



FEATURE PROCESSING & LEARNING

Bor-shen Lin at NTUST

DIMENSION REDUCTION: WHY

- Complexity depends on the number of input dimension
 - Images/speeches/documents are high-dimensional
- Simpler models are more robust on small data set
- Fewer features are easy for explanation and knowledge discovery
- Few dimensions are easy to plot and visualize.



OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- Transformer / GAN
- t-SNE



VECTOR QUANTIZATION

- Similar to K-Means Clustering
 - Minimization of square error
- Build a *codebook* that contains *codewords*
 - Select representatives for the vector space
 - Might induce quantization errors
- Vectors with continuous variables can be converted into discrete symbols
 - CHMM → DHMM
- Applied to data compression or statistical modeling



OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- t-SNE



PRINCIPAL COMPONENTS ANALYSIS

- Also known as *Karhunen-Loeve transformation (KLT)*
- Projection of \mathbf{x} on the direction \mathbf{w} is

$$\mathbf{z} = \mathbf{w}^t \mathbf{x}, \quad \mu_z = \mathbf{w}^t \boldsymbol{\mu}_x$$

$$\text{var}(z) = E((\mathbf{z} - \boldsymbol{\mu}_z)^2) = E((\mathbf{w}^t \mathbf{x} - \boldsymbol{\mu}_x)^2) = \mathbf{w}^t \boldsymbol{\Sigma} \mathbf{w}.$$
- To maximize $\text{var}(z)$ for $|\mathbf{w}| = 1$
 Lagrangian function $J(\mathbf{w}) = \mathbf{w}^t \boldsymbol{\Sigma} \mathbf{w} + \alpha(\mathbf{w}^t \mathbf{w} - 1)$

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \frac{\partial J}{\partial \mathbf{w}} = 2\boldsymbol{\Sigma} \mathbf{w} - 2\alpha \mathbf{w} = 0, \quad \boldsymbol{\Sigma} \mathbf{w} = \alpha \mathbf{w}$$
- \mathbf{w}^* is the *eigenvector* for $\boldsymbol{\Sigma}$ with corresponding *eigenvalue* α
- $\text{var}(z) = \mathbf{w}^t \boldsymbol{\Sigma} \mathbf{w} = \mathbf{w}^t \alpha \mathbf{w} = \alpha |\mathbf{w}|^2 = \alpha$
 \rightarrow Eigenvalue α equals to *the variance of z*.
- Principal components: with largest eigenvalues (variances)
 \rightarrow reserve discriminant capability for reduced dimensions

REALIZATION OF PCA

- Given the samples $X = \{\mathbf{x}_i\}$

$$\boldsymbol{\mu} \equiv E[\mathbf{x}] \cong \frac{1}{n} \sum_i \mathbf{x}_i$$

$$\boldsymbol{\Sigma} \equiv E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t] \cong \frac{1}{n} \sum_i (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^t$$

$$D = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_n^2 \end{bmatrix}$$

- Covariance matrix $\boldsymbol{\Sigma}$ can be diagonalized through **Eigen-decomposition**

$$\boldsymbol{\Sigma} = P D P^t \text{ where } P P^t = P^t P = I$$

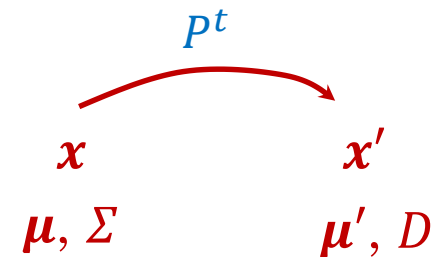
$$\begin{aligned} D &= P^t \boldsymbol{\Sigma} P = P^t E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t] P \\ &= E[P^t (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t P] \end{aligned}$$

- Let $\mathbf{x}' = P^t \mathbf{x}$ then

$$\boldsymbol{\mu}' = P^t \boldsymbol{\mu}$$

$$\boldsymbol{\Sigma}' = E[(\mathbf{x}' - \boldsymbol{\mu}')(\mathbf{x}' - \boldsymbol{\mu}')^t] = D$$

- Transform linearly the vectors to a vector space whose dimensions are mutually uncorrelated.

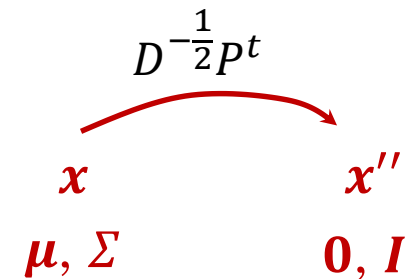


REALIZATION OF PCA

- Let $\mathbf{x}'' = D^{-\frac{1}{2}} P^t (\mathbf{x} - \boldsymbol{\mu})$ then
 $\boldsymbol{\mu}'' = 0$

$$\begin{aligned}\Sigma'' &= E \left[D^{-\frac{1}{2}} P^t (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^t P D^{-\frac{1}{2}} \right] \\ &= D^{-\frac{1}{2}} P^t E \left[(\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^t \right] P D^{-\frac{1}{2}} \\ &= D^{-\frac{1}{2}} P^t \Sigma P D^{-\frac{1}{2}} = D^{-\frac{1}{2}} D D^{-\frac{1}{2}} = I\end{aligned}$$

$$D = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_n^2 \end{bmatrix}$$



- Transformed linearly to a vector space of which each dimension has zero mean and unit variance.
- Pre-whitening* for images/speeches



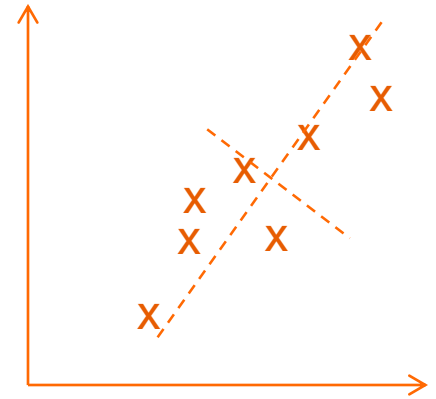
DISCUSSIONS ON PCA

- PCA: perform diagonalization & normalization on the covariance.
 - All dimensions become **uncorrelated** after being transformed
- $P^t \mathbf{x}$: representation of \mathbf{x} on eigenspace (orthogonal)
 - P is the rotary matrix that contains the eigenvectors and forms the basis of eigenspace
- $D^{-\frac{1}{2}}$: normalization by the standard deviations
 - containing the eigenvectors of Σ
 - D contains the eigenvalues of Σ
 - $\lambda_1 > \lambda_2 > \dots > 0$ (usually sorted in descending order)



DISCUSSIONS ON PCA

- Reduction of dimension
 - Eigenvalues are sorted and k largest eigenvalues with the corresponding eigenvectors are reserved (P_k)
 - Keep those dimensions with the largest variances (most relevant) that preserve a large portion total variance.
- All vectors can be projected onto the K-dimensional eigenspace
 - $\mathbf{x}' = P_k^t \mathbf{x}$ or $\mathbf{x}' = D_k^{-\frac{1}{2}} P_k^t \mathbf{x}$
- In eigenspace, distance or cosine similarity between any two vectors can be computed.



