

DEEP BOLTZMANN MACHINES

COMPLETELY UNDIRECTED. MULTIPLE HIDDEN LAYERS. IN EACH LAYER UNITS ARE MUTUALLY INDEPENDENT, CONDITIONED ON NEIGHBORING LAYERS. BINARY BUT VISIBLES MAY BE REALS

- $P(v, h^1, h^2, h^3) = \frac{1}{Z(\theta)} \exp(-E(v, h^1, h^2, h^3 | \theta))$ • $E(v, h^1, h^2, h^3 | \theta) = -v^T W^1 h^1 - h^1^T W^2 h^2 - h^2^T W^3 h^3$ • CONNECTIONS BETWEEN HIDDEN UNITS
- DBM LAYERS CAN BE ORGANIZED IN DIPARTITE GRAPHS. ODD/EVEN LAYERS. CONDITIONING ON ONE PARTITION MAKES UNITS OF THE OTHER BECOME CONDITIONALLY INDEPENDENT
- CONDITIONALS ON NEIGHBORING LAYERS ARE FULLY FACTORIAL → $P(h^1_i | v, h^2) = \text{SIGM}(v^T W^1_i + W^2_i h^2)$ $P(h^1 | v, h^2) = \prod P(h^1_i | v, h^2)$

PROPERTIES: • SIMPLER POSTERIOR THAN DBN → E1 OF UNITS WAS OTHERS ALWAYS. FIXED POINT OPTIMIZATION OF VARIATIONAL LOWER BOUND, PROPER MEAN FIBER

• INEFFICIENT SAMPLING BECAUSE MCMC AT EACH LAYER **WHAT DO?** GENERATIVE OR CLASSIFIER (+ MLP)

MEAN FIELD INFERENCE

POSTERIOR OVER ALL h IS COMPUTED $P(h | h^2 | v)$ BECAUSE INTERACTION WEIGHTS W^2 → MAKE h_1, h_2 MUTUALLY DEPENDENT GIVEN v

→ MEAN FIELD ON BECAUSE CONDITIONALS ARE FULLY FACTORIZED $q(h^1, h^2 | v) = \prod q(h^1_i | v) \prod q(h^2_i | v)$ FULLY FACTORED APPROXIMATION

→ MINIMIZE $KL(Q || P) = \sum_i Q(h^1, h^2 | v) \log \frac{Q(h^1, h^2 | v)}{P(h^1, h^2 | v)}$ • NOT NECESSARY TO SPECIFY FORM OF APPROXIMATION BUT HERE WE CAN AS A PRODUCT OF DETERMINISTICS

→ MINIMIZE THE VARIATIONAL / EVIDENCE LOWER BOUND (OR EQUIVALENTLY $KL(Q || P)$ EXPLICITLY) $L(q) = -\sum_i Q(h^1, h^2 | v) E(v, h^1, h^2 | \theta) - \log Z(\theta) + H(Q)$

→ SUBSTITUTE FORM OF Q → DERIVE WRT EACH h^i AND SET TO 0: $\begin{cases} \hat{h}^1_i = \text{SIGM}(\sum_j v_j W^1_{ji} + \sum_n W^2_{in} \hat{h}^2_n) \\ \hat{h}^2_n = \text{SIGM}(\sum_i W^2_{in} \hat{h}^1_i) \end{cases} \forall i, n$

→ ITERATE UPDATES UNTIL CONVERGENCE $\propto 10$ ITERATIONS

PARAMETER LEARNING

- EVALUATING PMF REQUIRES APPROXIMATE METHODS (AIS). TRAINING REQUIRES LOG PARTITION GRADIENT APPROXIMATIONS. • POSTERIOR ALSO INTERACTIVE → APPROX DURING SAMPLING TOO
- MF / CO WON'T WORK / OFFER SPEEDUPS → **USUALLY STOCHASTIC ML IS USED** NEGATIVE PHASE SAMPLES WITH Gibbs AUTOMATICALLY EVEN AM ODD LAYERS

→ **VARIATIONAL E-M** CAN BE SEEN AS

- **E:** OPTIMIZE $L(q, \theta)$ WRT VARIATIONAL PARAMS - EXPLORED ABOVE h^1, h^2
- **M:** OPTIMIZE $L(\theta, \theta)$ WRT θ → WE HAVE INTERACTIVE Z . WE ONLY TAKE SMALL STEP ALONG DIRECTION OF GRADIENT; NO $\nabla \theta = 0$ OFF THE BAT.

