

Unsupervised Learning: Dimensionality Reduction (PCA and other methods)

CS771: Introduction to Machine Learning

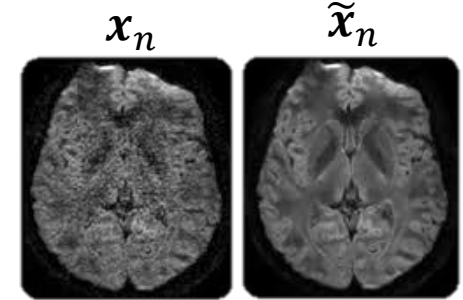
Piyush Rai

Dimensionality Reduction

- Goal: Reduce the dimensionality of each input $\mathbf{x}_n \in \mathbb{R}^D$

$\mathbf{z}_n \in \mathbb{R}^K$ ($K \ll D$) is a compressed version of \mathbf{x}_n

$$\mathbf{z}_n = f(\mathbf{x}_n)$$



- Also want to be able to (approximately) reconstruct \mathbf{x}_n from \mathbf{z}_n

Often $\tilde{\mathbf{x}}_n$ is a “cleaned” version of \mathbf{x}_n (the loss in information is often the noise/redundant information in \mathbf{x}_n)

$$\tilde{\mathbf{x}}_n = g(\mathbf{z}_n) = g(f(\mathbf{x}_n)) \approx \mathbf{x}_n$$

- Sometimes f is called “encoder” and g is called “decoder”. Can be linear/nonlinear
- These functions are learned by minimizing the distortion/reconstruction error of inputs

$$\mathcal{L} = \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 = \sum_{n=1}^N \|\mathbf{x}_n - g(f(\mathbf{x}_n))\|^2$$



Dimensionality Reduction

- Choosing f and g as linear transformations \mathbf{W}^T ($K \times D$) and \mathbf{W} , respectively

$$\mathcal{L} = \sum_{n=1}^N \|\mathbf{x}_n - g(f(\mathbf{x}_n))\|^2 = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{W}\mathbf{W}^T \mathbf{x}_n\|^2$$

$\tilde{\mathbf{x}}_n = \mathbf{W}\mathbf{z}_n$

- Minimizer of \mathcal{L} , if the K columns of \mathbf{W} are orthonormal, are top K eigenvectors of

$D \times D$ empirical
covariance
matrix of inputs

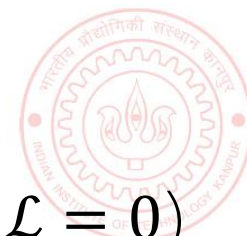
$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T = \frac{1}{N} \mathbf{X}_c^T \mathbf{X}_c$$

\mathbf{X}_c is the $N \times D$ matrix of inputs after centering each input (subtracting off the mean of inputs from each input)

- The matrix \mathbf{W} does a “linear projection” of each input $\mathbf{x}_n \in \mathbb{R}^D$ into a K dim space

- $\mathbf{z}_n = \mathbf{W}^T \mathbf{x}_n \in \mathbb{R}^K$ denotes this linear projection

- Note: If we use $K = D$ eigenvectors for \mathbf{W} , the reconstruction will be perfect ($\mathcal{L} = 0$)



Dimensionality Reduction

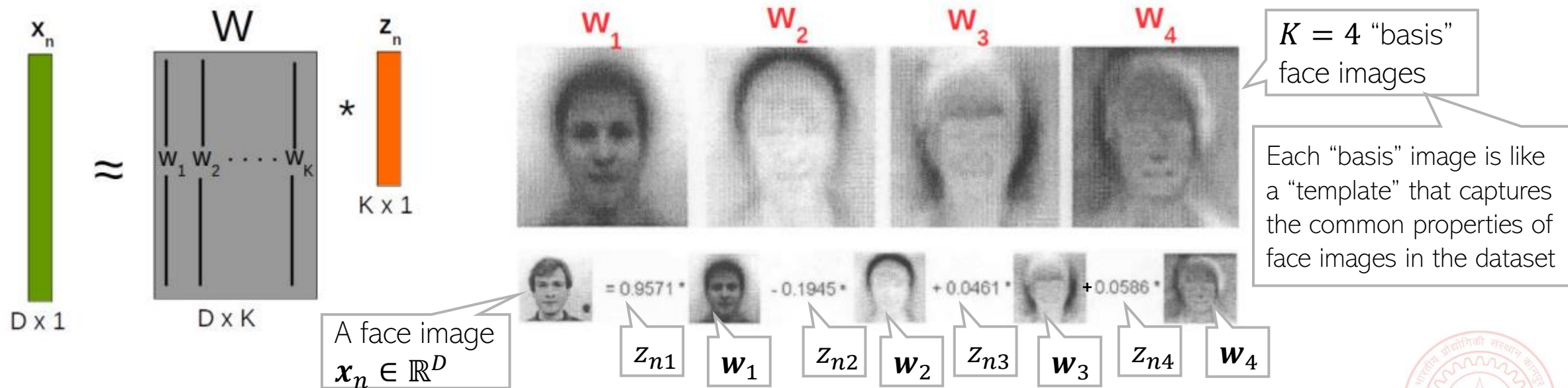
- Consider a linear model of the form

Not necessarily PCA where the columns of \mathbf{W} were orthonormal

$$\mathbf{x}_n \approx \tilde{\mathbf{x}}_n = \mathbf{W} \mathbf{z}_n = \sum_{k=1}^K z_{nk} \mathbf{w}_k$$

