# Dimensionality Reduction

**Data Intelligence and Learning (DIAL) Lab**

**Prof. Jongwuk Lee**

# Dimensionality Reduction Basics

# High-dimensional Data

➢ **High-resolution images**
  ◆ Thousands of pixels

# High-dimensional Data

➤ **News articles**

 ◆ Vocabulary of hundreds of thousands of words

# High-dimensional Data

➢ **3D brain imaging data**
  ◆ > 100 MB per scan

# What is Dimensionality Reduction?

➢ **Goal: Unsupervised learning methods for extracting hidden structures from high-dimensional datasets**

➢ **Assumption: Data objects lie on or near a low $d$-dimensional subspace.**

# Example of Dimensionality Reduction

➢ **How to reduce dimensionality without loss of information?**

➢ **Using different axis [2 3 0 0 0] and [0 0 2 4 2], data can be represented by two-dimensional space.**

**Five-dimensional space**

|        | Mon | Tue | Wed | Thu | Fri |
|--------|-----|-----|-----|-----|-----|
| **Alice** | 2   | 3   | 0   | 0   | 0   |
| **Bob**   | 4   | 6   | 0   | 0   | 0   |
| **Carol** | 6   | 9   | 0   | 0   | 0   |
| **David** | 0   | 0   | 2   | 4   | 2   |
| **Eve**   | 0   | 0   | 3   | 6   | 3   |
| **Frank** | 0   | 0   | 1   | 2   | 1   |

**Two-dimensional space**

|        | F1 | F2  |
|--------|-----|-----|
| **Alice** | 1   | 0   |
| **Bob**   | 2   | 0   |
| **Carol** | 3   | 0   |
| **David** | 0   | 1   |
| **Eve**   | 0   | 1.5 |
| **Frank** | 0   | 0.5 |

7

# Low-Dimensional Projection

➢ **Instead of picking a subset of the features, we obtain new features by combining existing features $\mathbf{x}_1, \ldots, \mathbf{x}_d$.**

➢ **We can reduce dimensions, i.e., $k < d$.**

$$\mathbf{z}_1 = \sum_{i=1}^{d} w_i^{(1)} \mathbf{x}_i$$

$$\vdots$$

$$\mathbf{z}_k = \sum_{i=1}^{d} w_i^{(k)} \mathbf{x}_i$$

➢ **New features are represented by linear combinations of old features.**

# Linear Dimensionality Reduction

➢ **Principle component analysis (PCA)**
- ◆ Finding directions of maximal variance

➢ **Independent component analysis (ICA)**
- ◆ Finding directions of maximal independence

➢ **Non-negative matrix factorization (NMF)**



PCA

ICA

# Non-linear Dimensionality Reduction

➢ **t-Distributed stochastic neighbor embedding (t-SNE)**

➢ **Uniform Manifold Approximation and Projection (UMAP)**



UMAP      t-SNE

| T-shirt/top | Shirt | Pullover | Coat | Dress | Sandal | Sneaker | Ankle boot | Trouser | Bag |

# Why Reducing Dimensions?

➢ **Discovering hidden correlations of features**

➢ **Removing redundant and noisy features**

➢ **Easier storage and processing of the data**

➢ **Interpretation and visualization**

$D = 3$
$d = 2$

# Review: Linear Algebra

# Rank of a Matrix

➤ **Q: What is the rank of a matrix $X$?**

➤ **A: Dimension of the vector space spanned by its columns**

- ◆ A: Number of **linearly independent** columns of **X**.

➤ **Example**

$$\mathbf{X} = \begin{bmatrix} 1 & -2 & 3 \\ 2 & -3 & 5 \\ 1 & 1 & 0 \end{bmatrix} \; has \; rank \; r = 2.$$

- ◆ The first two columns are **linearly independent**, so the rank is 2.
- ◆ The third column is equal to the sum of the first and second columns, so the rank must be less than 3.

# Rank is Dimensionality

➢ **Cloud of points 3D space:**
  - ◆ Each column means a point.
  - ◆ Think of point positions as a matrix:

$$\mathbf{X} = \begin{bmatrix} 1 & -2 & 3 \\ 2 & -3 & 5 \\ 1 & 1 & 0 \end{bmatrix}$$



$D = 3$
$d = 2$

➢ **Old basis vectors: [1 0 0] [0 1 0] [0 0 1]**

➢ **We can rewrite coordinates more efficiently!**

➢ **New basis vectors: [1 2 1] [-2 -3 1]**
  - ◆ **We reduced the number of coordinates!**

# Linear Independence

➢ **A set of vectors is said to be linearly dependent.**

  ◆ if one of the vectors in the set can be defined as a linear combination of the others.

