

VARIATIONAL INFERENCE

IDEA: PICK APPROXIMATION $Q(x)$, MAKE IT AS CLOSE AS POSSIBLE TO TRUE POSTERIOR $P^*(x) = P(x|D)$ | COMPUTE W/O FULL CONDITIONAL, AVERAGE OUT NEIGHBOURS.

→ MINIMIZE KL DIVERGENCE $KL(P^*||Q)$ OR ITS REVERSE $KL(Q||P^*)$, BECAUSE EXPECTATIONS W/ Q ARE TRACTABLE.

→ $J(Q) = \sum_x Q(x) \log \frac{Q(x)}{P^*(x)} = \sum_x Q(x) \log \frac{Q(x)}{P(x)} - \log Z = KL(Q||P) - \log Z$ BECAUSE Z IS CONSTANT

→ $J(Q) = KL(Q||P^*) - \log Z \geq -\log Z = -\log P(D)$, $J(Q)$ IS UPPER BOUND ON NLL. $Z = P(D)$, NORMALIZATION CONSTANT

→ ALTERNATIVE: $L(Q) = -J(Q) = \log P(D)$, MAXIMIZE LOWER BOUND ON DATA LL. ENERGY FUNCTIONAL RELATED TO EM! LOWER BOUNDS!

→ $J(Q) = E_Q[\log q(x)] + E_Q[-\log \tilde{p}(x)] = -H(Q) + E_Q[E(x)]$ • EXPECTED ENERGY - ENTROPY OF SYSTEM J = VARIATIONAL FREE ENERGY
 $L(E(x)) = -\log \tilde{p}(x)$

→ $J(Q) = E_Q[-\log P(D|x)] + KL(Q(x)||P(x))$ • EXPECTED NLL + PENALTY TERM

REVERSE KL $KL(Q||P) = \sum_x Q(x) \ln \frac{Q(x)}{P(x)}$ • INFORMATION, I-PROJECTION. ZERO-FORCING FOR Q UNDERESTIMATES P SUPPORT

FORWARDS KL $KL(P||Q) = \sum_x P(x) \ln \frac{P(x)}{Q(x)}$ • MOMENT, M-PROJECTION. ZERO-AVOIDING FOR Q . OVERESTIMATES P SUPPORT

• IF TRUE DISTRIBUTION MULTIMODAL → FORWARDS KL IS GENERALLY BAD IDEA

ALPHA DIVERGENCE: FAMILY OF DIVERGENCE MEASURES $D_\alpha(P||Q) = \frac{1}{1-\alpha^2} \left(1 - \int P(x)^{(1+\alpha)/2} Q(x)^{(1-\alpha)/2} dx \right)$ • 0 IFF $P=Q$
 • NOT SYMMETRIC → NOT A TRUE METRIC

$\alpha \rightarrow 1 = KL(P||Q)$
 $\alpha \rightarrow -1 = KL(Q||P)$
 $\alpha \rightarrow 0 = D_H(P||Q) = \int (P(x)^{1/2} Q(x)^{1/2})^2 dx \sim$ HELLINGER DISTANCE $\sqrt{D_H(P||Q)}$ IS TRUE DISTANCE METRIC

MEAN FIELD METHOD

ASSUMPTION: POSTERIOR IS FULLY-FACTORED APPROXIMATION $Q(x) = \prod_i Q(x_i)$ GOAL: MIN $KL(Q||P)$ OPTIMIZING OVER PARAMS OF EACH MARGINAL Q_i

- COORDINATE DESCENT, UPDATE $\log Q_i(x_i) = E_{-Q_i}[\log \tilde{p}(x)] + \text{CONST}$ EACH STEP. $\tilde{p}(x) = P(x, D)$ UNNORMALIZED POSTERIOR, E_{-Q_i} EXPECTATION OVER $f(x)$ WAS ALL VARS EXCEPT x_i
- WHEN UPDATING Q_i WE REASON IN TERMS OF J 'S MARGINAL BOUND, OTHERS ARE ABSORBED INTO CONSTANT.
- CAN BE USED FOR MANY MODELS
- UPDATE DERIVATION MINIMIZE $L(Q_i) \rightarrow L(Q_i) = -KL(Q_i||P_i) \rightarrow Q_i = f_i \rightarrow \log Q_i(x_i) = E_{-Q_i}[\log \tilde{p}(x)] + \text{CONST}$ • SIMILAR TO GIBBS SAMPLING BUT JOINTS MEAN MSGS.

MEAN FIELD FOR ISING MODEL

EXAMPLE → IMAGE DENSIFYING FACTORED APPROXIMATION IS $Q(x) = \prod_i Q(x_i, m_i)$, m_i MEAN VALUE

- $\log \tilde{p}(x) = x_i \sum_j W_{ij} x_j + L_i(x_i) + \text{CONST}$
- $Q_i(x_i) \propto \text{EXP} \left(x_i \sum_j W_{ij} m_j + L_i(x_i) \right) \rightarrow$ REPLACE NEIGHBOUR W/ AVE VALUES
- $m_i = \sum_j W_{ij} m_j$
- $Q_i = m_i + 0.5(L_i^+ - L_i^-)$ APPROX MARGINAL POSTERIOR
- $m_i = E_i[x_i] = \text{TANH}(\alpha_i) \rightarrow m_i = \text{TANH} \left(\sum_j W_{ij} m_j + 0.5(L_i^+ - L_i^-) \right)$ CAN TURN INTO FIXED POINT
- DAMPED UPDATES $m_i^t = (1-\lambda)m_i^{t-1} + \lambda \text{TANH} \left(\sum_j W_{ij} m_j^{t-1} + 0.5(L_i^+ - L_i^-) \right)$

STRUCTURED MEAN FIELD

EXPLOIT TRactable SUBSTRUCTURE IN PROBLEM. GROUP VAR SETS TOGETHER AND UPDATE SIMULTANEOUSLY.