\mathbf{w}_k is the k -th column of \mathbf{W}

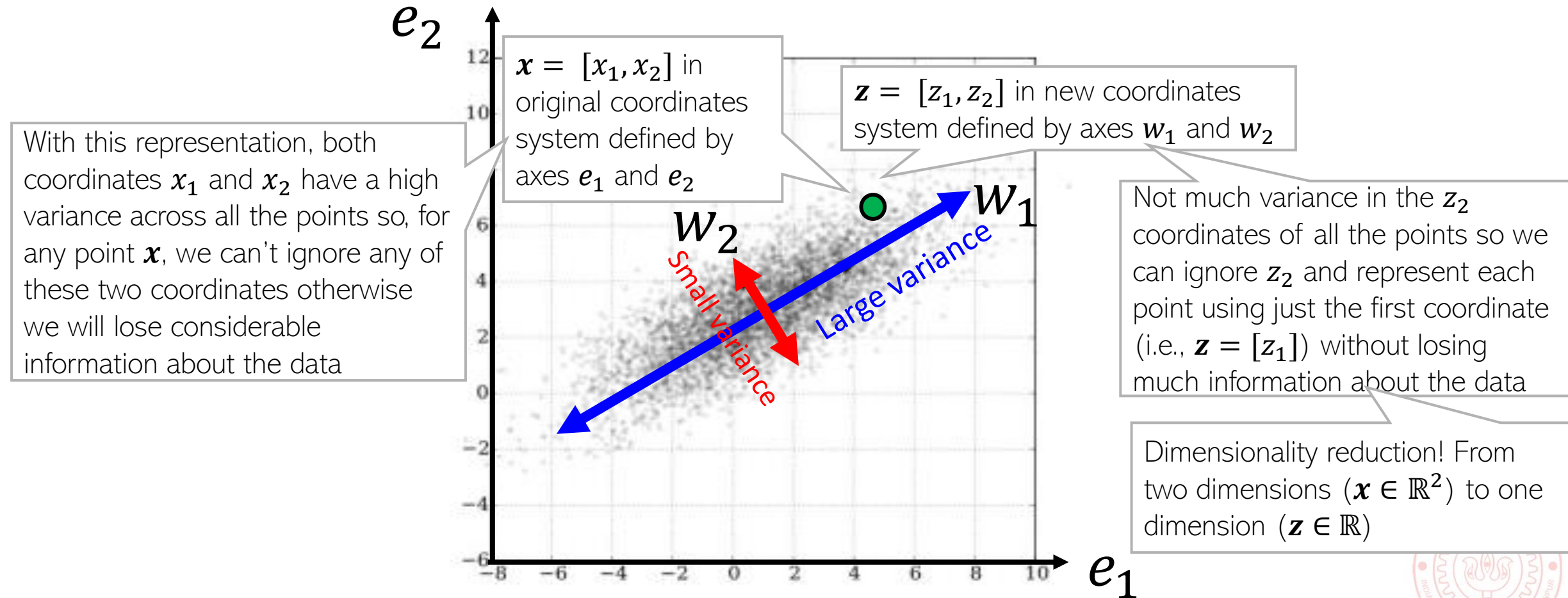
- Above means that each \mathbf{x}_n is approx a linear comb of K vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$



- In this example, $\mathbf{z}_n \in \mathbb{R}^K$ ($K = 4$) is a low-dim feature rep. for each image $\mathbf{x}_n \in \mathbb{R}^D$

Principal Component Analysis (PCA)

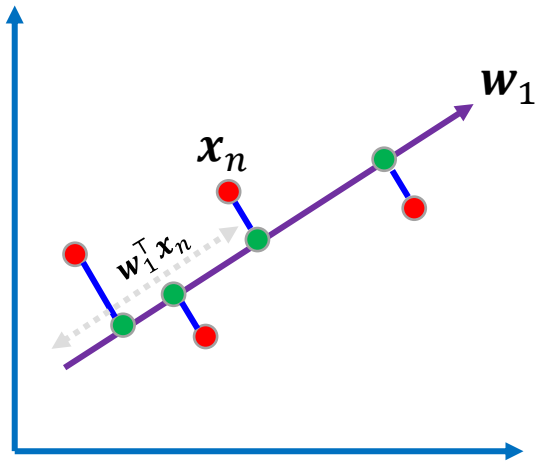
- PCA learns a different and more economical coordinate system to represent data



- Top eigenvectors of the covariance matrix of inputs give us the large variance directions

Finding Max. Variance Directions

- Consider projecting an input $\mathbf{x}_n \in \mathbb{R}^D$ along a direction $\mathbf{w}_1 \in \mathbb{R}^D$
- Projection/**embedding** of \mathbf{x}_n (red points below) will be $\mathbf{w}_1^\top \mathbf{x}_n$ (green pts below)



Mean of projections of all inputs:

$$\frac{1}{N} \sum_{n=1}^N \mathbf{w}_1^\top \mathbf{x}_n = \mathbf{w}_1^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right) = \mathbf{w}_1^\top \boldsymbol{\mu}$$

Variance of the projections:

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{w}_1^\top \mathbf{x}_n - \mathbf{w}_1^\top \boldsymbol{\mu})^2 = \frac{1}{N} \sum_{n=1}^N \{\mathbf{w}_1^\top (\mathbf{x}_n - \boldsymbol{\mu})\}^2 = \mathbf{w}_1^\top \mathbf{S} \mathbf{w}_1$$

\mathbf{S} is the $D \times D$ cov matrix of the data:

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top$$

- Want \mathbf{w}_1 such that variance $\mathbf{w}_1^\top \mathbf{S} \mathbf{w}_1$ is maximized

$$\operatorname{argmax}_{\mathbf{w}_1} \mathbf{w}_1^\top \mathbf{S} \mathbf{w}_1 \quad \text{s.t.} \quad \mathbf{w}_1^\top \mathbf{w}_1 = 1$$

Need this constraint otherwise the objective's max will be infinity

For already centered data, $\boldsymbol{\mu} = \mathbf{0}$ and
 $\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top$

Max. Variance Direction

Variance along the direction \mathbf{w}_1

- Our objective function was $\operatorname{argmax}_{\mathbf{w}_1} \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1$ s.t. $\mathbf{w}_1^T \mathbf{w}_1 = 1$
- Can construct a Lagrangian for this problem

$$\operatorname{argmax}_{\mathbf{w}_1} \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 + \lambda_1 (1 - \mathbf{w}_1^T \mathbf{w}_1)$$

- Taking derivative w.r.t. \mathbf{w}_1 and setting to zero gives $\mathbf{S} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1$
- Therefore \mathbf{w}_1 is an **eigenvector** of the cov matrix \mathbf{S} with eigenvalue λ_1
- Claim:** \mathbf{w}_1 is the eigenvector of \mathbf{S} with largest eigenvalue λ_1 . Note that

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1^T \mathbf{w}_1 = \lambda_1$$

- Thus variance $\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1$ will be max. if λ_1 is the largest eigenvalue (and \mathbf{w}_1 is the corresponding top eigenvector; also known as the first **Principal Component**)
- Other large variance directions can also be found likewise (with each being orthogonal to all others) using the eigendecomposition of cov matrix \mathbf{S} (this is PCA)

Note: Total variance of the data is equal to the sum of eigenvalues of \mathbf{S} , i.e., $\sum_{d=1}^D \lambda_d$

PCA would keep the top $K < D$ such directions of largest variances

Note: In general, \mathbf{S} will have D eigvecs



The PCA Algorithm

- Center the data (subtract the mean $\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$ from each data point)
- Compute the $D \times D$ covariance matrix \mathbf{S} using the centered data matrix \mathbf{X} as

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad (\text{Assuming } \mathbf{X} \text{ is arranged as } N \times D)$$

- Do an eigendecomposition of the covariance matrix \mathbf{S} (many methods exist)
- Take top $K < D$ leading eigenvectors $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K\}$ with eigvalues $\{\lambda_1, \lambda_2, \dots, \lambda_K\}$
- The K -dimensional projection/embedding of each input is

$$\mathbf{z}_n \approx \mathbf{W}_K^T \mathbf{x}_n$$

$\mathbf{W}_K = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ is the “projection matrix” of size $D \times K$

