

# UNDIRECTED GRAPHICAL MODELS - MARKOV RANDOM FIELDS

CAUSAL MRF: DAG MODEL W/ 2D LATTICE TOPOLOGY. ALSO MARKOV MESH. CI PROPERTIES ARE KINDA FUNNY

→ UNDIRECTED MODEL: MRF MARKOV NETWORK. MARKOV BLANKET IS NOW JUST SET OF NEIGHBOURS

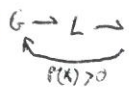
## PROS

- SYMMETRIC → BETTER FOR CERTAIN DOMAINS
- DISCRIMINATIVE UGM, CONDITIONAL RANDOM FIELDS,  $P(Y|X)$  BETTER THAN DGM.

## CONS

- PARAMS LESS EXPENSIVE, LESS MODULAR
- PARAM ESTIMATION MORE COMPUTATIONALLY EXPENSIVE

## CI PROPERTIES

- GLOBAL MARKOV PROPERTY:  $X_A \perp\!\!\!\perp X_B \mid X_C$  IFF  $C$  SEPARATES  $A, B$  IN GRAPH. REMOVE  $C$ ; NO PATHS CONNECT  $A, B$ .
- UNDIRECTED LOCAL MARKOV PROPERTY: MARKOV BLANKET OF NODE IS SET OF IMMEDIATE NEIGHBOURS
- PAIRWISE MARKOV PROPERTY: TWO NODES ARE CI IF NO DIRECT EDGE BETWEEN THEM.
- $G \rightarrow L \rightarrow P$ : ALL PROPERTIES ARE THE SAME. CAN RELY ONLY ON PAIRWISE FOR EASY THINGS  


## MORALIZATION

CONVERSION OF DGM TO UGM. CAN'T JUST DROP EDGE DIRECTION TO MAINTAIN CI STRUCTURE. CONNECT NODES THAT HAVE A COMMON CHILD TOO.  
→ STILL LOSES SOME INFORMATION. FORM ANCESTRAL GRAPH OF  $G$  WRT  $U = \{A \cup B \cup C\}$ . → REMOVES NODES NOT IN  $U$  AND NOT ANCESTORS OF  $U$ . MORALIZE ANCESTRAL GRAPH. THEN APPLY SEPARATION RULES FOR UGM. ALTERNATIVE TO D-SEPARATION

## UGM vs DGM

NEITHER IS MORE POWERFUL THAN OTHER. DGM AND UGM ARE PERFECT MAPS FOR DIFFERENT SETS OF DISTRIBUTIONS.  
DGM: V-STRUCTURE  $A \rightarrow C \leftarrow B$ .  $A \perp B$ ,  $A \not\perp B \mid C$ . NO UGM CAN PRECISELY REPRESENT THESE TWO STATEMENTS  
UGM: 4-CYCLE

UGM AND DGM BOTH PERFECT; ARE CHORDAL / DECOMPOSABLE GRAPHS. → IF WE COLLAPSE EACH MAXIMAL CLIQUE, THE GRAPH WILL NOW BE A TREE.

## MRF PARAMETERIZATION

NO TOPOLOGICAL ORDERING → NO CAN USE CHAIN RULE FOR  $P(Y)$ .

POTENTIAL FUNCTIONS OR FACTORS: ASSOCIATED WITH EACH MAXIMAL CLIQUE.  $\psi_c(y_c | \theta_c)$  CAN BE ANY NON-NEGATIVE CLIQUE  
JOINT DISTRIBUTION: PROPORTIONAL TO PRODUCT OF CLIQUE POTENTIALS.