EXAMPLE: FACTORIAL HMM

M CHAINS, LENGTH T, K STATES $P(x, y) = \prod_M \prod_T (x_{TM} | x_{T-1, M}) P(y_t | x_{tM})$. EACH CHAIN IS APPROX INDEPENDENT BUT COUPLED IN POSTERIOR DUE TO OBSERVATIONS. JUNCTION TREE IS $O(TM^2 K^2)$. MEAN FIELD IS $O(TM^2 K^2)$

• APPROXIMATION: PRODUCT OF CHAINS $Q(x_{TM} | x_{T-1, M}, \xi_{TM}) = \prod_{t=1}^T \left(\xi_{TM} \prod_{j=1}^M (A_{Mj})^{x_{T-1, Mj}} \right)^{x_{TM}}$

$E_Q(x) = \langle \text{NO WAY} \rangle$, $E(x, y)$ EXACT POSTERIOR

EACH CHAIN INDIVIDUALLY UPDATED WITH FWD-BWD

• OBJECTIVE: $KL(Q||P) = E[E] - E[E_Q] - \log Z_u + \log Z$

• UPDATES $\xi_{TM} = \exp(W_M^T \sum \tilde{y}_{tm} - \frac{1}{2} \delta_M)$

• ξ_{TM} IS LOCAL EVIDENCE, AGGREGATES OVER NEIGHBORING CHAINS.

$\delta_M = \text{DIAG}(W_M^T \Sigma^{-1} W_M)$

FWD-BWD UPDATE EACH CHAIN IN PARALLEL

$\tilde{y}_{tm} = y_t - \sum_M W_M E[x_{tM}]$

• $O(TM^2)$ FOR FULL UPDATE SWEEP.

VARIATIONAL BAYES

TO INFER MODEL PARAMETERS, NOT HIDDEN VARS. $P(D|D) \approx \prod_n q(\theta_n)$

• LOWER BOUND

IF PARAMS + LATENT INFERENCE \rightarrow VARIATIONAL BAYES EM

VB FOR UNIVARIATE GAUSSIAN

CONJUGATE PRIOR: $P(\mu, \lambda) = N(\mu | \mu_0, (\nu_0 \lambda)^{-1}) G(\lambda | a_0, b_0)$

APPROX FACTORED POSTERIOR: $Q(\mu, \lambda) = q_\mu(\mu) q_\lambda(\lambda)$

TARGET $\log \tilde{P}(\lambda, \mu) = \frac{N}{2} \log \lambda - \frac{\lambda}{2} \sum (x_i - \mu)^2 - \frac{\nu_0 \lambda}{2} (\mu - \mu_0)^2 + \frac{1}{2} \log(\nu_0 \lambda) + (a_0 - 1) \log \lambda - b_0 \lambda + \text{const}$

$q_\mu(\mu) = N(\mu | \mu_N, \lambda_N^{-1})$ OBTAINED FROM LOG = AVERAGE OVER λ

$q_\lambda(\lambda) = G(\lambda | a_N, b_N)$

• COMPUTE EXPECTATIONS, DERIVE EXPLICIT UPDATE FORMS.
 $\mu_N, \lambda_N, a_N, b_N$. a_N, μ_N CONSTANTS \rightarrow FIXED POINT UPDATES FOR OTHERS

• OBJECTIVE: MINIMIZE $L(Q)$, LOWER BOUND ON LOG MARINAL LL

$L(Q) \leq \log P(D) = \log \int P(D | \mu, \lambda) P(\mu, \lambda) d\mu d\lambda = \langle \text{TEDIOUS ALGEBRA} \rangle = \frac{1}{2} \log \frac{1}{\lambda_N} + \log \Gamma(a_N) - a_N \log b_N + \text{const}$
- PUSH LOG INTO EVA
- INTEGRALS BECOME EXPECTATIONS
- ...
MONOTONICALLY INCREASES WITH VARIATES

VB FOR LINEAR REGRESSION

PRIOR: $P(w, \alpha, \lambda) = N(w | 0, (\lambda \alpha)^{-1} I) G(\lambda | a_0^{\lambda}, b_0^{\lambda}) G(\alpha | a_0^{\alpha}, b_0^{\alpha})$

FACTORED POSTERIOR: $Q(w, \alpha, \lambda) = q(w, \lambda) q(\alpha)$

OPTIMAL POSTERIOR: $Q(w, \alpha, \lambda) = N(w | w_N, \lambda_N^{-1} V_N) G(\lambda | a_N^{\lambda}, b_N^{\lambda}) G(\alpha | a_N^{\alpha}, b_N^{\alpha})$ (MURPHY P 747)

• CAN FORMULATE WITH APD PRIORS

• ALTERNATE UPDATES $Q(w, \lambda)$ AND $Q(\alpha)$

• POSTERIOR PROPORTIONAL TO $P(y | x, D) = T(y | w_N^T x; \frac{b_N^{\lambda}}{a_N^{\lambda}} (1 + x^T V_N x), \frac{b_N^{\lambda}}{a_N^{\lambda}})$ ONCE PARAMS INFERRD

• EXACT MLL $P(D) = \int \int P(y | x, w, \lambda) P(w | \alpha) P(\lambda) dw d\alpha d\lambda$

• SIMILARITIES TO EMPIRICAL BAYES: MAX $\log P(D)$ WHILE VB MAXES LOWER BOUND ON IT

VARIATIONAL BAYES EM

- WITH LATENTS AND PARAMS $z_i \rightarrow x_i \leftarrow \theta$. MIXTURE MODELS, PCA, HMM. VBEM IS MORE BAYESIAN, MODELS UNCERTAINTY IN θ TOO. SAME COMPUTATIONAL COST AS REGULAR EM. **MEAN FIELD:** $P(\theta, z_{1:N}|D) \approx Q(\theta)Q(z) = Q(\theta)\prod Q(z_i)$
- VARIATIONAL E STEP:** UPDATES $Q(z_i|\theta)$. AVERAGES OVER PARAMETERS INSTEAD THAN PLUGGING IN $\hat{\theta}$ MAP ESTIMATES. PLUG-IN POSTERIOR MEAN THEN USE STUFF SUCH AS FWD-BWD
 - VARIATIONAL M STEP:** UPDATES $Q(\theta|D)$. UPDATES HYPERPARAMS USING EXPECTED SUFFICIENT STATISTICS. **VBEM = EM** $\rightarrow q(\theta|D) \approx \delta_{\hat{\theta}}(\theta)$
 - ADVANTAGE** MARGINALIZE PARAMS OUT \rightarrow CAN COMPUTE LOWER BOUND \rightarrow MODEL SELECTION