Note: Can decide how many eigvecs to use based on how much variance we want to capture (recall that each λ_k gives the variance in the k^{th} direction (and their sum is the total variance))



The Reconstruction Error View of PCA

- Representing a data point $\mathbf{x}_n = [x_{n1}, x_{n2}, \dots, x_{nD}]^T$ in the standard orthonormal basis $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_D\}$

x_{nd} is the coordinate of \mathbf{x}_n along the direction \mathbf{e}_d

$$\mathbf{x}_n = \sum_{d=1}^D x_{nd} \mathbf{e}_d$$

\mathbf{e}_d is a vector of all zeros except a single 1 at the d^{th} position. Also, $\mathbf{e}_d^T \mathbf{e}_{d'} = 0$ for $d \neq d'$

- Let's represent the same data point in a new orthonormal basis $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_D\}$

z_{nd} is the projection/coordinate of \mathbf{x}_n along the direction \mathbf{w}_d since $z_{nd} = \mathbf{w}_d^T \mathbf{x}_n = \mathbf{x}_n^T \mathbf{w}_d$ (verify)

$$\mathbf{x}_n = \sum_{d=1}^D z_{nd} \mathbf{w}_d$$

$\mathbf{z}_n = [z_{n1}, z_{n2}, \dots, z_{nD}]^T$ denotes the co-ordinates of \mathbf{x}_n in the new basis

- Ignoring directions along which projection z_{nd} is small, we can approximate \mathbf{x}_n as

$$\mathbf{x}_n \approx \hat{\mathbf{x}}_n = \sum_{d=1}^K z_{nd} \mathbf{w}_d = \sum_{d=1}^K (\mathbf{x}_n^T \mathbf{w}_d) \mathbf{w}_d = \sum_{d=1}^K (\mathbf{w}_d \mathbf{w}_d^T) \mathbf{x}_n$$

Note that $\|\mathbf{x}_n - \sum_{d=1}^K (\mathbf{w}_d \mathbf{w}_d^T) \mathbf{x}_n\|^2$ is the **reconstruction error** on \mathbf{x}_n . Would like it to minimize w.r.t. $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$

- Now \mathbf{x}_n is represented by $K < D$ dim. rep. $\mathbf{z}_n = [z_{n1}, z_{n2}, \dots, z_{nK}]$ and

Also, $\mathbf{x}_n \approx \mathbf{W}_K \mathbf{z}_n$

$$\mathbf{z}_n = \mathbf{W}_K^T \mathbf{x}_n$$

$\mathbf{W}_K = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ is the "projection matrix" of size $D \times K$



PCA Minimizes Reconstruction Error

- We plan to use only K directions $[\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ so would like them to be such that the total reconstruction error is minimized

$$\mathcal{L}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = \sum_{n=1}^N \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2 = \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{d=1}^K (\mathbf{w}_d \mathbf{w}_d^T) \mathbf{x}_n \right\|^2$$

Constant; doesn't
depend on the \mathbf{w}_d 's

Variance along \mathbf{w}_d

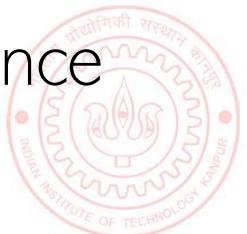
$$= C - \sum_{d=1}^K \mathbf{w}_d^T \mathbf{S} \mathbf{w}_d \text{ (verify)}$$

- Each optimal \mathbf{w}_d can be found by solving

$$\operatorname{argmin}_{\mathbf{w}_d} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = \operatorname{argmax}_{\mathbf{w}_d} \mathbf{w}_d^T \mathbf{S} \mathbf{w}_d$$

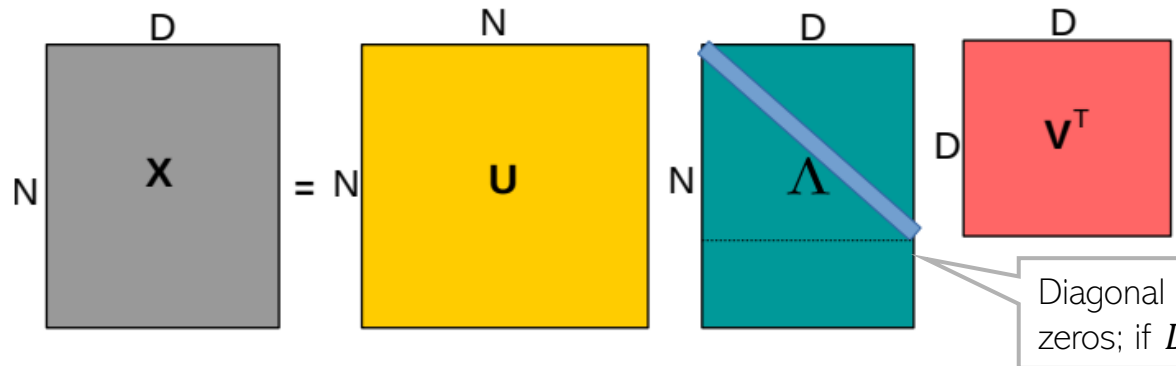
Subject to
 $\mathbf{w}_d^T \mathbf{w}_d = 1$

- Thus minimizing the reconstruction error is equivalent to maximizing variance
- The K directions can be found by solving the eigendecomposition of \mathbf{S}



Singular Value Decomposition (SVD)

- Any matrix \mathbf{X} of size $N \times D$ can be represented as the following decomposition



$$\mathbf{X} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T = \sum_{k=1}^{\min\{N,D\}} \lambda_k \mathbf{u}_k \mathbf{v}_k^T$$

Diagonal matrix. If $N > D$, last $D - N$ rows are all zeros; if $D > N$, last $D - N$ columns are all zeros

- $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N]$ is $N \times N$ matrix of **left singular vectors**, each $\mathbf{u}_n \in \mathbb{R}^N$
 - \mathbf{U} is also orthonormal ($\mathbf{u}_n^T \mathbf{u}_n = 1 \forall n$ and $\mathbf{u}_n^T \mathbf{u}_{n'} = 0 \forall n \neq n'$)
- $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D]$ is $D \times D$ matrix of **right singular vectors**, each $\mathbf{v}_d \in \mathbb{R}^D$
 - \mathbf{V} is also orthonormal ($\mathbf{v}_d^T \mathbf{v}_d = 1 \forall d$ and $\mathbf{v}_d^T \mathbf{v}_{d'} = 0 \forall d \neq d'$)
- $\mathbf{\Lambda}$ is $N \times D$ with only $\min(N, D)$ diagonal entries - **singular values**
- Note: If \mathbf{X} is symmetric then it is known as eigenvalue decomposition ($\mathbf{U} = \mathbf{V}$)