Hammersley-Clifford Theorem: A POSITIVE DISTRIBUTION  $P(Y) > 0$  SATISFIES CI PROPERTIES OF U-GRAPH  $G$  → IT CAN BE REPRESENTED AS PRODUCT OF POTENTIALS ONE PER MAX-CLIQUE

$$P(Y|\theta) = \frac{1}{Z(\theta)} \prod_c \psi_c(y_c | \theta_c), \quad Z(\theta) = \sum_Y \prod_c \psi_c(y_c | \theta_c)$$

PARTITION FUNCTION  
ENSURES SUM TO 1

## RELATION TO GIBBS DISTRIBUTION:

GIBBS IS  $P(Y|\theta) = \frac{1}{Z(\theta)} \cdot \exp(-\sum E(y_c | \theta_c))$ ,  $E(y_c) > 0$  IS ENERGY OF VARS IN CLIQUE. IF WE  $\psi_c(y_c | \theta_c) = \exp(-E(y_c | \theta_c))$ , IT BECOMES UGM

ENERGY-BASED MODELS! HIGH PROBABILITY STATES ↔ LOW ENERGY CONFIGURATIONS. PAIRWISE MRF: EDGE POTENTIALS:  $P(Y|\theta) \propto \prod_{st} \psi_{st}(y_s, y_t)$   
VS MAXCLIQUES

MAX ENTROPY/LOGLINEAR FORM OF POTENTIAL FUNCTIONS:  $\log \psi_c(y_c) = \phi_c(y_c)^T \theta_c$ ,  $\phi_c(y_c)$  IS FEATURE VECTOR, ONE PER EDGE/CLIQUE.

LOG PROBABILITY IS  $\log P(Y|\theta) = \sum_c \phi_c(y_c)^T \theta_c - Z(\theta)$ . IF WEIGHT PER FEATURE →  $\psi_{st}(y_s = u, y_t = v) = \exp([\theta_{st}^T \phi_{st}]_{uv}) = \exp(\theta_{st}(u, v))$  IS  $K \times K$ . TABULAR POTENTIAL IN LOGLINEAR FORM

## ISING MODEL

BORN FOR MODELING MAGNETS. 2D/3D LATTICE PAIRWISE CUQUE POTENTIAL NO CONNECTION,  $w_{st}=0$ .  $w$  SYMMETRIC. OFTEN  $w_{st} = 1$  CONSTANT. 0.370

UNNORMALIZED LOG PROB:  $\log \tilde{P}(y) = -\sum_{s,t} y_s w_{st} y_t = -\frac{1}{2} y^T W y$  0.320

A BIAS/EXTERNAL FIELD:  $\log \tilde{P}(y) = \frac{1}{2} y^T W y + b^T y$ . CAN WRITE IN FORM SIMILAR TO GAUSSIAN GM

BUT NP-HARD VS  $O(D^3)$

$\psi_{st}(y_s, y_t) = \begin{pmatrix} e^{w_{st}} & e^{-w_{st}} \\ e^{-w_{st}} & e^{w_{st}} \end{pmatrix}$   $w_{st}$  IS COUPLING STRENGTH OF S-T

NEIGHBORING SPINS LIKELY TO BE IN SAME STATE,  $+/-1$ . GROUP STATE DIFF FROM NEIGHBORS. FRUSTRATED SYSTEM. TWO MODES. ASSOCIATIVE. MULTIPLE MODES. NO ALL CONSTRAINTS AT SAME TIME. MARKOV NETWORK

## HOPFIELD NETWORKS

IS FULLY CONNECTED ISING MODEL WITH  $W=W^T$  SYMMETRIC. LEARNED FROM DATA. ASSOCIATIVE MEMORY / CONTENT ADDRESSABLE MEMORY, USED IN PATTERN COMPLETION TASKS. RETRIEVE EXAMPLE FROM MEMORY GIVEN ONLY A PIECE. NO EXACT INFERENCE  $\rightarrow$  COORD DESCENT WITH

ICM SETS EACH NODE TO MIN ENERGY STATE GIVEN ALL NEIGHBORS. CAN BE INTERPRETED AS RECURRENT NEURAL NET ITERATIVE CONDITIONAL MODELS

BOLTZMANN MACHINE IS ISING/HOPFIELD GENERALIZATION THAT ADDS HIDDEN NODES.

## POTTS MODEL

ISING GENERALIZATION WITH MULTIPLE DISCRETE STATES  $y_i = \{1, 2, \dots, K\}$   $\psi_{st} = \begin{pmatrix} e^{\beta} & 0 & 0 \\ 0 & e^{\beta} & 0 \\ 0 & 0 & e^{\beta} \end{pmatrix}$  PHASE TRANSITIONS: DIFFERENT VALUES OF  $\beta$  INDUCE DIFFERENT CLUSTERING BEHAVIORS. USED IN IMAGE SEGMENTATION AS PRIOR (NEIGHBORING PIXELS  $\rightarrow$  SAME LABEL LIKELY).