VBEM FOR MIXTURE OF GAUSSIANS

- EXACT PRIOR:** $P(\theta) = \text{Dir}(\pi|\alpha_0) \prod_n N(\mu_n|\mu_0, (\beta_0 \Lambda_n)^{-1}) W_n(\Lambda|L_0, \nu_0)$ ASSUME ALL PRIOR PARAMS ARE SAME FOR ALL CLUSTERS
- FACTORED POSTERIOR:** $Q(z, \theta) = Q(z|\theta)Q(\theta) = \left[\prod_i \text{Cat}(z_i|r_i) \right] \left[\text{Dir}(\pi, \alpha) \prod_n N(\mu_n|\mu_n, (\beta_n \Lambda_n)^{-1}) W_n(\Lambda_n|L_n, \nu_n) \right]$

E-STEPM-STEP
- ML LOWER BOUND** $L = \sum_z \int Q(z|\theta) \log \frac{P(x, z, \theta)}{Q(z, \theta)} d\theta \leq \log P(D)$
- POSTERIOR PREDICTIVE:** SUM OF WEIGHTED T_s
 $P(x|D) = \sum_z \int P(x|z, \theta) P(z|\theta) Q(\theta) d\theta = \sum_n \frac{\alpha_n}{\sum \alpha_n} T(x|\mu_n, \Lambda_n, \nu_n, 1-D)$
- MODEL SELECTION:** SELECT k . FIT SEVERAL MODELS AND COME TO LOWER BOUND OR ML. WATCH OUT FOR UNIDENTIFIABILITY, $k!$ EQUIVALENT MODES. $\log P(D|k) \approx \log(k!) + L(k)$
- SPARSITY:** FIT A SINGLE MODEL WITH LARGE k AND $\alpha_0 \ll 1$, ENCOURAGES SPARSE MIXING VECTOR. IN VBEM MIXING WEIGHTS ARE SUBJECT TO A PENALTY MORE SEVERE FOR SMALL CLUSTERS (FEW WEIGHTED COMPONENTS) \rightarrow THEY WILL EMPTY OUT OVER ITERATIONS. EFFICIENT WAY TO SEARCH OPTIMAL NO OF CLUSTERS: SPINES NEAR EDGES OF SIMPLEX

VARIATIONAL MESSAGE PASSING

GEN PURPOSE METHOD FOR DGM WHOSE CRD ARE IN EXP FAMILY AND PARENTS ARE CONJ. SWEEP OVER GRAPH ON 1 AT TIME UPDATE. **VIBES/MP ON FOR CONTINUOUS LATENTS** (NEEDN'T TOUCH!)

LOCAL VARIATIONAL BOUNDS

- WE NOW REPLACE A TERM IN THE JOINT WITH A SIMPLER ONE TO FACILITATE COMPUTING THE POSTERIOR. **LOCAL VARIATIONAL APPROXIMATION**
- GENERALLY NEEDED WHEN GAUSSIAN PRIOR X MULTINOMIAL LIKELIHOOD: MULTI-TASK LEARNING, DISCRETE FACTOR ANALYSIS, CORRELATE TOPIC MODEL
- GENERAL DIFFICULTY:** $LSE(\eta_i) = \log(1 + \sum e^{\eta_i \cdot m})$ \rightarrow COMES FROM GAUSSIAN PRIOR X MULTINOMIAL LIKELIHOOD. **LOG-SUM-EXP**. **LIKELIHOOD:** $P(y|x, w) = \prod \exp[y_i^T \eta_i - LSE(\eta_i)]$ MULTI-TASK LOGISTIC REGRESSION \leftarrow LOGISTIC REGRESSION
- BOHNING QUADRATIC BOUND:** $P(y_i|x_i, w) \geq f(x_i, \psi_i) N(\psi_i|x_i, w, A_i^{-1}) \rightarrow$ NOW EASY TO COMPUTE **APPROX POSTERIOR** $Q(w) = N(\mu_w, \Sigma_w)$
- COMPUTE $L(\hat{a})$; INTRODUCE BOUNDS $\rightarrow L(\hat{a}) \geq [\text{VERY LONG EXPRESSION}] \rightarrow$ USE COORDINATE ASCENT \rightarrow UPDATE VARIATIONAL LIKELIHOOD \rightarrow UPDATE VARIATIONAL POSTERIOR μ_w, Σ_w

SIGMOID FUNCTION BOUND:

- U BOUND.** HAS ADAPTIVE CURVATURE TERM, OR STILL BOHNING BOUND
- $\log(1 + e^\eta) \leq \frac{1}{2} a(\xi) \eta^2 - b(\xi) \eta + c(\xi)$
 - HAS ADAPTIVE CURVATURE

$\log(1 + e^\eta) \leq \frac{1}{2} a \eta^2 - b \eta + c$
 - CONSTANT CURVATURE

a, b, c , ME STUFF.

OTHER BOUNDS, EVEN MORE MESSY:

- PRODUCT OF SIGMOIDS** (MURPHY P. 762-3)
 - JENSEN'S INEQUALITY** ALL ON LSE FUNCTION
 - MULTIVARIATE DELTA**
- * SOMETIMES WE DO VARIATIONAL INFERENCE WITH UPPER BOUNDS INSTEAD *

MOAR VARIATIONAL INFERENCE

Q NOW ISN'T FACTORIZED, NOT EVEN GLOBALLY VALID JOINT. JUST LOCALLY CONSISTENT → DIST OF TWO ADJACENT NODES MESSAGES WITH CORRESPONDING MESSAGES

LOOPY BELIEF PROPAGATION

IDEA: DISREGARD LOOPS, APPLY BELIEF PROPAGATION UNTIL 'CONVERGENCE'!

