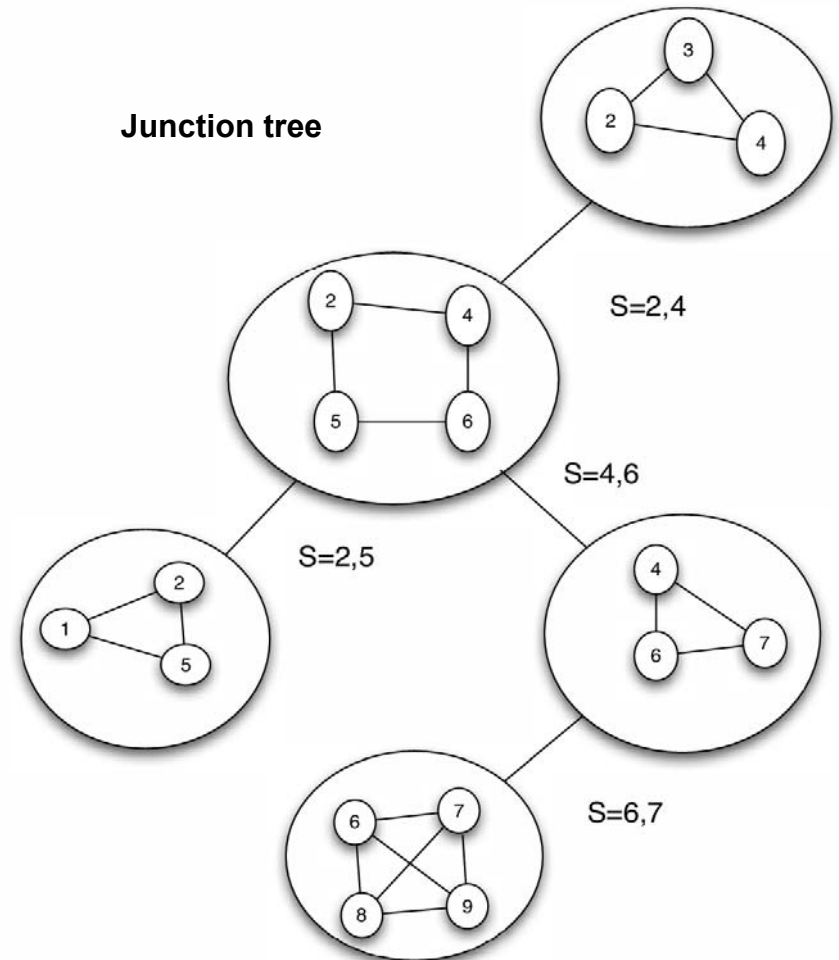
Decompose with
Separator 4,6

Complete decomposition

**Junction tree**

S=2,4

S=4,6

S=2,5

S=6,7

Decomposable graphs:
PC=clique (maximal complete subgraph)

Non-decomposable:
PC=anything goes (e.g., big chains, cycles, …

# Gaussian Graphical Models

Precision $\mathbf{K} = \mathbf{\Sigma}^{-1}$ on a graph $G$

$$\mathrm{E}(x_i | x_{-i}) = \sum_j \ (-\mathrm{K}_{ij} / \mathrm{K}_{ii}) x_j$$

Edges ~ $|\mathrm{K}_{ij}| > 0$

Priors:
- over graphs G ...
- then non-zero elements of **K** on G
- **sparsity inducing**

Inference:

"covariance selection"

Computation: finding graphs: MCMC, stochastic search, annealing
- Dimension!
- Distributed/cluster computation

# Graph, Model and Prior Decompositions: Dealing with Dimension

Dimension 'reduces' via graph decompositions:

$$\mathrm{p}(x\,|\,G,K) = \prod_{PC} \mathrm{p}(x_{PC}\,|\,K_{PC}) / \prod_{S} \mathrm{p}(x_S\,|\,K_S)$$

- PC : prime components
- S : separators – complete subgraphs separating PCs

Large, realistic sparse graphs:  massive decomposition

Decomposable graphs:  PC=clique

Non-decomposable:  Anything goes (e.g., big chains)

General graphical model

# Priors over Covariance/Precision Matrix

Likelihood:

$$\mathrm{p}(x \mid G, K) = \prod_{PC} \mathrm{p}(x_{PC} \mid K_{PC}) / \prod_{S} \mathrm{p}(x_S \mid K_S)$$

