

LATENT LINEAR MODELS

• LIMITATION OF MIXTURE MODELS \rightarrow CATEGORICAL HIDDEN VARIABLES

• IN LLM HIDDEN VARS ARE REAL VALUES $z_i \in \mathbb{R}^L$

FACTOR ANALYSIS

• PRIOR $P(z_i) = N(z_i | \mu_0, \Sigma_0)$ GAUSSIAN

• LIKELIHOOD ALSO GAUSSIAN $P(x_i | z_i) = N(Wz_i + \mu, \Psi)$

- W $D \times L$ FACTOR LOADING MATRIX

- Ψ $D \times D$ COVARIANCE MATRIX, DIAGONAL BECAUSE WE WANT TO EXPLAIN CORRELATION AND NOT BAKE IT IN

- IF $\Psi = \sigma^2 I \rightarrow$ PPCA

• CAN BE SEEN AS WAY TO SPECIFY JOINT DENSITY WITH LOW NUMBER OF PARAMS; $\text{COV}[x|\theta] = WW^T + \Psi \approx \text{COV}[x]$
 $O(LD)$ PARAMS VS $O(D^2)$ FOR FULL GAUSSIAN AND $O(D)$ FOR DIAGONAL COVARIANCE

• HOPEFULLY LATENT FACTOR z REVEAL INTERESTING STUFF \rightarrow INFERENCE. POSTERIOR OVER z 'S. $P(z_i | x_i, \theta) = N(z_i | m_i, \Sigma)$

LATENT FACTORS $m_i = \Sigma_i (W^T \Psi^{-1} (x_i - \mu) + \Sigma_0^{-1} \mu_0)$

• LINEAR COMBINATIONS / PROJECTIONS OF ORIG. FEATURES IN LOWER DIMENSIONAL SPACE

• LATENT FACTORS ARE UNIDENTIFIABLE; CANNOT UNIQUELY IDENTIFY W TOO MANY DOF. $D + LD + L(L-1)/2 \leq D(D+1)/2$

- BUT STILL ROTATIONAL AMBIGUITY.

FIXES:

$$L_{\max} = L D + 0.5(1 - \sqrt{1 + 8D})$$

- W FORCED ORTHONORMAL: COLUMNS ORDERED BY DECREASING VARIANCE \rightarrow PCA. NOT NEC MORE INTERPRETABLE BUT UNIQUE

- W FORCED LOWER TRIANGULAR! 1ST VISIBLE FEATURE \rightarrow 1ST LATENT FACTOR, 2ND V.F \rightarrow 1ST, 2ND LF, ...
 FIRST L VISIBLES AFFECT L.F. \rightarrow FORMER VARIABLES. CAREFUL.

- SPARSITY PROMOTING PRIORS ON W : L_1 REGULARIZATION \rightarrow SPARSE FACTOR ANALYSIS

- INFORMATIVE ROTATION MATRIX! FIND R THAT MODIFY W TO BE SPARSE \rightarrow VARIMAX

- NON-GAUSSIAN PRIORS FOR LF! CAN ACHIEVE IDENTIFIABILITY \rightarrow ICA

MIXTURES OF FA

APPROXIMATES DATA CURVED MANIFOLD W/ PIECEWISE LINEAR MANIFOLDS. LOW-RANK APPROX OF MIXTURE OF GAUSSIANS. $O(KLD)$ VS $O(KD^2)$ PARAMS REDUCES OVERFITTING.

• LATENT INDICATOR TO PICK THE MIXTURE COMPONENT TO USE. WITH $P(x_i | z_i, \alpha_i = u, \theta) = N(x_i | \mu_u + W_u z_i, \Psi)$

$$P(z_i | \theta) = N(z_i | 0, I) \quad P(u_i | \theta) = \text{CAT}(g_i | \pi)$$

EM FOR FA

• WADA STRAIGHTFORWARD E-STEP: $R_{ic} = P(q=c | x_i, \theta) = \pi N(x_i | \mu_c, W_c W_c^T + \Psi)$ M-STEP $\rightarrow \hat{W}_c, \hat{\Psi}, \hat{\pi}$