- ON PAIRWISE MODELS (BINARY ~~FAIR~~ GRAPHS), WORKS WELL! BECAUSE THEY'RE TREE-LIKE LOCALLY. CYCLE IS $\log n$ LONG.
- ON FACTOR GRAPHS

FACTOR GRAPH: WAY TO UNIFY DGM AND VGM. UNDIRECTED BIPARTITE GRAPH WITH ROUND NODES, VARIABLES, AND SQUARE NODES, FACTORS. EDGES FROM EACH VAR TO FACTOR THAT MENTION IT. IN VGM → FACTOR = POTENTIALS. IN DGM FACTORS → CPDs

$V \rightarrow F$ MSG: $m_{x \rightarrow f}(x) = \prod_{h \in \text{NEIGH}(x) \setminus \{f\}} m_{h \rightarrow x}(x)$

$F \rightarrow V$ MSG: $m_{f \rightarrow x}(x) = \sum_y f(x,y) \prod_{y \in \text{NEIGH}(f) \setminus \{x\}} m_{y \rightarrow f}(y)$

$BEL(x) \propto \prod_{f \in \text{NEIGH}(x)} m_{f \rightarrow x}(x)$

CONVERGENCE

NO GUARANTEES LBP WILL CONVERGE. INVESTIGATE WITH COMPUTATION TREES. T STEPS OF LBP → EXACT COMPUTATION IN TREE OF HEIGHT T+1 IF EDGE STRENGTHS ARE SUFF. WEAK → INFLUENCE DIMINISHES OVER DIFFERENT PASSES

DAMPING: SEND OUT DAMPED MESSAGES $M_{TS}^n(x_i) = \lambda M_{TS}(x_i) + (1-\lambda) M_{TS}^{n-1}(x_i)$, USUALLY $\lambda \sim 0.5$

MESSAGE SCHEDULING:

- SYNCHRONOUS: ARBORS/UPDATE ALL IN PARALLEL
- ASYNCHRONOUS: COMPUTE USING NEW MSG FROM EARLIER IN THE ORDERING, OLD FROM LATER IN THE ORDERING.

HOW TO ORDER? FIXED, RANDOM, PICKING SPANNING TREES AND SWEEP ONE AT A TIME; ADAPTIVELY → BASED ON DIFF FROM PREVIOUS VALUES, HIGHEST FIRST.

TREE REPARAMETERIZATION

RESIDUAL BELIEF PROPAGATION

ACCURACY! SINGLE LOOP → EXACT MAP ESTIMATES, ELSE NOT BUT ERROR CAN BE BOUNDED. GAUSSIAN MODELS → IF CONVERGES, MEANS EXACT, VARIATIONAL VARIANCES

SPEEDUP TRICKS:

- LARGE STATE SPACES: EACH LBP MSG IS $O(u^d)$. MANY STATES (IE, 256 PER VISION) IS TOO MUCH. FFT MAKES IT $O(u \log u)$ BECAUSE MESSAGES ARE JUST CONVOLUTIONS. $\psi_{ST}(x_s, x_t) = \psi(x_s - x_t)$. ALSO DISTANCE TRANSFORM IF APPROPRIATE. $O(K)$
- MULTI-SCALE: 2D BITMAPS, VISION, INITIALIZE ACTUAL GRID BY COMPUTING VALS ON A STACK OF COARSE GRIDS.
- CASCADE: FILTER ON SPEED/ACCURACY TRADEOFF, PRUNE IMPROBABLE STATES.

LBP FROM VARIATIONAL POV

WHY IS LBP VARIATIONAL?

VGM IN EXP FAMILY REPRESENTATION: $P(x|\theta) = \frac{1}{Z(\theta)} \exp(-E(x))$, $E(x) = -\theta^T \phi(x)$. θ NODE/EDGE PARAMS. ϕ EDGE IMAGERS (SUFF STATS)

$M(\theta) = E[\phi(x)]$ VECTOR MESSAGES. COMPREHENSIVELY CHARACTERIZES $P(x|\theta)$. MEAN PARAMS.

MARGINAL POLYTOPE: $M(\theta)$, SPACE OF ALLOWABLE μ VECTORS. SET OF ALL MEAN PARAMS. GENERATED FROM VALID PROBABILITY DISTRIBUTION. OBTAINED BY TAKING CONVEX COMBINATION OF $\phi(x) \rightarrow$ IS CONVEX HULL OF FEATURE SET. $M(\theta) = \text{CONV}\{\phi(x_i), \phi(x_j)\}$. IT DEFINES A VOLUME, OR INTERSECTION OF HALF-PLANES.

INFERENCE AS VARIATIONAL PROBLEM: V.I. FINDS Q MAXIMIZING ENERGY FUNCTIONAL $L(q) = E_q[\log \tilde{P}(x)] + H(q) \leq \log Z$
 $\rightarrow \log \tilde{P} = \theta^T \phi(x)$, $q = P \rightarrow \max_{\mu \in M(\theta)} \theta^T \mu + H(\mu)$ μ IS JOINT DISTRIBUTION OVER ALL STATE CONFIGURATION. ENTROPY! YAM! $\mu = E_P[\phi(x)]$
 $\max_{\mu \in M(\theta)} \theta^T \mu + H(\mu) = \log Z(\theta)$

• **MEAN FIELD AS VARIATIONAL:** NATURAL PARAMS ASSOCIATED W/ SUFF STATS & QUISIDE CLASS $\rightarrow 0$

$M_F(\mathcal{G}) = \{ \mu \in \mathcal{R}^D : \mu = E_\theta[\phi(x)] \text{ for } \theta \in \Sigma \}$ IS INNER APPROXIMATION OF MARGINAL POLYTOPE $M_F \subseteq M_G$

IS NON-CONVEX \rightarrow MULTIPLE LOCAL OPTIMA

MEAN FIELD ENERGY FUNCTIONAL: $\max_{\mu \in M_F(\mathcal{G})} \theta^T \mu + H(\mu) \leq \log Z(\theta)$

• MF MAXIMIZES A CONCAVE OBJECTIVE OVER A NON-CONVEX SET

PSEUDO-MARGINALS

• **LBP AS VARIATIONAL:** $M(\mathcal{G})$ HERE IS EXPONENTIALLY LARGE. WE RELAX CONSTRAINTS TO LOCAL CONSISTENCY: $\mu \rightarrow \gamma$ • $\sum_{x_S} \gamma_S(x_S) = 1$

• $L(\mathcal{G}) = \{ \gamma \geq 0 : \text{CONSTRAINTS HOLD} \}$ IS CONVEX OUTER APPROXIMATION ON $M(\mathcal{G})$, $M(\mathcal{G}) \subseteq L(\mathcal{G})$ • $\sum_{x_T} \gamma_{ST}(x_S, x_T) = \gamma_S(x_S)$