PRACTICALITIES: TRAINING DBM STARTING FROM RANDOM INIT IS USUALLY FAILURE. → FAIL TO REPRESENT DISTRIBUTION → NO BETTER LIKELIHOOD THAN SHALLOW RBM **NO IDEA WHY**

INITIALIZATION FROM PREBURNED CONFIGURATION DOES WORK INSTEAD. MAYBE BECAUSE LEARNING RATE FOR GRADIENT DESCENT DOESN'T PLAY WELL WITH NO OF GRADIENTS IN NEGATIVE PHASE. NOT GOOD MIXING. OR MAYBE BECAUSE ILL-CONDITIONED HESSIAN.

LAYERWISE PRETRAINING EACH LAYER TRAINED AS RBM IN ISOLATION. 7 RBM TRAINED ON SAMPLES DRAWN FROM $L-1$ POSTERIOR COMBINED INTO A DBM. PERSISTENT CD TRAINING SML TOO

BUT WATCH OUT WITH WEIGHTS $\times 2/\sqrt{2}$ WHEN RBM → DBM. ALSO REMOVE y IN MLP

JOINT TRAINING TO MAKE PERFORMANCE TRACKING EASIER. (HYPERPARAMS INFO)

• **(CENTRED) DBM:** REPARAMETRIZATION OF THE MODEL TO MAKE HESSIAN BETTER CONDITIONED → NO MORE PRETRAINING. GREAT FOR GENERATIVE USE CASES. NOT SO GREAT IN CLASSIFICATION

$$E(s, w; b) = -\frac{1}{2} s^T W s - s^T b \rightarrow E(s, w; b) = -\frac{1}{2} (s - \mu)^T W (s - \mu) - (s - \mu)^T b$$

W: WEIGHTS; b: BIAS; s: UNITS. $s - \mu \approx 0$ AT BEGINNING AFFECTS SML DYNAMICS

• **MULTI-PREDICTION DBM (MP-DBM)** VIEWES MEAN FIELD EQUATIONS AS DEFINING FAMILY OF RECURRENT NETS FOR APPROX SOLVING EVIDENCE INFERENCE PROBLEM

• TRAIN TO MAKE RMN OBTAIN ACCURATE ANSWERS TO INFERENCE PROBLEM. SAMPLE TRAIN EXAMPLE, SAMPLE INPUT SUBSET FOR INFERENCE NETWORK, TRAIN INFERENCE NETWORK TO PREDICT VALUES OF REMAINING UNITS. BALANCE THROUGH INFERENCE GRAPH. LOSS IS NOT LOWER BOUND BUT APPROX CONDITIONAL THAT INFERENCE NET IMPOSES ON MISSING VALUES. TRAINS MODELS AS THEY'RE USED. BETTER CLASSIFICATION PERFORMANCE. BETTER MEAN FIELD. LOSS EXACT GRADIENT, WRT SML IS APPROX GRADIENT

GROUPOUT BY MAKING UNITS INACTIVE/NOT GIVEN/ACTIVE BECOME TARGETS

BOLTZMANN MACHINES FOR REAL DATA (MISSING:

MCBPM, MPOT)

SPIKE-AM-SLAP RBM

- TWO SETS OF HIDDEN UNITS: SPIKE (BINARY) SLAP (REAL). MEAN OF $P(V|H) = (h\sigma_s)W^T$. SPIKE = DETERMINE COMPONENT PRESENCE SLAP = COMPONENT INTENSITY
- MODEL INPUT COVARIANCE. NO NEED FOR MATRIX INVERSION: CD AND PERSISTENT CD WORK.

ISSUE: SOME PARAM SETTINGS RESULT IN NON POSITIVE DEFINITE COV MATRIX, INTEGRAL DIVERGES. • EXTENDED TO CONVOLUTION, HIGHER ORDER INTERACTION AND AVG-POOLING OF LAY FEATURES

CONVOLUTIONAL BOLTZMANN MACHINES

- IN STANDARD CONVNETS POOLING REDUCES SPATIAL INPUT SIZE AT EACH LAYER, WHAT DO IN ENERGY MODEL? PROBABILISTIC MAX POOLING: CONSTRAIN DETECTOR UNITS SO ONLY 1 AS AT MOST ACTIVE AT A TIME. POOLING ON IFF DETECTOR IS ON, NO ACTIVES → ENERGY 0
- DRAWBACK: DETECTORS MUTUALLY EXCLUSIVE, NON OVERLAPPING POOLING
- CBM CAN DO IMAGE IMITATION BUT NO BETTER THAN SIMILAR MODELS FOR PRETRAINING

(MISSING: CONDITIONAL BM, RNN-RBM, DISCRIMINATIVE BM, ...)