Conjugate prior:           Hyper Wishart (Roverato 2002)

$$\mathrm{p}(K \mid G) \propto \prod_{PC} \mathrm{p}(K_{PC}) / \prod_{S} \mathrm{p}(K_S)$$

"Local Hyper-Wishart Prior"

Wishart on cliques

Constrained Wishart on non-complete prime components

# Graph Decomposition
## - Stochastic Computational Strategies -

Wander around G space:  evaluate p(x|G),  then p(G|x)
(marginalised over parameters **K**)

$$\mathrm{p}(x\,|\,G) = \prod_{PC} \mathrm{m}_{PC} \Big/ \prod_{S} \mathrm{m}_{S}$$

$\mathrm{m}_S$:  Wishart normalising constants

$\mathrm{m}_{PC}$:

   — PC complete: analytic (inverse Wishart normalizing constants)
     (e.g., decomposable: Giudici & Green 99; Wong & Carter 02)

   – PC incomplete: (hard) integral over constrained Wishart
      Monte Carlo (Atay-Kayis & Massam 02 … 04)

(Roverato 2002, Scand J Stat)

# Priors on Graphs

Uniform prior ?

Unrestricted graphs:

$E(\#edges)=p(p-1)/4$

Decomposable cases

# Sparsity Priors on Graphs

Edge inclusion prob.

$\beta = 2/(p-1)$

Unrestricted graphs:
mode of p edges

Decomposable cases

# Local Computations on Graph Space

Current graph $G$ : $p(x|G)$
$G$ 1-edge neighbours $G*$ : $p(x|G*)$

Major efficiencies if decomposable:
- change 2 cliques at most
- efficiently check decomposability

(Giudici & Green 99)

Non-decomposable? Anything can happen

Practical relevance of decomposability?

Most published examples: p=4-12

Challenges: efficiency, dimension

# MCMC on Graph Space

Gibbs:  random 1-edge moves

      Decomposable models: Giudici & Green 99, Wong & Carter 02

      Conditional posterior for edge in/out

Metropolis Hasting

      random choice of add/delete,

      random 1-edge move  (Jones et al, Duke team, 03)

Challenges:  computations in non-decomposable cases

> Monte Carlo (Atay-Kayis & Massam 02/04) hugely challenging

Issues: Horrible/impossible as p increases

# Annealed Stochastic Search

Wander around G space: find "high probability" graphs

- 1-edge different graphs ("neighbours")
  - store top h
- Select "next" :    $p(G|x)^a$
- repeat

Parallelisable

Multiple "interesting" regions of model space

# 12 Node Example



- Non-decomposable

  A non-complete prime component

- n=250

- 10,000 SS steps

- 10,000x66 MH steps

- a=1

# 12 Node Example

| Method | Time (s) | Top log posterior | Evaluations to top grph | Time to top graph |
|---|---|---|---|---|
| MH decom-posable | 36 | -2591.18 | 912 | 1 |
| SS decom-posable | 183 | -2591.18 | 792 | 2 |
| MH un-restricted | 15,220 | -2590.94 | 415 | 2 |
| SS un-restricted | 2773 | -2590.94 | 13266 | 5 |

True graph has lower marginal likelihood than either of these

# 12 Node Example "Top Graphs"



Decomposable

Unrestricted

# 15 Node Example



- Decomposable
- n=250
- 10,000 SS steps
- 10,000x105 MH steps
- a=1

# 15 Node Example

| Method | Time (s) | Top log posterior | Evaluations to top grph | Time to top graph |
|---|---|---|---|---|
| MH decom-posable | 93 | 15633.76 | 349,484 | 36 |
| SS decom-posable | 234 | 15633.76 | 33,495 | 9 |
| MH un-restricted | 513,077 | 15633.83 | 666, 425 | 309, 222 |
| SS un-restricted | 5930 | 15636.33 | 82845 | 112 |

True graph has lower marginal likelihood than any of these

# 15 Node Example "Top Graphs"



Decomposable

Unrestricted

# 150 Node Example

- gene expression – breast cancer : ER

- (O)Estrogen receptor pathway

- n=49

- 25,000 SS steps

- 25,000x11,175 MH steps

- annealing?