➢ **How to check the linear independence?**

  ◆ $n$ vectors are **linearly independent** if and only if the determinant of the matrix formed by taking the vectors as its columns is **non-zero**.

$$A = \begin{bmatrix} 1 & -3 \\ 1 & 2 \end{bmatrix} \quad \Longrightarrow \quad \det(A) = 1 \cdot 2 - 1 \cdot (-3) = 5 \neq 0$$

  ◆ Because the determinant is non-zero, two vectors $(1, 1)$ and $(-3, 2)$ are linearly independent.

# Orthogonality

➤ **Two vectors $x$ and $y$ are orthogonal if their inner product $\langle x, y \rangle$ is zero.**
  ◆ $(1, 3, 2)^\mathrm{T}, (3, -1, 0)^\mathrm{T}, (1, 3, -5)^\mathrm{T}$ are orthogonal to each other.

➤ **If two vectors are orthogonal, then they are linearly independent.**
  ◆ However, the inverse relationship does not hold.

➤ **Two vectors are orthonormal if they are orthogonal and unit vectors.**

# Eigenvalue and Eigenvector

> ## What is an eigenvector?

- ◆ A non-zero vector **v** whose direction does not change when a linear transformation **A** is applied to it
- ◆ An eigenvalue $\lambda$ is a scalar associated with the eigenvector **v**.

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$



**The red arrow changes direction, but the blue arrow does not.**

**The blue arrow is an eigenvector of this mapping** because it doesn't change direction, and since its length is unchanged, its eigenvalue is 1.

# How to Compute Eigenvalues?

➢ **Given** $X = \begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix}$**, compute eigenvalues and eigenvectors.**

- $\det(X - \lambda I) = \begin{vmatrix} 2 - \lambda & 2 \\ 5 & -1 - \lambda \end{vmatrix}$

$$= (2 - \lambda)(-1 - \lambda) - 10 = \lambda^2 - \lambda - 12 = (\lambda - 4)(\lambda + 3) = 0$$

- When $\lambda = 4$, $\begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 4 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow x_1 = x_2$

- When $\lambda = -3$, $\begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = -3 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow x_1 = -\frac{2}{5} x_2$

# How to Compute Eigenvalues?

➢ **Given X** $= \begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix}$, **compute eigenvalues and eigenvectors.**

◆ When $\lambda = 4$, $\begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 4 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow x_1 = x_2$

◆ When $\lambda = -3$, $\begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = -3 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow x_1 = -\frac{2}{5} x_2$



**The direction of the red arrow is not changed.**

# Principle Component Analysis (PCA)

# Motivation

➢ An **orthogonal linear transformation** that transfers the data to a **new coordinate system**

# Maximizing the Variance

➢ Which maximizes the **variance**?

➢ Which minimizes **the sum of projected distances**?

# Maximizing the Variance

➢ **PCA is the <span style="color:red">orthogonal projection</span> of the data onto a lower-dimension linear space that**

  ◆ Maximizes the variance of projected data (**red line**).
  ◆ Minimizes the mean squared distance between data point and projections (**sum of blue lines**).

➢ **Data in two-dimensional space is projected into a new one-dimensional space.**

# Principle Component Analysis (PCA)

➢ **PCA Vectors originate from the center of mass.**

➢ **The first principal component: points in the direction of the largest variance.**

➢ **Each subsequent principal component**
- ◆ is **orthogonal to the previous ones**, and
- ◆ points in the directions of the **largest variance of the residual subspace**.

# Example

# Multivariate Data

➤ **Let $\mathcal{D} = \{\mathbf{x}^{(1)}, \cdots, \mathbf{x}^{(n)}\}$ be a training dataset.**

➤ **Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a data matrix.**

- ◆ $n$ samples/instances/examples
- ◆ $d$ features/attributes

$$\mathbf{X} = \begin{bmatrix} \left[\mathbf{x}^{(1)}\right]^{\mathrm{T}} \\ \left[\mathbf{x}^{(2)}\right]^{\mathrm{T}} \\ \vdots \\ \left[\mathbf{x}^{(n)}\right]^{\mathrm{T}} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_d^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$$

\# of samples

\# of features

# Zero-Centered Normalization

➢ **We assume that data is <span style="color:red">centered</span>.**

$$\mu = \frac{1}{n}\sum_{i=1}^{n} \mathbf{x}^{(i)} = 0$$

➢ **Q: What if our data is not centered?**

➢ **A: Subtract the sample mean as <span style="color:red">pre-processing</span>.**

# Sample Covariance Matrix

➢ **How to compute a covariance matrix?**

$$\Sigma_{jk} = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_j^{(i)} - \mu_j)(\mathbf{x}_k^{(i)} - \mu_k)$$

➢ **Since the data matrix is <span style="color:red">centered</span>, we rewrite it as:**

$$\Sigma = Cov(\mathbf{X}) = \frac{1}{n}(\mathbf{X} - \boldsymbol{\mu})^{\mathrm{T}}(\mathbf{X} - \boldsymbol{\mu}) \quad \Longrightarrow \quad \Sigma = \frac{1}{n}\mathbf{X}^{\mathrm{T}}\mathbf{X}$$

Note: $\boldsymbol{\mu}$ is a zero vector.