DISCUSSIONS ON PCA

- When $\mathbf{x}' = D^{-\frac{1}{2}} P^t \mathbf{x}$, all dimensions in eigenspace are uncorrelated.

- $$\begin{aligned} \langle \mathbf{x}', \mathbf{y}' \rangle &= \mathbf{x}'^t \mathbf{y}' = (D^{-\frac{1}{2}} P^t \mathbf{x})^t D^{-\frac{1}{2}} P^t \mathbf{y} \\ &= \mathbf{x}^t P D^{-\frac{1}{2}} D^{-\frac{1}{2}} P^t \mathbf{y} = \mathbf{x}^t P D^{-1} P^t \mathbf{y} \\ &= \mathbf{x}^t \Sigma^{-1} \mathbf{y} \neq \mathbf{x}^t \mathbf{y} \end{aligned}$$

$$d(\mathbf{x}', \mathbf{y}') \equiv \langle \mathbf{x}' - \mathbf{y}', \mathbf{x}' - \mathbf{y}' \rangle = (\mathbf{x} - \mathbf{y})^t \Sigma^{-1} (\mathbf{x} - \mathbf{y})$$

$$d(\mathbf{x}, \mathbf{y}) \equiv \langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle = (\mathbf{x} - \mathbf{y})^t (\mathbf{x} - \mathbf{y})$$

$$d(\mathbf{x}', \mathbf{y}') \neq d(\mathbf{x}, \mathbf{y})$$

- Definition of inner product in original space should consider covariance Σ with correlated dimensions.

- Gaussian: $p(\mathbf{x}) = \frac{1}{(2\pi|\Sigma|)^{\frac{D}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^t \Sigma^{-1} (\mathbf{x}-\mu)}$



PRINCIPAL COMPONENT ANALYSIS (PCA)

- $\mathbf{z} = D^{-\frac{1}{2}}P^t(\mathbf{x} - \boldsymbol{\mu}_x)$, then
 - $\boldsymbol{\mu}_z = 0, \Sigma_z = I$
 - All dimensions of variable z are independent statistically with zero mean and unit variance
- $\mathbf{x}' = D_k^{-\frac{1}{2}}P_k^t\mathbf{x}$ is a k -dimensional sub-space of the eigenspace (k eigenvectors with k eigenvalues)
 - The k dimensions with the highest variances (eigenvalues) are reserved while the others are discarded.
 - A high-dimensional vector can be reduced to a vector with dimension k . (e.g. image feature)
 - All dimensions become uncorrelated statistically after the conversion.

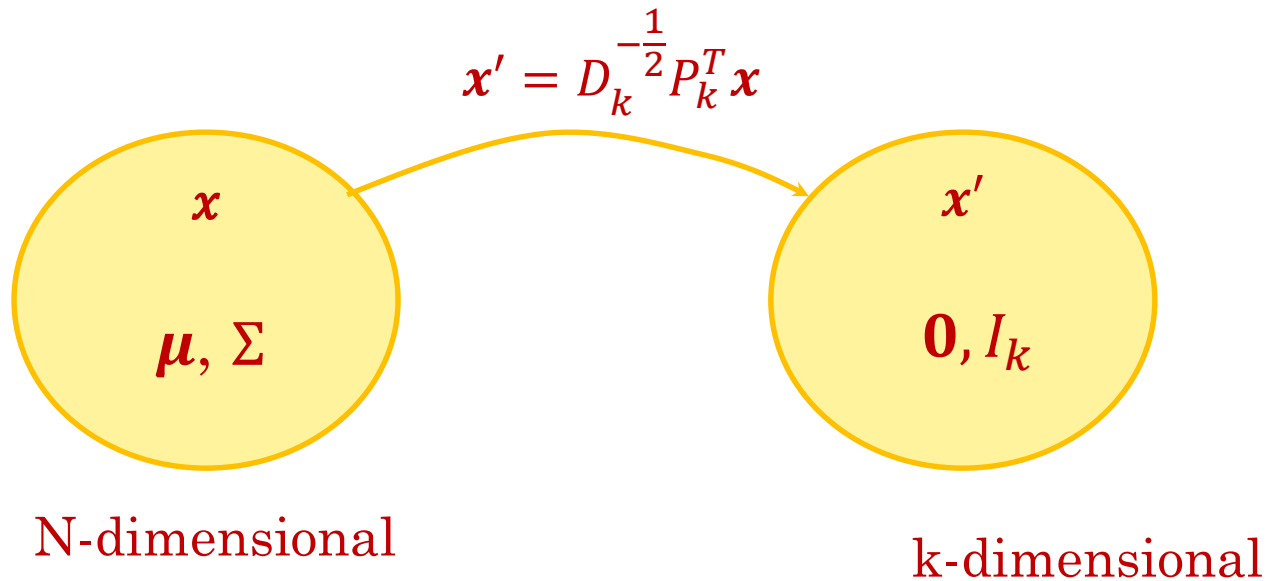


PRINCIPAL COMPONENT ANALYSIS

- $x' = D_k^{-\frac{1}{2}} P_k^t x$: transformed to a k-dimensional orthonormal space (covariance $\Sigma_k'^{-1} = I_k$)
 - $\langle x', y' \rangle = x'^t y' = x^t P_k D_k^{-1} P_k^t y = x^t \Sigma_k^{-1} y$.
- Similarity or distance measures in transformed space
 - Norm-2: $|x'|^2 = \langle x', x' \rangle = x^t \Sigma_k^{-1} x$
 - Distance: $d(x', y') = |x' - y'|^2 = (x - y)^t \Sigma_k^{-1} (x - y)$
 - $\cos(x', y') = \frac{\langle x', y' \rangle}{|x'| |y'|} = \frac{x^t \Sigma_k^{-1} y}{((x^t \Sigma_k^{-1} x)(y^t \Sigma_k^{-1} y))^{\frac{1}{2}}}$
- Note $\langle x, y \rangle \equiv x^t y$, then $\langle x', y' \rangle \neq \langle x, y \rangle$ even if $k=n$
 - $\langle x, y \rangle \equiv x^t y = \sum x_i y_i$ appropriate for orthonormal space
- If $\langle x, y \rangle$ is redefined as $x^t \Sigma^{-1} y$ (instead of $x^t y$),
 - $\langle x, y \rangle = x^t \Sigma^{-1} y \neq x^t y$ for non-orthonormal space (normalized by Σ^{-1})
 - $\langle x', y' \rangle = x'^t \Sigma_k'^{-1} y' = x'^t I_k y' = x'^t y'$ for orthonormal space



