

AUTO ENCODERS

UNSUPERVISED METHODS LEARNING NONLINEAR PARAMETRIC DATA TRANSFORMATION  $x \rightarrow h$ , AS WELL AS  $h \rightarrow x$ . RELATED TO LINEAR FACTOR MODELS, IE PCA IS LINEAR AUTOENCODER

• AUTOENCODERS ARE NEURAL NETS TRYING TO COPY THEIR INPUT TO THEIR OUTPUT  
• GENERAL STRUCTURE: INPUT  $x$ , ENCODER  $f$ , CODE/REPRESENTATION  $h = f(x)$ , DECODER  $g$ , OUTPUT/RECONSTRUCTION  $\hat{x} = g(f(x))$ , LOSS  $\mathcal{L} = \mathcal{L}(x, \hat{x})$

ISSUE: HOW TO MAKE THE AE NOT LEARN THE IDENTITY FCN?

• HISTORICALLY: MAKE THE DIMENSION OF  $h$  SMALLER THAN THAT OF  $x$ , BOTTLENECK. MORE RECENTLY:  $h$  EVEN LARGER THAN  $x$  (OVERCOMPLETE) BUT OTHER RECONSTRUCTIONS IN  $\rightarrow$  BECAUSE COMPLEX DISTRIBUTIONS REQUIRE LARGE MODEL CAPACITY / POWER

- SPARSITY

$|h| > |x|$  BUT MOST ELEMENTS ARE AT OR CLOSE TO 0  $\left( \frac{\partial h_i}{\partial x} \approx 0 \right)$ . MOTIVATED BY MANIFOLD LEARNING

• SPARSE CODING: UNSUPERVISED LEARNING / FEATURE INFERENCE. ACTUALLY A LFM  $\rightarrow$  ENCODER IS NOT PARAMETRIC, BUT CODE IS OBTAINED THROUGH ITERATIVE OPTIMIZATION. RELATED ALSO TO GRAPHICAL MODELS.  $\hat{h} = f(x) = \underset{h}{\text{ARGMIN}} \mathcal{L}(g(h), x) + \lambda \Omega(h)$ .  $L_1$  PENALTY OR OTHERS.

• PREDICTIVE SPARSE DECOMPOSITION: COMBINES OPTIMIZED + PARAMETRIC ENCODER

• SPARSE AUTOENCODERS STANDARD AUTOENCODER + SPARSITY PENALTY.  $L_1$ , STUDENT-T PENALTY  $\sum \log(1 + d^2 h^2)$ ,  $W-L$  PENALTY  $-\sum (t \cdot \log h_i + (1-t) \log(1-h_i))$   $t$  IS DIST

• CONTRACTIVE AUTOENCODERS PENALTY  $\left\| \frac{\partial h}{\partial x} \right\|_F^2$ ; SUM OF SQUARED NORMS OF  $\frac{\partial h_i}{\partial x}$ . ENCOURAGES CONTRACTIVE (SMALL IN ALL DIRECTIONS) MAPPING

- ROBUSTNESS

INJECT NOISE IN INPUT OR HIDDEN LAYERS, ASK TO RECONSTRUCT CLEAN INPUT. • CONNECTION BETWEEN CONTRACTIVE AND DENOISING AUTOENCODERS.  $\rightarrow$  SINCE  $x$  AND  $x + \epsilon$  ARE TO YIELD SAME OUTPUT  $\rightarrow$  INSENSITIVITY TO CHANGES IN ALL DIRECTIONS  $\in$ ; SMALL DERIVATIVES.