# Formulating PCA

> **The matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ is projected into a vector $\mathbf{v} \in \mathbb{R}^{d \times 1}$.**

Assume that $\mathbf{v}$ is the unit vector, i.e., $\mathbf{v}^{\mathbf{T}}\mathbf{v} = 1$.

$$\mathbf{v}^* = \underset{\mathbf{v}:\,\|\mathbf{v}\|^2=1}{\operatorname{argmax}} Var(\mathbf{X}\mathbf{v}) = \underset{\mathbf{v}:\,\|\mathbf{v}\|^2=1}{\operatorname{argmax}} \frac{1}{n}\sum_{i=1}^{n}\left(\mathbf{x}^{(i)}\mathbf{v}\right)^{\mathbf{T}}\left(\mathbf{x}^{(i)}\mathbf{v}\right)$$

Data are zero-centered.

$$= \underset{\mathbf{v}:\,\|\mathbf{v}\|^2=1}{\operatorname{argmax}} \frac{1}{n}\sum_{i=1}^{n}\mathbf{v}^{\mathbf{T}}\left(\mathbf{x}^{(i)}\right)^{\mathbf{T}}\mathbf{x}^{(i)}\mathbf{v} = \underset{\mathbf{v}:\,\|\mathbf{v}\|^2=1}{\operatorname{argmax}} \mathbf{v}^{\mathbf{T}}\Sigma\mathbf{v}$$

Introduce a Lagrange multiplier for the equality constraint $\|\mathbf{v}\|^2 = 1$.

$$\Sigma = \frac{1}{n}\sum_{i=1}^{n}\mathbf{x}^{\mathbf{T}}\mathbf{x}$$

$$\mathcal{L}(v, \lambda) = \mathbf{v}^{\mathbf{T}}\Sigma\mathbf{v} - \lambda\left(\mathbf{v}^{\mathbf{T}}\mathbf{v} - 1\right)$$

# Formulating PCA

➢ **Taking the derivative of the Lagrangian and setting it to zero,**

$$\mathcal{L}(v, \lambda) = \mathbf{v}^{\mathrm{T}}\mathbf{\Sigma}\mathbf{v} - \lambda(\mathbf{v}^{\mathrm{T}}\mathbf{v} - 1)$$

$$\frac{d}{d\mathbf{v}}\left(\mathbf{v}^{\mathrm{T}}\mathbf{\Sigma}\mathbf{v} - \lambda(\mathbf{v}^{\mathrm{T}}\mathbf{v} - 1)\right) = 0$$

➡ $$\mathbf{\Sigma}\mathbf{v} - \lambda\mathbf{v} = 0 \iff \mathbf{\Sigma}\mathbf{v} = \lambda\mathbf{v}$$

➢ **We can compute the PCs as the eigenvector of $\Sigma$.**

➢ **Recall: For a square matrix $\mathbf{A}$, the vector $\mathbf{v}$ is an eigenvector if and only if there exist eigenvalue $\lambda$ such that $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$.**

# How to Choose PCs?

➢ **The eigenvalue $\lambda$ denotes the amount of variability captured along that dimension.**

➢ **For eigenvalues $\lambda_1 \geq \lambda_2 \geq \cdots$,**

- The 1st PC is the eigenvector $\mathbf{v}_1$ of the covariance matrix $\mathbf{X}^{\mathrm{T}}\mathbf{X}$ associated with the largest eigenvalue.

- The 2nd PC is the eigenvector $\mathbf{v}_2$ of the covariance matrix $\mathbf{X}^{\mathrm{T}}\mathbf{X}$ associated with the second largest eigenvalue.

- And so on …

# How to Choose PCs?

➢ **For $d$ $(d < n)$ original dimensions, sample covariance matrix is $d \times d$, and has up to $d$ eigenvectors.**

➢ **Q: How do we perform dimensionality reduction?**

➢ **A: Can ignore the components of lesser significance.**

- ◆ May lose some information, but we lose insignificant information if the eigenvalues are small.