• MISSING DATA: USE MAP OR BAYESIAN INFERENCE

PRINCIPAL COMPONENT ANALYSIS

ORTHOGONALITY = $W^T W = I$

- FA WHERE $\Psi = \sigma^2 I$ AND W ORTHOGONAL $\therefore \sigma^2 = 0$ CLASSICAL PCA
- $\sigma^2 > 0$ PROBABILISTIC/SENSIBLE PCA

PCA

DEFINITION: WE WANT ORTHOGONAL SET OF L LINEAR BASIS VECTORS $w_i \in \mathbb{R}^D$ AND SCORES $z_i \in \mathbb{R}^L$ SO TO MINIMIZE RECONSTRUCTION ERROR $J(W, Z) = \frac{1}{N} \sum \|x_i - \hat{x}_i\|^2$ WHERE $\hat{x}_i = W z_i$ IS CONSTRAINED ORTHOGONAL

$\rightarrow J(W, Z) = \|X - W Z^T\|_F^2$, $Z = N \times L$ WITH z_i IN ROWS; $\|A\|_F = \sqrt{\sum \sum a_{ij}^2} = \sqrt{\text{TR}(A A^T)}$ FROBENIUS NORM

SOLUTION: $\hat{W}_L = V_L$, V_L CONTAINS L EIGENVECTORS WITH LARGEST EIGENVALUES OF EMPIRICAL COVARIANCE MATRIX $\hat{\Sigma} = \frac{1}{N} \sum x_i x_i^T$
OPTIMAL LOW-DIMENSIONAL DATA ENCODING IS $\hat{z}_i = W^T x_i$, ORTHOGONAL PROJECTION OF DATA ONTO EIGENVECTORS' COLSPACE

- PRINCIPAL COMPONENTS: PRINCIPAL DIRECTIONS ALONG WHICH DATA SHOWS MAXIMAL VARIANCE \rightarrow STANDARDIZE DATA TO AVOID ISSUES (DUE TO SCALE, ETC)
- PRINCIPAL COMPONENTS IS EIGENSTUFF (FACES/DIGITS), LIN-COMBS OF ORIGINAL DIMS \rightarrow USE CORRELATION MATRIX VS COVARIANCE MATRIX
- ALSO CENTER THE DATA

PROOF: DERIVATIVE OF J WRT $W \rightarrow$ OPTIMAL WEIGHTS ARE \perp PROJ ON PRINCIPAL DIRECTIONS \rightarrow ARG MIN $J(W) = \text{ARGMAX VAR}[Z]$
DERIVATIVE OF PROJECTION OF VARIANCE $= 0 \rightarrow \hat{\Sigma} w_i = \lambda_i w_i \rightarrow$ EIGENVECTOR OF COV. MATRIX \therefore VARIANCE = $w_i^T \hat{\Sigma} w_i = \lambda_i$

ALT! USING SVD \rightarrow OK FOR NON-SQUARE MATRICES, SINGULAR VECTORS = EIGENVECTORS OF $\hat{\Sigma}$
SVD $X = U S V^T$, $X^T X = V S^T U^T U S V^T = V S^2 V^T$, $(X^T X) V = V D \rightarrow$ RIGHT SINGULAR VECTORS OF X ARE TO λ_i OF $\hat{\Sigma}$
 $(X^T X) U = U D$

PPCA

LL MAXIMA: $\hat{W} = V (\Lambda - \sigma^2 I)^{1/2} R$, R ORTHOGONAL, V COL EIGENVECTORS, Λ DIAGONAL EIGENVALUES $R \rightarrow I$ $\hat{\sigma}^2 = \frac{1}{D-L} \sum_{i=1}^L \lambda_i$ AVG VARIANCE DISCARDED DIMENSIONS

- $\sigma^2 \rightarrow 0 \rightarrow$ PCA
- $\sigma^2 > 0 \rightarrow$ POSTERIOR MEAN $\hat{F} = \hat{W}^T \hat{W} + \sigma^2 I$ IS NOT ORTHOGONAL PROJECTION \rightarrow RECONSTRUCTION ERROR CLOSER TO DATA MEAN
- CAN BAYESIAN METHODS • OPTIMIZE PARAMS, INTEGRATE LATENTS • USEFUL WHEN $|D| \gg$