DIRECTED GENERATIVE NETS

UNTIL RECENTLY IN DL ONLY UNDIRECTED GENERATIVE MODELS WERE POPULAR (RBM AND FRIENDS)

MISSING: DIFFERENTIABLE GENERATION NETS, VARIATIONAL AUTOENCODER, CONVOLUTIONAL GENERATIVE NETWORK

GENERATIVE ADVERSARIAL NETWORKS

DIFFERENTIABLE MAPPINGS FROM INPUT NOISE TO SAMPLES RESEMBLING DATA. SIMILAR TO VARIATIONAL AUTOENCODERS BUT DIFFERENT TRAINING PROCEDURE, NO INFERENCE NETWORK

- BASED ON GAME THEORY GENERATOR NETWORK: TRAINED TO MAP INPUT NOISE z TO SAMPLES x . $y(z)$ GENERATIVE MODEL. $P(z)$ NOT LEARNED.
- g DEFINES CONDITIONAL $P(x|z) = N(x|y(z), 1/\beta)$ DISCRIMINATOR NETWORK: ESTIMATE PROBABILITY OF x TO HAVE BEEN SAMPLED FROM THE DATA VS THE MODEL
- TRAINING: D MAXIMIZES, G MINIMIZES $V(x, d)$ VALUE FUN OF D BEING CORRECT $y^* = \arg \min_G \max_D V(y, d)$. $V(y, d) = E_{x \sim \text{DATA}} [\log d(x)] + E_{y \sim \text{MODEL}} [\log (1 - d(x))]$
- SIMPLY TRAIN V /BACKPROP, NO APPROXIMATE INFERENCE REQUIRED

* GAUSSIAN-BERNOULLI RBM

BINARY HIDDEN UNITS, REAL VISIBL UNITS. CONDITIONAL OVER VISIBL IS GAUSSIAN WITH MEAN IS FUNCTION OF HIDDEN. DIFFERENT PARAMETRIZATIONS. COVARIANCE/PRECISION FORM.

- $P(V|h) \propto N(V|Wh, \beta^{-1})$, $\log N(V|Wh, \beta^{-1}) = -\frac{1}{2}(V-Wh)^T \beta (V-Wh) + f(\beta)$ WE HAVE SOME CHOICE, NOT FULL COVARIANCE MATRIX BECAUSE ELSE WE USE LINEAR FACTOR MODEL, NO CONNECTION BETWEEN EDGES
- USUALLY DIAGONAL β WITH VARI. BIAS OR NOT

OTHER BOLTZMANN MACHINES

- MANY OTHER VARIANTS. • CONDITIONAL RBM WITH COV. TERMS IN ENERGY FUN. • RNN-RBM → RNN FEEDS RBM PARAMS AT EACH TIMESTEP
- OTHER TRAINING CRITERIONS IE TRAIN ON $\log P(y|v)$ • TRAIN WITH HIGHER (72) ORDER INTERACTIONS

*

SIGMOID BELIEF NETS

DBN WITH SPECIFIC TYPE OF CONDITIONAL VECTOR OF BINARY s . $P(s) = \text{SIGM}\left(\sum_i W_{si}s_i + b_s\right)$ s INFLUENCED BY ANCESTORS. MANY LAYERS. ANCESTRAL SAMPLING

SIMILAR TO DBN BUT SAMPLES ARE INITIALLY INDEPENDENT. SAMPLING THE VISIBL IS EASY. OTHER OPERATIONS NOT SO MUCH. INFERENCE INTRACTABLE

RECENT BREAKTHROUGHS: IMPORTANCE SAMPLING, WAVE-SURF, HELMHOLTZ MACHINES MADE IT FEASIBLE, GOOD PERFORMANCE. INFERENCE WITH SPECIAL LOCAL BOUND ON INFERENCE NETWORK. SBN WITHOUT LATENTS IS A TYPE OF AUTO-REGRESSIVE NET