PCA FOR DIMENSION REDUCTION



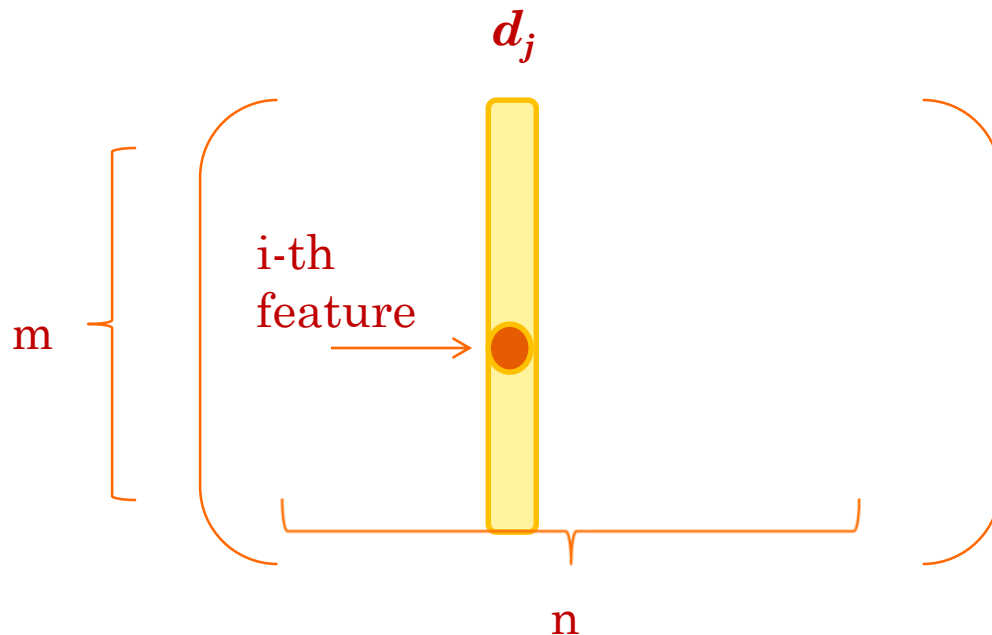
OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- t-SNE



LATENT SEMANTIC INDEXING (LSI)

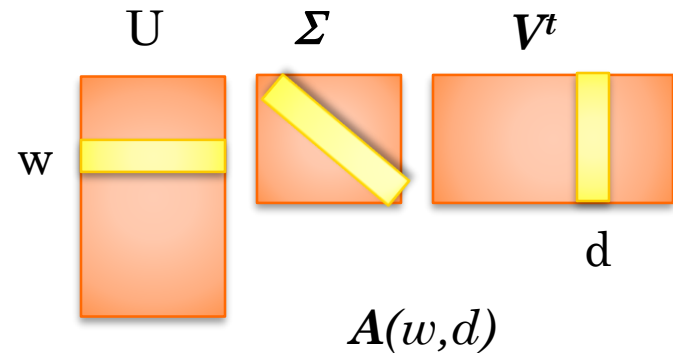
- Assume there are n vectors (e.g. documents) each containing m features (terms).
- The n vectors, $\{ \mathbf{d}_j \}$, can be represented as an $m \times n$ matrix, \mathbf{A} , which is called the *term-document matrix*.
- j -th column vector \mathbf{d}_j for matrix $\mathbf{A} \sim j$ -th document.



$$\mathbf{A}_{m \times n} = \mathbf{U}_{m \times r} \boldsymbol{\Sigma}_{r \times r} \mathbf{V}_{r \times n}^t$$



SINGULAR VALUE DECOMPOSITION (SVD)



- $A_{m \times n} = U_{m \times r} \Sigma_{r \times r} V_{r \times n}^t$
- U : right singular vectors
- V : left singular vectors
- $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$ where $\sigma_i \geq 0$
 σ_i : singular values (square root of eigen-values)
- Select the largest k among r singular values (corresponding to the highest variances)
- $A' = U_k \Sigma_k V_k^t$ (reduce to k -dimensional vector)
- Document vectors: in V_k^t
- Term vectors are in U_k
- Query vector: $q_k^t = q^t U_k^t \Sigma_k^{-1}$



- Convert document/term into eigenspace in which all dimensions are statistically independent.
 - In LSI, dimensions with largest variances (σ_i) will be kept while the rest dimensions are eliminated so as to reduce the computational complexity.
 - In eigenspace, terms and documents are represented in the same form. Therefore, the similarity between documents or that between terms can be computed.
 - In original document space, the cosine similarity between, say, *tour* and *travel* is 0. In eigenspace, their similarity might not be 0.
 - Search, clustering, classification, and GMM training can be performed in eigenspace.



APPLICABILITY OF LSI

- It is possible to defined **similarity between terms** (or abstract concepts)
 - $\text{similarity}(\text{tour} \sim \text{travel}) ? \text{similarity}(\text{tour}, \text{visa})?$
 - Computed not by *word senses*, but by the *usage contexts* (the documents) of the words
- Generalization
 - Matrix A is NOT limited to *term-document* relationships. It can be generalized for **arbitrary M-to-M mappings** i.e. **any tabular data in attribute-value form**
 - SVD may be performed, and so do search, clustering, classification, statistical models (e.g. GMM) etc.
 - *topic/author* \rightarrow similarity between authors/topics
user/items \rightarrow clustering products/users, recommendation
album/style \rightarrow distance between albums/styles

OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- t-SNE



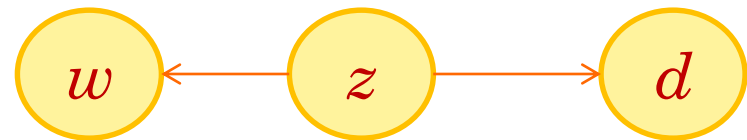
PROBABILISTIC LATENT SEMANTIC INDEXING (PLSI)