EM FOR PCA

HAVE SVD PERCEP • E STEP: $\bar{z} = (W^T W)^{-1} W^T X$ ANALOGY: E STEPS MOVES POINTS TO ORTHOGONAL, M STEP 'ROTATES ROD/LINE'
• M STEP: $W = \bar{X} \bar{z}^T (\bar{z} \bar{z}^T)^{-1} \rightarrow$ LINE LINEAR REGRESSION WITH EXPECTED LATENT VALUES INSTEAD OF OBSERVED INPUTS

- CONVERGES TO W IS SAME LINEAR SUBSPACE OF L EIGENVECTORS BUT HAS TO BE ORTHOGONALIZED AN ORDERED
- EM FASTER ESP $N, D \gg L$, DOMINATED BY E STEP, $O(TLND)$. EIGENVECTOR METHODS ARE $O(\min(ND^2, DN^2))$. EMROS METHOD COMPARABLE TO EM
- EM IS ONLINE-FRIENDLY
- EM CAN HANDLE MISSING DATA
- EM CAN BE USED FOR MIXTURE MODELS
- EM CAN BE UPGRADED TO VARIATIONAL METHODS FOR GREAT JUSTICE

FA/PPCA MODEL SELECTION

HOW TO PICK OPTIMAL NUMBER OF L

- SIMPLE: BIC, VARIATIONAL LOWER BOUNDS, CROSS VALIDATION \rightarrow BUT EXPENSIVE.

- USUALLY EXHAUSTIVE SEARCH OVER VALUES OF L

- AUTOMATIC RELEVANCY DETERMINATION (ARD) + EM TO PRUNE OUT IRRELEVANT WEIGHTS

- STILL SEARCH OVER $K \rightarrow$ BIRTH/DEATH MOVES, STOCHASTIC SAMPLING OF MODEL SPACES, GRASS SAMPLING + MCMC POORS

PCA MODEL SELECTION

NON PROBABILISTIC \rightarrow ABOVE METHODS NOT OK

- APPROX 4 RECONSTRUCTION ERROR: $E(D, L) = \frac{1}{|D|} \sum ||x_i - \hat{x}_i||^2$, $\hat{x}_i = Wz_i + \mu$, $z_i = W^T(x_i - \mu)$

- RESIDUAL ERROR WITH L TERMS: $E(\text{Dimin}_L) = \sum_{j=L+1}^D \lambda_j$, SUM OF DISCARDED EIGENVALUES

- \rightarrow CAN PLOT RETAINED EIGENVALUES / SCORE PLOT

FRACTION OF EXPLAINED VARIANCE:

$$\frac{\sum_{j=1}^L \lambda_j}{\sum_{j=1}^{L_{\max}} \lambda_j}$$

- OBS: ERROR KEEPS GOING DOWN ON TEST SET TOO \rightarrow PCA IS NOT GENERATIVE MODEL BUT COMPRESSION TECH!

MORE DIMENSIONS \rightarrow MORE ACCURATE APPROXIMATION NO MATTER WHAT

SOLUTION: PROFILE LIKELIHOOD $\ell(L) = \sum \log N(\lambda_u | \mu_1(L), \sigma^2(L)) + \sum \log N(\lambda_u | \mu_2(L), \sigma^2(L))$

~~PARITION~~ MODEL OF SIZE K , ERROR λ_u , ERROR SO THAT $\lambda_1 \geq \dots \geq \lambda_{L_{\max}}$. THRESHOLD L . PARTITION $u \leq L, u > L$

CHANGE POINT MODEL $\lambda_u \sim N(\mu_1, \sigma^2)$, $\lambda_u \sim N(\mu_2, \sigma^2)$, σ SAME. FIT FOR $L=1:L_{\max}$, MLE, POOLED VARIANCE ESTIMATE

$$L^* = \text{ARGMAX } \ell(L)$$

CATEGORICAL PCA

OBSERVED DATA IS CATEGORICAL. EACH y GENERATED FROM LATENT VAR $z, \in \mathbb{R}^L$ W/ GAUSSIAN PRIOR, SOFTMAXED.

