

GRAPHICAL MODEL STRUCTURE LEARNING

- WE WANT $P(G|D)$, G IS GRAPH STRUCTURE AS $V \times V$ ADJ MATRIX. PROBLEM EXPONENTIAL IN NUMBER OF NODES $O(2^{V(V-1)/2})$
FULL POSTERIOR IS PROHIBITIVELY LARGE.

FOR KNOWLEDGE DISCOVERY

JUST GRAPH TOPOLOGY. POSTEDGE EDGE MAX/MIN $P(G_{ST}=1/D)$. EDGE 'THICKNESS' IS CONFIDENCE VALUE

FOR DENSITY ESTIMATION

- MAP & MAPIT $\hat{G} \in \text{argmax}_G \tau(G|D)$. HEURISTICS BECAUSE EXPTIME. FOR TREES IS ON TO FIND GLOBAL OPTIMUM EXACTLY.

- CONSIDER WHETHER A LATENT VARIABLE MODEL WOULD BE APPROPRIATE BECAUSE IT IS A LOT EASIER

FOR KNOWLEDGE DISCOVERY

QUICK AND DIRTY. NO JOINT LOSS \rightarrow NO REGULARIZATION. NO CAP GUARANTEES OF FIT. GOOD FOR VISUALIZATION.

RELEVANCE NETWORKS: PLOTS PAIRWISE MUTUAL INFORMATION. WE DRAW EDGE IF $I(X_i, X_j) > \text{THRESHOLD}$. IF GAUSSIAN, THIS IS Σ . CONTINUOUS GRAPH.
 → ALSO PLOT EDGES WHERE CONDITIONAL INTERDEPENDENCE.

DEPENDENCY NETWORKS! FITS D. SPARSE FULL CONDITIONALS $P(X_i | X_{-i})$. CHOSEN VARS ARE INPUTS PLOTTED FOR MODE. CPD FITTED USING ANY REG/CLASS METHOD.
→ CAN BE USED FOR INFERENCE VIA Gibbs. OK-ish IF NOT MUCH MISSING DATA. USEFUL FOR DATA IMPUTATION / INITIALIZATION

LEARNING TREES

BECAUSE MUCH EASIER THAN ANY GRAPH. AND SUPPORT EXACT INFERENCE

- SAME NUM OF PARAMS FOR DIRECTED OR UNDIRECTED TREES, UNIL + FOR STRUCTURE LEARNING. DIRECTED + FOR PARAMETER LEARNING

TREE LOG-LIKELIHOOD:

$$\frac{\log p(\theta|\tau)}{N} = \sum_{t \in V} \sum_u f_{\text{emp}}(x_t = u) \log f_{\text{emp}}(x_t = u) + \sum_{s, \tau} 1(x_s, x_\tau) \hat{\theta}_{s, \tau}$$

DOES NOT DEPEND ON TOPOLOGY

→ IGNORED WHEN LEARNING STRUCTURE

MUTUAL INFO GIVEN
EMPIRICAL DISTR.

MAX LL: MAXIMUM WEIGHT SPANNING TREE, EDGE WEIGHTS ARE MUTUAL INFO
→ CHOI-LIU ALGORITHM

- USE ALGO TO FIND MAX SPANNING TREE. $O(E \lg V)$

- NO OVERFIT \rightarrow ALL TREES HAVE SAME NUM OF PARAMS

IF WE WANT FOREST BECAUSE BP INFERENCE IS MUCH FASTER. NO CAN USE MLE (IT DOESN'T OMIT EDGES). USE **MARGINAL LIKELIHOOD** OR OTHER PENALIZED MEASURE (BIC) $\left\{ \begin{aligned} \log p(D|T) &= \sum \text{SCORE}(N_i, p_i(N_i)) - N \text{ ME COUNTS.} \\ \log p(D|T) &= \sum w_{pq}(T)_i + \sum_j \text{SCORE}(t|U) \end{aligned} \right.$

- MOST PROBABLE TREE IS MAXIMAL BRANCHING OF MATCHING WEIGHTED DIRECTED GRAPH
- **MODIFIED MST**

MIXTURES OF TREES

→ **FIT USING EM** : E COMPUTES RESPONSIBILITIES FOR EACH TREE

→ FIT USING EM :

- E COMPUTES RESPONSIBILITIES FOR EACH CLUSTER FOR EACH DATAPoint.
- M USES WEIGHTED CHAIN-LINK

LEARNING DAG

AND BAYESIAN STRUCTURAL LEARNING: NO HIDDEN VARS.