EMERGING GRAPHICAL MODEL IS MIX OF DIRECTED AND UNDIRECTED EDGES  $\rightarrow$  CHAIN GRAPH. IS 2D ANALOG OF HMM, PARTIALLY OBSERVED MRF BUT INFERENCE MUCH HARDER.

OK FOR REGULARIZING SUPERVISED PROBLEMS, NOT ACCURATE ENOUGH FOR UNSUPERVISED SEGMENTATION.

## GAUSSIAN MRFs

UNDIRECTED GGM  $\begin{cases} P(y|t) = \prod \psi_{st}(y_s, y_t) \prod \psi_t(y_t) \\ \psi_{st}(y_s, y_t) = \exp(-\frac{1}{2} y_s \Lambda_{st} y_t) \\ \psi_t(y_t) = \exp(-\frac{1}{2} \Lambda_{tt} y_t^2 + \eta_t y_t) \end{cases}$  JOINT:  $P(y|t) = \exp[\eta^T y - \frac{1}{2} y^T \Lambda y] \rightarrow$  MVN IN INFORMATION FORM!  $\Lambda_{st}=0 \iff y_s \perp y_t | y_{st}$

$\bullet$  CORRESPOND TO SPARSE PRECISION MATRICES!

$\bullet$  D-GGM ARE SPARSE FACTORIZATIONS OF COVARIANCE MATRICES

$\bullet$  COMBINE D + U GGM:  $P(y_t | y_{t-1}, y_{t-2}, t) = N(y_t | \Lambda_1 y_{t-1} + \Lambda_2 y_{t-2}, \Sigma)$  VECTOR AUTO-REGRESSIVE PROCESS OF ORDER 2. IE TIME-SERIES ASPECT IS MODELED WITH DAG. ZEROS IN TRANSITION MATRICES ARE NO DIRECTED EDGES FROM  $y_{t-1}, y_{t-2}$  INTO  $y_t$ . ZEROS IN  $\Sigma^{-1}$  ARE NO UNDIRECTED EDGES IN  $y_t$ .

$\bullet$  BIDIRECTED GRAPH: USED FOR SPARSE COVARIANCE MATRICES, NON CONNECTED NODES ARE UNCONDITIONALLY INDEPENDENT.

CAN BE CONVERTED INTO A DAG WITH LATENT VARIABLES. EACH DIR EDGE BECOMES A HIDDEN, CONFOUNDER VARIABLE

$\bullet$  BIDIR + DIR: DIRECTED MIXED GRAPHICAL MODEL

## MARCOV LOGIC NETWORKS

FULL FIRST ORDER LOGIC AND FORMAL RULES APPROACH. USED UNROLLED GGM.  $\rightarrow$  MARCOV LOGIC NETWORK

$\bullet$  REWRITE RULES IN CNF/CNJSAL FORM. RESTRICT LANGUAGE TO HORN CLAUSES TO MAKE INFERENCE DECIDABLE

$\bullet$  ATTACH WEIGHT TO RULES, DEFINING CUQUE POTENTIALS  $\psi_c(x_c) = \exp(w_c \phi_c(x_c))$ ,  $\phi_c$  LOGICAL EVALUATES CONJ.

$\bullet$  CONSTRUCT GROUND NETWORK BY CREATING RANDOM BINARY VARS AND WIRE ACCORDING TO CLAUSES

LEARNING FOR MRF

ML AND MAP ESTIMATION EXPENSIVE SO RARELY DONE

WITH GRADIENT METHODS (FULLY VISIBLE)

LOG LINEAR FORM  $P(y|\theta) = \frac{1}{Z(\theta)} \text{EXP}(\sum_i \theta_i^T \phi(y))$ ,  $l(\theta) = \frac{1}{N} \log P(y|\theta)$ ,  $\frac{\partial l}{\partial \theta_i} = \left[ \frac{1}{N} \sum \phi_i(y_i) \right] - E[\phi_i(y)]$

CLAMPED FORM      UNCLAMPED FORM  
CONTINUOUS      → INFERENCE ONCE PER GRADIENT STEP. TRAINING IS SLOWER THAN DGM

THEY ARE IN EXPONENTIAL FAMILY, SO CONVEX WRT  $\theta$ , GRADIENT.

