

GRAPHICAL MODELS

WITHIN DL, GMs ARE A 'MODELING TOOL'. **TASKS**: DENSITY ESTIMATION, DENOISING, IMPUTATION, SAMPLING

TABULAR, FULL CASE NOT VIABLE: K^N PARAMS, STORAGE COST, CURSE OF DIMENSIONALITY, COST OF SAMPLING, COST OF INFERENCE

• **DIRECTED MODELS** \rightarrow CONDITIONAL DISTRIBUTIONS

$$P(x) = \prod_i P(x_i | \text{PAG}(x)) \quad \bullet \quad \underbrace{O(K^N)}_{\text{TAB}} \rightarrow \underbrace{O(N^M)}_{\text{DAG}}, \quad M \text{ IS MAX NO VARS APPEARING IN A CONDITIONAL}$$

• **UNDIRECTED MODELS** \rightarrow FACTORS, CLIQUE POTENTIALS $\phi(C)$

• $\tilde{P}(x) = \prod_C \phi(C)$ UNNORMALIZED DISTRIBUTION • **PARTITION FCN** $Z = \int \tilde{P}(x) dx = 1 \rightarrow$ NORMALIZED DISTRIBUTION $P(x) = \frac{1}{Z} \tilde{P}(x)$

• **ENERGY FORMULATION**

$\tilde{P}(x) = \exp(-E(x))$ ALLOW NOT TO EXPLICITLY CONSTRAINT POTENTIAL/FACTOR FCN \rightarrow **BOLTZMANN DISTRIBUTION**, HENCE 'BOLTZMANN MACHINES'

• MAKS EASY COMPUTATIONS: $\exp(a) \exp(b) = \exp(a+b)$

• **D-SEPARATION & MORALIZATION**

• **FACTOR GRAPHS**

UNDIRECTED MODELS WITH 'FACTOR' NODES EXPLICITLY REPRESENTING SCOPES OF EACH ϕ . FACTORS OF UNKNOWN ALIBED DISTRIBUTION, RESOLVE AMBIGUITY

• **LATENT VARIABLES VS STRUCTURE LEARNING**

USING LATENTS TO MODEL STRUCTURE AVOIDS THE NEED TO PERFORM DISCRETE SEARCHES AND MULTIPLE ROUNDS OF TRAINING, IMPOSITION OF INTERACTIONS BETWEEN VISIBLE VARS. \rightarrow CAN ALSO USE LATENT VARS FOR FEATURE LEARNING

• **APPROXIMATE INFERENCE**

ANA VARIATIONAL OR SAMPLING APPROXIMATIONS **REPARAMETRIZATION TRICK** ALLOWS TO COMPUTE GRADIENTS THROUGH h_s AND DO BACKPROP IS BACKPROPAGATION THROUGH SAMPLING. $\rightarrow h \sim P(h|\theta) \rightarrow h = f(\theta, \eta)$ f CONTINUOUS, η NOISE. $\rightarrow \frac{\partial}{\partial \theta} \int L(h) P(h|\theta) dh = \frac{\partial}{\partial \theta} \int L(f(\theta, \eta)) P(\eta) d\eta$

$$\rightarrow \frac{\partial}{\partial \theta} \int L(f(\theta, \eta)) P(\eta) d\eta$$

GRAPHICAL MODELS IN DEEP LEARNING

EMPHASIS ON LATENT VARS, LOTS OF LATENT VARS. $|h| > |x|$. LATENTS NOT DESIGNED TO HAVE SPECIFIC MEANING. TRAINING SORTS IT OUT.

ALSO, IN DL WE RARELY CARE TO OBTAIN EXACT INFERENCE. WE USE MIN AMOUNT OF INFO AND APPROXIMATE IT AS QUICKLY AS POSSIBLE.

MODEL POWER / CAPACITY RAISED UNTIL IT 'BREAKS'

RESTRICTED BOLTZMANN MACHINE

ENERGY MODEL WITH BINARY VISIBLES AND HIDDEN $E(v, h) = -b^T v - c^T h - v^T W h$. NO DIRECT INTERACTIONS BETWEEN v_i OR h_i

$$P(h|v) = \prod_i P(h_i|v) \quad P(v|h) = \prod_i P(v_i|h) \quad P(h_i = 1|v) = \sigma(v^T W_{\cdot i} + b_i) \quad \rightarrow \text{EFFICIENT BLOCK GIBBS SAMPLING} \text{ IN ON } v \text{ SIMULTANEOUS}$$

\rightarrow CAN TRAIN W/ APPROX OF $\nabla_{\theta} \log Z$

$$\rightarrow \frac{\partial}{\partial W_{ij}} \mathbb{E}_{v, h} [E(v, h)] = -v_i h_j$$

Monte Carlo Methods

FOR SAMPLING FROM GRAPHICAL/ENERGY MODELS.

- **ANCESTRAL SAMPLING** SAMPLE IN TOPOLOGICAL ORDER, CONDITIONING ON PARENTS
- **MALLOV CHAINS** EQUILIBRIUM DISTRIBUTION, DETAILED BALANCE, IT IS EIGENVECTOR OF T WITH $\lambda=1$ AND IT'S THE LARGEST AND ONLY ONE WITH $\lambda=1$
→ BURNIN, MIXING-RATE