Low-Rank Approximation via SVD

- If we just use the top $K < \min\{N, D\}$ singular values, we get a rank- K SVD

Diagram illustrating the rank- K SVD approximation of matrix \mathbf{X} (size $N \times D$). The approximation is shown as:

$$\mathbf{X} \approx \mathbf{\hat{X}} = \sum_{k=1}^K \lambda_k \mathbf{u}_k \mathbf{v}_k^T = \mathbf{U}_K \mathbf{\Lambda}_K \mathbf{V}_K^T$$

The diagram shows the matrix \mathbf{X} (size $N \times D$) being approximated by the product of three matrices: \mathbf{U}_K (size $N \times K$), $\mathbf{\Lambda}_K$ (size $K \times K$), and \mathbf{V}_K^T (size $K \times D$). The matrix $\mathbf{\Lambda}_K$ is shown as a diagonal matrix. A callout points to the term $\mathbf{u}_k \mathbf{v}_k^T$ in the summation, stating "This is a rank-1 matrix".

- Above SVD approx. can be shown to minimize the reconstruction error $\|\mathbf{X} - \mathbf{\hat{X}}\|$
 - Fact: SVD gives the **best rank- K approximation** of a matrix
- PCA is done by doing SVD on the covariance matrix \mathbf{S} (left and right singular vectors are the same and become eigenvectors, singular values become eigenvalues)

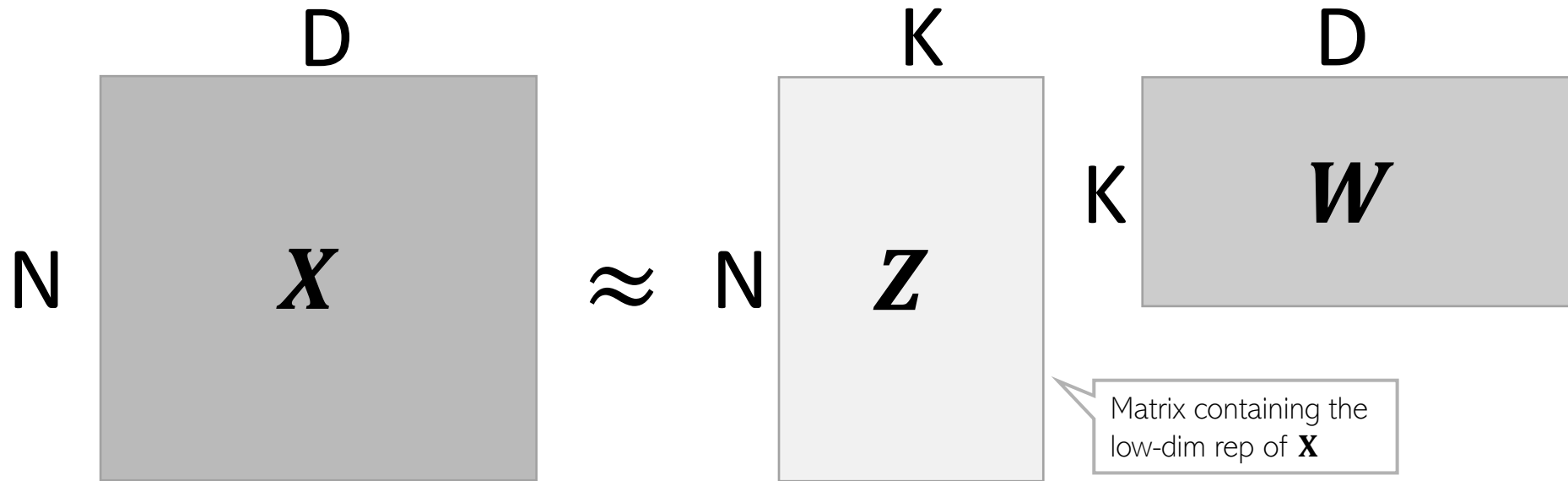


Dimensionality Reduction: Beyond PCA



Dim-Red as Matrix Factorization

- If we don't care about the orthonormality constraints on \mathbf{W} , then dim-red can also be achieved by solving a matrix factorization problem on the data matrix \mathbf{X}



$$\{\hat{\mathbf{Z}}, \hat{\mathbf{W}}\} = \operatorname{argmin}_{\mathbf{Z}, \mathbf{W}} \|\mathbf{X} - \mathbf{Z}\mathbf{W}\|^2$$

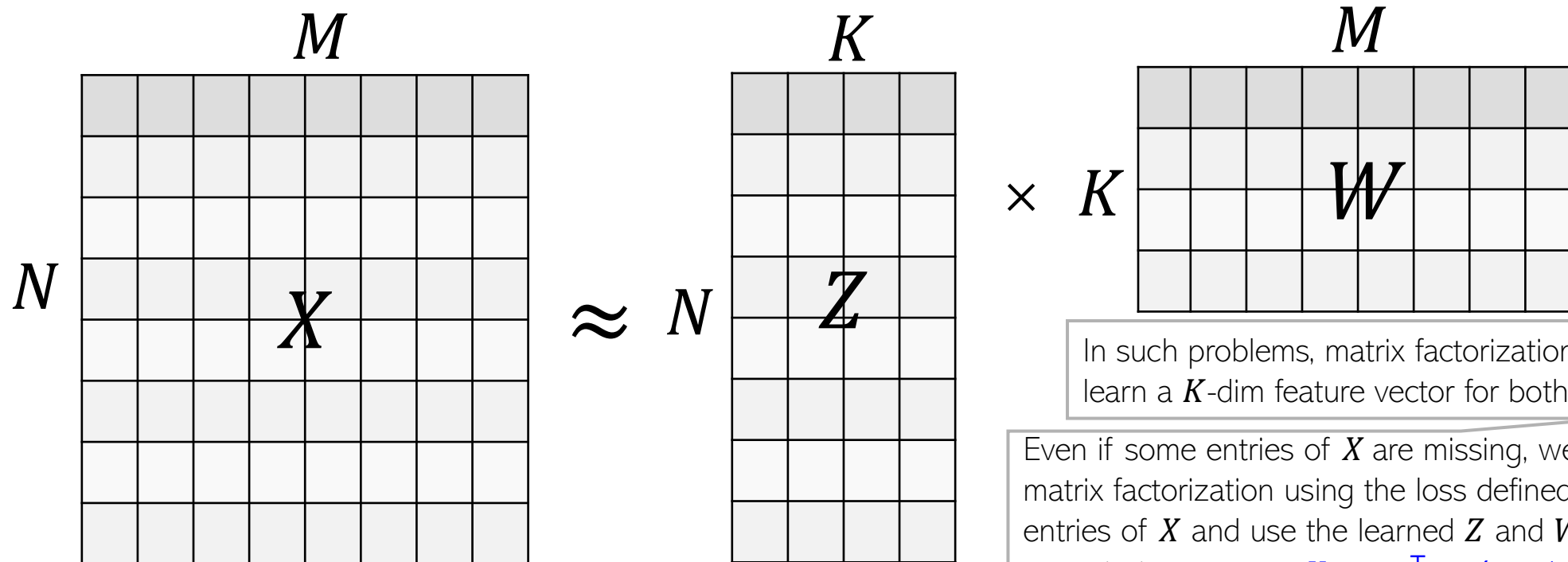
If $K < \min\{D, N\}$, such a factorization gives a low-rank approximation of the data matrix \mathbf{X}

- Can solve such problems using ALT-OPT
- Can impose various constraints on \mathbf{Z} and \mathbf{W} (e.g., sparsity, non-negativity, etc)



Matrix Factorization is a very useful method!