• REGULARIZED REPRESENTATIONS

CAN BE UNDERSTOOD AS PUTTING A PRIOR ON  $h \rightarrow -\log P(h)$ . IS A DATA-DEPENDENT PRIOR, ON ACTUAL SAMPLED VALUES FROM POSTERIOR/ENCODER, IN TURN BASED ON INPUT  $x$ . INDIRECT PARAM REGULARIZATION

• DEPTH

BOTH ENCODER AND DECODER BENEFIT FROM THE ADDITION OF DEPTH. UNIV. APPROX. THEOREM. MORE LAYERS  $\rightarrow$  MORE CONSTRAINTS. DEPTH ALSO REDUCES COST OF EVALUATION / TRAINING. BETTER COMPRESSION RATES, TRAINING: GROSSLY PRETRAINING OF DEEP AE BY TRAINING STACK OF SHALLOW MODELS

• RECONSTRUCTION

RAW RECONSTRUCTION LOSS MAY NOT BE APPROPRIATE, IE  $x$  IS DISCRETE OR NOT WELL APPROXED BY GAUSSIAN  $\rightarrow$  DEFINE  $\mathcal{L}$  AS NEGATIVE LOG LIKELIHOOD OVER TARGETS RANDOM VARS.  $\rightarrow$  OUTPUT IS NOW PROBABILITY OF RECONSTRUCTION  $P(x|h)$ ; MODELING UNCERTAINTY AS WELL AS EXPECTATION

•  $\rightarrow$  ENCODING DISTRIBUTION  $Q(h|x) \rightarrow$  DECODING DISTRIBUTION  $P(x|h)$  •  $h$  IS NOW A LATENT VARIABLE

• IN SIMPLE CASES (GAUSSIAN, BERNULLI)  $P$  FACTORIZES:  $P(x|h) = \prod P(x_i|h)$  • RBM CAN BE SEEN AS AUTOENCODER WHERE  $P(x|h) = Q(h|x)$ . UNLIKE JOINT WITH  $P, Q$  AS CONDITIONALS, SPECIAL CASE

LINEAR FACTOR MODELS

CAN BE SEEN AS AUTOENCODERS WHERE  $h$  GENERATES  $x$  VIA LINEAR TRANSFORMATION •  $h \sim P(h)$ ,  $x = Wh + b + \epsilon$ .  $\epsilon$  NOISE GAUSSIAN, DIAGONAL. DIFFERENT LFM HAVE DIFFERENT FORMS FOR PRIOR AND NOISE. FACTOR ANALYSIS: PRIOR IS  $h \sim N(0, I)$  UNIT VARIANCE MVN,  $x_i$  ASSUMED CONDITIONALLY INDEPENDENT GIVEN  $h$ , NOISE COMING FROM DIAGONAL COV MATRIX  $\rightarrow x$  IS STILL A MVN  $x \sim N(b, WW^T + \Psi)$  PPCA: CONDITIONAL VARIABLES NOW EQUAL  $x \sim N(b, WW^T + \sigma^2 I)$ , COVARIANCE IS NOW MOSTLY CAPTURED BY  $h$ , ITSELF UP TO SMALL RESIDUAL RECONSTRUCTION ERROR  $\sigma^2$ .  $\sigma^2 \rightarrow 0$  WE GET NORMAL PCA. REDUCTION TO  $\perp$  PROJECTION

ICA: ALSO A LFM. NONGAUSSIAN PROJECTIONS / PRIORS, BUT STILL FACTORIZED  $P(h) = \prod P(h_i)$ . BECAUSE IF GAUSSIAN WE CAN'T IDENTIFY / RECOVER / DENTANGLE FACTORS

DOMINANT FORM OF NONGAUSSIANITY IN REAL DATA IS DUE TO SPARSITY.  $P(h)$  NONPARAMETRIC

SPARSE CODING:

LATENT VAR PRIOR IS PARAMETRIC, STILL NONGAUSSIAN • FACTORIZED UNIFORM  $P(h) = \prod \frac{1}{2} e^{-|h_i|}$  • FACTORIZED T-STUDENT  $P(h) \propto \prod \frac{1}{1 + \frac{h_i^2}{2}}$  AND GAUSSIAN RECONSTRUCTION NOISE. NORMALLY DONE WITH MAP INFERENCE  $\rightarrow$  INDUCE SPARSITY AND IMPROVEMENT OF FACTORS • ENCODER IS APPROX INFERENCE ALGORITHM. \*