- d document, w word, z latent variable

$$d \in D = \{d_1, d_2, \dots, d_N\}$$

$$w \in W = \{w_1, w_2, \dots, w_M\}$$

$$z \in Z = \{z_1, z_2, \dots, z_K\}$$



- To maximize $L = \sum_{d,w} p(d,w) \log(p(d,w))$

based on Expectation Maximization

$$\text{E-step: } p(z|d,w) = \frac{p(z,d,w)}{p(d,w)} = \frac{p(z)p(w|z)p(d|z)}{\sum_{z'} p(z')p(w|z')p(d|z')}$$

$$p(d,w) = \sum_{z'} p(z')p(w|z')p(d|z') \rightarrow L$$

$$\text{M-step: } p'(w|z) = \frac{p(z,w)}{p(z)} = \frac{\sum_d n(d,w)p(z|d,w)}{\sum_{d,w'} n(d,w')p(z|d,w')}$$

$$p'(d|z) = \frac{p(z,d)}{p(z)} = \frac{\sum_w n(d,w)p(z|d,w)}{\sum_{d',w} n(d',w)p(z|d',w)}$$

$$p'(z) = \frac{\sum_{d,w} n(d,w)p(z|d,w)}{\sum_{d,w} n(d,w)}$$



TRAINING PROCEDURE

- Parameters for estimation
 - $p(z) \sim \Sigma_{k,k} = p(z_k)$ $K \times K$ diagonal matrix
 - $p(w | z) \sim U_{k,j} = p(w_j | z_k)$ $M \times K$ matrix
 - $p(d | z) \sim V_{j,k} = p(d_j | z_k)$ $N \times K$ matrix
- Training
 - Random initial values
 - Reestimation of parameters based on EM algorithm
- Estimation of $p(d | z)$ for new document d
 - Compute $p(d, w, z)$ for d with $p(z)$, $p(w | z)$ fixed
 - $p(z, d) = \sum_w p(z, w, d)$, $p(d | z) = p(z, d) / p(z) \rightarrow$ iteration
- After converging, compute $p(z | d) = p(z, d) / \sum_z p(z', d)$
 - $p(z_k | d)$'s form a K -dimensional vector (summed to 1)



DISCUSSIONS ON PLSI

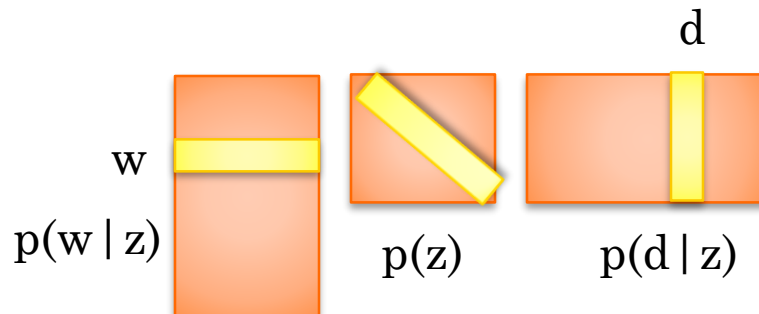
Comparison with LSI

$$p(d, w) = \sum_z p(z)p(w | z)p(d | z)$$

$U\Sigma V^t$ is a $M \times N$ matrix \sim term-doc matrix

Interpretation

- U (for $p(w | z)$) and V (for $p(d | z)$) correspond to the eigenvectors by SVD in LSI
- Σ (for $p(z)$) corresponds to the singular values in LSI



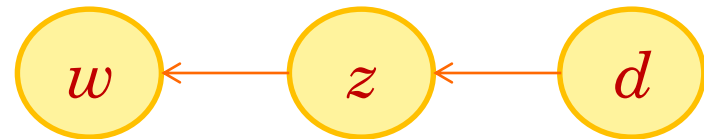
PLSI – DIFFERENT FORMULATION

- document, w word, z latent variable

$$d \in D = \{d_1, d_2, \dots, d_N\}$$

$$w \in W = \{w_1, w_2, \dots, w_M\}$$

$$z \in Z = \{z_1, z_2, \dots, z_K\}$$



- To maximize $L = \sum_{d,w,z} p(d, w, z) \log(p(d, w, z))$
based on Expectation Maximization

$$\text{E-step: } p(z|d, w) = \frac{p(z, w|d)}{p(w|d)} = \frac{p(z|d)p(w|z)}{\sum_{z'} p(z'|d)p(w|z')}$$

$$p(d, w) = p(d) \sum_{z'} p(z'|d)p(w|z')$$

$$\text{M-step: } p'(w|z) = \frac{p(z, w)}{p(z)} = \frac{\sum_d n(d, w)p(z|d, w)}{\sum_{d, w'} n(d, w')p(z|d, w')}$$

$$p'(z|d) = \frac{p(z, d)}{p(d)} = \frac{\sum_w n(d, w)p(z|d, w)}{\sum_w n(d, w)}$$



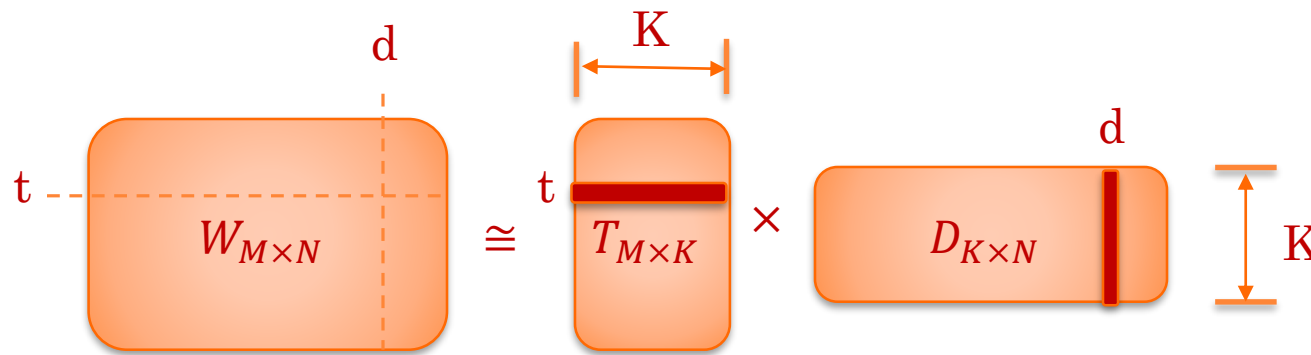
OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- t-SNE