- In many problems, we are given co-occurrence data in form of an $N \times M$ matrix X
- Data consists of relationship b/w two sets of entities containing N and M entities
- Each entry X_{ij} denotes how many times the pair (i, j) co-occurs, e.g.,
 - Number of times a document i (total N docs) contains word j of a vocabulary (total M words)
 - Rating user i gave to item (or movie) j on a shopping (or movie streaming) website



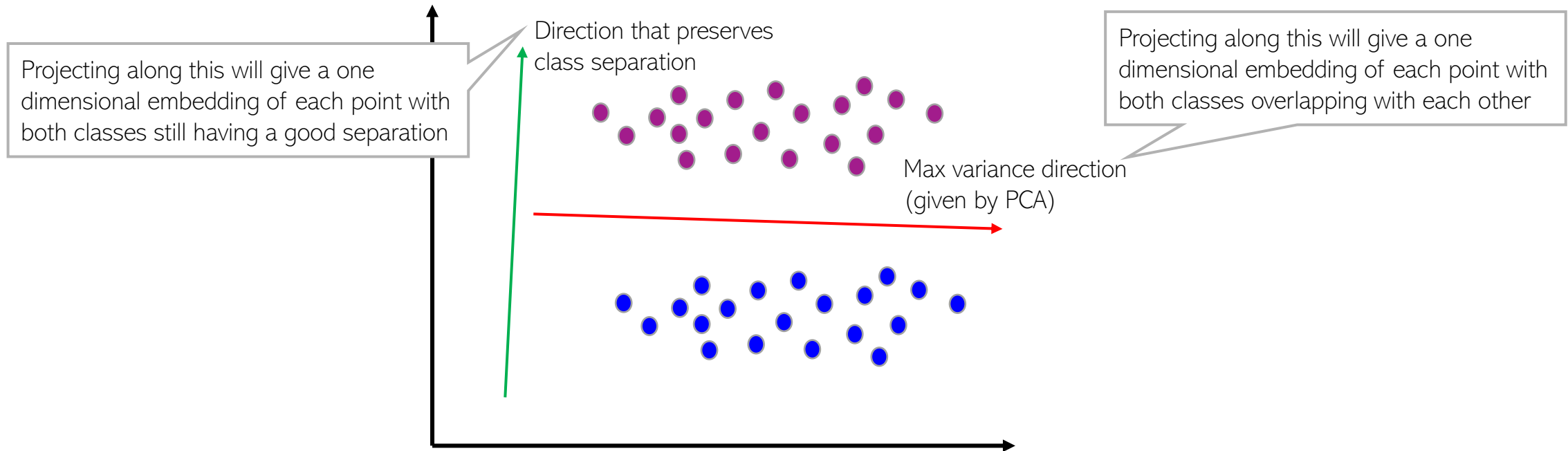
In such problems, matrix factorization can be used to learn a K -dim feature vector for both set of entities

Even if some entries of X are missing, we can still do matrix factorization using the loss defined on the given entries of X and use the learned Z and W to predict any missing entry as $X_{ij} \approx \mathbf{z}_i^T \mathbf{w}_j$ (matrix completion)



Supervised Dimensionality Reduction

- Maximum variance directions may not be aligned with class separation directions (focusing only on variance/reconstruction error of the inputs \mathbf{x}_n , is not always ideal)



- Be careful when using methods like PCA for supervised learning problems
- A better option would be to find projection directions such that after projection
 - Points within the same class are close (low intra-class variance)
 - Points from different classes are well separated (the class means are far apart)



Dim. Reduction by Preserving Pairwise Distances

- PCA/SVD etc assume we are given points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ as vectors (e.g., in D dim)
- Often the data is given in form of **distances** d_{ij} between \mathbf{x}_i and \mathbf{x}_j ($i, j = 1, 2, \dots, N$)
- Would like to project data such that pairwise distances between points are preserved

$$\hat{\mathbf{Z}} = \arg \min_{\mathbf{Z}} \mathcal{L}(\mathbf{Z}) = \arg \min_{\mathbf{Z}} \sum_{i,j=1}^N (d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|)^2$$

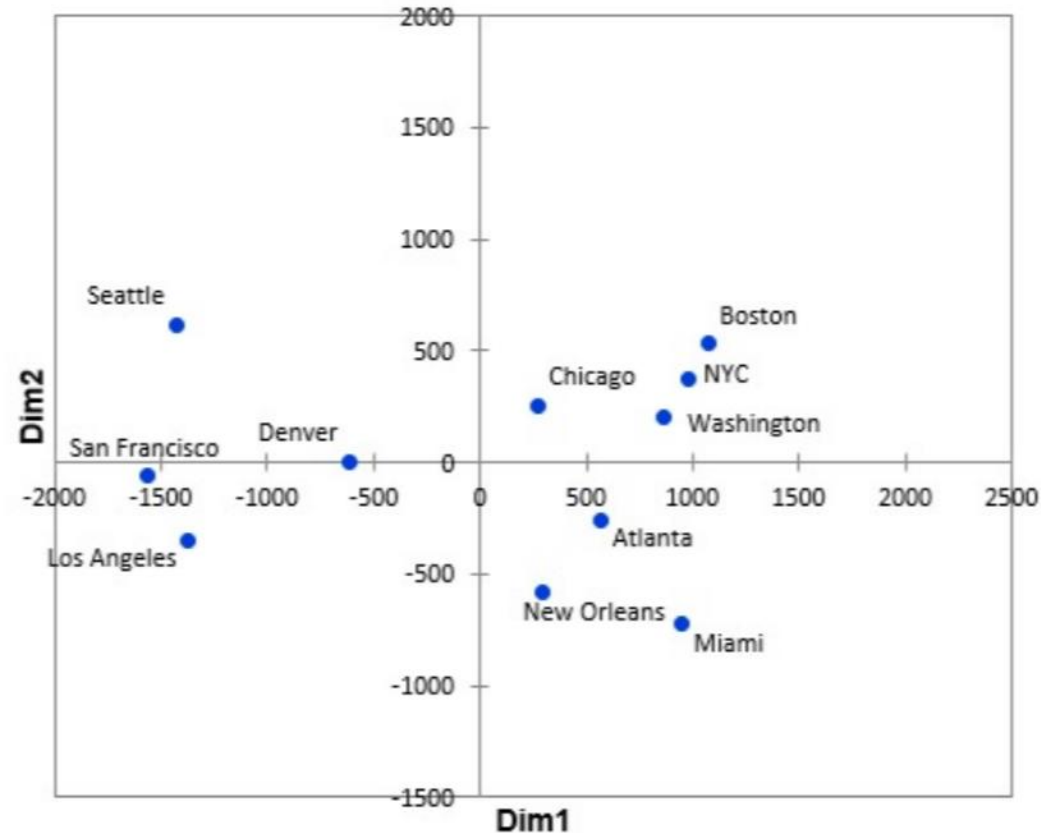
\mathbf{z}_i and \mathbf{z}_j denote low-dim embeddings/projections of \mathbf{x}_i and \mathbf{x}_j , respectively