$\frac{\partial l}{\partial \theta_i} = E_{P(y|\theta)}[\phi_i(y)] - E_{P(y|\theta)}[\phi_i(y)]$  EXPECTED FEATURE VECTOR ACCORDING TO EMPIRICAL DISTRIBUTION - MODEL'S EXPECTATION OF F.V. → MOMENT MATCHING

WITH MISSING DATA/HIDDEN VARS

$P(y, h|\theta) = \frac{1}{Z(\theta)} \text{EXP}(\sum_i \theta_i^T \phi_i(h, y))$  → GRADIENT OF LL IS EXPECTED FEATURES WHERE WE CLAMP  $y_i$  AND AVERAGE OVER  $h$

$\frac{\partial l}{\partial \theta_i} = \frac{1}{N} \sum \{ E[\phi_i(h, y_i)|\theta] - E[\phi_i(h, y_i)|\theta] \}$  1ST TERM CLAMPS VISIBL 2ND TERM LEAVES THEM FREE. WE MINIMIZE OVER  $h_i$

ALTERNATIVE USE EM WITH GRADIENT METHODS AT M STEP

APPROXIMATE METHODS

NO CLOSED FORM SOLUTION → GRADIENT OPTIMIZERS → REQUIRES INFERENCE → INFERENCE INTRACTABLE → LEARNING INTRACTABLE

PSEUDO-LIKELIHOOD

MINIMIZE THAT INSTEAD OF MLE.  $Q_{PL}(\theta) = \frac{1}{N} \sum \log P(y_i|\theta)$  IS PRODUCT OF FULL CONDITIONALS/COMPOSITE LIKELIHOOD. FOR GAUSSIAN MRF IS AS ML FASTER BECAUSE EACH FULL CONDITIONAL ONLY REQUIRE STATES OF SIBLING NODES TO NORMALIZE. NODE GIVEN ALL ITS NEIGHBOURS. HARD TO APPLY WHEN HIDDEN VARS. STILL VERY MUCH FASTER

STOCHASTIC MLE

FOR PARTIALLY OBSERVED MODELS. USE MCMC TO GENERATE SAMPLES FOR GRADIENT DESCENT. TRICKS TO MAKE MCMC FASTER! START AT PREV VAL AND TAKE A FEW STEPS.

FEATURE INDUCTION

UNSUPERVISED WAY TO LEARN FEATURES. START WITH BASE SET THEN GREEDILY ADD BETTER FOUND ONES. LINE GRAPHICAL MODEL STRUCTURE LEARNING BUT MORE FINE GRAINED.

$P(y) = \frac{1}{Z} \text{EXP}(\theta_1 \phi_1(y) + \theta_2 \phi_2(y))$

ITERATIVE PROPORTIONAL FITTING

PAIRWISE MRF, TABLE POTENTIAL,  $\psi_{ij} = \text{EXP}(\theta_{ij}^T \text{STATS})$  FEATURE VECTORS ARE INDICATORS. COUNTS.  $P(y_i) = P(y_i|\theta)$  AT OPTIMUM. FOR DECOMPOSABLE GRAPHS  $P(y|\theta) = \prod_i \psi_i(y_i)$  ELSE  $\psi_i^{\text{TM}}(y_i) = \psi_i(y_i) \times \frac{P_{\text{sum}}(y_i)}{P(y_i|\psi^t)}$  COORDINATE ASCENT. IS FIXED POINT ALGORITHM FOR MOMENT MATCHING CONSTRAINTS. IF DECOMPOSABLE CONVERGES IN SINGLE ITERATION.

CAN BE MADE FASTER → EFFICIENT IFF, OR 'PARALLELIZED' WRT PARAMETERS UPDATES. CONTRAST TO OTHER FEATURE SCHEMES → ITERATIVE SCALING. CLOSED FORM SOLUTIONS

DECOMPOSABLE GRAPH: TREE-LIKE GRAPHS WHERE EITHER UGM OR DGM ARE ON.