Variance (%) = ratio of variance along given principal component to total variance of all principal components

# Example of PCA

# Zero-Centered Normalization

> **Step 1: Given a data, we subtract the average.**

| A1 | A2 |
|---|---|
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

| A1 | A2 |
|---|---|
| 0.69 | 0.49 |
| -1.31 | -1.21 |
| 0.39 | 0.99 |
| 0.09 | 0.29 |
| 1.29 | 1.09 |
| 0.49 | 0.79 |
| 0.19 | -0.31 |
| -0.81 | -0.81 |
| -0.31 | -0.31 |
| -0.71 | -1.01 |

# Zero-Centered Normalization

➢ **Step 1: Given a data, we subtract the average.**

zero-centered normalization

# Computing a Covariance Matrix

> **Step 2: Calculate the covariance matrix.**
> - $cov = \begin{pmatrix} .6166 & .6154 \\ .6154 & .7166 \end{pmatrix}$

$$\Sigma = \frac{1}{n}\mathbf{X}^{\mathrm{T}}\mathbf{X}$$

http://www.visiondummy.com/2014/04/geometric-interpretation-covariance-matrix/

# Computing Eigenvectors

➢ **Step 3: Calculate the eigenvectors and eigenvalues for the covariance matrix.**

◆ $eigenvalues = \begin{pmatrix} 1.2840 \\ 0.0491 \end{pmatrix}$

◆ $eigenvectors = \begin{pmatrix} -0.6779 & -0.7352 \\ -0.7352 & 0.6779 \end{pmatrix}$

# Transforming into New Features

➢ **Construct a new feature vector $f = (\mathbf{v}_1, \dots, \mathbf{v}_d)$.**

$$new\ feature\ vector = \begin{pmatrix} -0.6779 & -0.7352 \\ -0.7352 & 0.6779 \end{pmatrix}$$

➢ **For approximation, it is possible to reduce the feature vector on one-dimensional space.**

$$new\ feature\ vector = \begin{pmatrix} -0.6779 \\ -0.7352 \end{pmatrix}$$

# Transforming into New Features

> **Deriving a new data set**
> - **Final data = row data adjust * new feature vector**

| A1 | A2 |
|---|---|
| 0.69 | 0.49 |
| -1.31 | -1.21 |
| 0.39 | 0.99 |
| 0.09 | 0.29 |
| 1.29 | 1.09 |
| 0.49 | 0.79 |
| 0.19 | -0.31 |
| -0.81 | -0.81 |
| -0.31 | -0.31 |
| -0.71 | -1.01 |

$$\cdot \begin{pmatrix} -0.6779 & -0.7352 \\ -0.7352 & 0.6779 \end{pmatrix} =$$

**Represent data into new transformed space.**

| A1' | A2' |
|---|---|
| -0.82797 | -0.17512 |
| 1.77758 | 0.142857 |
| -0.9922 | 0.384375 |
| -0.27421 | 0.130417 |
| -1.6758 | -0.2095 |
| -0.91295 | 0.175282 |
| 0.099109 | -0.34982 |
| 1.144572 | 0.046417 |
| 0.438046 | 0.017765 |
| 1.223821 | -0.16268 |

# Transforming into New Features

> **Deriving a new data set**

- ◆ **Final data = row data adjust * new feature vector**

| A1 | A2 |
|-----|-----|
| 0.69 | 0.49 |
| -1.31 | -1.21 |
| 0.39 | 0.99 |
| 0.09 | 0.29 |
| 1.29 | 1.09 |
| 0.49 | 0.79 |
| 0.19 | -0.31 |
| -0.81 | -0.81 |
| -0.31 | -0.31 |
| -0.71 | -1.01 |

$$\cdot \begin{pmatrix} -0.6779 \\ -0.7352 \end{pmatrix} =$$

| A1' |
|-----|
| -0.82797 |
| 1.77758 |
| -0.9922 |
| -0.27421 |
| -1.6758 |
| -0.91295 |
| 0.099109 |
| 1.144572 |
| 0.438046 |
| 1.223821 |

# Transforming into New Features

➢ **Transform original data into new data space using chosen components.**

◆ **Final data = row data adjust * new feature vector**

| A1' | A2' |
|---|---|
| -0.82797 | -0.17512 |
| 1.77758 | 0.142857 |
| -0.9922 | 0.384375 |
| -0.27421 | 0.130417 |
| -1.6758 | -0.2095 |
| -0.91295 | 0.175282 |
| 0.099109 | -0.34982 |
| 1.144572 | 0.046417 |
| 0.438046 | 0.017765 |
| 1.223821 | -0.16268 |

# Applications of PCA

# Face Recognition

- **We want to identify a specific person based on facial images.**
  - It should be robust to glasses, lighting, etc.
  - However, 256 x 256 images are high-dimensional data.

- **Solution: Build one PCA database for the whole dataset and then classify based on the weights.**

# Eigenfaces: Applying PCA

➢ **Example data set: Images of faces**

- ◆ Eigenface approach
- ◆ Each face **x** is 256 x 256 pixels
- ◆ **x** is a 64K dimensional vector.

➢ **Form** $\mathbf{X} \in \mathbb{R}^{n \times 64K}$ **centered data**

- ◆ Let $n$ be the number of samples.