SUM TO 1, MARGINALIZATION

• IF GRAPH IS A TREE $M(\mathcal{G}) = L(\mathcal{G})$

• **ENTROPY:** $H(\mu) = \sum_{S \in V} H_S(\mu_S) - \sum_{S \in E} I_{ST}(\mu_{ST})$, EXACT FOR TREES

• **BETHE APPROX:** $H_{\text{BETHE}}(\gamma) = \sum_S H_S(\gamma_S) - \sum_{ST} I_{ST}(\gamma_{ST}) \rightarrow$ SCREW IT, USE IT EVEN IF GRAPH IS NOT TREE

• **BETHE FREE ENERGY:** $F_{\text{BETHE}}(\gamma) = -[\theta^T \gamma + H_{\text{BETHE}}(\gamma)]$

• **LBP OBJECTIVE:** $\min_{\gamma \in L(\mathcal{G})} F_{\text{BETHE}}(\gamma) = \max_{\gamma \in L(\mathcal{G})} \theta^T \gamma + H_{\text{BETHE}}(\gamma)$ • OPTIMIZE NON-CONCAVE OBJECTIVE OVER CONVEX SET
MULTIPLE LOCAL OPTIMA, APPROXES $\log Z(\theta)$, EXACT IF TREE

• ANY FIXED POINT OF LBP ALGO IS STATIONARY POINT FOR CONSTRAINED OBJECTIVE

MF VS LBP

LBP XACT FOR TREES, MF NOT

LBP OPTIMIZES W/ NODE, EDGE MARGINALS; MF ONLY NODE \rightarrow LBP MORE ACCURATE

IF EDGE MARGINALS FACTORIZE \rightarrow SAME TREE ENERGY APPROXIMATION

MF OBJECTIVE HAS MORE LOCAL OPTIMA \rightarrow HARDER; CAN IF STARTS FROM UNIFORM MSG BUT OK IF WE INIT MF WITH LBP MARGINALS

MF USEFUL BECAUSE GIVES A LOWER BOUND.

INIT MF WITH UNIFORM/RANDOM IS WAY BAD \rightarrow INIT WITH LBP MARGINALS

GENERALIZED BP

CLUSTER VARIATIONAL METHOD: CLUSTER TOGETHER NODES FORMING A TIGHT LOOP. \rightarrow HYPER EDGES BETWEEN SETS OF VERICES
REPRESENT USING POSETS. IF HYPEREDGE SIZE = TREEWIDTH \rightarrow GRAPH IS TREE; METHOD IS EXACT.

ENTROPY: $H_{\text{MINUCHI}}(\gamma) = \sum_{Y \in \mathcal{E}} C(Y) H_Y(\gamma_Y)$. H ENTROPY OF JOINT OF VERTICES IN SET Y . $C(Y)$ IS OVERCOUNTING NUMBER

ENERGY FUNCTIONAL: $F_{\text{MINUCHI}}(\gamma) = -[\theta^T \gamma + H_{\text{MINUCHI}}(\gamma)]$ **VARIATIONAL PROBLEM:** $\min_{\gamma \in L(\mathcal{G})} F_{\text{MINUCHI}}(\gamma) = \max_{\gamma \in L(\mathcal{G})} \theta^T \gamma + H_{\text{MINUCHI}}(\gamma)$

• OBJECTIVE NOT CONCAVE. GENERALIZED BP ALGO \rightarrow MORE ACCURATE THAN LBP, INCREASED COMPUTATIONAL COST

CONVEX BP

- EJS A CONCAVE OBJECTIVE ON A CONVEX SET, TREES OR PUNAR GRAPHS
- WORK WITH SUBMODELS $F \subseteq G$. (TREES OR PUNAR GRAPHS). • FEWER CONSTRAINTS \rightarrow ENTROPY IS HIGHER THAN G . $H(M, \rho)$ IS CONCAVE WRT M
- ENERGY $F_{\text{CONVEX}}(M, \rho) = -[M^T \theta + H(M, \rho)]$. WE HAVE ENTROPY UPPER BOUND
- OBJECTIVE! $\min_{\tau \in L(G, F)} F_{\text{CONVEX}}(\tau, \rho) = \max_{\tau \in L(G, F)} \tau^T \theta + H(\tau, \rho)$
- POLYTOPE! SO THAT PROJECTION OF τ ON G IS ON PROJ OF M ON F .
 $L(G, F) = \{\tau \in \mathbb{R}^E: \tau(F) \in M(F) \forall F \in F\}$
- TREE- REWEIGHTED BP:
 - CONSIDER SET OF ALL SPANNING TREES OF A GRAPH
 - FOR SINGLE NODES $\rho = 1$, FOR EDGES IS EDGE APPEARANCE PROBABILITY
 - IN THIS CASE $L(G, F) = L(G)$
 - GENERALLY DOES NOT CONVERGE, USE DAMPING OR DOUBLE-LOOP UPDATES
- OPTIMIZATION PROBLEM! $\max_{\tau \in L(G)} \left\{ \tau^T \theta + \sum_{s \in V} H_s(\tau_s) - \sum_{s, t \in E(G)} \rho_{st} \log(\tau_{st}) \right\}$ SAME AS LBP BUT FOR ρ_s
- ALGO! TRBP MESSAGE $T \rightarrow S$ IS FCN OF ALL MSG FROM $V \rightarrow T + S \rightarrow T$. MAIN DIFFERENCE IS STILL THE ρ WEIGHTS
 IF $\rho_{ST} = 1 \forall S, T$ IS STANDARD LBP BUT $\rho_{ST} = 1 \iff$ ORIGINAL GRAPH IS ALREADY A TREE.

