# DEEP LEARNING

- IS STUFF WITH $\geq 1$ HIDDEN LAYERS

## DEEP GENERATIVE MODELS

USED UNSUPERVISED, TO LABEL DATA FOR FURTHER STAGES.

**DIRECTED:** DEEP DIRECTED NETWORKS. ALL NODES BINARY + ALL CPD LOGISTIC $\longrightarrow$ SIGMOID BELIEF NETWORK INFERENCE INTRACTABLE. POSTERIOR OF HIDDENS IS CORRELATED. FAST MEAN FIELD IS INACCURATE. MCMC SLOW BECAUSE CORRELATIONS

**UNDIRECTED:** STACK RBM ON TOP OF EACHOTHER. DEEP BOLTZMANN MACHINE. CAN DO EFFICIENT BLOCK (LAYER) GIBBS OR MEAN FIELD. (CI ABOVE | BELOW) NODES
DIFFICULT TO TRAIN DUE TO PARTITION. DO IT GREEDILY.

**MIXED:** DEEP BELIEF NETWORKS DIRECTED EXCEPT AT THE TOP. TOP ACTS AS ASSOCIATIVE MEMORY, REST GENERATES OUTPUT. PRIOR MULT LIKELIHOOD = FACTORED POSTERIOR
WE CAN INFER POSTERIOR EXACTLY LIKE IN RBM. FULLY FACTORIZED. $P(h_i | w_i, v_i)$. THIS ONLY HAPPENS IF PRIOR $P(h_i | w_i)$ IS COMPLEMENTARY
RBM ACTS AS SUCH. IF MORE LEVELS OR WEIGHTS NOT TIED $\longrightarrow$ NOT EXACT ANYMORE BUT VARIATIONAL LOWER BOUND. APPROXIMATE INFERENCE
CANNOT INTO TOP-DOWN INFERENCE, ONLY FEEDFORWARD

## GREEDY DBN TRAINING

- FIT A RBM
- UNROLL INTO DBN WITH TWO HIDDENS, UNTIE WEIGHTS
- FIT A 2ND RBM WITH ACTIVATION OF $h_1$ HIDDEN UNITS AS INPUT. $\longrightarrow$ BETTER PRIOR FOR $P(h_1 | w_1)$
- RINSE, REPEAT
- REFINE WEIGHTS WITH BACKFITTING: DO UPWARDS SAMPLING PASS TO TOP. DO GIBBS IN TOP RBM / CD UPDATE. DOWNWARDS ANCESTRAL SAMPLING [APPROX POSTERIOR]
  UPDATE LOGISTIC CPD PARAMS

## DEEP NEURAL NETWORKS

- IS MULTILAYER PERCEPTRON, FEEDFORWARD. • DIFFICULT TO TRAIN BECAUSE VANISHING GRADIENT AND MANY PLATEAUS. GPU + 2ND ORDER BACKPROP METHODS ARE USED.
  $\longrightarrow$ IDEA: GENERATIVE PRE-TRAINING: INIT PARAMS WITH UNSUPERVISED LEARNING. MODEL LEARNS TO MODEL ITS OWN INPUT FEATURE VECTOR, BUT TOUGH STILL
  'DATA INDUCED REGULARIZER', HELPS BACKPROP A LOT, FIND GOOD GENERALIZING LOCAL MINIMA.

- **DEEP AUTOENCODER**
  UNSUPERVISED ANN USED FOR DIM. REDUCTION AND FEATURE DISCOVERY. TRAINED TO PREDICT INPUT ITSELF. BOTTLENECK IN HIDDEN LAYERS TO PREVENT LEARNING
  IDENTITY FUNCTION. • LINEAR, SHALLOW AUTOENCODERS ARE EQUIVALENT TO PCA, SAME FIRST k PRINCIPAL COMPONENTS
- DIRECT BACKPROP TRAINING DOES NOT WORK WELL BECAUSE V.G. $\longrightarrow$ TRAIN SOME RBM, USE THEIR WEIGHTS TO INIT AUTOENCODER, FINETUNE W/ BACKPROP

- **STACKED DENOISING AUTOENCODERS**
  NO BOTTLENECK $\longrightarrow$ OTHER TRICKS TO PREVENT LEARNING IDENTITY. • IMPOSE SPARSITY CONSTRAINTS ON HIDDEN ACTIVATIONS
  $\longrightarrow$ ADD NOISE TO INPUT. SIMILAR TO APPROX ML TRAINING. • CAN BE STACKED, FINETUNING W/ BACKPROP. LIKE ANN

## APPLICATIONS

- HANDWRITTEN DIGIT CLASSIFICATION. MNIST. DBN. SOFTMAX CLASSIF. SEMINAL RESULT.

- VISUALIZATION, FEATURE DISCOVERY WITH DEEP AUTOENCODERS. 2D BOTTLENECK. SEMANTIC TOPIC ANALYSIS. NO LABELS BUT MORE HUMAN-FRIENDLY RESULTS THAN LDA/LSA

- SEMANTIC HASHING: BINARY, LOW DIM REPRESENTATION IN AUTOENCODER BOTTLENECK. USE LEARNED REPRESENTATION AS HASH KEYS. WIN.

- 1D CONVNETS - AUDIO: CONVOLUTION AUTOMATICALLY RESULTS IN PARAMETER TYING. MAXPOOLING: LOCAL MAX OVER FILTERED RESPONSE. FOR INVARIANCE AND SPEEDUP.
  NOISY-OR CPD TO ALLOW BACKWARDS INFO FLOW.

- 2D CONVNETS - IMAGES: STRAIGHTFORWARD EXTENSION. HIERARCHICAL FEATURES. WIN. SPLIT UP R,G,B CHANNELS.