- Basically, if d_{ij} is large (resp. small), would like $\|\mathbf{z}_i - \mathbf{z}_j\|$ to be large (resp. small)
- **Multi-dimensional Scaling (MDS)** is one such algorithm
- Note: If d_{ij} is the Euclidean distance, MDS is equivalent to PCA



MDS: An Example

- Result of applying MDS (with $K = 2$) on pairwise distances between some US cities



- Here MDS produces 2D embedding of each city such that geographically close cities are also close in 2D embedding space

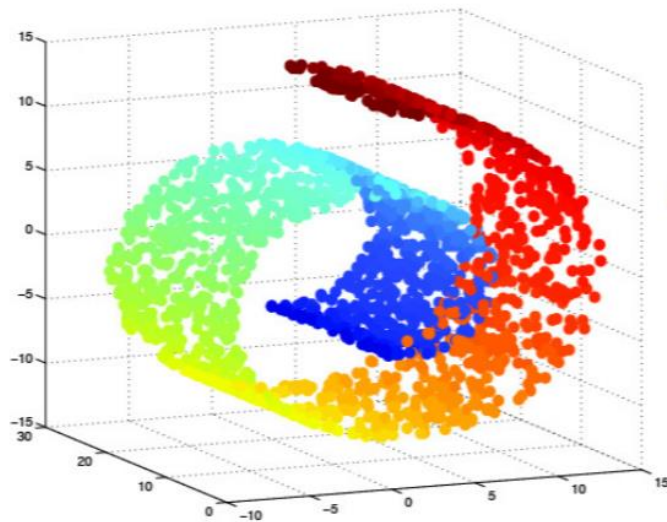


Nonlinear Dimensionality Reduction

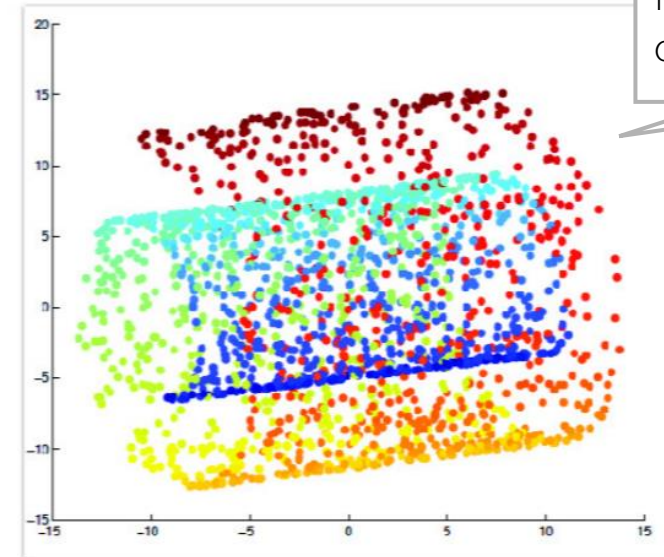


Beyond Linear Projections

- Consider the swiss-roll dataset (points lying close to a manifold)



PCA (Linear Projection)



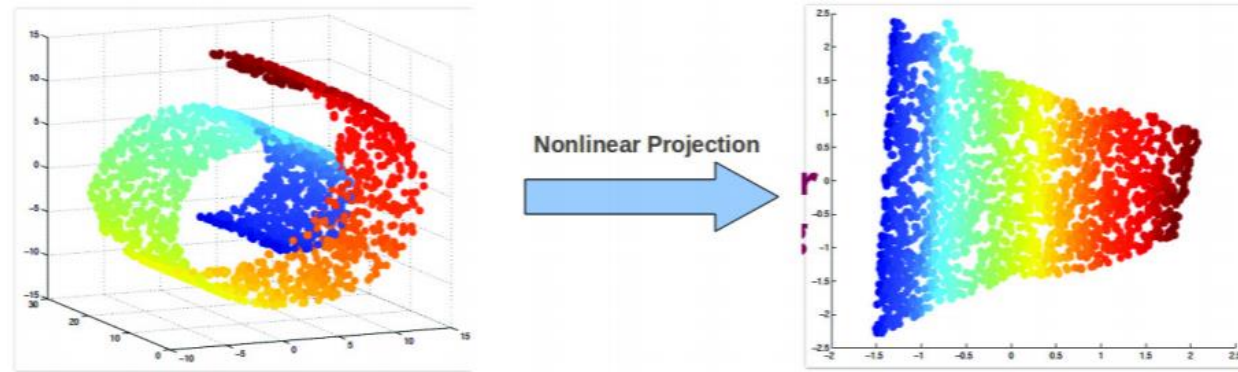
Relative positions of points destroyed after the projection

- Linear projection methods (e.g., PCA) can't capture intrinsic nonlinearities
 - Maximum variance directions may not be the most interesting ones



Nonlinear Dimensionality Reduction

- We want to learn **nonlinear** low-dim projection



Relative positions of points preserved after the projection

- Some ways of doing this
 - Nonlinearize a linear dimensionality reduction method. E.g.:
 - Cluster data and apply linear PCA within each cluster (**mixture of PCA**)
 - **Kernel** PCA (nonlinear PCA)
 - Using **manifold based methods** that intrinsically preserve nonlinear geometry, e.g.,
 - Locally Linear Embedding (LLE), Isomap
 - Maximum Variance Unfolding
 - Laplacian Eigenmap, and others such as SNE/tSNE, etc.
- .. or use unsupervised deep learning techniques (later)