VLL LOSSES

RECONSTRUCTION ERROR IS  $\mathcal{L} = -\log P(x|h) \rightarrow$  OBJECTIVE  $\mathcal{L} = -\log P(x|g(f(x)))$  TELLS US WHICH LOSS FUN DEPENDS ON INPUT TYPE. • REAL, UNBOUNDED  $\rightarrow$  BOUND FROM,  $P(x|h)$  IS

• BIT VECTOR  $\rightarrow$  CROSS ENTROPY LOSS  $P(x|h) = \prod P(x_i|h)$  BERNULLI. DECODER TRAINING  $\rightarrow$  ESTIMATING  $P(x|h)$ , AND WE CAN TALK ABOUT IMPLICIT OR EXPLICIT  $P(x)$  NORMAL

\* SPARSE CODING PROS:

MINIMIZES REC. ERROR AND LOG PRIOR BETTER THAN PARAMETRIC. DOES EXPLAINING AWAY, INHIBITING HIDDEN FACTORS

SPARSE CODING CONS:

INFERENCE ON  $h|x$  CAN BE LONGER THAN PARAMETRIC, RESULTING ENCODER COULD BE NONSMOOTH AND/OR NONLINEAR, MAKES IT FOR DIFFICULT GENERALIZATION

## SPARSE AUTOENCODERS

- WRT SPARSE CODING, SPARSE AUTOENCODERS HAVE EXPLICIT PARAMETRIC ENCODER. •  $L = -\log P(x|g(h)) + \Omega(h)$  SPARSITY PUNISH/REGULARIZER. LAPLACE PRIOR  $\rightarrow$  L1 PENALTY
- $-\log P(h) = \sum_i \log \lambda/2 + \lambda|h_i| = \text{CONST} + \Omega(h)$  • STUDENT-T PRIOR  $\rightarrow \Omega(h) = \sum_i \frac{v+1}{2} \log \left(1 + \frac{h_i^2}{v}\right)$  • INITIALLY SPARSITY WAS CONNECTED TO ENERGY PARTITION FUNCTION GRADIENTS, CONSIDERING RECONSTRUCTION ERROR A PROXY FOR ENERGY  $\rightarrow$  REGULARIZER PREVENTS IT FROM BRING LOW EVERYWHERE, JUST ON TRAINING SET
- LATER, OTHER NON-PROBABILISTIC SPARSITY PENALTIES SUCCESSFULLY INTRODUCED, IE CROSS-ENTROPIES ON ULS BETWEEN BERNULLI  $p=h_i$  AND BERNULLI  $p=(0,0.5)$  AND WEAKY THRESHOLD
- ACTUAL ZEROS IN  $h \rightarrow$  USE RELU UNITS AS ENCODER OUTPUTS, PRIOR PUSHED TO 0  $\rightarrow$  WE CAN CONTROL HOW MANY ZS
- REGULARIZER CORRESPONDS TO LOG PRIOR ON LATENT VARS | REPRESENTATION, NOT THE DATA, PREFERENCE OVER FUNCTIONS OF DATA, NOT OF THE PARAMETERS DATA-DEPENDENT REGULARIZER

## PREDICTIVE SPARSE DECOMPOSITION

- COMBINES SPARSE CODING AND PARAMETRIC ENCODER. REPRESENTATION IS FREE VARIABLE. •  $L = \text{ARGMIN}_h (\|x - g(h)\|^2 + \lambda|h|_1 + \gamma\|h - f(x)\|^2)$
- SPARSE CODING + CRITERION OPTIMIZING SPARSE  $h$  (AFTER INFERENCE) TO BE CLOSE TO ENCODER OUTPUT  $f(x)$  • ITERATIVE OPTIMIZATION PERFORMED ON ENCODER OUTPUT  $\rightarrow$  FEW ITERATIONS
- 2  $g, f$  UPDATED TOWARDS MINIMIZATION. • VARIATIONAL/EM-LIKE APPROACH GRADIENT DESCENT  $\rightarrow$   $h$  OPTIMIZED
- $f$  IS PARAMETRIC APPROXIMATION TO NON-PARAMETRIC SPARSE CODING 'ENCODER' • ITERATIVE OPTIMIZATION ONLY USED AT TRAIN-TIME. LEARNED  $f$  CAN BE USED FOR INITIALIZATION OF DEEP MODELS
- CAN STACK GREEDILY