## CONDITIONAL RANDOM FIELDS

IS MRF WHERE ALL CLIQUE POTENTIALS ARE CONDITIONED ON INPUT FEATURES, CAN BE SEEN AS **STRUCTURED OUTPUT** EXTENSION OF LOGISTIC REGRESSION.  
ADVANTAGE OF DISCRIMINATIVE VS GENERATIVE. NO NEED TO MODEL STUFF WE ALWAYS SEE. CAN MAKE POTENTIALS DATA-DEPENDANT AND EV. TURN OFF.  
DISADVANTAGE: REQUIRES LABELED TRAINING DATA, SLOWER. SAME AS LOGISTIC REGRESSION VS NAIVE BAYES, 'LIVE HAM BUT DISCRIMINATIVE'!

$$P(y|x, w) = \frac{1}{Z(x, w)} \prod_c \psi_c(y_c | x, w) \quad \bullet \text{ LOG-LINEAR POTENTIALS } \psi_c(y_c | x, w) = \exp(w_c^T \phi(x, y_c))$$

$\phi(x, y)$  IS FEATURE VECTOR FROM GLOBAL INPUTS  $x$  AND LOCAL LABELS  $y_c$

SOME FEATURES ARE GLOBALS, OTHERS ARE LOCAL ON THE NODE.

**MAXIMUM ENTROPY MARKOV MODEL (MEMM):** 'REVERSE  $y \rightarrow x$  MARKOV ON HMM AND PUT  $x_2$  ON TOP OF  $y_1$ '

$$P(y|x, w) = \prod_t P(y_t | y_{t-1}, x, w)$$

STATE TRANSITION PROBABILITIES ARE CONDITIONED ON INPUTS.

**LABEL BIAS:** LOCAL FEATURES AT  $t$  DO NOT INFLUENCE PRIOR STATES. INFORMATION DOESN'T FLOW BACKWARDS. ISSUES WITH IE DISAMBIGUATION OF SENTENCES.

**CHAIN-STRUCTURED CRF:** UNDIRECTED HMM, PUT  $x_2$  ON TOP. NOW STUFF IS GLOBALLY NORMALIZED. LOCAL FACTORS NEED NOT SUM TO 1.  
NO VALID PROBABILITY UNTIL WE'VE SEEN WHOLE SENTENCE; THEN WE CAN NORMALIZE.  $\rightarrow$  CRF NOT OK FOR ONLINE INFERENCE

## APPLICATIONS

**HANDWRITING RECOGNITION:** LOCAL LETTERS VERY AMBIGUOUS, CONTEXT HELPS. NODE POTENTIALS ARE PROBABILISTIC DISCRIMINATIVE CLASSIFIERS, (NN, RVM), EDGE POTENTIALS ARE BIGRAM MODELS

**NOVN PHRASE CHUNKING:** PHRASE SEGMENTATION. BIO TAGGING. BIO ARE STATES. CONSTRUCTED FEATURES.  $L_1$  REGULARIZATION. FWD-BWD FOR INFERENCE.

**NAMED-ENTITY RECOGNITION:** SMALL SCALE TASK EXTRACTION. IE WORDS FOR PLACES, NOUNS, LOCATIONS, ETC... CONSIDERS LONG-RANGE CORRELATIONS BETWEEN WORDS  $\rightarrow$  **SHIP-CHAIN CRF**

**NATURAL-LANGUAGE PARSING:** PROBABILISTIC CONTEXT-FREE GRAMMARS. <sup>GENERATIVE</sup> PROBABILITY OF SEQUENCE IS SUM OF ALL TREES THAT GENERATE IT.  
DISCRIMINATIVE VARIANTS ENCODE PAIRS OF TREE  $y$  GIVEN SEQUENCE OF WORDS  $x$ .  $P(y|x) \propto \exp(w^T \phi(x, y))$

**HIERARCHICAL CLASSIFICATION:** MULTI-CLASS CLASSIFICATION WITH LABEL TAXONOMY SPECIFYING HIERARCHY OF CLASSES. POSITION IN H. ENCODED VIA BINARY VECTOR. COMBINED  $w$ /FEATURES WITH TENSOR PRODUCT.