Will look at KPCA, LLE, SNE/tSNE



Kernel PCA

- Recall PCA: Given N observations $\mathbf{x}_n \in \mathbb{R}^D$, $n = 1, 2, \dots, N$,

$D \times D$ cov matrix
assuming centered data

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top$$

D eigenvectors of \mathbf{S}

$$\mathbf{S} \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad \forall i = 1, \dots, D$$

- Assume a kernel k with associated M dimensional nonlinear map ϕ

$M \times M$ cov matrix assuming
centered data in the kernel-
induced feature space

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top$$

M eigenvectors of \mathbf{C}

$$\mathbf{C} \mathbf{v}_i = \lambda_i \mathbf{v}_i \quad \forall i = 1, \dots, M$$

- Would like to do it without computing \mathbf{C} and the mappings $\phi(\mathbf{x}_n)$'s since M can be very large (even infinite, e.g., when using an RBF kernel)
- Boils down to doing eigendecomposition of the $N \times N$ kernel matrix \mathbf{K} (PRML 12.3)
 - Can verify that each \mathbf{v}_i above can be written as a lin-comb of the inputs: $\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n)$
 - Can show that finding $\mathbf{a}_i = [a_{i1}, a_{i2}, \dots, a_{iN}]$ reduces to solving an eigendecomposition of \mathbf{K}
 - Note: Due to req. of centering, we work with a centered kernel matrix $\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$

$N \times N$ matrix of all 1s

Locally Linear Embedding

Several non-lin dim-red
algos use this idea

Essentially, neighbourhood
preservation, but only local

23

- Basic idea: If two points are **local neighbors** in the original space then they should be local neighbors in the projected space too
- Given N observations $\mathbf{x}_n \in \mathbb{R}^D$, $n = 1, 2, \dots, N$, LLE is formulated as

Solve this to learn weights W_{ij} such that each point \mathbf{x}_i can be written as a weighted linear combination of its local neighbors in the original feature space

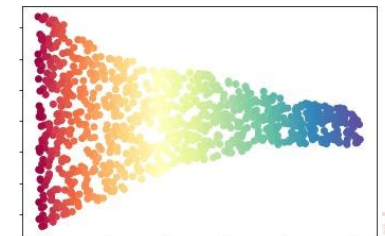
$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}(i)} W_{ij} \mathbf{x}_j \right\|^2$$

$\mathcal{N}(i)$ denotes the local neighbors (a predefined number, say K , of them) of point \mathbf{x}_i

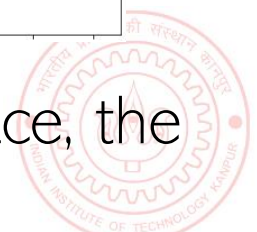
- For each point $\mathbf{x}_n \in \mathbb{R}^D$, LLE learns $\mathbf{z}_n \in \mathbb{R}^K$, $n = 1, 2, \dots, N$ such that the same neighborhood structure exists in low-dim space too

Requires solving an
eigenvalue problem

$$\hat{\mathbf{Z}} = \arg \min_{\mathbf{Z}} \sum_{i=1}^N \left\| \mathbf{z}_i - \sum_{j \in \mathcal{N}(i)} W_{ij} \mathbf{z}_j \right\|^2$$



- Basically, if point \mathbf{x}_i can be reconstructed from its neighbors in the original space, the same weights W_{ij} should be able to reconstruct \mathbf{z}_i in the new space too



SNE and t-SNE

Thus very useful if we want to visualize some high-dim data in two or three dims

- Also nonlin. dim-red methods, especially suited for projecting to 2D or 3D
- SNE stands for **Stochastic Neighbor Embedding** (Hinton and Roweis, 2002)
- Uses the idea of preserving **probabilistically defined neighborhoods**
- SNE, for each point \mathbf{x}_i , defines the probability of a point \mathbf{x}_j being its neighbor as

Neighbor probability in the original space

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma^2)}$$

Neighbor probability in the projected/embedding space

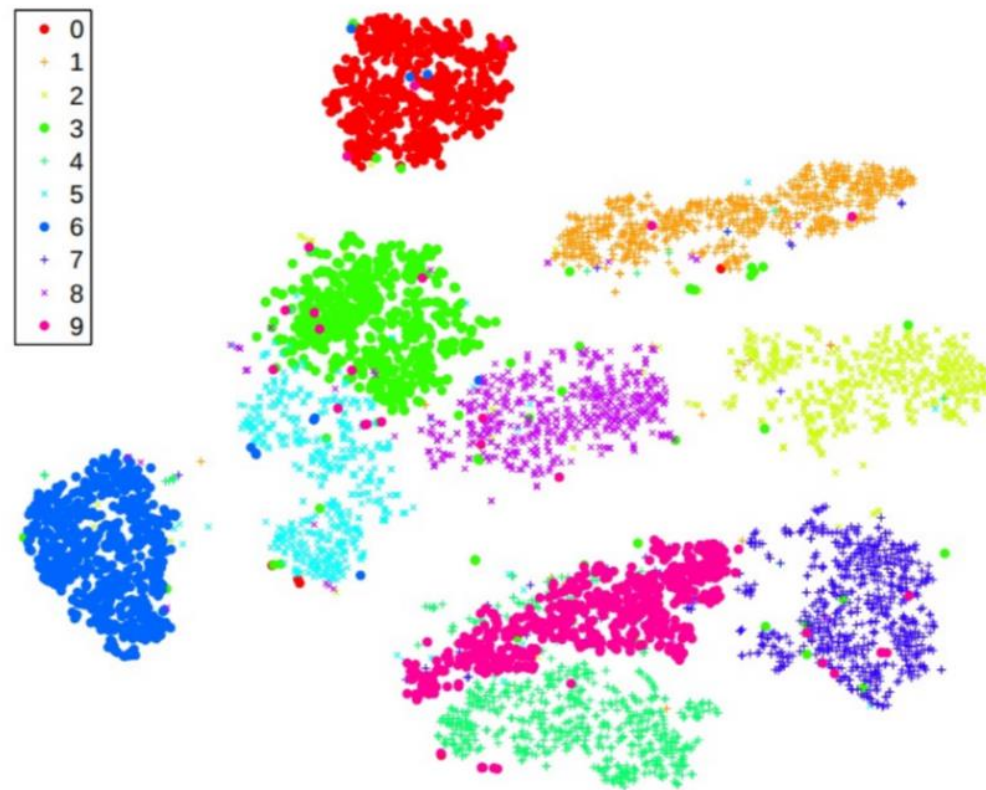
$$q_{j|i} = \frac{\exp(-\|\mathbf{z}_i - \mathbf{z}_j\|^2 / 2\sigma^2)}{\sum_{k \neq i} \exp(-\|\mathbf{z}_i - \mathbf{z}_k\|^2 / 2\sigma^2)}$$