- unrestricted: accuracy of MC evaluations

# 150 Node Example

| Method | Time (hrs) | Top log posterior | Evaluations to top grph | Time to top graph |
|---|---|---|---|---|
| MH decomposable | 18.02 | -9417.97 | 100,467k | 6.51 |
| SS decomposable | 0.03 | -9260.35 | 1699k | 0.03 |
| ⭐SS unrestricted | **6.29** | **-9227.68** | **44.7k** | **3.39** |

⭐ **Starting from the 'best' decomposable  graph, and is a (local) mode**

# Scaling-Up?

# Scaling Up: p>15
# High-Dimensional Sparse Models

- Build directed graphical models
- Induce conditional independence graph
- Priors on directed (acyclic) graphs

Compositional Networks: Parallel Regressions

$$\mathrm{p}(x \mid K) = \prod_{i=1}^{p} \mathrm{p}(x_i \mid x_{cne(i)}, \theta_i), \qquad cne(i) \subseteq \{(i+1) : p\}$$

Regression parameters functions of **K**
Triangular array: Order matters

# EGFR

## Brain cancer gene expression
## Duke Keck Center for Neurooncogenomics



DAG to G: moralise

# Compositional Nets: Priors & Computation

- Normal/inverse gamma priors ~ Wisharts

- Sparsity: Regression variable inclusion/exclusion priors

- Stochastic model search via sets of regressions:
    - parallel MCMC for regression search
    - select from multiple 'neighbouring' graphs
    - 'local' exploration near 'good' graphs
    - 'local' MCMC ?

- Evaluate relative posterior prob on directed graph

# Shotgun Stochastic Search on DAGs

- Current DAG: ordering, edges, posterior probability
- Local shuffle of ordering: switch two neighbours
- Local recomputation of regression search, sampling
- Compute posterior prob on new DAG

- Optimisation steps to initialise and update orderings

  Concept: "explanatory" genes low in ordering

- Gibbs sampling for full conditional regressions
  - to initialise
  - and reduce candidate predictors for each x

- Shotgun search for regressions (forthcoming)

Code:  C++/MPI Beowulf cluster implementation
HdBCS (Dobra, Duke)

# Breast Cancer Gene Expression

- n=158 tumour samples

- p=12558

- 48x2cpu cluster

- Summarisation of high prob graphs?

Code:  C++/MPI Beowulf cluster implementation
HdBCS: Adrian Dobra, Duke

# Beowulf cluster is needed



**Zillions**

# p=12558
# Adjacency Matrix of An "Interesting" Graph

p=12558: ER Genes Subgraph

# p=12558: Breast Cancer Gene Expression

# Neighbour & Path Exploration

# Exploring Graphs: Gene Discovery & Annotation



http://dig.cgt.duke.edu

Large-scale graphical model search and evaluation

Inference on large, sparse inverse covariance matrix

(Dobra et al, JMVA 2004)          graphexplore.cagp.duke.edu

## Some key refs:

- Giudici and Green 1999, Biometrika

  MCMC in decomposable graphical models

- Roverato 2002, Scand J Stat

  HIW priors on graphical models

- Wong and Carter 2002, tech report, Hong Kong Univ

  covariance selection

- Atay-Kayis and Massam 2002, tech report, York Univ

  Monte Carlo evaluation of marginal likelihoods

- Jones, West et al 2003, tech report, SAMSI & Duke Univ

  stochastic computation and search

- Dobra, West et al 2004, J Multivariate Anal

  initial DAG based approach, stochastic search & gene expression studies

# Graphical Models – Some Current Foci

Weighting paths between two nodes
(Beatrix Jones & MW 2004)

SSS: "Shotgun" Stochastic Search – annealing & rapid
search over Zillions of models in regression
(Chris Hans, Adrian Dobra & MW 2004)

Latent graphs: measurement error

Data exploration (Duke team, DIG paper)
- One pathway, multiple data sets
- Transcription factor binding sites
- Graphs from literature data

Software and visualization tools (see links)

Adrian Dobra                              Beatrix Jones
Carlos Carvalho                          Chris Hans

Institute of Statistics and Decision Sciences
&
Computational and Applied Genomics Program

Quanli Wang      Joseph Nevins      Guang Yao

Duke University