MARNOV EQUIVALENCY: IF THEY ENCODE SAME SET OF CI ASSUMPTIONS

- WHEN LEARNING DAG STRUCTURE \rightarrow UP TO MARKOV EQUIVALENCY

ESSENTIAL GRAPH F PATTERN

DAG SOME EDGES ARE REVERSIBLE | UNDIRECTED, OTHERS ARE DIRECTED, (COMPELLED)
SAME PATTERN, SAME V-STRUCTURES

EXACT INFERENCE

• LIMITATION $P(D|G, \theta) = \prod_{i=1}^V \prod_{j=1}^C \prod_{k=1}^H \theta_{TCk}^{N_{TCk}}$ IN STATE OF NODE, C STATE OF PARENTS

- ML WILL YIELD FULLY CONNECTED GRAPHS (B/C MAX PPMs) \rightarrow MAX MARGINAL LIKELIHOOD $P(\text{Dig})$, WE NEED POCS ON PPMs

$$\bullet P(D|G) = \prod_{t=1}^T \text{score}(N_t, p_t(r))$$

$$\left\{ \begin{aligned} \text{score}(N_t, p_t(r)) &= \frac{C_t}{\prod_{c=1}^C} \frac{B(N_{tc} + \alpha_{tc})}{B(\alpha_{tc})} \end{aligned} \right.$$

ASSUME: GLOBAL PRIOR PARAM IMPERFORMANCE

$$P(\theta) = \prod_{c=1}^C P(\theta_c)$$

FROM FOR EACH ROW
MUST BE DIRECTLY

$$P(\theta_{Tc}) = \text{Dir}(\theta_{Tc} | \alpha_{Tc})$$

$$N_{TC} = \sum N_{TCU}, \quad \alpha_{TC} = \sum \alpha_{TCU}; \quad N_{T, PA(T)} \text{ is VECTOR OF COUNTS (SUFF. STATS)}$$

- MARGINAL WELFARE DECOMPOSES / FACTORIZES OVER GRANT

HOW TO SET α

WHY NOT JEFFREYS? \rightarrow VIOLATES LIKELIHOOD EQUIVALENCE (SAME MARG. EQUIVALENCE \rightarrow SAME MARG. LIKELIHOOD), ONLY DIRICHLET WORKS

- **BOE**: $\alpha_{TCH} = \alpha P_0(x_T = h, x_{PA(T)} = c)$, $\alpha > 0$. P_0 IS PRIOR JOINT
- **BOEU**: N HAS SAME PARENTS IN $G_1, G_2 \iff P(\theta_T | G_1) = P(\theta_T | G_2)$
 P_0 ASSUMED UNIFORM $\alpha_{TCH} = \frac{\alpha}{N_T(T)}$

OBS! DAGS W/O HIDDEN MAY STILL HAVE LATENT HIDDEN IN DATA, CAREFUL WHEN INTERPRETING CAUSALLY

- **M2 ALGORITHM** WHEN NODES ARE TOTALLY ORDERED, ENUMERATE OVER ALL POSSIBLE ANCESTORS SUBSET AND COMPUTE MARGINAL LIKELIHOODS.
- **GAUSSIAN CPD** (VS TABULAR) \rightarrow **CONDITIONAL GAUSSIAN DAG** CAN MIX GAUSSIAN + DISCRETE NODES. BIC APPROX OF MARGINAL LIKELIHOOD.

$$P(D|G) = \sum_{\tau} \log P(D_{\tau} | \hat{\theta}_{\tau}) - \frac{N_T(\tau)}{2} \log N$$

LARGE SCALE APPROXIMATIONS

- **HYPOTHESIS SPACE IS HUGE** MOTHEAFUENING LARGE! $f(D) = \sum_{i=1}^D (-1)^{i+1} \binom{D}{i} 2^{i(D-i)} f(D-i)$ WHAT DO?
- **MAP** \rightarrow DYNAMIC PROGRAMMING WORKS UP TO 16 NODES. THEN GREEDY HILL-CLIMBING. CAN INIT WITH BEST RANDOM TREE (FOUND ANALYTICALLY). RANDOM RESTARTS, ETC.
- **OTHER STUFF** \rightarrow IE KNOWLEDGE DISCOVERY. SAMPLE FROM POSTERIOR AND COMPUTE EMPIRICAL COUNTS. USE MH SAMPLING, COLLAPSED.

DAG WITH LATENT VARS

- **MARGINAL LIKELIHOOD**: $P(D|G) = \sum_{\theta} \int P(D, h | \theta, G) P(\theta | G) d\theta$ INTRACTABLE, USE DETERMINISTIC APPROXIMATIONS
- **BIC**: $BIC(G) = \log P(D | \hat{\theta}, G) - \frac{\log N}{2} \dim(G)$ \rightarrow EM
- **CHEBSEMAN-SZURZ**: COMPUTE MAP PARAMS $\hat{\theta}$, $\bar{D} = D(\hat{\theta})$, $P(D|G) \approx P(\bar{D}|G)$, DOES NOT SUM FOR h . $\log P(D|G) \approx \log P(\bar{D}|G) + \log P(D|\hat{\theta}, G) - \log P(\bar{D}|\hat{\theta}, G)$
- **VARIATIONAL BAYES EM**: $P(\theta, z_{1:N} | D) = Q(\theta) \prod Q(z_i)$