## DENOISING AUTOENCODERS

- NO EXPLICIT CONSTRAINTS ON DIMENSION OR SPARSITY OF  $h$ .  $|h|$  HERE NOT RESTRICTED TO  $|x|$ . • CORRUPTION PROCESS  $C(\tilde{x}|x)$ , DAE LEARNS  $P(x|\tilde{x})$  RECONSTRUCTION DISTRIBUTION
- WE CAN PERFORM GRADIENT-BASED APPROX. MINIMIZATION ON  $KL - \log(P(x|h))$  WITH BACKPROP. • DAE DOES SGD ON  $-E_{x \sim Q(x)} E_{\tilde{x} \sim C(\tilde{x}|x)} \log P(x|g(f(\tilde{x})))$
- DAE TRAINING MAKES THEM LEARN A VECTOR FIELD  $(g(f(x)) - x)$  WHICH ESTIMATES SCORE, GRADIENT FIELD  $\frac{\partial \log Q(x)}{\partial x} \rightarrow$  DENOISING IS SCORE MATCHING REGULARIZED
- CONNECTION BETWEEN RDM AND AUTOENCODERS • SCORE MATCHING PROVEN FOR GAUSSIAN CORRUPTION AND RECONSTRUCTION DISTRIBUTIONS L CONSISTENT ESTIMATOR
- $\frac{g(f(x)) - x}{\sigma^2}$  IS CONSISTENT ESTIMATOR FOR  $\frac{\partial \log Q(x)}{\partial x}$ , GRADIENT OF ENERGY FUNCTION/LOG-DENSITY, LEADS TO POINT TOWARDS LOWER ENERGY/HIGHER PROB. REGIONS
- RECONSTRUCTION ERROR NORM COULD ALSO BE HIGH WHILE PROBABILITY IS LOW. WATCH OUT.

## CONTRACTIVE AUTOENCODERS

- INTRODUCE AN EXPLICIT REGULARIZER ON CODE  $h=f(x)$ , ENCOURAGING DERIVATIVES TO BE AS SMALL AS POSSIBLE.  $\Omega(f) = \left\| \frac{\partial f(x)}{\partial x} \right\|_F^2$  SQUARED FROBENIUS NORM (SUM SQ. ELEMENTS)
- OF ENCODER JACOBIAN MATRIX • CAE LEARNS TO CONTRACT ENCODER. OPPOSED TO RECONSTRUCTION ERROR  $\rightarrow$  DERIVATIVES TINY EXCEPT WHERE NOISY TO RECONSTRUCT TRAINING EXAMPLES  $\rightarrow$  DIRECTIONS TANGENT TO MANIFOLD WHERE DATA IS CONCENTRATED. PROMOTES STRONG INVARIANCE TO DIRECTION ORTHOGONAL TO MANIFOLD, CHECK SINGULAR VALUE SPECTRUM OF JACOBIAN. • CAE CONCENTRATE/COMPRESS REPRESENTATION SENSITIVITY IN VERY LITTLE DIMS, BETTER THAN VARIATIONAL AE  $\rightarrow$  SINGULAR VALUES OF INPUT-CODING MAPPING

- ISSUES: - DV FOR SHALLOW, EXPENSIVE FOR DEEP STACKS  $\rightarrow$  PREFER IN SINGLE CAES AND STACK. - CONTRACTION PENALTY ON  $f$  USELESS IF  $g$  COMPENSATES/EXPANDS  $\rightarrow$  TIE THE WEIGHTS  $\rightarrow W_g = W_f^T$  ON SOMETHING