- SNE ensures that neighbourhood distributions in both spaces are as close as possible
 - This is ensured by minimizing their total mismatch (KL divergence) $\mathcal{L} = \sum_{i=1}^N \sum_{j=1}^N p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$
- t-SNE (van der Maaten and Hinton, 2008) offers a couple of improvements to SNE
 - Learns \mathbf{z}_i 's by minimizing **symmetric KL divergence**
 - Uses **Student-t distribution** instead of Gaussian for defining $q_{j|i}$



SNE and t-SNE

- Especially useful for visualizing data by projecting into 2D or 3D



Result of visualizing MNIST digits data in 2D (Figure from van der Maaten and Hinton, 2008)



Word Embeddings: Dim-Reduction for Words

Or sentences, paragraphs, documents, etc
which are basically a set of words

- Feature representation/embeddings of words are very useful in many applications
- Naively we can have a one-hot vector of size V for each word (where V is the vocab size)

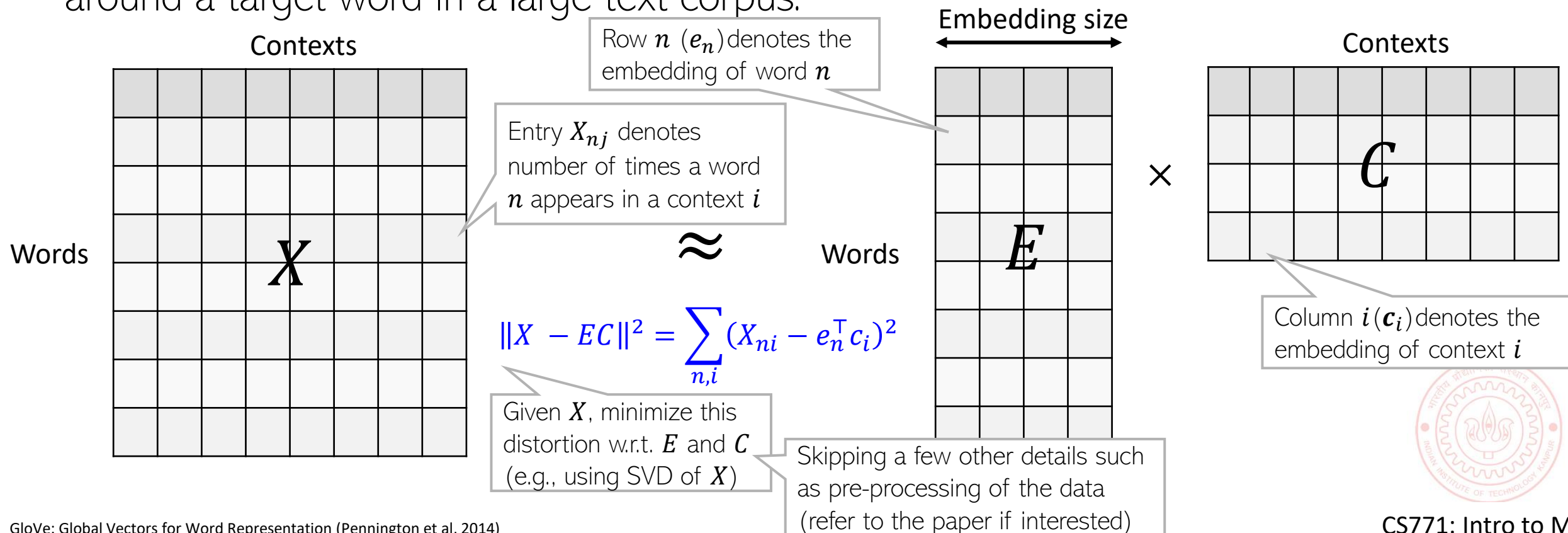


- One-hot representation of a word has two main issues
 - Very high dimensionality (V) for each word
 - One-hot vector does not capture word semantics (any pair of words will have zero similarity)
- Desirable: Learning low-dim **word embeddings** that capture the meaning/semantics
- We want embedding of each word n to be low-dimensional vector $\mathbf{e}_n \in \mathbb{R}^K$
 - If two words n and n' are semantically similar (dissimilar), we want \mathbf{e}_n and $\mathbf{e}_{n'}$ to be close (far)
- Many methods to learn word embeddings (e.g., [Glove](#) and [Word2Vec](#))



GloVe

- GloVe (Global Vectors for Word Representation) is a linear word embedding method
- Based on matrix factorization of a word-context **co-occurrence matrix**
- In GloVe, the context consists of words that co-occur within a specified **context window** around a target word in a large text corpus.



Word2Vec

- A deep neural network based nonlinear word embedding method
- Usually learned using one of the following two objectives
 - Skip-gram
 - Continuous bag of words (CBOW)
- Skip-gram: Probability of a context i occurring around a word n

Conditional probability which can be estimated from training data

$$p(i|n) = \frac{\exp(\mathbf{c}_i^T \mathbf{e}_n)}{\sum_i \exp(\mathbf{c}_i^T \mathbf{e}_n)}$$

Embeddings are learned by optimizing a neural network based loss function which makes the difference b/w LHS and RHS small

- CBOW: Probability of word n occurring given a context window, e.g., k previous and k next words

Conditional probability which can be estimated from training data

$$p(n|n-k:n+k) = \frac{\exp(\mathbf{e}_n^T \mathbf{c}_n)}{\sum_n \exp(\mathbf{e}_n^T \mathbf{c}_n)}$$

Sum/average of the embeddings of the context window for word n

Embeddings are learned by optimizing a neural network based loss function which makes the difference b/w LHS and RHS small



Dimensionality Reduction: Out-of-sample Embedding

- Some dim-red methods can only compute the embedding of the training data
- Given N training samples $\{x_1, x_2, \dots, x_N\}$ they will give their embedding $\{z_1, z_2, \dots, z_N\}$
- However, given a new point x_* (not in the training samples), they can't produce its embedding z_* easily
 - Thus **no easy way of getting “out-of-sample” embedding**
- Some of the nonlinear dim-red methods like LLE, SNE, KPCA, etc have this limitation
 - Reason: They don't learn an explicit encoder and directly optimize for $\{z_n\}_{n=1}^N$ given $\{x_n\}_{n=1}^N$
 - To get “out-of-sample” embeddings, these methods require some modifications*
- But many other methods do explicitly learn a mapping $z = f(x)$ in form of an “encoder” f that can give z_* for any new x_* as well (such methods are more useful)
 - For PCA, the $D \times K$ projection matrix W_K is this encoder function and $z_* = W_K^T x_*$
 - Neural network based **autoencoders** can also do this (will see them later)



*Out-of-Sample Extensions for LLE, Isomap, MDS, Eigenmaps, and Spectral Clustering (Bengio et al, 2003)