EXPECTATION PROPAGATION

- BP WITH APPROXIMATED MSGS. GENERALIZES ASSUMED DENSITY FILTERING, APPROX POSTERIOR AT EACH TIME USING ASSUMED FUNCTIONAL FORM. EXTENDS ADF
- CUTTER PROBLEM: INFERRING UNKNOWN VECTOR x , WHEN OBSERVATION MODEL IS MIXTURE OF TWO GAUSSIANS, ONE AT x , ONE AT 0 .
- $P(y|x) = (1-w)N(y|x, 1) + wN(y|0, \alpha)$: • WITH FIXED PARS: GET EXPONENTIAL FORM
 - INFERENCE SPACE $M(\phi, \Phi)$ IS SET OF MEAN PARAMS REALIZABLE BY ANY PROB. DISTRIBUTION 'SEEN' THROUGH SUFFICIENT STATISTICS
 \rightarrow IS INTRACTABLE, 2^N MODES.
 - IDEA: USE $\tilde{\Phi}$ APPROXIMATED DISTRIBUTIONS, INCORPORATE INTRACTABLE TERMS AND WORK ITERATIVELY $P(x|\theta, \tilde{\theta}_1) \propto f_0(x) \exp(\theta^T \phi(x)) \exp(\tilde{\theta}_1^T \tilde{\Phi}(x))$
 - $\rightarrow P(x|\theta, \tilde{\theta}_1) = \exp(-\frac{1}{2} x^T \Sigma^{-1} x) [wN(y|0, \alpha) + (1-w)N(y|x, 1)]$ NOW INTRACTABLE!!
 - \rightarrow APPROX $M(\phi, \Phi)$ WITH $L(\phi, \tilde{\Phi})$, WHICH STILL CONVEX
 - \rightarrow VARIATIONAL PROBLEM! $\max_{(\tau, \tilde{\tau}) \in L(\phi, \tilde{\Phi})} \tau^T \theta + \tilde{\tau}^T \tilde{\theta} + H_{EP}(\tau, \tilde{\tau})$ ENTROPY! $H_{EP}(\tau, \tilde{\tau}) = H(\tau) + \sum_1 [H(\tau, \tilde{\tau}_1) - H(\tau)]$
 - BETTER ACCURACY THAN VARIATIONES OR MCMC PER UNIT OF CPU TIME
 - LBP IS SPECIAL CASE OF EP WHERE BASE DISTRIBUTION HAS NODE MARKOVALS AND INTRACTABLE TERM ARE EDGE POTENTIALS
 - APPLICATION: XBOX TRUESKILL
 - OPTIMIZE VIA MOMENT-MATCHING:
 - PICK A TERM 1 ,
 - COMPUTE Q_1 (REMOVE OLD APPROX $\tilde{\Phi}_1$) \rightarrow FACTOR
 - COMPUTE NEW Q WITH MOMENT MATCHING. ~~THE~~ SOLVING VAR. OBJECTIVE
 - COMPUTE NEW FACTOR MESSAGE.
 - AT (IF) CONVERGENCE: APPROX MARGINAL LIKELIHOOD $P(D) \propto \prod f_i(x) dx$

MAP STATE ESTIMATION

- FINDING MOST PROBABLE CONFIGURATION OF VARIABLES FOR DISCRETE-STATE GM. IF TREEWIDTH LOW \rightarrow CAN USE JUNCTION TREE.
- $x^* = \arg \max_{x \in X} \theta^T \phi(x)$, θ = NODE + FACTOR POTENTIALS.
- LINEAR PROGRAMMING RELAXATION: $\max_{MEM(G)} \theta^T \mu \leq \max_{\tau \in L(G)} \theta^T \tau$ L IS CONVEX OUTER BOUND OF POLYTOPE
- OPTIMIZE WITH DISTRIBUTED MESSAGE-PASSING
 - OBS! OBJECTIVE IS SAME AS STD VARIATIONAL OBJECTIVE WITHOUT ENTROPY
 - TERM \rightarrow ZERO TEMPERATURE LIMIT: DISTRIBUTION HAS ALL ITS MASS CONCENTRATED ON MODE. BETHE OBJECTIVE DOES NOT WORK, BECAUSE NOT CONCAVE. USE TRBP, WITH SPECIAL PARALLEL SCHEDULING ALWAYS CONVERGES

GRAPH CUTS

FIND MAP STATE ESTIMATES / MIN ENERGY CONFIGURATIONS USING MAX FLOW / MIN CUT ALGORITHMS ON GRAPHS.

- ADD SOURCE S AND SINK T TO GRAPH.

- COMPUTE MINIMUM $S-T$ CUT: PARTITION OF S, T CONNECTED NODES MINIMIZING SUM OF EDGE COSTS BETWEEN NODES ON DIFFERENT SIDES OF PARTITION

→ EQUIVALENT TO MINIMIZING ENERGY

→ SUBMODULAR ENERGIES: SUM OF DIAGONAL ENERGIES \leq SUM OF OFF-DIAGONAL ENERGIES. E.G. STD ISING MODEL WITH $\lambda > 0$, EXACT MAP.

→ NONBINARY METRIC MRF: REQUIRES PAIRWISE ENERGIES FORM A METRIC, METRIC MRF PICKS ONE STATE, EACH VAR STAYS SAME OR MOVES TO Q . Q -EXPANSION: STRONG LOCAL OPTIMUM AT EVERY STEP, BETTER THAN GREEDY SEARCH

$Q-P$ SWAP: VMS FLIP STATE IF Δ ENERGY

GRAPH CUTS VS BP

TRAP AND CUTS RUDE. VANILLA BP NOT SO MUCH. CUTS ARE FASTEST BY LITTLE MARGIN.

DUAL DECOMPOSITION:

$P^* = \max \sum \theta_i(x_i) + \sum \theta_f(x_f)$ WE CAN LOCALLY OPTIMIZE EACH FACTOR BUT COMBINATION MAKES THINGS UNTRACTABLE.

IDEA: OPTIMIZE EACH TERM INDIVIDUALLY, THEN INTRODUCE CONSTRAINTS FORCING STUFF TO AGREE.

- LAGRANGIAN MULTIPLIERS EVERYWHERE $\delta \rightarrow L(\delta) = \max_x L(\delta, x, x')$. IS DUAL OF LP RELAXATION. ALLOWS TO MIX AND MATCH DIFFERENT OPTIMIZATION ALGOS

- GRADIENT DESCENT: MUST USE SUBGRADIENT BECAUSE $L(\delta)$ IS CONVEX BUT NONDIFFERENTIABLE AT POINTS. PARALLEL UPDATES.

- COORD DESCENT: UPDATE ALL δ VECTOR AT ONCE. WITH MAX PRODUCT UNFOLD PROGRAMMING, USUALLY FASTER THAN GRAD. NO CONVERGENCE GUARANTEES THO.