**PROTEIN SIDE-CHAIN PREDICTION:** SHIP-CHAIN MODEL FOR PROTEIN STRUCTURE. PREDICT ANGLES GIVEN AMINO-ACID SEQUENCE, ENERGY FUNCTION TO BE MINIMIZED COMES FROM CHEMICAL-PHYSICAL MODEL

**STEREO-VISION:** LOW LEVEL DEPTH ESTIMATION GIVEN TWO DIFFERENT IMAGES. CRF. NODE POTENTIAL IS DISPARITY OF MATCHING PIXELS.  
GAUSSIAN POTENTIALS  $\rightarrow$  FUZZY BORDERS. TRUNCATED GAUSSIAN POTENTIALS  $\rightarrow$  PRESERVE DISCONTINUITIES / EDGES.  
**MEINIC CRF**  $\rightarrow$  DISCRETIZE VARIABLES FOR FAST COMPUTATION. POTENTIALS FORM A MATRIX

## TRAINING

$$\frac{\partial \ell}{\partial w_c} = \frac{1}{N} \sum_i [\phi_c(y_i, x_i) - E[\phi_c(y, x_i)]]$$

INFERENCE FOR EVERY TRAINING CASE IN EACH GRADIENT STEP.  $O(N)$  TIMES SLOWER THAN MRF.

**PARAMETER TYING** TO ENSURE DISTRIBUTIONS OF SIMILAR SIZE. **REGULARIZATION** TO PREVENT OVERFITTING:  $L_0, L_1, L_2$

$$\ell'(w) = \frac{1}{N} \sum_i \log p(y_i | x_i, w) - \lambda_1 \|w\|_1 - \lambda_2 \|w\|_2^2$$

## STRUCTURAL SVM

STRUCTURAL OUTPUT CLASSIFIERS LEVERAGING ON FAST MAP SOLVERS.

- $R_{MAP} = -\log p(w) - \sum_{i=1}^N \log p(y_i | x_i, w)$  **IDEA!** MINIMIZE POSTERIOR EXPECTED LOSS ON TRAINING SETS  $\rightarrow$  BOUND ON EXPECTED LOSS.  $\rightarrow$  SET SPHERICAL GAUSSIAN PRIOR

$\rightarrow R_{SUM}(w) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N [\max\{0, 1 - y_i w^T x_i\}]$  OPTIMIZES UPPER BOUND ON BAYESIAN OBJECTIVE

ALSO (MAINLY) DERIVED NON PROBABILISTICALLY WITH MARGIN FORMULATION. SAME OBJECTIVE

- QUADRATIC PROGRAMS BUT EXPONENTIALLY MANY CONSTRAINTS. IF LOSS AND FEATURE VECTOR DECOMPOSE WITH GRAPHICAL MODEL CAN REDUCE TO POLYNOMIAL. ELSE CUTTING PLANE OR STOCHASTIC SUBGRADIENT

## CUTTING PLANE

ITERATIVELY FIND MOST VIOLATED CONSTRAINT. IF  $\gamma$  MARGIN, ADD IT TO SET OF WORKING CONSTRAINT. WHEN ALL ARE, WE HAVE A TOLERANCE OF  $\epsilon$ . FTIME.

- HAS LOSS AUGMENTED DECODING, IS EFFICIENT
- CAN DO BETTER  $\rightarrow$  LINEAR TIME USING KENNEL TRICK QP IN  $O(N)$ !!

## STRUCTURED PERCEPTION

$\hat{y} = \text{ARGMAX } f(y|x)$  WITH EG VITERBI. IF  $\hat{y} \neq y$  UPDATE WEIGHTS WITH  $w_{n+1} = w_n + \phi(y, x) - \phi(\hat{y}, x)$

ONLINE



## STOCHASTIC SUBGRADIENT DESCENT

PEWASOS

$g(w) = \sum_{i=1}^N \phi(x_i, \hat{y}_i) - \phi(x_i, y_i) + \lambda \|w\|$   $w_{n+1} = w_n - \eta_n g(w_n)$   $\eta$  STEPSIZE. IS GENERALIZATION OF PERCEPTION

