

APPROXIMATE INFERENCE

PROBLEM: COMPUTING P(h|v), IF MULTIPLE LAYERS EXACT INFERENCE REQUIRES EXPONENTIAL TIME, OBSTACLE IS LATENT INTERACTIONS IN GRAPHICAL MODEL OR EXPANDING AWAY EFFECTS

INFERENCE AS OPTIMIZATION

EXACT INFERENCE CAN BE INTERPRETED AS OPTIMIZATION PROBLEM. log p(v, theta) DIFFICULT -> LET'S GET A LOWER BOUND ON IT

- log p(v) = L(q) + KL(q||p) -> L(v; q, theta) = log p(v, theta) - D_KL(q(h)||p(h|v, theta))
- LET'S FIND q MAXIMIZING L. BUT q BUT IMPROVED OPTIMIZATION OR RESTRICTED q NO OPTIMAL OPTIMIZATION.

DECOMPOSE KL DO MATH

$L(v, \theta, q) = E_{h \sim q} [\log p(h, v)] + H(q)$

ENERGY ENTROPY

EXPECTATION-MAXIMIZATION

CAN BE FRAMED AS COORDINATE ASCENT MAXIMIZATION FOR L. E STEP WRT q, M STEP WRT theta.

- DUE EXACT INFERENCE BUT IS 'APPROXIMATE' M STEP PRODUCES 'GAP' BETWEEN L AND log p, BUT LATER E-STEP CLOSES IT FULLY.
- SOME TIMES WE JUMP TO EXACT SOLUTION IN 1 ITERATION

SPARSE CODING/MAP

FRAMING OF SPARSE CODING AS PROBABILISTIC MODEL. MAP(DUMMIE): $h^* = \text{ARGMAX}_h p(h|v)$ -> SEEN AS MAXIMIZATION OF L UNDER q WHERE q IS DIRAC DISTRIBUTION

- EXTRACTION OF CODES -> OPTIMIZATION OF W FOR OPTIMAL RECONSTRUCTION ERROR EQUIVALENT TO MAXIMIZING L WRT theta OBTAINED FROM MAP, EM w/ DIRAC POSTERIOR
- WE MAXIMIZE A BOUND ON TRUE LIKELIHOOD USING EXACT MAP INFERENCE

SEQUENCE MODELING

WITH GRAPHICAL MODELS. ALTERNATIVE TO RNN. IF MODELS HAVE MARKOV STRUCTURE -> STUFF IS EASY, EXACT INFERENCE. IF RECURSIVE -> APPROXIMATE

HMMs

- MARGINALIZATION SUM OVER ALL COMPLETE SOURCE TO SINK PATHS OF PRODUCT OF EXP SCORES $m(i) = \sum_{\text{PATH } a \in \text{PATH}} \prod e^a$ - ON TRELLIS DIAGRAM -> m IS LIKELIHOOD
- INFERENCE $\pi(i) = \text{ARGMAX}_{\text{PATH } b \in \text{PATH}} \sum_{a \in \text{PATH}} a$ MUST PRODUCE PATH. $v(i) = \text{MAX}_{\text{PATH } b \in \text{PATH}} \sum_{a \in \text{PATH}} a$ LOG-SCORE
- LET'S EXPLOIT DYNAMIC PROGRAMMING. MARGINALIZE VIA FWD-BWD:
$$\begin{cases} m(i) = \sum m(i^n) \\ m(i^n) = \sum_{n' \in \text{PRED}(i)} m(i^{n'}) e^{a_{n', i}} \end{cases}$$
- TRANSITION PROBABILITIES, EMISSION PROBABILITIES...
- OVERALL LIKELIHOOD $P(x_1 \dots x_N) = \sum_{s_1} \prod_{t=1}^N P(x_t | s_t) P(s_t | s_{t-1})$
- LOGSCORES: $a_{mn} = \log P(x_t | s_t=1) + \log P(s_t=1 | s_{t-1}=n)$ GRADIENT OPTIMIZATION WORKS BUT EM IS FASTER
- TOO MANY STATES IE N-GRAM MODELS -> VITERBI BREAKS DOWN, TOO EXPENSIVE -> DO BEAM SEARCH

MAP VIA VITERBI:

$$\begin{cases} v(i) = \max_{\text{NORMAL}} v(i^n) \\ v(i^n) = \max_{m \in \text{PRED}(i)} v(i^m) + a_{m, i} \end{cases}$$

$N^k \leftarrow \text{MAX}_{\text{FINAL}} v(i^k)$

$N^k \leftarrow \text{MAX}_{\text{PRE}} v(i^m) + a_{m, i}$

$\pi(i)$ RECONSTRUCT UP ARGMAXES

CAN FORMULATE DISCRIMINATIVE LIKELIHOOD TOO -> $P(y \dots | x \dots) = \frac{P(x \dots | y \dots) P(y \dots)}{P(x \dots)}$