PLUG-IN
INFERENCE
PRIOR LIKELIHOOD
- **STRUCTURAL EM**: FILL-IN DATA ONCE, USE IT TO EVALUATE NEIGHBOR SCORES. GOOD APPROX FOR DIFFERENCE IN MARGINAL LIKELIHOOD BETWEEN MODELS.
 \rightarrow OR FOR FINDING NEIGHBORS

HOW TO DISCOVER HIDDEN VARS? LOOK FOR SIGNS IN THE DATA, IE CLUSTERS. INTRODUCE HIDDEN VARS. SEE HOW IT PERFORMS. HIDDEN MAY HAVE HIERARCHY \rightarrow AUTOCORRELATION. DATA MAY BE CONSTRAINED AT USAGES OR NOT. WORKS EXCELLENTLY IF TRUE GRAPH IS A TREE.

STRUCTURAL EQUATION MODELS (SEM): IS DIRECTED MIXED GRAPH WHERE CPD ARE GAUSSIAN. PATH DIAGRAMS.

$$X = WX + \mu + \epsilon, \quad P(x) = N(\mu, \Sigma), \quad \Sigma = (I - W)^{-1} \Psi (I - W)^{-T}$$

W LOWER TRIANGULAR \rightarrow ACYCLIC. Ψ NOT DIAGONAL \rightarrow DIRECTED EDGES, CORRELATIONS
 RELATED TO FACTOR ANALYSIS, WHERE NOISE HAS FULL COVARIANCE MATRIX

CAUSAL DAG

STRONGER THAN ASSOCIATIVE CLAIMS. $A \rightarrow B$ IS NOW A DIRECTLY CAUSES B. CAUSAL MARKOV ASSUMPTION. NO CONFOUNDERS

- ASSIGNMENT \neq OBSERVATION. DO NOTATION.
- GRAPH SURGERY: JOINT IS USUAL DAG, CUT ARCS COMING INTO NODES INTERVENED
- GRAPH SURGERY, THEN CAN COMPUTE $P(X_i | DO(X_j))$
- DIFFERENT CAUSAL ASSUMPTIONS \rightarrow DIFFERENT CAUSAL CONCLUSIONS. (SIMPSON'S PARADOX)

LEARN FROM OBSERVATIONAL DATA: LEARN A PDAG. IF WE KNOW TRUE DAG \rightarrow COMPUTE CAUSAL EFFECTS. IF WE DON'T KNOW TRUE DAG \rightarrow COMPUTE FOR

LEARN FROM INTERVENTION DATA: SKIP CASES WHERE INTERVENTION HAPPENED BEFORE LEARNING θ .
OR AUGMENT NORMAL DAG WITH NODES FOR INTERVENTIONAL FACTORS

[IDA]

ALL DAGS IN CLASS
OR LOCAL NEIGHBORHOOD

LEARNING U-GGM

- EASIER THAN DAG \rightarrow NO CYCLE PROBLEM. • HARDER THAN DAG \rightarrow LIKELIHOOD DOES NOT DECOMPOSE. NO GREEDY SEARCH/MCMC
- MLE $\ell(\Omega) = \log \det(\Omega) - \text{TR}(S\Omega)$. $\Omega = \Sigma^{-1}$, S = EMPIRICAL COVARIANCE MATRIX $\rightarrow \nabla \ell(\Omega) = \Omega^{-1} - S$: COVARIANCE ESTIMATION
- 1:1 CORRESPONDENCE BETWEEN ZEROS IN Ω AND ABSENT EDGES IN GRAPH.
- LET'S USE A SPARSE INVOLVING OBJECTIVE: $J(\Omega) = -\log \det(\Omega) + \text{TR}(S\Omega) + \lambda \|\Omega\|_1$ GRAPHICAL LASSO!!! CONVEX BUT NON-SMOOTH BECAUSE ℓ_1
- COMPUTING EXACTLY $P(G|D)$. POSTERIOR INFERENCE IN SPACE OF GRAPHS. CAN DO IF GRAPH IS DECOMPOSABLE.
NOT DECOMPOSABLE \rightarrow V. HARD. MODIFIED GRADIENT DESCENT + DIAGONAL LAPLACE PRIOR
- NON GAUSSIAN STILL CONTINUOUS DATA: USE D MONOTONIC TRANSFORMATIONS f_j SO THAT RESULTING DATA IS JOINTLY GAUSSIAN NONPARAMETRIC DISTRIBUTION

DISCRETE UGM

HARDEST. STOCHASTIC LOCAL SEARCH IS NOT TRACTABLE. BUT STUFF HAS BEEN TRIED

GLASSO: FOR DISCRETE CRF/MRF. MODIFIED OBJECTIVE. V. COSTLY. APPROX FOR EVERYTHING.

LEARNING SPARSE MODELS MAY NOT BE ENOUGH BECAUSE IF TREewidth IS LARGE, INFERENCE IS INTRACTABLE. \rightarrow BOUND THE TREewidth

- THIN JUNCTION TREES, OR LOW CIRCUIT COMPLEXITY