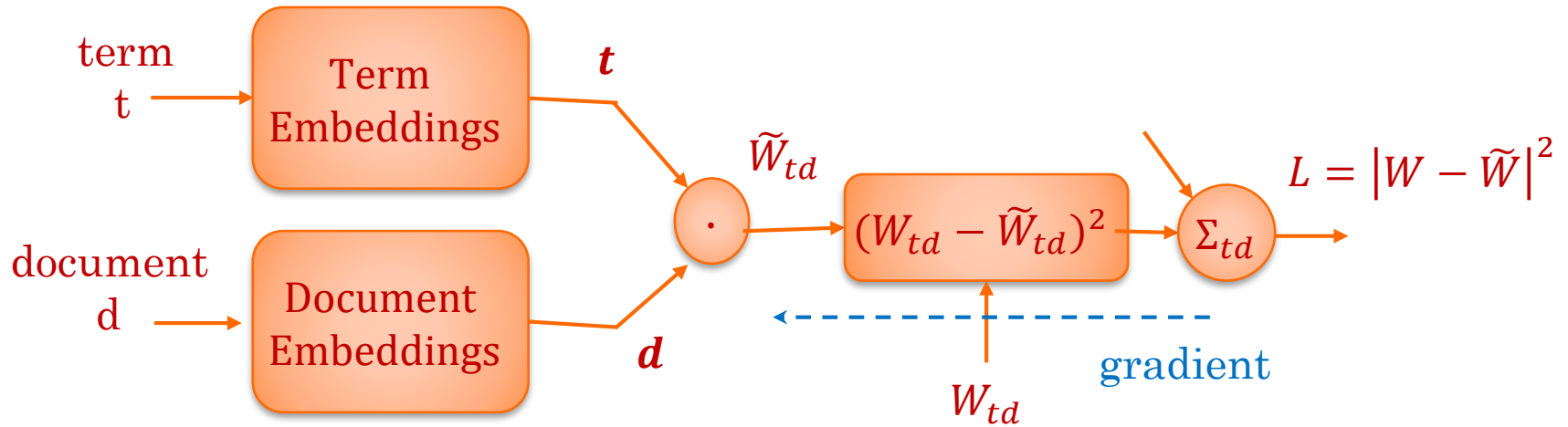


MATRIX FACTORIZATION (MF)

- A is a term-document (or user-item) matrix may be represented with a generative model of latent space.
 - Factorized as the multiplication of two matrix
 - Term matrix T and document matrix D are learnable
 - $\tilde{W}_{td} = \sum_k T_{tk} D_{kd}$ is the inner product of T_t and D_d
 - All vectors are of dimension k.



MF BY GRADIENT DESCENT



- Embeddings: a set of learnable k -dimensional vectors
- \mathbf{t} and \mathbf{d} are K -dimensional and $\mathbf{t} \cdot \mathbf{d} = \sum_k \mathbf{t}(k)\mathbf{d}(k)$
- $\frac{\partial L}{\partial \tilde{W}_{td}} = 2(W_{td} - \tilde{W}_{td})(-1) = -2(W_{td} - \tilde{W}_{td})$
- $\nabla_{\mathbf{d}} L = \frac{\partial L}{\partial \tilde{W}_{td}} \frac{\partial \tilde{W}_{td}}{\partial \mathbf{d}} = 2(W_{td} - \tilde{W}_{td})\mathbf{t} \rightarrow \Delta \mathbf{d} = -\varepsilon \nabla_{\mathbf{d}} L$
- $\nabla_{\mathbf{t}} L = \frac{\partial L}{\partial \tilde{W}_{td}} \frac{\partial \tilde{W}_{td}}{\partial \mathbf{t}} = 2(W_{td} - \tilde{W}_{td})\mathbf{d} \rightarrow \Delta \mathbf{t} = -\varepsilon \nabla_{\mathbf{t}} L$



DISCUSSIONS

- Notice MF does not guarantee the latent vectors to be positive or summed to one (like PLSI).
 - Every k-D vector belongs to R^K
- Similar computation graph could also be used to learn latent vectors with the constraint of PLSI based on gradient descent (without using singular value decomposition)
 - May use softmax to make an array of probabilities (non-negative) summed to 1.
 - e.g. add softmax to make $|\mathbf{d}| = 1, |\mathbf{t}| = 1, \sum_t \tilde{W}_{td} = 1$



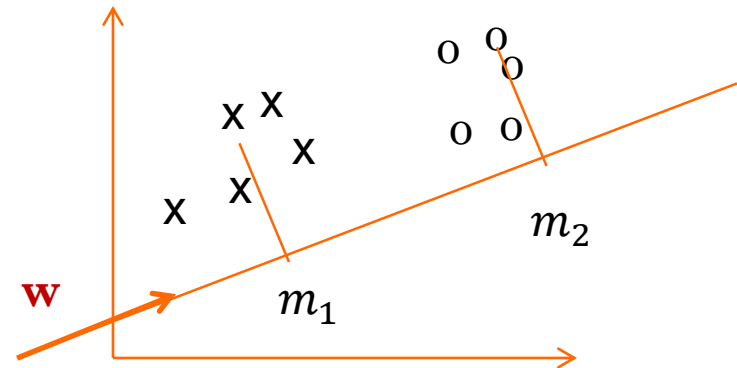
OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- t-SNE



LINEAR DISCRIMINANT ANALYSIS (LDA)

- Reduction of dimension for training high-dimensional vectors such that the low-dimensional feature is discriminant for classification
- Binary classification (C_1, C_2) & $K = 2$
- Conversion of feature $z = \mathbf{w}^t \mathbf{x}$
find optimal direction \mathbf{w}
- $\mathbf{m}_i = \frac{1}{|C_i|} \sum_{\mathbf{x}_t \in C_i} \mathbf{x}_t$
- $m_i = \mathbf{w}^t \mathbf{m}_i$
- $s_i^2 = \sum_{\mathbf{x}_t \in C_i} (\mathbf{w}^t \mathbf{x}_t - m_i)^2$
- $J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$



LINEAR DISCRIMINANT ANALYSIS (LDA)

- $(m_1 - m_2)^2 = (w^t m_1 - w^t m_2)^2$
 $= w^t (m_1 - m_2) (m_1 - m_2)^t w$
 $= w^t S_B w$
- Between-class scatter matrix
 $S_B = (m_1 - m_2)(m_1 - m_2)^t$
- $s_i^2 = \sum_{x_t \in C_i} (w^t x_t - m_i)^2 = \sum_{x_t \in C_i} (w^t (x_t - m_i))^2$
 $= \sum_{x_t \in C_i} w^t (x_t - m_i) (x_t - m_i)^t w = w^t S_i w$
 $S_i \equiv \sum_{x_t \in C_i} (x_t - m_i)(x_t - m_i)^t$
 $s_1^2 + s_2^2 = w^t (S_1 + S_2) w = w^t S_w w$
- Within-class scatter matrix
 $S_w = S_1 + S_2$



LINEAR DISCRIMINANT ANALYSIS (LDA)

- Maximization of $J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{S}_B \mathbf{w}}{\mathbf{w}^t \mathbf{S}_w \mathbf{w}}$

$$\mathbf{w}^* = c \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \text{ (for } K = 1 \text{)}$$

is called Fisher's linear discriminant

- A special case

When $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$, $\mathbf{w}^* = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$.