CONDITIONAL RANDOM FIELDS

UNDIRECTED MODELS TRAINED TO MAXIMIZE JOINT P(y|x)

GRAPH TRANSFORMER: COMBINATIONAL MODEL TRANSFORMING A WEIGHTED DAG INTO ANOTHER. SET OF IN WEIGHTS -> SET OF OUT WEIGHTS. CRF ARE GRAPH TRANSFORMERS

- CRF IS MRF BUT WITH POTENTIALS ARE CONDITIONED, PARAMETERIZED, ON INPUTS
- 2 WOULD BE IMTRACTABLE, BUT MARKOV PROPERTY MAKES IT POSSIBLE TO USE DYNAMIC PROGRAMMING AND TRACK IT
- NO OF SUBSUMS FOR INFERENCE SCALES WITH DEGREE
- FOR MAP: SUM-PRODUCT -> MAX-SUM

NEURAL NETWORKS AND SEARCH COMBO

ANN + HMM/SEARCH IS OLD IDEA. SPEECH & HANDWRITING RECOGNITION. MARGINALIZATION AND INFERENCE (VIA DYNPROG) FOR TEMPORALLY STRUCTURED OUTPUTS CAN BE APPLIED ALSO WHEN LOGSCORES ARE LEARNED LINEAR FUNCTIONS (LIKE IN ANN) -> GRAPH TRANSFORMER

EXAMPLE: SEGMENTATION -> RECOGNITION -> SPEECH -> CONTEXT AWARENESS GRAPH TRANSFORMER NETWORK

BEAM SEARCH: IE WHEN NO OF NODES GROWS EXPONENTIALLY IN SEQUENCE LENGTH. VITERBI NO GOOD.

- BREAK NODES INTO GROUPS OF COMPATIBLE NODES -> SAME PATH LENGTH
- SEQUENTIALLY PROCESS GROUPS, KEEPING ONLY A GROUP SUBSET AT STEP t SUBSET IS BEAM. BASED ON S_{t-1}. EACH NODE HAS APPROX OF MAX TOTAL LOG-SCORE OF PATH
- S_t OBTAINED FROM FOLLOWING ARCS IN S_{t-1}, SOME RESULTING GROUPS WRT NEXT STEP ESTIMATION. KEEP THE BEST k.

BEAM OFTEN LACKS IN DIVERSITY. IS GROSSLY WRT t

VARIATIONAL INFERENCE

- THE COMPLICATED DISTRIBUTIONS ARE THE CONDITIONALS $p(v|h)$ POSTERIORS. CORRESPOND TO THE COMPLEX GRAPHS OVER THE HIDDEN UNITS \rightarrow BRUTE FORCE LOG
- BUT IF $p(v|h)$ AND $E[p]$ ARE MESSY UP • APPROXIMATE p WITH $q(h)$, SIMPLER. • IS MAXIMIZATION OF A BOUND \rightarrow VARIATIONAL ELIMINATION AND BOUND.

CALCULUS OF VARIATIONS

FUNCTIONAL: FCN OF FCN. MINIMIZATION OF FUNCTIONALS. FUNCTIONAL DERIVATIVES $\frac{\delta}{\delta f(x)}$ • $\frac{\delta}{\delta f(x)} \int g(f(x), x) dx = \frac{\partial}{\partial y} g(f(x), x)$

EXAMPLE: GAUSSIAN AS MAX-ENTROPY DISTRIBUTION

$H[p] = -E_x \log p(x) = -\int p(x) \log p(x) dx$ • WE NEED CONSTRAINTS TO FIX INTEGRAL TO 1, BOUND THE VARIANCE, SPECIFY SHIFT. LAGRANGIAN FUNCTIONAL.

$$L(p) = \lambda_1 \left(\int p(x) dx - 1 \right) + \lambda_2 (E[x] - \mu) + \lambda_3 (E[(x-\mu)^2] - \sigma^2) + H[p] = \dots$$

MINIMIZE FCN DERIVATIVE

$$\forall x \frac{\delta}{\delta p(x)} L = \lambda_1 + \lambda_2 x + \lambda_3 (x-\mu)^2 - 1 - \log p = 0$$

$$\rightarrow p(x) = \exp(-\lambda_1 - \lambda_2 x + \lambda_3 (x-\mu)^2 + 1) \xrightarrow{\text{SET CONSTRAINTS } \lambda} \lambda_1 = \log \sigma \sqrt{2\pi}, \lambda_2 = 0, \lambda_3 = 1/2\sigma^2 \rightarrow p(x) = N(x | \mu, \sigma^2)$$