## LATENT STRUCTURAL SVM

HAS HIDDEN VARIABLES. CAN 'TURN' CRF INTO SSUM FORM/LATENT SVM BUT OBJECTIVE NOT CONVEX **CCCP** PROCEDURE: MINIMIZES  $f(w) - g(w)$  WHERE  $f, g$  CONVEX

# EXACT INFERENCE FOR GRAPHICAL MODELS

GENERALIZE STUFF LIKE FWD-BWD OR VAMP TO ARBITRARY DIRECTED/UNDIRECTED GRAPHS

## BELIEF PROPAGATION FOR TREES (SUM-PRODUCT ALGORITHM)

GENERALIZES FWD-BWD, AND SUM-PRODUCT ALGORITHM

• **SEQUENTIAL PROTOCOL:**  $\psi_s$  IS NODE EVIDENCE  $\psi_{st}$  IS EDGE POTENTIAL. PICK 'ARBITRARY ROOT'. LEAVES  $\rightarrow$  ROOT + ROOT  $\rightarrow$  LEAVES PHASES.

**L  $\rightarrow$  R:** MULTIPLY MSG BY EDGE POTENTIALS. AT NODE BELIEF STATE IS  $BEL_t^-(x_t) = P(x_t | V_t^-) = \frac{1}{2} \psi_t(x_t) \prod m_{c \rightarrow t}^-(x_t)$

**R  $\rightarrow$  L:** ALL INFO PARENT HAS RECEIVED, EXCEPT THAT FROM SPECIFIC CHILD.  $m_{t \rightarrow s}^+ = P(x_s | V_t^+) = \sum_{x_t} \psi_{st}(x_s, x_t) \frac{BEL_t^-(x_t)}{m_{s \rightarrow t}^-(x_t)}$   
OR WE MULTIPLY ALL BUT ONE MSGS.

• T-D MSG CONDITIONAL PROBABILITIES POSTERIOR (DEPEND ON B-U) - MARKOV BLANKET

• SUM PRODUCT: T-D DO NOT DEPEND ON B-U? CONDITIONAL LIKELIHOODS

## • **PARALLEL PROTOCOL:**

CAN USE IN NON-TREES. **SYSTOLIC ARRAY**. EACH NODE ABSORBS MSG FROM ITS NEIGHBORS IN PARALLEL AND THEN SENDS MSG TO ALL NEIGHBORS STILL IN PARALLEL BY MULTIPLYING RECEIVED ONES (MINUS ONE FROM RECIPIENT) BY EDGE POTENTIAL. CONVERGES AFTER  $O(G)$  STEPS, GRAPH DIAMETER LINEAR

## • **GAUSSIAN BP**

$P(x|v)$  IS JOINTLY GAUSSIAN. GAUSSIAN PAIRWISE MRF. NODES, EDGE POTENTIALS ARE GAUSSIANS  $\rightarrow$  ALL MSGS AND MARGINALS ARE TOO. ALGORITHM TRICKS MAKE IT COMPUTE IN  $O(D)$  VS  $O(D^3)$ .

$$\psi_t(x_t) = \exp\left(-\frac{1}{2} A_{tt} x_t^2 + b_t x_t\right)$$

$$\psi_{st}(x_s, x_t) = \exp\left(-\frac{1}{2} x_s A_{st} x_t\right)$$

## • **MAX-PRODUCT BP**

REPLACE  $\sum$  WITH MAX  $\rightarrow$  LOCAL MAP ESTIMATES. ONLY NOT BE GLOBALLY CONSISTENT  $\rightarrow$  VITERBI

## • **SAMPLING**

CAN SAMPLE FROM TREE MODEL WITH FWD/BWD SAMPLING ALGOS

## • **POSTERIORS ON SETS OF VARIABLES**

COMBINE  $x_t, x_{t+1}$  INTO 'MEGA NODE'. COMBINATION  $\Psi$  MATRICES.