DIFFERENTIABLE GENERATOR NETS

IS FAMILY OF VARIATIONAL AUTOENCODER. WE HAVE A GENERATOR NETWORK. TRANSFORM SAMPLES OF LATENTS z TO SAMPLES OF x OR DISTRIBUTION OF SAMPLES OVER x VIA DIFFERENTIABLE FUNCTION $g(z, \theta)$. THE GENERATOR NETWORK. • NET ARCHITECTURES GIVES FAMILY OF DISTRIBUTIONS TO SAMPLE FROM, PARAMS SELECT THE SPECIFIC DISTRIBUTION

- GENERATIVE MODELING IS (OBN) MORE DIFFICULT: BUT NONOBYVUS COMPUTATIONAL POWER IS SUFFICIENT.

→ DIFFICULTY IS IN TRAINING WHEN VALUES OF z FOR EACH x ARE NOT FIXED AND UNKNOWN REPRODUCING

*2 [NEXT PAGE]

AUTOREGRESSIVE NETWORKS

SIMILAR TO RECURRENT NETS BUT NO MORE PARAMETER SHARING ACROSS TIME, ELEMENTS NOT A FUNCTION EQUIVARIANT STRUCTURE, BUT ARBITRARY TUPLE.

LOGISTIC A/R NETS: NO HIDDEN UNITS, NO SHARING. $P(x_t | x_{1:t-1})$ PARAMETERIZED AS LOGISTIC REGRESSION, FIXED CAPACITY MODEL. CONTINUOUS VARIABLES \rightarrow LINEAR A/R MODEL.

NEURAL A/R NETS: MADE TO AVOID COO ISSUES IN NONPARAMETRIC GRAPHICAL MODELS. ESTIMATION OF CONDITIONAL PROBABILITIES W/O EXPONENTIAL NO. OF PARAMS LEFT-TO-RIGHT CONNECTIVITY. HIDDEN LAYERS FEATURES FOR x_t AND REUSED FOR x_{t+n} . ALL UNITS FOR t DEPEND ONLY FROM x_1, \dots, x_t . MULTITASK / TRANSFER LEARNING.

• NN OUTPUTS PREDICT PARAMETERS OF CONDITIONAL OF x_t

NADE: LINEAR NEURAL A/R NET BUT HAS WEIGHT SHARING. ALL $x_t \rightarrow h_t, i=1..T$ HAVE SAME WEIGHT MATRIX. RATIONALE IS SIMILARITY TO MEAN-FIELD INFERENCE IN RBM \rightarrow **RNADE:** FOR CONTINUOUS PROB DISTRIBUTIONS: MODELED AS GAUSSIAN MIXTURE, • GETTING RID OF CHOOSING ARBITRARILY ORDER: RANDOMLY SAMPLE ANY ORDER & TELL TO HIDDEN WHAT IS BEING OBSERVED, AND WHAT IS TO BE PREDICTED / MISSING \rightarrow ANY INFERENCE EFFICIENTLY.

ENSEMBLE ON ORDER: $P_{ENS}(x) = \frac{1}{M} \sum_{i=1}^M P(x|O_i)$. DOESN'T SCALE WELL FOR DEEP MODELS

AUTOENCODERS AS GENERATIVE MODELS

DECODING AE CAN BE SAMPLED FROM WITH MCMC FOR GAUSSIAN RBM. CONTRACTIVE AE ESTIMATE TANGENT MANIFOLD \rightarrow TO SAMPLE DU ENCODE / DECODE AND INJECT NOISE.

DAE MARKOV CHAIN: FROM PREVIOUS $x \rightarrow ((\hat{x}|x) \bullet h = f(\hat{x}) \bullet$ DECODE $w = g(h)$ OF $P(x|w = g(h)) = P(x|\hat{x}) \bullet$ SAMPLE NEXT STATE x FROM $P(x|w = g(h)) = P(x|\hat{x})$

• IF AE IS CONSISTENT ESTIMATOR WRT TRUE CONDITIONAL \rightarrow STATIONARY DISTRIBUTION OF SUCH MARKOV CHAIN IS A CONSISTENT ESTIMATOR FOR x

CONDITIONAL SAMPLING: CLAMP OBSERVED UNITS x_t , ONLY RESAMPLE FREE UNITS $x_{t+1}|x_t$ AND SAMPLE LATENT VARS

WALK BACK TRAINING: SPEEDS UP CONVERGENCE FOR GENERATIVE DAE **IDEA:** MULTIPLE ENCODE-DECODE STEPS FROM CHAIN INITIALIZED AT TRAINING EXAMPLE, WITH PENALTIES ON THE LAST RECONSTRUCTIONS. LIVE IN CONSTRAINED DIVERGENCE, BETTER AT REMOVING SPURIOUS MODES