➢ **Compute** $\mathbf{\Sigma} = \mathbf{X}^{\mathrm{T}}\mathbf{X}.$

➢ **Problem: The** $64K \times 64K$ **matrix is too HUGE!**

# Computational Complexity

➢ **Suppose that we have $n$ samples.**

◆ Eigenfaces: $n = 500$ faces, each of size $d = 64K$

➢ **Given $d \times d$ covariance matrix $\Sigma$, we have to compute**

◆ All $d$ eigenvectors/eigenvalues in O($d^3$)
◆ First $k$ eigenvectors/eigenvalues in O($kd^2$)

➢ **However, if $d = 64K$, it is too EXPENSIVE!**

# A Clever Idea

➢ **Note that** $n \ll 64K$**.**

  ◆ Use $\mathbf{L} = \mathbf{XX}^{\mathrm{T}}$ instead of $\Sigma = \mathbf{X}^{\mathrm{T}}\mathbf{X}$

➢ **If $\mathbf{v}$ is eigenvector of $\mathbf{L}$, then $\mathbf{X}^{\mathrm{T}}\mathbf{v}$ is eigenvector of $\Sigma$.**

➢ **Proof:**

$$\mathbf{Lv} = \lambda \mathbf{v}$$

$$\mathbf{XX}^{\mathrm{T}}\mathbf{v} = \lambda \mathbf{v}$$

$$\mathbf{X}^{\mathrm{T}}(\mathbf{XX}^{\mathrm{T}}\mathbf{v}) = \mathbf{X}^{\mathrm{T}}(\lambda \mathbf{v})$$

$$(\mathbf{X}^{\mathrm{T}}\mathbf{X})\mathbf{X}^{\mathrm{T}}\mathbf{v} = \lambda(\mathbf{X}^{\mathrm{T}}\mathbf{v})$$

$$\Sigma \mathbf{X}^{\mathrm{T}}\mathbf{v} = \lambda(\mathbf{X}^{\mathrm{T}}\mathbf{v})$$

Instead of computing the eigenvector of $\Sigma$ directly, compute the eigenvector of $\Sigma$ using the eigenvector of $\mathbf{L}$.

# Face Recognition with EigenFaces

> Using PCA, we represent eigenfaces from original images.

# Face Recognition with EigenFaces

➢ **Finding the most similar face from databases**

# Autoencoders

# What are Autoencoders (AE)?

➤ **It is the feed-forward neural network trained to reconstruct its input at the output layer**

- ◆ It has a **bottleneck layer**.
- ◆ The key purpose of AE is **dimensionality reduction**.
- ◆ This idea is also used for **learning generative data models**.



**Original data**  **Reconstructed data**

**Encoding**  **Decoding**

# What are Autoencoders (AE)?

➢ **Encoder-decoder architecture**

- ◆ **Encoder**: $\mathbf{z} = f(\mathbf{Wx}), \mathbf{W} \in \mathbb{R}^{k \times d}$
- ◆ **Decoder**: $\hat{\mathbf{x}} = g(\mathbf{W'z}), \mathbf{W'} \in \mathbb{R}^{d \times k}$
- ◆ Often, use the tied weight, i.e., $\mathbf{W}^{\mathrm{T}} = \mathbf{W'}$.

$$\hat{\mathbf{x}} = f\big(g(\mathbf{x})\big)$$

**Input layer**

**Hidden layer**

**Output layer**

$f(\mathbf{Wx})$

$g(\mathbf{W'z})$

$\mathbf{z}$

$\mathbf{X}$

$\tilde{\mathbf{x}}$

# Encoders in AE

> **Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data**

**Features z**

Originally: Linear + nonlinearity
Later: Deep, fully-connected
Later: CNN layers

**Encoder**

**Input data x**

# Decoders in AE

➢ **How to learn feature representation?**

- ◆ Train such that features can be used to **reconstruct** original data.
- ◆ Encoding itself, i.e., "**autoencoding**"

**Reconstructed input data $\hat{x}$**



Usually, **z** is smaller than **x**.
(dimensionality reduction)

**Decoder**

**Features z**

Originally: Linear + nonlinearity
Later: Deep, fully-connected
Later: CNN layers

**Encoder**

**Input data x**

# Feature Representation

> **How to learn feature representation?**

- Train such that features can be used to **reconstruct** original data.
- Encoding itself, i.e., "**autoencoding**"

**Reconstructed input data $\hat{\mathbf{x}}$**

**Decoder**

**Does not use labels!**

$$\|\mathbf{x} - \hat{\mathbf{x}}\|^2$$

**Features z**

**Encoder**

**Input data $\mathbf{x}$**

# Linear Autoencoders

➢ **The simplest autoencoder has one hidden layer, linear activations, and squared error loss.**

- This network computes $\tilde{\mathbf{x}} = \mathbf{W}'\mathbf{W}\mathbf{x}$, which is a linear function.
- If $k \geq d$, we can choose $\mathbf{W}$ and $\mathbf{W}'$ such that $\mathbf{W}'\mathbf{W}$ is the identity matrix. This is not useful.