## VARIABLE ELIMINATION ALGORITHM

WORKS ON ANY KIND OF GRAPH

• IF DGM  $\rightarrow$  MINIMIZE INTO UGM FIRST

• ASSIGN POTENTIALS TO CPD, ALREADY LOCALLY NORMALIZED.  $Z=1$

• ENUMERATING ON ALL POSSIBLE ASSIGNMENTS IS  $O(2^n)$   $\rightarrow$  PUSH SUMS INSIDE PRODUCTS

• CREATE TEMPORARY FACTORS BY MULT, THEN MARGINALIZE SUMMED-OVER VAR. PROCESS ~~LEFT TO RIGHT~~ **RIGHT TO LEFT**

• THIS COMPUTES ANY MARGINAL

• IF WE WANT CONDITIONAL  $\rightarrow$  RATIO OF TWO MARGINALS NORMALIZED BY EVIDENCE / CLAMPING VARS BY THEIR OBSERVED VALUES

• VE IS NON-SEQUENTIAL DYNAMIC PROGRAMMING.

• CAN APPLY TO OBTAIN MAP VIA  $x^* = \underset{x}{\text{ARGMAX}} \prod \psi_c(x_c)$  BY REPLACING SUMS WITH MAX (+TRACEBACK STEP)

• VE WORKS ON ANY COMMUTATIVE SEMI-RING! SET WITH + AND  $\times$  OPERATIONS ASSOCIATIVE, COMMUTATIVE, WITH IDENTITY, AND DISTRIBUTE LAW  
RUNNING TIME IS EXPONENTIAL IN SIZE OF LARGEST FACTOR. **ELIMINATION ORDER** MATTERS A LOT. BECAUSE OF CYCLES CREATED. MINIMIZE THE TREEWIDTH  $W = \min \max |C| - 1$ . GENERALLY NP-HARD. CHAINS AND TREES ARE EASY THO.  $O(VN^2)$  FOR CHAINS;  $O(\min\{M, N\})$  FOR WIDGES

• **WEAKNESS:** NOT EFFICIENT FOR COMPUTING MULTIPLE QUERIES CONDITIONAL ON SAME EVIDENCE W/ FWD-BWD  $O(K^2T)$   
DOES NOT REUSE/REUSE PREVIOUS MESSAGES

## JUNCTION-TREE ALGORITHM

GENERALIZES EFFICIENCY OF BP TO ARBITRARY GRAPHS.

- MCM-RUN VE, ADDING FILL-IN EDGES  $\rightarrow$  **CHORDAL GRAPH** (EACH UNBI CYCLE HAS CHORDS CONNECTING NON-ADJACENT COMPONENTS)
- FIND MAXIMAL CLIQUES  $\rightarrow$  EASY FOR CHORDALS.
- ARRANGE THE CLIQUES INTO THE **JUNCTION TREE**  $\rightarrow$  **RUNNING INTERSECTION PROPERTY** AT LEAST 1 VAR SHARED BETWEEN ADJACENT NODES.
- **BP TO JTREE**: EXACT VALUES OF  $P(x_i | v)$  FOR CLIQUE NODES  $\rightarrow$  EXTRACT NODE/EDGE MARGINALS FOR ORIGINAL MODEL BY MARGINALIZATION

### MESSAGE PASSING!

CONCEPTUALLY SAME AS BP

ADJUST POTENTIALS BY MULTIPLYING CLIQUE POTENTIALS AND SEPARATING SET POTENTIALS.

**BOTTOM  $\rightarrow$  UP**: MARGINALIZING OUT THE SEPARATING SET, ALREADY KNOWS ABOUT IT

**TOP  $\rightarrow$  BOTTOM**: ALWAYS DIVIDING OR LEAVING OUT NODE.

**COMPLEXITY**: IF NODES DISCRETE W/  $n$  STATES.  $O(|C| \cdot K^{w+1})$  SPACE AND TIME. FOR GGM WE USE INFORMATION FORM AND IS  $O(|C| \cdot w^3)$  TIME AND  $O(|C| \cdot w^2)$  TIME. FOR CHAINS IS SAME AS HMM SMOOTHING

**GENERALIZATIONS**: MAP WITH MAX-PRODUCT.  $N$ -MOST PROBABLE CONFIGURATIONS. POSTERIOR SAMPLES. SOLVING CSP/SAT PROBLEMS. LOGICAL REASONING.

**FINAL REMARK**: STUFF IS EXPONENTIAL IN TREE-WIDTH. APPROXIMATE INFERENCE FTW