GENERATIVE STOCHASTIC NETWORKS

ARE GENERALIZATIONS OF DAE, INCLUDE HIDDEN VARIABLES IN GENERATIVE MC, TWO CONDITIONALS: $P(x_u|h_u)$ RECONSTRUCTION DISTRIBUTION

• $P(h_u|h_{u-1}, x_{u-1})$ LATENT STATE UPDATE **GSN AM DAE** MODEL THE GENERATIVE PROCESS ITSELF, NOT THE JOINT DISTRIBUTION OF h AND v . IF IT EXISTS IT'S IMPLICIT AND ITS THE STATIONARY MC DISTRIBUTION

• TRAIN WITH RECONSTRUCTION LOG-PROBS ON VISIBLES, WALKBACK, ETC...

DISCRIMINANT GSN

LET'S USE GSN TO OPTIMIZE $P(y|x)$. BACKPROP LOG-PROBS OVER OUTPUT VARS, KEEPING INPUTS FIXED: STRUCTURED OUTPUT \rightarrow MODEL'S CHAIN OVER OUTPUT, INPUTS ARE CONDITIONAL

*2[PREVIOUS PAGE]

GENERATIVE ADVERSARIAL NETWORKS

IS DIFFERENT GENERATION NETWORK. GAME-THEORETIC SETTING, GEN NETWORK COMPETES AGAINST ADVERSARY \rightarrow **DISCRIMINATION NETWORK** HAS TO DISTINGUISH TRAINING SAMPLES FROM GENERATED ONE

USUALLY ZERO-SUM GAME WITH PAYOFF $V(\theta^g, \theta^d) = E_{x \sim P_{DATA}} \log d(x) + E_{x \sim P_{MODEL}} \log(1 - d(x))$. AT CONVERGENCE SAMPLES ARE INDISTINGUISHABLE

• DISCRIMINATION GETS $V(\theta^g, \theta^d)$, GENERATOR GETS $-V(\theta^g, \theta^d)$ • NO GUARANTEES FOR CONVERGENCE, SADDLE POINTS, PAYOFFS EQUIVARIANT MAY NOT BE MINIM FOR V

• CAN TRAIN CONDITIONAL GAN TO SAMPLE FROM $P(x|y) \rightarrow$ IMPROV • **UNUSUAL PROPERTY:** CAN FIT PROBS ASSIGNING ZERO MASS TO CERTAIN POINTS, BECAUSE IT DOESN'T MAXIMIZE LOG-PROBS OF POINTS BUT FITS A MANIFOLD

GENERATIVE MOMENT-MATCHING NETWORKS

ONLY GENERATOR NETWORK NO INFERENCE / DISCRIMINATION. MOMENT-MATCHING OF TRAINING SAMPLES STATISTICS AM GENERATED SAMPLES STATISTICS

• **COST FCN:** MAXIMUM MEAN DISCREPANCY. WHEEL TRUCK, SO PERMUT FCN, TO IMPLICITLY MATCH 1ST MOMENT IN D -DIM SPACE. WORKS ON BATCH OBVIOUSLY

CONVOLUTIONAL GENERATIVE NETWORKS

CONVOLUTIONAL STRUCTURE BUT ADD INFORMATION THROUGH LAYERS INSTEAD OF REMOVING IT. POOLING IS NOT INVARIANT \rightarrow SO INCREASE FEATURE MAPS SPATIAL SIZE