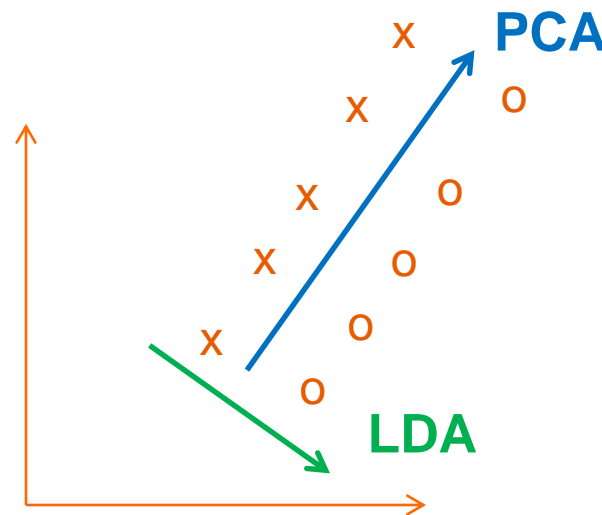
LDA GENERALIZED FOR $K > 2$

- $\mathbf{z} = \mathbf{w}^t \mathbf{x}$ (reduced dimension > 1)
- $S_w = \sum_{i=1}^K S_i$
- $S_B = \sum_{i=1}^K |C_i| (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^t$
 $\mathbf{m} = \sum_{i=1}^K \mathbf{m}_i$
- Maximization of $J(\mathbf{w}) = \frac{|\mathbf{w}^t S_B \mathbf{w}|}{|\mathbf{w}^t S_w \mathbf{w}|}$
 \mathbf{w}^* : eigenvectors of $S_w^{-1} S_B$ with largest eigenvalues.



LDA vs. PCA

- LDA: for classification (with class labels)
 - To find the dimensions for discriminating the classes
- PCA: dimensions with maximum variance
(class labels are not required or provided)



OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- t-SNE

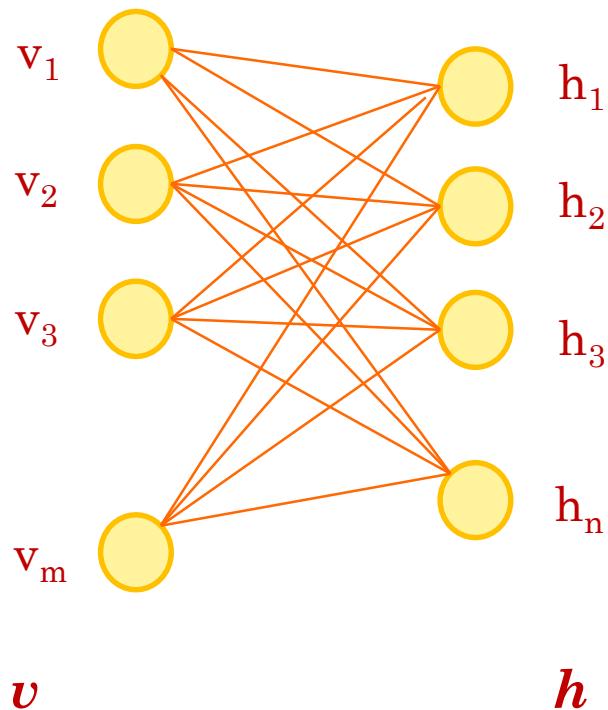


FEATURE LEARNING USING ANN

- Conventionally, features for classification need to be determined in advance based on domain knowledge
 - Gaussian pyramid for image recognition
 - MFCC for speech recognition
- In classification task, the input features could be **pre-trained** automatically using ANN (pre-training)
 - Restricted Boltzman Machine (RBM)
 - Auto Encoder (AE)
- The input could also be the raw data (images/speeches), and the features can be trained and optimized jointly with the classifier in a deep neural network.
 - Convolutional neural network (CNN) for image classification
 - DNN-HMM for speech recognition



RESTRICTED BOLTZMAN MACHINE



- Model: $p(\mathbf{v}, \mathbf{h}) \propto e^{-E(\mathbf{v}, \mathbf{h})}$
- $E(\mathbf{v}, \mathbf{h}) = -\mathbf{h}^t W \mathbf{v} - \mathbf{b}^t \mathbf{v} - \mathbf{c}^t \mathbf{h}$
 $= -\sum_{i,j} w_{ij} h_i v_j - \sum_{j=1}^m b_j v_j - \sum_{i=1}^n c_i h_i$
- $p(\mathbf{v}) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}) = \frac{\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}}{Z}$
- $Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$
 - normalizing the probability weights such that $\sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}) = 1$
- Objective: $\sum_{\mathbf{v}} \log p(\mathbf{v})$
 - $p(\mathbf{v})$ fits for all the \mathbf{v} 's

