

# DIRECTED GRAPHICAL MODELS

COMPACTLY REPRESENT JOINT DISTRIBUTION  $P(X|D)$ , HOW TO INFERR GIVEN SET OF VARS GIVEN MOTHER, HOW TO LEARN DISTRIBUTION PARAMS

CHAIN RULE:  $P(X_1:V) = P(x_1)P(x_2|x_1)P(x_3|x_2, x_1) \dots$  **CONDITIONAL PROBABILITY TABLES:**  $O(U^V)$  SPACE, A LOT!

**CONDITIONAL PROBABILITY DISTRIBUTION:**  $O(U^2 V^2)$  SPACE, BETTER! BUT STILL DEPENDS ON ALL PREVIOUS VARS

**CONDITIONAL INDEPENDENCE / MARKOV ASSUMPTION** CONDITIONAL JOINT AS PRODUCT OF MARGINALS  $P(X_1:V) = P(x_1) \prod_t P(x_t | x_{t-1})$

**GRAPHICAL MODEL:** MODELS CI IN  $> 1D$ , LACK OF EDGES  $\rightarrow$  CI

-(GRAPH) NEIGHBOURS  $\rightarrow$  SET OF DIRECTLY CONNECTED NODES

-(GRAPH) CLIQUE  $\rightarrow$  SET OF NODES ALL NEIGHBOURS OF EACH OTHER

**MAROV CHAIN:** INITIAL DISTRIBUTION + STATE TRANSITION MATRIX

**DIRECTED GRAPHICAL MODEL:** A GM WHOSE GRAPH IS A DAG. ALSO BAYESIAN NETWORKS / BELIEF NETWORKS / CAUSAL NETWORKS

$\rightarrow$  **ORDERED MARKOV PROPERTY** A NODE ONLY DEPENDS ON ITS IMMEDIATE PARENTS  $X_S \perp X_{PRED(S)} | X_{PA(S)}$  EACH  $P(x_t | x_{pa(t)})$  IS CPD

$\rightarrow$  **GENERALLY:**  $P(1:V|E) = \prod_t P(x_t | x_{pa(t)})$ ,  $O(VK^F)$  SPACE,  $F =$  PARENTS  $U =$  STATES

## NAIVE BAYES

CAN BE SEEN AS DGM. FEATURES ARE ASSUMED CI FROM LABELS  $P(Y|X) = P(Y) \prod_i P(x_i | Y)$

**TREE-AUGMENTED NAIVE BAYES:** CAPTURES CORRELATION BETWEEN FEATURES, GRAPH IS A TREE, CAN FIND OPTIMAL STRUCTURE, CAN HANDLE MISSING FEATURES

## MAROV AND HIDDEN MAROV MODELS

• **2ND ORDER MAROV CHAIN:**  $P(1:T) = P(x_1, x_2)P(x_3|x_1, x_2)P(x_4|x_2, x_3) = P(x_1, x_2) \prod_t P(x_t | x_{t-1}, x_{t-2})$

• **HIDDEN MAROV MODEL:** HIDDEN VARIABLES  $z_t$ , OBSERVED VARIABLES  $x_t$

$\rightarrow$  **TRANSITION MODEL**  $P(z_t | z_{t-1})$   $\rightarrow$  **OBSERVATION MODEL**  $P(x_t | z_t)$   $\rightarrow$  **STATE ESTIMATION**  $P(z_t | x_{1:t}, \theta)$  ESTIMATE HIDDEN STATE

• IF EXTEND TOO MUCH NUMBER OF PARAMS BLOWS UP

## NOISY-OR MODEL / SIGMOID BELIEF NET

IF A PARENT ON  $\rightarrow$  CHILD USUALLY ON, BUT OCCASIONALLY PARENT  $\rightarrow$  CHILD FAILS  $\theta_{st} = 1 - q_{st}$  FROM FAILURE

FROM CHILD = OFF  $P(V_t = 0 | h) = \prod \theta_{st}^{I(h_i=1)}$  • **DUMMY LEAK NOISE:** CHILD ON IF ALL PARENTS OFF  $\theta_{0t}$

• IF  $w_{st} = \log(\theta_{st})$   $P(V_t = 1 | h) = 1 - \exp(-w_{0t} - \sum h_i w_{st})$  SIMILAR TO LOGISTIC REGRESSION