HOPFIELD NETWORKS

- FULLY CONNECTED ISING WITH SYMMETRIC WEIGHTS • PERFORMS AS ASSOCIATIVE MEMORY, CONTENT ADDRESSABLE • DENOISING / RECONSTRUCTION TASKS
- $S_i = \text{SIGN}\left(\sum_j W_{ij} S_j\right)$ THRESHOLD AT 0, $+1/-1$ ENCODING. $W_{ij} = \frac{1}{N} \sum_{n=1}^N S_i(n) S_j(n)$ ← TRAINING / UPDATE, SYNCH / ASYNCH. UNTIL CONVERGENCE
- INFERENCE: PRESENT INPUT UNTIL CONVERGENCE, READ S_i AS OUTPUT
- ENERGY FUNCTION: $H = -\frac{1}{2} \sum_i \sum_j W_{ij} S_i S_j = -\frac{1}{2} S^T W S$. DECODES 'AMOUNT OF MISMATCH' • $H(S) - H(S') = -\frac{1}{2} (S - S')^T \sum_j W_{ij} S_j$
- CAPACITY: $N \approx 0.138 d$ PATTERNS
- CONTINUOUS HOPFIELD NET: FOR IE GRAYSCALE IMAGES: $+1/-1$ OUTPUTS VIA SIGMOID, $P(x|W) = \frac{1}{Z(W)} \exp\left[\frac{1}{2} x^T W x\right]$ IS BOLTZMANN MACHINE

VARIATIONAL AUTOENCODER

MADE TO WORK ON: • INTERMEDIATE MANUAL LIKELIHOOD, POSTERIOR, INTERMEDIATE MP VARIATIONAL DATES • LARGE DATASET

ASSUMPTION: $p_\theta(z)$, $p_\theta(z|x)$ ARE PARAMETRIC AND DIFFERENTIABLE ALMOST EVERYWHERE WRT θ, z ,

WHAT IT SOLVES: 1- ML, MAP FOR θ 2- POSTERIOR INFERENCE OF $z|x, \theta$ 3- EFFICIENT MANUAL INFERENCE OF x

$q_\phi(z|x)$: VARIATIONAL APPROXIMATION, ENCODER z : CODE

$p_\theta(x|z)$: DECODER

$$= L(\theta, \phi, x) =$$

$$\text{LOWER BOUND: } \log p_\theta(x) \geq E_{q_\phi(z|x)} \left[\underbrace{-\log q_\phi(z|x)}_A + \log p_\theta(x|z) \right] = -KL(q_\phi(z|x) || p_\theta(z)) + \underbrace{E_{q_\phi(z|x)} [\log p_\theta(x|z)]}_B$$

• DIFFERENTIATE AND OPTIMIZE WRT $\theta, \phi \rightarrow$ MC ESTIMATION IS NO, HIGH VARIANCE

REPARAMETRIZATION!

$z \sim q_\phi(z|x) \rightarrow q_\phi(\epsilon, x)$ TRANSFORMATION + NOISE $\epsilon \sim p(\epsilon)$ MONTE-CARLO ESTIMATE OF $f(z)$ WAS $q_\phi(z|x) \rightarrow E_{q_\phi}[f(z)] = E_{p(\epsilon)}[f(q_\phi(\epsilon, x))] \approx \frac{1}{L} \sum_{\epsilon} f(q_\phi(\epsilon, x))$

SGVB ESTIMATOR: $\hat{L}^A(\theta, \phi, x) = \frac{1}{L} \sum_{\epsilon} \log p_\theta(x|z) - \log q_\phi(z|x)$, $z = q_\phi(\epsilon, x)$ a

$\hat{L}^B(\theta, \phi, x) = -KL(q_\phi(z|x) || p_\theta(z)) + \frac{1}{L} \sum_{\epsilon} (\log p_\theta(x|z))$ b WHEN KL IS COMPUTED ITSELF. KL: RECOMPUTED ON ϕ ; MAKES IT EASIER TO GRAD P

AEVB ALGORITHM:

- θ, ϕ INIT
- PICK MINIBATCH
- SAMPLE NOISE ϵ

E • $g = \nabla_{\theta, \phi} \hat{L}(\theta, \phi, x; \epsilon)$ GRADIENTS WRT MINIBATCH

M • UPDATE ϕ, θ WITH GRAD DESCENT NOW \rightarrow SGD, ADAM, ...

E_q - RECONSTRUCTION ERROR

NOISE 'INJECTED' IN SAMPLES FROM APPROXIMATE POSTERIOR

REPARAMETRIZATION TRICK!