# RESTRICTED BOLTZMANN MACHINES

IS UGM. • PAIRWISE MRF WITH HIDDEN AND VISIBLE NODES, → INFERENCE INTRACTABLE • **RESTRICTION:** NO CONNECTIONS BETWEEN NODE IN SAME LAYER

• $P(h|v,\theta) = \frac{1}{Z(\theta)} \prod^R \prod^H \psi_{RH}(V_R, H_H)$   • IS SPECIAL CASE OF **PRODUCT OF EXPERTS**, WORKS BETTER BECAUSE YIELDS SHARPER DISTRIBUTIONS,

CONSTRAINTS SATISFIED MORE EASILY. ADDING EXPERTS INSTEAD ONLY MAKES IT BROADER.

• **DISTRIBUTED ENCODING**, MANY UNITS GENERATE OUTPUT. VS LOCALIST ENCODING.   • **MAIN DIFF.** HIDDEN VARS ARE CI GIVEN VISIBLES. POSTERIOR FACTORIZES

RBM VS       $P(h|V,\theta) = \prod_H P(h_H|V,\theta)$ → EACH $h_H$ IN PARALLEL LIKE ANN

2 LAYER DGM

## TYPES OF RBM: DIFFERENT PAIRWISE POTENTIAL FCNS

– **BINARY** BINARY HIDDEN / BINARY VISIBLES. JOINT: $P(V,H|\theta) = \frac{1}{Z(\theta)} EXP(-E(V,H|\theta))$ POSTERIOR: $P(h|v,\theta) = \prod^H BER(h_H|SIGM(w^Tv))$ VARS/DATA

   • $E[h|V,\theta] = SIGM(w^Tv)$   • **W** GENERATIVE WEIGHTS      DATA#VARS   $P(v|h,\theta) = \prod_R BER(v_R|SIGM(w^Th))$

   • $E[v|h,\theta] = SIGM(Wh)$   • **$W^T$** RECOGNITION WEIGHTS   • HIDDEN NODE H ACTIVATES PROPORTIONALLY TO HOW MUCH V LOOKS LIKE $W_H$ → FF-ANN LIKE BEHAVIOR

– **GAUSSIAN CATEGORICAL:** USES 1-OF-C ENCODING. C NO. OF STATES FOR EACH V IR $P(v_R|h,\theta) = CAT(\cdots)$ ; $P(h_H=1|v,\theta) = SIGM(\cdots)$

– **GAUSSIAN:** HANDLES REAL-VALUED DATA. $P(v_R|h,\theta) = N(v_R|B_R + \sum_H w_{RH} h_H, 1)$ , $P(h_H=1,\theta) = SIGM(\cdots)$

– **HIDDEN GAUSSIANS:** LATENT GAUSSIAN + VISIBLE GAUSSIAN → UNDIRECTED FACTOR ANALYSIS. SAME AS DIRECTED.

   LATENT GAUSSIAN + CAT OBSERVED → UNDIRECTED CAT PCA. MEH PERFORMANCE.

## LEARNING

USE A SGD METHOD. AND $\ell_2$ REGULARIZATION. HERE WEIGHT DECAY. **GRADIENT:** $\frac{\partial}{\partial w_{RH}} \ell(\theta) = \frac{1}{N} E[v_R h_H|V,\theta] - E[v_R h_H|\theta]$   • MODEL XPECTATIONS

• **EXPECTATIONS:** • DO BLOCK GIBBS SAMPLING ON JOINT $P(v|h,\theta)$. WAIT FOR CHAIN BURN-IN.

CLAMPED PHASE / UNCLAMPED PHASE — EMPIRICAL.     EMPIRICAL XPECTATIONS = LEARNING STOPS

   • MEAN FIELD WORKS REALLY BAD.

   • **CONTRASTIVE DIVERGENCE:** $\nabla_W \ell \propto E[vh^T|v_i] - E_\alpha[vh^T]$   SIMILAR TO SGD?

   ⌐→ N UP-DOWN GIBBS PASSES INITIALIZED AT DATA VECTOR.

   FANTASY DATA: V' IS ATTEMPT AT RECONSTRUCTION → **AUTOENCODER-LIKE BEHAVIOR**

   • **PERSISTENT CD:** LIKE STOCHASTIC MAXIMUM LIKELIHOOD. INIT WEIGHTS. INIT CHAINS. MF UPDATES. MCMC UPDATES. PARAM UPDATES.

      ⌐→ v   ⌐→ h

## APPLICATIONS

• LANGUAGE MODELING INSTEAD OF LDA. SOFTMAX AT THE END.  UNIGRAM MODEL < LDA < RBM.

• COLLABORATIVE FILTERING

• BUILDING BLOCK FOR DEEP MODELS

## OBS! WAKE-SLEEP ALGORITHM

• IS GENERAL APPROACH USED BY ~~CONTRASTIVE DIVERGENCE~~. USABLE ON OTHER GENERATIVE MODELS.

• GENERATIVE + RECOGNITION WEIGHTS

• $\log P(D|G) \geq \log P(D|G) - KL(Q(h|D,R)||P(h|D,G)) =$ DECREASING FREE ENERGY INCREASES LOWER BOUND → INCREASES LL → **WAKE PHASE**  TRAIN $G$;

   $= -F(D,R,G)$ FREE ENERGY

• USE **REVERSE KL** FOR SLEEP PHASE: $KL(P(h|D,G)||Q(h|D,g))$ → SAMPLE FROM HIDDENS, UPDATE $R_s$

• FREE ENERGY $\propto$ SCORE  $F(D,R,G) = \#\left(-\sum_h Q(h|D,R) \log P(h|D,G)\right) - \left(-\sum_h Q(h|D,R) \log Q(h|D,R)\right)$