$$\mathcal{L}(\mathbf{x}, \tilde{\mathbf{x}}) = \|\mathbf{x} - \tilde{\mathbf{x}}\|^2$$
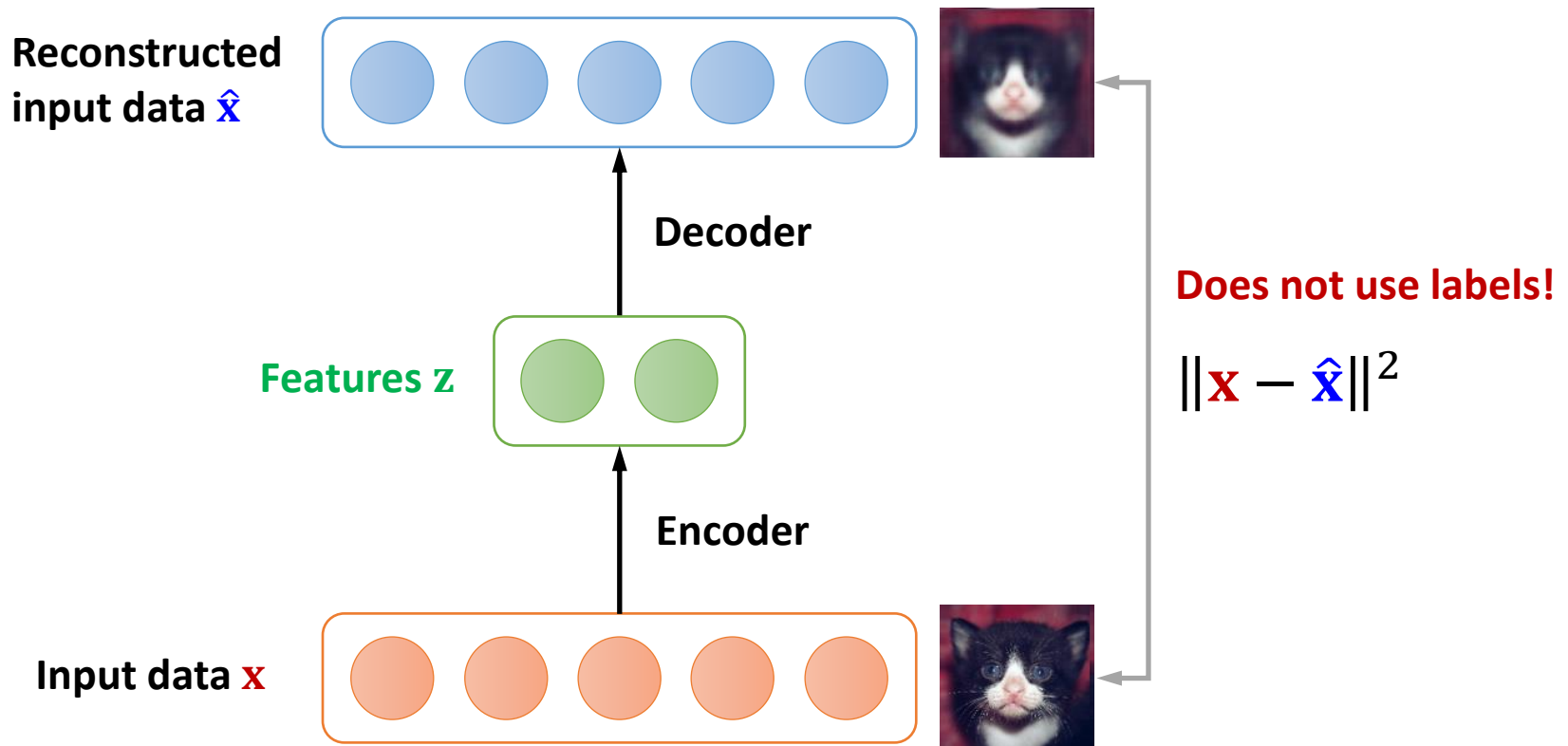
$\tilde{\mathbf{x}}$  | $d$ units |

$\mathbf{W}'$   **Decoder**

$\mathbf{z}$  | $k$ units |

$\mathbf{W}$   **Encoder**

$\mathbf{x}$  | $d$ units |

➢ **When $k < d$, $\mathbf{W}$ maps $\mathbf{x}$ to a $k$-dimensional space, it makes dimensionality reduction.**

# Linear Autoencoders

➤ **The output must lie in a $k$-dimensional subspace spanned by the columns of $\mathbf{W}'$.**

- ◆ The best possible k-dimensional subspace is the PCA subspace regarding the reconstruction error.