- $P(z_i) = N(0, I)$

- FITTED WITH EM

- $P(y_i | z_i, \theta) = \prod_{r=1}^R \text{CAT}(y_{ir} | S(W_r^T z_i + W_{0r}))$

- USED TO VISUALIZE HIGH-DIM CATEGORICAL DATA

PCA FOR PAIRED / MULTI-VIEW DATA

WHEN RELATED DATASETS WE WANT TO COMBINE IN LOW-D FOR IE, PREDICTION. DATA FUSION

EASY TO DO PAIRS \rightarrow DATA SETS. FITTED USING EM OR GIBBS SAMPLING. DISCRETE/COUNT DATA \rightarrow EXP. FAMILY RESPONSE INSTEAD OF GAUSS \rightarrow APPROXIMATE INFERENCE (MCMC)

SUPERVISED PCA / LATENT FACTOR REGRESSION

LIKE PCA BUT y_i IS TAKEN INTO ACCOUNT WHEN LEARNING LOW DIM. EMBEDDING

$P(z_i) = N(0, I_0)$ • JOINTLY GAUSSIAN $y_i | x_i \sim N(x_i^T W_i \sigma_y^2 + W_i^T C W_i y_i)$, $W = \Psi^{-1} W_x C W_y$, $\Psi = \sigma_x^2 I_0$, $C^{-1} = I + W_x^T \Psi^{-1} W_x$
 $P(y_i | z_i) = N(W_y^T z_i + \mu_y, \sigma_y^2)$ • DEPENDENCE OF PRIOR FOR W ON X ARISES FROM W IS DERIVED FROM JOINT MODEL OF X, Y
 $P(x_i | z_i) = N(W_x z_i + \mu_x, \sigma_x^2 I_0)$

DISCRETE VARS \rightarrow GAUSS \rightarrow EXP FAMILY. NO MORE CLOSED-FORM CONDITIONALS. BUT STILL INFORMATION BOTTLENECK.

MULTILABEL: y_i IS VECTOR OF RESPONSES, COVAR. FUSION

FIND ENCODING DISTRIBUTION $P(z|x)$ SO TO

MINIMIZE $I(x, z) - \beta I(x, y)$; $\beta > 0$

DISCRIMINATIVE SUPERVISED PCA

DIFFERENT WEIGHTS TO INPUTS x_i AND OUTPUTS y_i WEIGHTED OBJECTIVE $\ell(\theta) = \prod P(y_i | \eta_{i,y})^{\alpha_y} P(x_i | \eta_{i,x})^{\alpha_x}$ $\eta_{i,m} = W_m z_i$

IN CASE OF GAUSSIAN α_m CONTROLS NOISE VARIANCE. HARD TO ESTIMATE α_m BECAUSE IT CHANGES LIKELIHOOD NORMALIZATION

PARTIAL LEAST SQUARES

MORE 'DISCRIMINATIVE' SUPERVISED PCA. ALLOWS SOME INPUT COVARIANCE TO BE EXPLAINED BY THEIR OWN SUBSPACE z_i^x WHILE REST OF I/O COV IS z_i^s

$P(z_i) = N(z_i^s | 0, I_{L_s}) N(z_i^x | 0, I_{L_x})$ • $P(v_i | \theta) = N(v_i | \mu, W W^T + \sigma^2 I)$ • L HAS TO BE PICKED 'LARGE ENOUGH'
 $P(y_i | z_i) = N(W_y z_i^s + \mu_y, \sigma_y^2 I_{D_y})$ • PROJECTS PREDICTED VARS y AND OBS x TO NEW SPACE
 $P(x_i | z_i) = N(W_x z_i^s + B_x z_i^x + \mu_x, \sigma_x^2 I_{D_x})$

CANONICAL CORRELATION ANALYSIS (CCA)

LIKE SYMMETRIC UNSUPERVISED VERSION OF PARTIAL LEAST SQUARES. EACH OBS HAS ITS OWN PRIVATE SUBSPACE + A SHARED ONE. z_i^x, z_i^y, z_i^s

$P(z_i) = N(z_i^s | 0, I_{L_s}) N(z_i^x | 0, I_{L_x}) N(z_i^y | 0, I_{L_y})$ • $P(v_i | \theta) = N(v_i | \mu, W W^T + \sigma^2 I_0)$ • CAN SPARSIFY VIA AND
 $P(x_i | z_i) = N(x_i | B_x z_i^s + W_x z_i^x + \mu_x, \sigma_x^2 I_{D_x})$ • CAN MLE VIA EM \rightarrow EQUIV TO NON PROBABILISTIC
 $P(y_i | z_i) = N(y_i | B_y z_i^s + W_y z_i^y + \mu_y, \sigma_y^2 I_{D_y})$ • CAN GENERALIZE TO $m > 2$ OBSERVED VARIABLES
 • CAN BAYESIAN INFERENCE
 • CAN CREATE MIXTURES OF CCA

FINDS LIN. COMB OF x, y OF MAX CORRELATION

INDEPENDENT COMPONENT ANALYSIS - ICA

BLIND SIGNAL / SOURCE SEPARATION. RECONSTITUTE ORIGINAL SIGNAL(S) WHERE MANY LATENT SOURCES GET LINEARLY MIXED TOGETHER.