## DIRECTED GAUSSIAN GRAPHICAL MODELS

DGM WHERE ALL VARS ARE REAL VALUED, CPD  $P(x_t | x_{pa(t)}) = N(x_t | \mu_t + W^T x_{pa(t)}, \sigma_t^2)$  **LINEAR GAUSSIAN CPD**

**GAUSSIAN BAYES NET:**  $P(x) = N(x | \mu, \Sigma)$   $\mu: x_t = \mu_t + \sum w_{ts} (x_s - \mu_s) + \sigma_t^2 z_t$   $z_t \sim N(0, 1)$ ,  $\sigma_t$  CONDITIONAL STANDARD OF  $x_t$  GIVEN PARENTS

$\rightarrow \mu = (\mu_1, \dots, \mu_D)$

$w_{ts}$  EDGE WEIGHT,  $\mu_t$  LOCAL MEAN

$\Sigma: (x - \mu) = W(x - \mu) + S z$ ,  $S = \text{DIAG}(\sigma)$ ;  $E = S z =$  VECTOR OF NOISE TERMS  $\rightarrow x - \mu = U \cdot S \cdot z$ ,  $U = (I - W)^{-1}$  **CHOLESKY DECOMPOSITION**

$\rightarrow \Sigma = U \cdot S^2 \cdot U^T$

# INFERENCE

WE HAVE  $p(x_{1:n} | \theta)$  AND  $\theta$  ARE UNKNOWN. VISIBLE VARS  $X_V$  HIDDEN VARS  $X_H$ . COMPUTING THE POSTERIOR OF THE UNKNOWNS GIVEN THE KNOWNS

$$p(x_h | x_v, \theta) = \frac{p(x_h, x_v | \theta)}{p(x_v | \theta)} = \frac{p(x_h, x_v | \theta)}{\sum_{x_h'} p(x_h', x_v | \theta)}$$

• WE CONDITION ON DATA BY CLAMPING ~~EVIDENCE~~ VISIBLE VARS TO OBSERVED VALUES  $x_v$  AND NORMALIZE ON DATA LIKELIHOOD / EVIDENCE

• IF ONLY SOME OF HIDDEN VARS ARE INTERESTING, OTHERS ARE MARGINALIZED OUT

$$p(x_q | x_v, \theta) = \sum_{x_n} p(x_q, x_n | x_v, \theta) \quad x_q = \text{INTERESTING} \quad x_n = \text{NUISANCE}$$

## LEARNING

$$\hat{\theta} = \underset{\theta}{\text{ARGMAX}} \sum_{i=1}^n \log p(x_{i:V} | \theta) + \log p(\theta)$$

MAP ESTIMATE OF PARAMS GIVEN DATA

• UNIFORM PRIOR: MLE

• IF BAYESIAN: PARAMS ARE ALSO UNKNOWN VARS. NO INFERENCE/LEARNING DISTINCTION

WE ADD PARAMS TO GRAPH AS NODES, CONDITION ON  $\emptyset$ , THEN INFER VALUES.

HOWEVER PARAM VARS DO NOT GROW WITH DATA, WE CAN POINT ESTIMATE THEM.

• PLATE NOTATION FOR CONVENIENT DRAWING, NESTED PLATES, ROLLED/UNROLLED REPRESENTATIONS

WHEN DATA IS COMPLETE:

- LIKELIHOOD:  $p(D | \theta) = \prod_{t=1}^T p(\theta_t | \theta)$  DATA ASSOCIATED WITH NODE FAMILY ( $n$  PARENTS). LIKELIHOOD DECOMPOSES ACCORDING TO GRAPH STRUCTURE

- PRIOR:  $p(\theta) = \prod_{t=1}^T p(\theta_t)$

- POSTERIOR:  $p(\theta | D) = \prod_{t=1}^T p(\theta_t | D_t) p(\theta_t)$  EACH CDD IS INDEPENDENT. • FACTORIZED PRIOR + FACTORIZED LIKELIHOOD  $\rightarrow$  FACTORIZED POSTERIOR

MISSING DATA/LATENT VARIABLES: LIKELIHOOD NO LONGER FACTORIZES. NO LONGER CONVEX. ONLY LOCAL ML/MAP ESTIMATES. APPROX. INFERENCE

## CI PROPERTIES

$I(G)$  = SET OF CI STATEMENTS ENCODED IN GRAPH.  $G$  IS **I-MAP** FOR  $P$  IFF  $I(G) \subseteq I(P)$ . IF GRAPH DOESN'T ASSESS FALSE STATEMENTS ABOUT  $P$ .