➤ **The autoencoder can achieve this by setting $\mathbf{W}' = \mathbf{W}^{\mathrm{T}}$.**

➤ **The optimal weights for a linear autoencoder are principal components!**

$$\min_{\mathbf{W}} \left\| \mathbf{X} - \mathbf{W}^{\mathrm{T}} \mathbf{W} \mathbf{X} \right\|^2$$

$\tilde{\mathbf{x}}$ | $d$ units

$\mathbf{W}' = \mathbf{W}^{\mathrm{T}}$  **Decoder**

$\mathbf{z}$ | $k$ units

$\mathbf{W}$  **Encoder**

$\mathbf{x}$ | $d$ units

# Example of Linear Autoencoders

➢ **Original image**

# Example of Linear Autoencoders

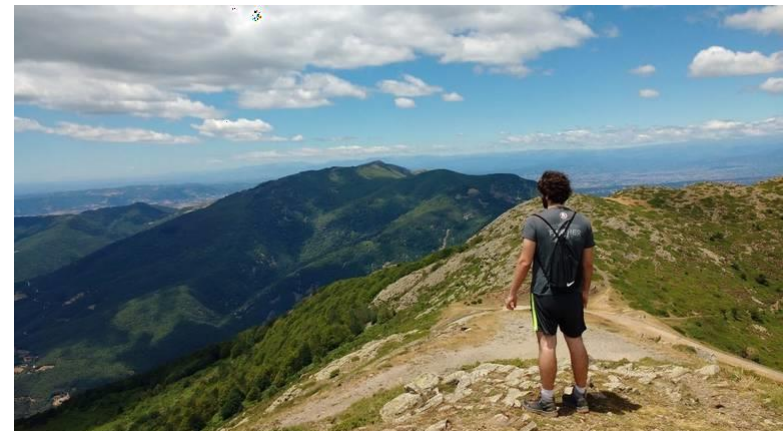**Image reconstruction using 2 principal components**



**Image reconstruction using 30 principal components**



**Image reconstruction using 200 principal components**



**Image reconstruction using 300 principal components**

# Nonlinear Autoencoders

➢ **Deep nonlinear autoencoders learn to project the data not onto a subspace but a nonlinear manifold.**

◆ This manifold is the image of the decoder.

➢ **This is called nonlinear dimensionality reduction.**

➢ **Nonlinear autoencoders can learn more powerful hidden vectors compared with linear autoencoders (PCA).**
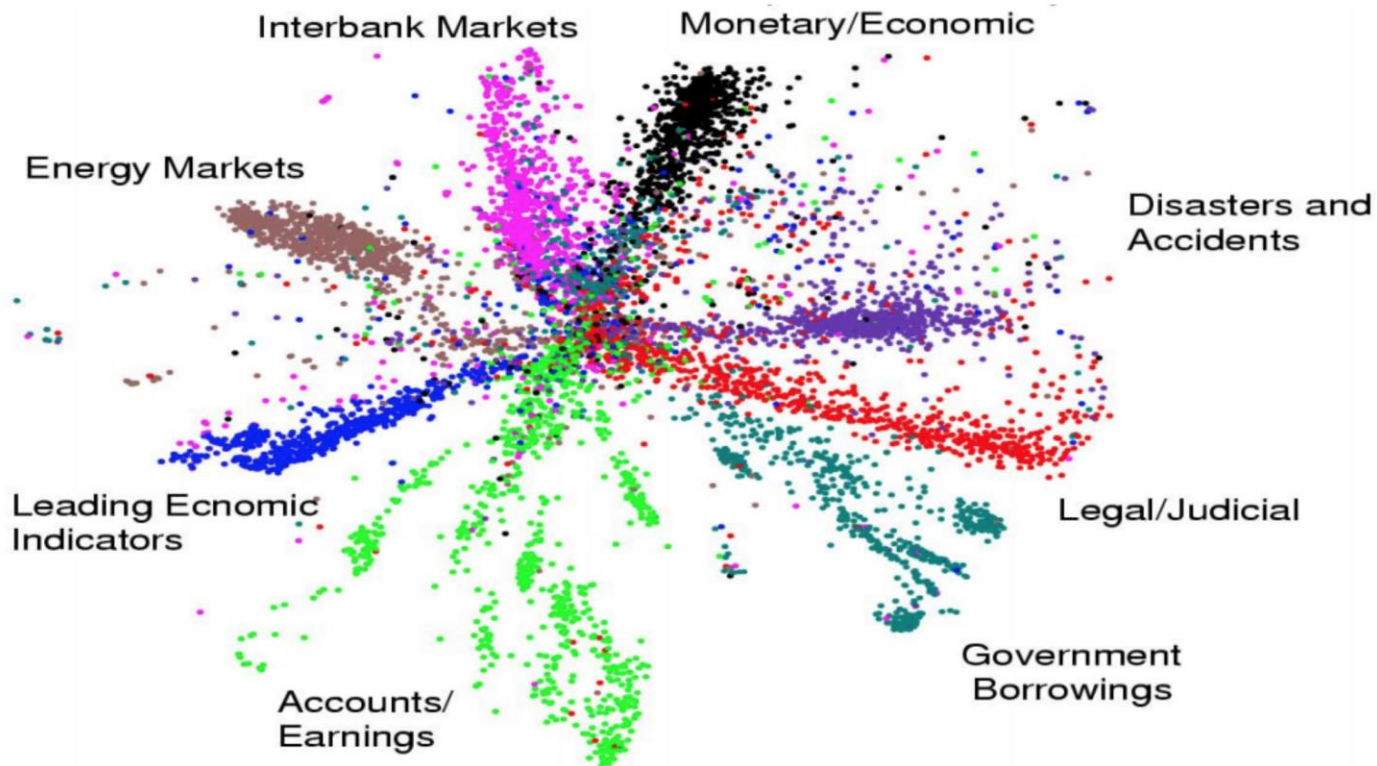