• $x_t = Wz_t + \epsilon_t$; W IS $D \times L$ MIXING MATRIX, $\epsilon_t \sim N(0, \psi)$. EACH TIMEPOINT IS INDEPENDENT OBSERVATION. • WE WANT TO INFER $p(z_t | x_t, 0)$.

• $L = D$ (SOURCES = SENSORS) $\rightarrow W$ IS SQUARE • $|\psi| = 0$ FOR SIMPLICITY

PRIOR: ANY NON-GAUSSIAN $p(z_t) = \prod_1^L p_i(z_{t,i})$; CONSTRAINED VARIANCE TO 1. NO GAUSSIAN BECAUSE DOES NOT ALLOW UNIQUE RECOVERY OF SOURCES

• PCA LIKELIHOOD INVARIANT TO ORTHOGONAL TRANSFORMATION OF LIKELIHOOD. • PCA RECOVERS BEST LINEAR SUBSPACE, NOT THE SIGNALS.

• IN SYMMETRIC GAUSSIAN POSTERIOR NO INFORMATION TO TELL US ANGLE WE NEED TO ROTATE. • PCA WHITENS DATA, SOLVES HALF THE PROBLEM.

• ICA IDENTIFIES ROTATION. • ESTIMATE W AND p_i 'S

• SQUARE $W \rightarrow$ ON CLOSED FORM UNIQUE SIGNALS

• NON-SQUARE $W \rightarrow$ NO UNIQUE SIGNALS BUT WE DO POSTERIOR $p(z_t | x_t, \hat{W})$

ESTIMATE: MLE

FOR NOISE-FREE SQUARE W 'S, ALREADY WHITENED WITH PCA. BECAUSE NOISE-FREE, WHITENED AND CENTERED DATA.

• W MUST BE ORTHOGONAL BECAUSE $E[x x^T] = I$, AND $\text{COV}[x] = W W^T$

$\rightarrow D(D-1)/2$ PARAMS AS OPPOSED TO D^2

• W GENERATIVE WEIGHTS • V RECOGNITION WEIGHTS $= W^{-1}$

$$p(x) = p(z) = p(z_t) / |\det(W^{-1})| = p(z_t) | \det(V) | \rightarrow \text{LL: } \frac{1}{T} \log p(D|V) = \underbrace{\log |\det(V)|}_{\text{CONSTANT}} + \frac{1}{T} \sum_{t=1}^T \log p_i(V^T x_t) \rightarrow \text{NLL}(V) = \sum_{t=1}^T E[G_i(z_t)]$$

$z_t = V^T x_t, G_i(z) = -\log p_i(z)$

• MINIMIZE SO TO ROWS OF V ORTHOGONAL AND UNIT NORM \rightarrow ORTHOGONAL

- CAN GRADIENT DESCENT BUT SLOW

- FASTICA

- EM

FASTICA: FASTICA

APPROXIMATE NEWTON METHOD FOR FITTING ICA. • ASSUME 1 LATENT FACTOR AND ALL SOURCE DISTRIBUTIONS KNOWN ARE SAME $\rightarrow G = -\log p(z)$

• $g(z) = \frac{d}{dz} G(z)$, GRADIENT, HESSIAN WRT W • NEWTON UPDATE: $V^* = E[x g(V^T x)] - E[g'(V^T x)] V$ • IN PRACTICE WE REPLACE EXPECTATIONS WITH MC ESTIMATES FROM TRAINING SET

• MAXIMIZES NONGAUSSIANITY OF THE PROJECTION

\rightarrow PROJECT BACK ON CONSTRAINT

$$\text{SURFACE } V^{\text{NEW}} = \frac{V^*}{\|V^*\|}$$

• CONVERGENCE: $|V^T V^{\text{NEW}}| \rightarrow 1$

• NON-CONVEX OBJECTIVE; X PLOTTABLE TO LEARN MULTIPLE FEATURES, SEQUENTIALLY OR IN PARALLEL. (ON BECAUSE UNLIKE PCA 1ST FEATURE IS NOT MORE IMPORTANT THAN 2ND)