• A FULLY CONNECTED GRAPH IS IMAP FOR ALL DISTRIBUTIONS

## D-SEPARATION

- A PATH IS D-SEPARATED BY A SET OF NODES  $E$  (FOR EVIDENCE) IFF EITHER

•  $P$  CONTAINS A CHAIN WHERE  $MIDDLE \in E \rightarrow S \rightarrow M \rightarrow T$

•  $P$  CONTAINS TEND/FOUR WITH  $M \in E$  IS VERTEX

•  $P$  CONTAINS A COLLIDER/V-STRUCT WHERE  $M$  IS VERTEX,  $M \notin E$ , AND ANY DESCENDANT OF  $M$

- A SET  $A$  IS D-SEPARATED FROM SET  $B$  BY SET  $E$  IFF EACH <sup>UNION</sup> PATH FROM  $\forall a \in A$  TO  $\forall b \in B$  IS D-SEP BY  $E$

$$X_A \perp_G X_B \mid X_E \iff A \text{ IS D-SEP FROM } B \text{ GIVEN } E$$

C.I.  $\rightarrow$  D-SEP  
DIAGNOSE GLOBAL  
MARKOV PROPERTY

- CONDITIONING ON BOTTOM OF V/COLLIDER MAKES PARENTS DEPENDANT  $\rightarrow$  EXPLAINING AWAY/INTER-CAUSAL REASONING/BERKSON'S PDOX  
IE IF YOU OBSERVE THE SUM OF TWO VALUES, KNOWING ONE GIVES YOU THE OTHER.

ALGO: BAYES BALLS

DIRECTED LOCAL MARKOV PROPERTY:  $t \perp \text{ND}(t) \mid \text{PA}(t)$  NO: NON-DESCENDANTS

- THESE PROPS ARE EQUIVALENT.

ORDERED MARKOV PROPERTY:  $t \perp \text{POSD}(t) \mid \text{PA}(t)$

## MARKOV BLANKET

SET OF NODES READING NODE  $t$  CI FROM ALL OTHER NODES IN GRAPH. MB IS PARENTS, CHILDREN, AND CO-PARENTS (NODES PARENTS OF ITS CHILDREN)

## FULL CONDITIONAL

FULL PRODUCT OF CPD WITH  $x_t$  IN THEIR SCOPE  $P(x_t | x_{-t}) \propto P(x_t | x_{pa(t)}) \prod_{s \in children(t)} P(x_s | x_{pa(s)})$

## INFLUENCE/DECISION DIAGRAM

IT'S A DAG + UTILITY NODES, DECISION NODES. OIL-WIND-CONTROL PROBLEM.

— EXPECTED UTILITIES — POSTERIORES — EXHAUST THEM FOR ALL SCENARIOS → COMPUTE OPTIMAL POLICIES

— MAX EXPECTED UTILITY:  $MEU = \sum_s P(s) EU(d^*(s) | s)$

— VALUE OF PERFECT INFORMATION:  $\Delta$  UTILITY DUE TO GAINING MORE INFO. CAN HAVE COSTS, USED TO SET POLICIES/THRESHOLDS

VPI:  $MEU(I+T \rightarrow D) - MEU(I)$

## POMDP - PARTIALLY OBSERVED MARKOV DECISION PROCESS

IT'S A HMM WITH AUGMENTED ACTION AND REWARD NODES → USED FOR PERCEPTION - ACTION CYCLES IN INTELLIGENT AGENTS

MDP - MARKOV DECISION PROCESS. LIKE POMDP BUT FULLY OBSERVED. EASIER TO SOLVE

↳ STATE NOT FULLY OBSERVED  
CHOOSE ACTIONS BASED ON BELIEF STATE  
 $P(z_t | x_t, a_t, t)$