• IS ALTERNATE SAMPLE GENERATION METHOD FOR $q(z|x)$ • MAKES MONTECARLO EXPECTATION DIFFERENTIABLE WRT ϕ

TYPICALLY $z \sim N(\mu, \sigma^2) \rightarrow z = \mu + \sigma \epsilon$, $\epsilon \sim N(0, 1)$ OTHERS: INVERSE CDF, CONJUGATE PRIORS, COMPOSITES

VARIATIONAL AUTOENCODER

• THE ENCODER IS A NEURAL NETWORK

- $p_\theta(z) = N(z; 0, I)$ ISOTROPIC MVN - $p_\theta(x|z)$: MVN OR BERNULLI, THEIR PARAMS COMPUTED FROM z WITH A MLP NETWORK

- APPROX $\log q_\phi(z|x) = \log N(z; \mu, \sigma^2 I)$ MVN DIAGONAL COVARIANCES, μ, σ ARE OUTPUTS OF MLP, AND VARIATIONAL ϕ

$$- L(\theta, \phi, x) = \frac{1}{2} \sum_i (1 + \log((\sigma_i)^2) - \mu_i^2 - \sigma_i^2) + \frac{1}{L} \sum_{\epsilon} \log p_\theta(x|z)$$

RELATED TO: WAKEFIELD, STOCHASTIC VARIATIONAL INFERENCE, PROBABILISTIC MODELS ARE FUNDAMENTALLY AUTOENCODERS

STOCHASTIC BACKPROPAGATION

DEEP LATENT GAUSSIAN MODELS: GENERAL CLASS OF DEEP MODELS, GAUSSIAN LATENTS AT EACH LAYER

GENERATIVE USE: TOP \rightarrow BOTTOM, PERMITS LAYERS ABOVE WITH GAUSSIAN NOISE \Rightarrow GENERATE INPUT DATA PARAMS θ , $P(\theta) = N(\theta|0, KI)$ GAUSSIAN PRIOR

$$\begin{cases} P(v|h) = P(v|h, \theta) P(h|\theta) \prod_{i=1}^{L-1} P_i(h_i|h_{i-1}, \theta) & \bullet \xi(0,1) \text{ IS IID MVN} \\ P(v|\xi(0,1)) = P(v|h_1(\xi_0, \dots, \xi_{L-1}), \theta) P(\theta) \prod_{i=1}^L N(\xi_i|0, I) \end{cases}$$

GAUSSIAN BACKPROP:

$$\begin{cases} \nabla_M E_{N(M,C)}[\ell(\xi)] = E_{N(M,C)}[\nabla_{\xi} \ell(\xi)] & \bullet \text{ BECAUSE EXP FAMILY} & \theta_g = M, R \\ \nabla_C E_{N(M,C)}[\ell(\xi)] = \frac{1}{2} E_{N(M,C)}[\nabla_{\xi\xi}^2 \ell(\xi)] & \text{FULL: } \nabla_{\theta} E_{N(M,C)}[\ell(\xi)] = E_{N(M,C)}\left[y^T \frac{\partial M}{\partial \theta} + \frac{1}{2} \text{TR}\left(H \frac{\partial C}{\partial \theta}\right)\right] & \text{IS 2ND ORDER} \end{cases}$$

WITH REPARAM TRICK: $\nabla_{\theta} E_{N(M,C)}[\ell(\xi)] = \nabla_{\theta} E_{N(\mu, \Sigma)}[\ell(\mu + R\epsilon)] = E_{N(0,1)}[\epsilon y^T]$, $y = \nabla \ell$ AT $\mu + R\epsilon$, $y = \mu + R\epsilon$, $R = RR^T$

INFERENCE

LOWER BOUND: $L(V) = -\log P(V) \leq \mathbb{H}(Q(\xi) || P(\xi)) - E_Q[\log P(V|\xi, \theta)]$, $Q(\xi|V)$ GAUSSIAN FACTORIZED ACROSS LAYERS (NOT WITHIN LAYER)
V INPUT DATA

• RECOGNITION DIFFERENT FROM GENERATIVE MODEL

$$Q(\xi|V, \theta^R) = \prod_{i=1}^N \prod_{l=1}^L N(\xi_{n,l} | \mu_L(V_n), \Sigma_L(V_n)) \quad \begin{matrix} \text{RECOGNITION MODEL} \\ \mu, C \text{ FROM NEURAL} \\ \text{NETS} = \theta^R \end{matrix}$$

GRADIENTS

$$\begin{cases} \nabla_{\theta} F(V) = -E_Q[\nabla_{\theta} \log P(V|h)] + \frac{1}{N} \theta^y & \bullet \text{ GOOD DESCENT WITH RMSPROP} \\ \nabla_{\theta^R} F(V) = \nabla_M F(V)^T \frac{\partial M}{\partial \theta^R} + \text{TR}\left[\nabla_R F(V) \frac{\partial R}{\partial \theta^R}\right] \end{cases}$$