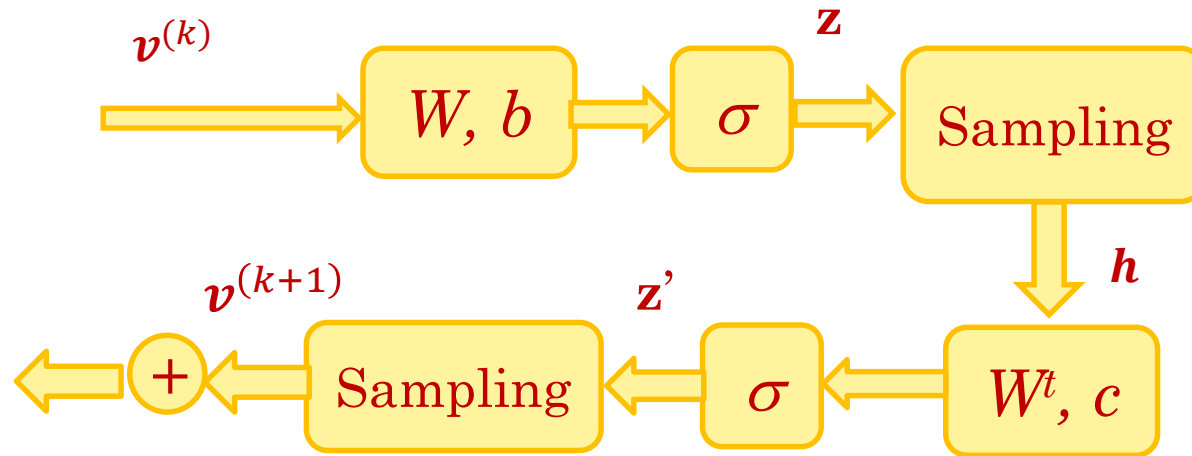


UNSUPERVISED LEARNING OF RBM

- Stochastic process
 - $z_i = P(H_i = 1|\mathbf{v}) = \sigma\left(\sum_{j=1}^m w_{ij}v_j + c_i\right)$
 - $z_j' = P(V_j = 1|\mathbf{h}) = \sigma(\sum_{i=1}^n w_{ij}h_i + b_j)$
- k-step contrastive divergence
 - $dw_{ij} = p(H_i = 1|\mathbf{v}^{(0)})v_j^{(0)} - p(H_i = 1|\mathbf{v}^{(k)})v_j^{(k)}$
 - $db_j = v_j^{(0)} - v_j^{(k)}$
 - $dc_i = h_i^{(0)} - v_i^{(k)}$
- High-level features \mathbf{h} 's can be learned from original features \mathbf{v} 's.
- Multiple features can be learned and stacked layer by layer.



COMPUTATIONAL GRAPH OF RBM



- \mathbf{v} and \mathbf{h} are binary vectors (0/1)
- \mathbf{z} and $\mathbf{z}' : z_i, z_i' \in [0,1]$
- W, b , and c are learnable parameters
- Uncertainty is included through random sampling



K-STEP CONTRASTIVE DIVERGENCE

Algorithm 1. k -step contrastive divergence

Input: RBM $(V_1, \dots, V_m, H_1, \dots, H_n)$, training batch S

Output: gradient approximation Δw_{ij} , Δb_j and Δc_i for $i = 1, \dots, n$,
 $j = 1, \dots, m$

```

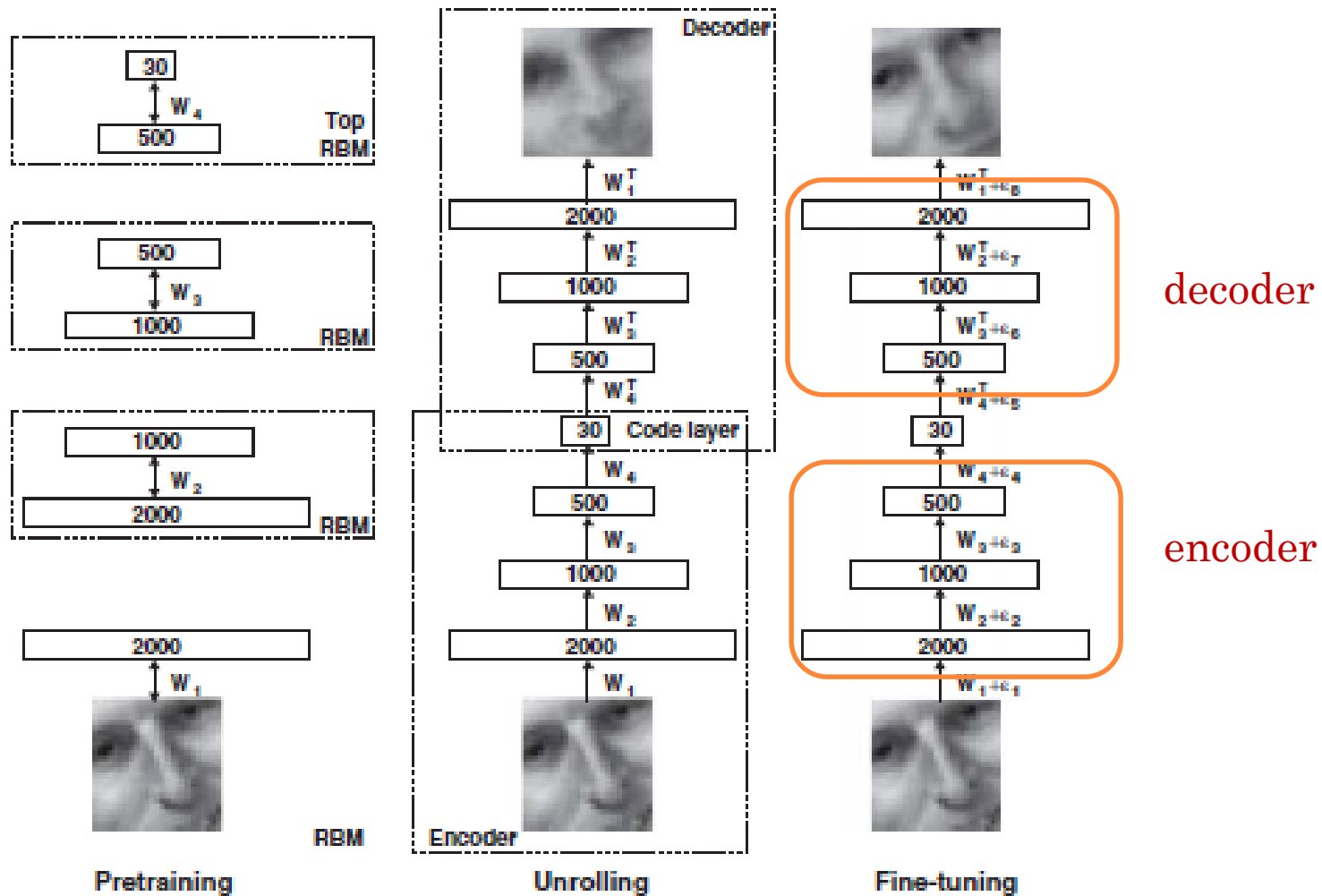
1  init  $\Delta w_{ij} = \Delta b_j = \Delta c_i = 0$  for  $i = 1, \dots, n, j = 1, \dots, m$ 
2  forall the  $v \in S$  do
3       $v^{(0)} \leftarrow v$ 
4      for  $t = 0, \dots, k - 1$  do
5          for  $i = 1, \dots, n$  do sample  $h_i^{(t)} \sim p(h_i | v^{(t)})$ 
6          for  $j = 1, \dots, m$  do sample  $v_j^{(t+1)} \sim p(v_j | h^{(t)})$ 
7      for  $i = 1, \dots, n, j = 1, \dots, m$  do
8           $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1 | v^{(0)}) \cdot v_j^{(0)} - p(H_i = 1 | v^{(k)}) \cdot v_j^{(k)}$ 
9           $\Delta b_j \leftarrow \Delta b_j + v_j^{(0)} - v_j^{(k)}$ 
10          $\Delta c_i \leftarrow \Delta c_i + p(H_i = 1 | v^{(0)}) - p(H_i = 1 | v^{(k)})$ 

```

- Cited from A. Fischer etc.