- X^* RECOVERY: GENERALLY NP-HARD, OR IF LOCALLY DECODABLE \rightarrow EACH θ_i^{δ} HAS UNIQUE MAXIMUM x_i^* , LP RELAXATION IS UNIQUE

VARIATIONAL INFERENCE - FROM BISHOP

LOG-MARGINAL

BASE = FUNCTIONALS. IE, ENTROPY OF PROB DISR IS A FUNCTIONAL. $\log p(x) = L(q) + KL(q||p)$, $L(q) = \int Q(z) \cdot \log \frac{p(x|z)}{Q(z)}$ $KL = - \int Q(z) \cdot \log \frac{p(z)}{Q(z)}$

MEAN FIELD: IN $q_i(z_i) = E_{i \neq i} [\ln p(x_i, z)] + \text{CONST}$ ←

LOWER BOUND

LOG OF OPTIMAL FOR FACTOR Q IS EXP OF JOINT OVER ALL OTHERS

EXPECTATION PROPAGATION: USES REVERSE WL - $KL(p||q) = - \int p(z) \left[\sum_i \log q_i(z_i) \right] dz + \text{CONST}$

LOWER BOUND: USEFUL TO MONITOR THE PROGRESS

VARIATIONAL MESSAGE PASSING: IN DAG. $p(x) = \prod p(x_i | p_{\pi_i})$, $q(x) = \prod q_i(x_i)$, $\ln q_i(x_i) = E \left[\sum_j \ln p(x_j | p_{\pi_j}) \right] + \text{CONST}$

MARGINAL DENSITY

LOCAL VARIATIONAL METHODS: PRE ON LOGISTIC SEMIOT

EXPECTATION PROPAGATION: // USES MOMENT MATCHING

VI - OTHER FORMULATIONS

EVIDENCE LOWER BOUND

- JENSEN'S INEQUALITY: f CONCAVE $\rightarrow f(E[x]) \geq E[f(x)]$
- JENSEN'S ON OBSERVATION LOG PROBS: $\log p(x) = \log \int_z p(x, z) = \log \int_z p(x, z) \frac{Q(z)}{Q(z)} = \log E_Q \left[\frac{p(x, z)}{Q(z)} \right] \geq E_Q \left[\log p(x, z) \right] - E_Q \left[\log Q(z) \right]$
EVIDENCE LOWER BOUND \rightarrow MAXIMIZE IT, ELBO IS ALSO FREE ENERGY

$\overbrace{\text{EXPECTED LOG JOINT}}^{+} \quad \overbrace{H(q)}^{+} \text{ ENTROPY}$
- RELATION TO WL
 $KL(Q(z)||P(z|x)) = E_Q \left[\log \frac{Q(z)}{P(z|x)} \right] = E_Q [\log Q(z)] - E_Q [\log P(z|x)] = E_Q [\log Q(z)] - E_Q [\log p(z, x)] + \log p(x) = -ELBO + \log p(x)$
 $\rightarrow \log p(x)$ DOES NOT DEPEND ON $Q \rightarrow$ MAXIMIZE ELBO = MINIMIZE WL
 $p(z|x) = \frac{p(z, x)}{p(x)}$
- $\ln(p(x)) = KL(Q(z)||P(z|x)) + F(z, x)$
 $\text{ELBO / FREE ENERGY} = E_Q [\log p(x, z)] + H(q)$
 $X \text{ LOG JOINT} \quad \text{ENTROPY} = - E_Q [\log Q(z)]$

IN MEAN FIELD

$q(z) = \prod q_i(z_i) \rightarrow F(z, x) = -KL(Q||\text{Exp}(\log p(z, x))_{Q_{-i}}) + C \rightarrow$ Gibbs LINE SCHEME APPROX AT A TIME \rightarrow VARIATIONAL EM

- EXPLOIT $\log p(x|z) - F(q, x) = KL(Q||P)$
 - E MAX F w.r.t q , θ FIXED
 - M MAX F w.r.t θ , q FIXED

VARIATIONAL BAYES \rightarrow EVIDENCE \rightarrow MODEL SELECTION

JENSEN

- PUT A LOWER BOUND ON MARGINAL LIKELIHOOD
- FACTORIZE Q , $\log p(y|M) - F(q(x), q(\theta), y) = KL(Q||P)$
- MODEL SELECTION METRICS: BIC, AIC, ANNEALED IMPORTANCE SAMPLING, CONTINUATION METHOD FOR BRIDGE OF $Q(x)$

FOR CONJUGATE - EXPONENTIAL MODELS

- JOINT IS IN EXPONENTIAL FAMILY $p(x, y | \theta) = p(x, y | \eta) \exp[\phi(\theta)^T \psi(x, y)]$
 η = NO OF OBSERVATIONS, ψ = VALUE OF OBSERVATIONS
 - CONJUGATE PRIOR $p(\theta | \eta, \psi) = h(\eta, \psi) y(\theta)^\eta \exp[\phi(\theta)^T \psi]$
 η = NO OF OBSERVATIONS, ψ = VALUE OF OBSERVATIONS
 - IS E-M VANILLA IF $q(\theta) = \delta(\theta - \theta^*)$
- GENERAL VARIATIONAL FORM

 - E $Q_x^{(t+1)}(x) = p(x | y, \bar{\theta}^{(t)})$
 - M $Q_\theta^{(t+1)}(\theta) = \text{Exp} \left[\int Q_x^{(t+1)}(x) \cdot \log p(x, y, \theta) dx \right]$

VI - RECONSTRUCTION FORMULATION

- $\log p(y) = \int p(y|z) p(z) dz \xrightarrow{\text{JENSEN}} \log p(y) \geq \int q(z) \log \left(p(y|z) \frac{p(z)}{q(z)} \right) dz = \int q(z) \log p(y|z) - \int q(z) \log \frac{q(z)}{p(z)} = E_{q(z)} [\log p(y|z)] - \text{KL}(q(z) || p(z))$

RECONSTRUCTION COST
PENALTY
- ENCODER VIEW
 - ENCODER $q(z|y) \sim$ VARIATIONAL DISTRIBUTION
 - DECODER $p(y|z) \sim$ LIKELIHOOD
- $E[q(z)] [\log p(y|z)]$ DATA CODE LENGTH UNDER q
- $\text{KL}(q(z) || p(z)) =$ HYPOTHESIS CODE, PENALTY $= \Omega(z, y)$