\downarrow 1D V_1 - PROJECT OUT

• IF $G(z)$ NOT KNOWN, WHAT DO? SIGNAL POWER

- NOT GAUSSIAN - STH LEPTOKURTIC 'SUPER-GAUSSIAN' BECAUSE BRAINZ

• LAPLACE $\log p(z) = -\sqrt{2}|z| - \log(\sqrt{2})$ MEANS VAR 1 | NOT DIFF. 1E

- NOT CRITICAL XACT SHAPE: $G(z) = \sqrt{2} \log \cosh(z)$

• LOGISTIC $\log p(z) = -2 \log \cosh\left(\frac{\pi}{2\sqrt{3}} z\right) - \log \frac{4\sqrt{3}}{\pi}$

ESTIMATE: EM

\uparrow MIXTURE OF GAUSSIANS

NO USE PARTICULAR FORM BUT FLEXY NONPARAMETRIC ESTIMATOR \rightarrow GMM! CAN DO VIA $E[z_t | x_t, \theta]$ VIA SUMMING OVER K^L COMPONENTS OF q_t FACTORS

THEN ESTIMATE ALL SOURCES IN PARALLEL VIA FITTING A GMM TO $E[z_t]$ \rightarrow MARGINALS $p_i(z_i) \rightarrow$ ICA FOR $W \rightarrow$ RINSE REPEAT

ESTIMATE: OTHERS

COOL BUT EQUIVALENT TO MLE

• **MAX NONGAUSSIANITY:** NEGENTROPY(z) $\approx H(N(\mu, \sigma^2)) - H(z)$. ASS GAUSSIAN IS MAXENTROPY, MAXIMIZING NEGENTROPY TAKES ME FAR.

$$J(V) = \sum \text{NEGENTROPY}(z_i) = \sum H(N(\mu_i, \sigma_i^2)) - H(z_i). \text{ IF } V \text{ ORTHOGONAL AND WHITENED DATA COV IS } I$$

$$J(V) = \sum E[\log p(z_i)] + \text{CONST} \rightarrow \text{EQUIV TO LL}$$

• **MINIMIZING MUTUAL INFORMATION:** MULTI-INFORMATION $I(z) = KL(p(z) || \prod p(z_i)) = \sum H(z_i) - H(z) \rightarrow I(z) = \sum H(z_i)$

\rightarrow SAME AS MAX. NEGENTROPY

• **MAXIMIZING MUTUAL INFORMATION:** INFO MAX PRINCIPLE: MAX INFORMATION FLOW THROUGH A SYSTEM

$$\text{MAXIMIZE MUTUAL INFO BETWEEN } Y \text{ INTERNAL REP AND } X \text{ INPUT. } I(x, y) = H(y) - H(y|x)$$

\rightarrow SAME AS ML

LOCALLY LINEAR EMBEDDING

• IDEA: MINIMIZE RECONSTRUCTION ERROR BY MAKING SMALL PATCHES WHERE THERE IS NONLINEARITY IN DATA.

• CHOOSE ON POINT NEIGHBOURHOOD, COMPUTE WEIGHTS FOR POINTS AS LINEARS OF THEIR NEIGHBOURS

• THEN FIND LOW-DIM EMBEDDING OF POINTS USING LINEARS. • UNUSUAL STUFF

• EVENS OF QUADRATIC FORM MATRIX

ISOMAP

VARIANT OF MULTIDIMENSIONAL SCALING. DIMENSIONALITY REDUCTION TECHNIQUE ATTEMPTING TO PRESERVE DISTANCES IN LOW-DIM SPACES. ASSUMES EUCLIDEAN/EUCLIDEAN SPACE.

ISOMAP: CALCULATE DISTANCES. FIND NEIGHBOURS (IE u_1, u_2, \dots), CONSTRUCT NEIGHBOURHOOD GRAPH. ESTIMATE GEODESICS BY FINDING SHORTEST PATHS (DIJIKSTRA, ...), APPLY MDS. THIS INSERTS NONLINEAR MANIFOLD STRUCTURE INTO THE PROBLEM.

• T-SNE