Real data

30-D deep AE

30-D PCA

# Example of Nonlinear Autoencoders

> **We use a two-dimensional autoencoder representation of newsgroup articles.**
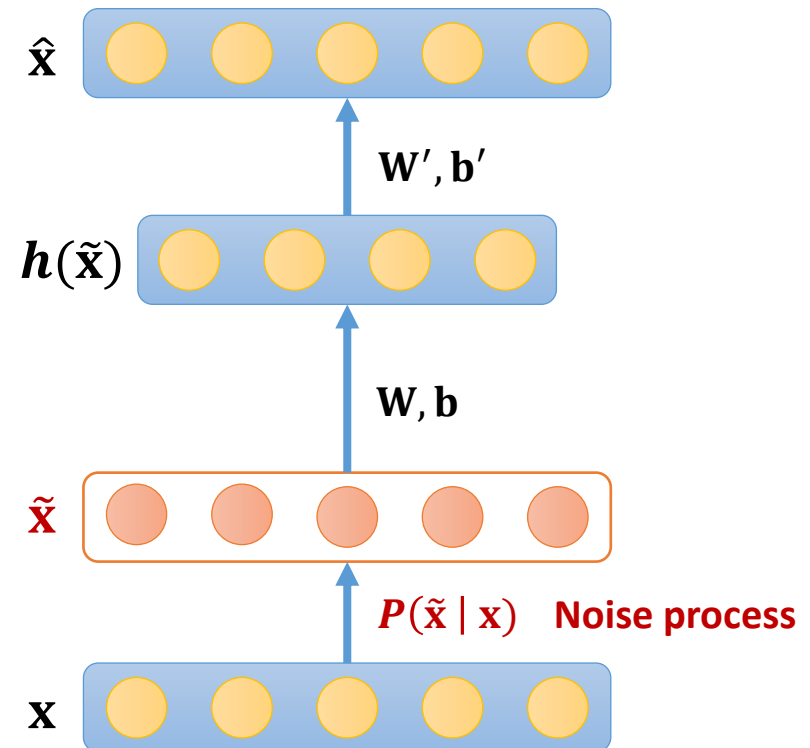> - Color-coded by topic without giving the labels

# Denoising Autoencoders (DAE)

➢ **Idea: Representation should be robust to noises.**

    ◆ Random assignment of subset of inputs to 0 with a probability

    ◆ Similar to dropouts on the input layer

➢ **Reconstruct $\hat{\mathbf{x}}$ from corrupted input $\tilde{\mathbf{x}}$**

➢ **Loss function compares reconstructed input $\hat{\mathbf{x}}$ with the noiseless input $\mathbf{x}$**

$\hat{\mathbf{x}}$

$\mathbf{W}', \mathbf{b}'$

$h(\tilde{\mathbf{x}})$

$\mathbf{W}, \mathbf{b}$

$\tilde{\mathbf{x}}$

$P(\tilde{\mathbf{x}} \mid \mathbf{x})$    **Noise process**

$\mathbf{x}$

# Denoising Autoencoders (DAE)

➢ **Add some noise to the input $\tilde{x} = x + \epsilon$.**
  - ◆ Ask the autoencoder to reconstruct the original input.
  - ◆ Learn more generalized representation.



**Reconstructed input data $\hat{x}$**

**Decoder**

**Features z**

**Encoder**

**Corrupted input data $\tilde{x}$**

# Properties of Autoencoders

➤ **Autoencoders are used for visualization**

➤ **Autoencoders are used for <span style="color:red">data compression.</span>**
- ◆ It is similar to PCA.
- ◆ In fact, linear autoencoders learn the same subspace as PCA.

➤ **Autoencoders can be used for <span style="color:red">pre-training</span>**
- ◆ No need for labeled data

# Q&A