VARIATIONAL INFERENCE

VARIATIONAL BAYES: ALL IS A DISTRIBUTION

- FREE-FORM: SOLVES FOR THE EXACT DISTRIBUTION $q(z)$
- FIXED FORM EXPLICIT FORM OF DISTRIBUTION $q(z) = f(z, \phi)$
- $\frac{\int F(y, q)}{\int q(z)} = 0$ ST $\int q(z) dz = 1 \rightarrow q(z) \propto p(z) \exp(\log p(y|z, \theta))$

VARIATIONAL EM

FOR $1 \dots N$ REPEAT UNTIL CONVERGES

- $E \phi \propto \nabla_{\phi} F(y, q) = \nabla_{\phi} E_{q(z)} [\log p(y|z)] - \nabla_{\phi} \text{KL}(q(z) || p(z))$ VARIATIONAL PARAMS
- $\theta \propto \nabla_{\theta} F(y, q) = \frac{1}{N} E_{q(z)} [\nabla_{\theta} \log p(y|z, \theta)]$ MODEL PARAMS

STOCHASTIC VARIATIONAL INFERENCE

LIKE EM, BUT ON DATA MINIBATCH

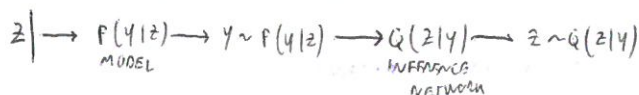
DOUBLY STOCHASTIC VI

LIKE EM BUT KL AM EXPECTED LL HARD TO COMPUTE / GRADIENT \rightarrow USE MONTECARLO

- $M \quad z_n \sim q(z_n | y_n), \quad \phi_n \propto \nabla_{\phi} \frac{1}{S} \sum_{s=1}^S \left[\log p(y_n | z_n(\phi)) - \log \frac{q(z_n(\phi) | y_n)}{p(z_n)} \right]$

AMORTIZED VI

INFERENCE NETWORK $\rightarrow q$ IS ENCODER / INVERSE MODEL, PARAMS OF q ARE GLOBAL ON TEST AND TRAIN POINTS. COST OF INFERENCE AMORTIZED ON ALL DATA. JOINT OPTIMIZATION OF MODEL AND VARIATIONAL PARAMS. NO ALLOCATION.



HOW TO GRADIENT

a. LOCAL VARIATIONAL METHODS / APPROXIMATIONS

b. STOCHASTIC BACKPROPAGATION [REPARAMETRIZATION TRICK]

$$\nabla_z E_{q(z)} [f(z)] \rightarrow z \sim N(\mu, \sigma^2), z = \mu + \sigma \epsilon, \epsilon \sim N(0, 1) \rightarrow \nabla_{\xi} E_{N(0,1)} [f(\mu + \sigma \epsilon)] \rightarrow E_{N(0,1)} [\nabla_{\xi} f(\mu + \sigma \epsilon)]$$

NO NEED FOR LOWER BOUNDS, LOW VARIANCE, CAN USE MANY DISTRIBUTIONS

c. MONTE CARLO VARIATE ESTIMATORS

SCORE FUN: $\nabla_{\xi} \log q_{\xi}(z|x) = \frac{\nabla_{\xi} q_{\xi}(z|x)}{q_{\xi}(z|x)}$

PROXIM $\nabla_{\phi} E_{q_{\phi}(z)} [\log p(y|z)]$

MCCV $E_{q_{\phi}(z)} [\log p(y|z) - c] \nabla_{\phi} \log q(z|y)$

\rightarrow CONTROL VARIABLE FOR ESTIMATION VARIANCE CONTROL

Stochastic Variational Inference (PAPER)

- NATURAL GRADIENT $\theta^{t+1} = \theta^t + \epsilon \nabla_{\theta} \ell(\theta)$ • IN DISTRIBUTION SPACE EVOLUTION DISTANCE OF PARAMS IS NOT REPRESENTATIVE OF DISTRIBUTION DISSIMILARITY
- SYMMETRIZED KL $D_{KL}(\theta, \theta') = E_{\theta} \left[\log \frac{q(y|\theta)}{q(y|\theta')} \right] + E_{\theta'} \left[\log \frac{q(y|\theta')}{q(y|\theta)} \right]$
- RIEMANNIAN METRIC TO MAP THE TWO SPACES \rightarrow LINEAR TRANSFORMATION OF θ MAKING DISTANCE $\theta =$ DISTANCE D_{KL} . $d\theta^T G(\theta) d\theta = D_{KL}^{sym}(\theta, \theta + d\theta)$
 $\rightarrow \nabla_{\theta}^T \ell(\theta) = G(\theta)^{-1} \nabla \ell(\theta)$, $G(\theta) =$ FISHER INFO OF $q(\theta) = E_{\theta} [\nabla_{\theta} \log q(y|\theta) (\nabla_{\theta} \log q(y|\theta))^T]$
- $D_{KL}(\theta, \theta + d\theta) = d\theta^T G(\theta) d\theta$ Φ IN EXP FAMILY: $G(\theta) = \nabla_{\theta}^2 \alpha_{\eta}(\theta)$ HESSIAN OF LOG-NORMALIZER FOR $[u]$ IS COV MATRIX OF SUFFICIENT STATISTICS $[u]$
- $\nabla_{\theta}^T L = E_{\phi} [\eta_y(x, z, a)] - \theta$ • CAN DO PARALLEL! • IS PROJECTIONS GRADIENT
- $\nabla_{\phi}^T L = E_{\phi} [\eta_l(x_N, z_N, y)] - \phi$ • η NATURAL PARAMETERS, COMPOSITE COMBINATOR

ALGORITHM!

- SAMPLE DATAPOINT x_i
- $\phi = E_{\theta^{t-1}} [\eta_y(x_i, z_i, y)]$ GLOBAL VAR PARAMS
- $\hat{\theta} = E_{\phi} [\eta_l(x_i, z_i, y)]$ INTERMEDIATE LOCAL PARAMS
- $\theta^t = (1 - \epsilon_t) \lambda^{t-1} + \epsilon_t \hat{\theta}$ FINAL LOCAL PARAMS

VARIANTS

- MINIBATCHES $\theta^t = (1 - \epsilon_t) \theta^{t-1} + \frac{\epsilon_t}{N} \sum \hat{\theta}_t$ - FOR EMPIRICAL BAYES