STACKED RBMs (AUTO-ENCODER)



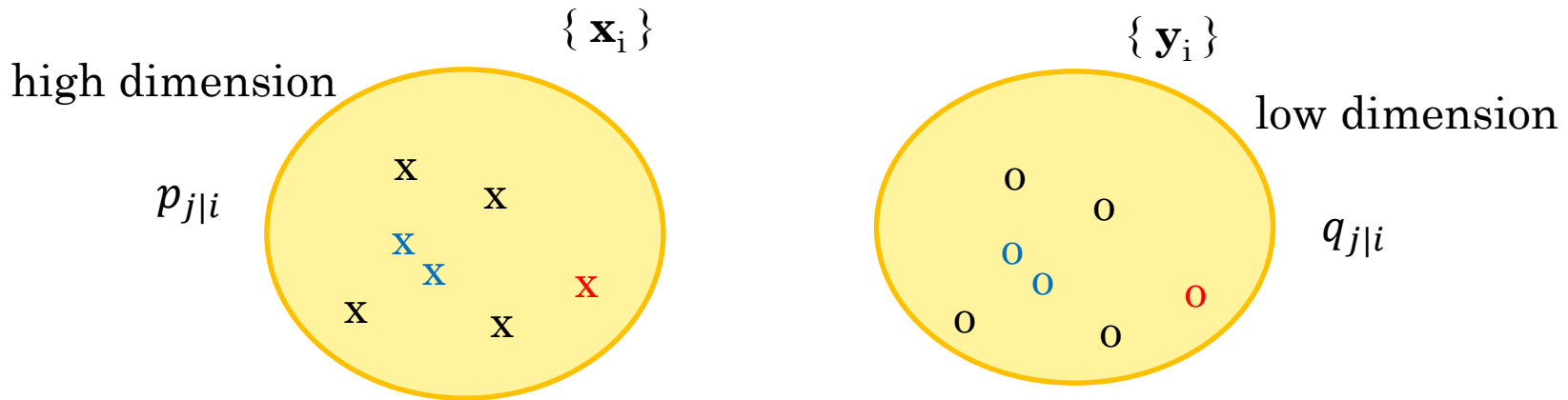
- Unsupervised learning (cited from Hinton etc.)

OUTLINES

- Vector Quantization
- Principal Component Analysis
- Latent Semantic Indexing
- Probabilistic Latent Semantic Indexing
- Matrix Factorization
- Linear Discriminant Analysis
- Restricted Boltzman Machines / Auto Encoder
- t-SNE



T-SNE



- $$p_{j|i} = \frac{\exp\left(-\frac{|\mathbf{x}_i - \mathbf{x}_j|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{|\mathbf{x}_i - \mathbf{x}_k|^2}{2\sigma_i^2}\right)}, \quad q_{j|i} = \frac{\exp(-|\mathbf{y}_i - \mathbf{y}_j|^2)}{\sum_{k \neq i} \exp(-|\mathbf{y}_i - \mathbf{y}_k|^2)}$$
- Thinking: observe points j 's from point i .
- KL-Divergence as cost $C = \sum_{i \neq j} p_{j|i} \log \left(\frac{p_{j|i}}{q_{j|i}} \right)$
- $$\frac{\partial C}{\partial \mathbf{y}_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(\mathbf{y}_i - \mathbf{y}_j)$$
- Asymmetric: $p_{i|j} \neq p_{j|i}, q_{i|j} \neq q_{j|i}$



SYMMETRIC SNE

- $p_{ij} = \frac{\exp\left(-\frac{|x_i - x_j|^2}{2\sigma^2}\right)}{\sum_{k \neq l} \exp\left(-\frac{|x_k - x_l|^2}{2\sigma^2}\right)}, q_{ij} = \frac{\exp(-|y_i - y_j|^2)}{\sum_{k \neq l} \exp(-|y_l - y_k|^2)}$
- $C = \sum_{i \neq j} p_{ij} \log \left(\frac{p_{ij}}{q_{ij}} \right)$
- $\frac{\partial C}{\partial y_i} \equiv \nabla_{y_i} C = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j)$
- Using student t-distribution (Cauchy distribution) to alleviate crowding problem, $q_{ij} = \frac{(1 + |y_i - y_j|^2)^{-1}}{\sum_{k \neq l} (1 + |y_k - y_l|^2)^{-1}}$
- Thinking: try to make
 - q_{ij} approaches p_{ij} after dimension reduction
 - Invariable for the relative distances among points
 - Gaussian: minus log probability \sim distance



EXPECTATION *MINIMIZATION*

- Build $\{ \mathbf{y}_i \}$ randomly
- Iteration
 1. Compute the pairwise probabilities p_{ij} and q_{ij}
 2. E-step: compute the objective KLD. If the KLD converges, then terminate.
 3. M-step: update \mathbf{y}_i by gradient descent, that is, $d\mathbf{y}_i = -\epsilon \nabla_{\mathbf{y}_i} C$, or, $\mathbf{y}_i' = \mathbf{y}_i - \epsilon \nabla_{\mathbf{y}_i} C$
 4. Repeat 1-3 till convergence.
- After the training, for a high-dimensional new point \mathbf{x} , the same procedure can be run with the \mathbf{y}_i 's fixed to obtain the corresponding.
 - Transformation is NOT in a closed form (e.g. PCA)



FEATURE LEARNING

	Name	Concept	Algorithm
PCA	Principal component analysis	Reserve those dimensions with highest variances	Eigen Decomposition
LSI/LSA	Latent semantic indexing	Reserve those dimensions with highest variances	Singular Value Decomposition
PLSI/PLSA	Probabilistic LSI	Topic model	EM/Gradient descent
MF	Matrix Factorization	Latent vectors	Gradient descent
t-SNE	t-Distributed Stochastic Neighbor Embedding	Dim low-D space similar to that in high-D space	Gradient descent
RBM / Auto-encoder	Restricted Boltzman Machine	Unsupervised learning (fine tuning)	Gradient descent

SOME REFERENCES

- Asja Fischer and Christian Igel, *An Introduction to Restricted Boltzman Machines*, LNCS 7441, pp.14-36, 2012. Springer-Verlag.
- G.E. Hinton and R. R. Salak Hutdinov, *Reducing the Dimensionality of Data with Neural Network*, Vol. 313, Science, 2006.
- Laurens van der Maaten, Geoffrey Hinton and Yoshua Bengio, *Visualizing Data using t-SNE*, Journal of Machine Learning Research 9 (2008) 2579-2605, 2008.

