



أكاديمية كاوست
KAUST ACADEMY

Day 4: Unsupervised Learning Summary

Note: Do Not depend entirely on this and study from the official slides

Contributed by: Hassan Mohammed Nasr

- In unsupervised learning, there is no label in the data, the algorithm should find hidden pattern and structure in the data.
- We use unsupervised learning as most of the data is unlabeled, and labeling them requires human effort, time and expertise.

Unsupervised Applications

Clustering

- Group similar data points together

Dimensionality Reduction

- Compress high dimensional data to lower dimensional

Generation

- Create New, realistic data samples

Anomaly Detection

- Identify unusual data points

1 - Clustering: K-Means

- Grouping data into K distinct groups (clusters) based on distance
 1. **Initialize** : Randomly place K clusters centers.
 2. **Assign**: for each data point, assign to the nearest center
 3. **Update step**: move each cluster's centroid to the cluster's mean
 4. **Repeat** 2 and 3 until assignments don't change or max iterations reached
- K-means assumes Spherical clusters, sensitive to outliers, sensitive to random initialization, bad results for high dimensional data, and need to specify K (number of clusters) in advance (we can use elbow method to approximate the number of clusters)

2- Dimensionality Reduction: PCA

- High dimensional data is sparse, noisy and hard to visualize
- **PCA is a dimensionality reduction** algorithm that finds the most important directions in your data by creating new features (principal components) that
 1. Capture maximum variance in the data
 2. Uncorrelated with each other
 3. Ranked by importance

How Does PCA work?

1. Center Data (zero mean)	2. Covariance Matrix	3. Eigenvalues and Eigen Vectors	4. Select and project
<ul style="list-style-type: none">• Shift data so they have zero mean	<ul style="list-style-type: none">• Calculate how features vary with respect to each other	<ul style="list-style-type: none">• Eigenvalues: The new direction• Eigenvector: The amount of Variance	<ul style="list-style-type: none">• Select top K eigenvectors

PCA are linear only, sensitive to scale (must standardize first), non interpretable and outliers can distort PCA



Most data are high dimensional that have nonlinear structure that can be hard to visualize and plot for exploration



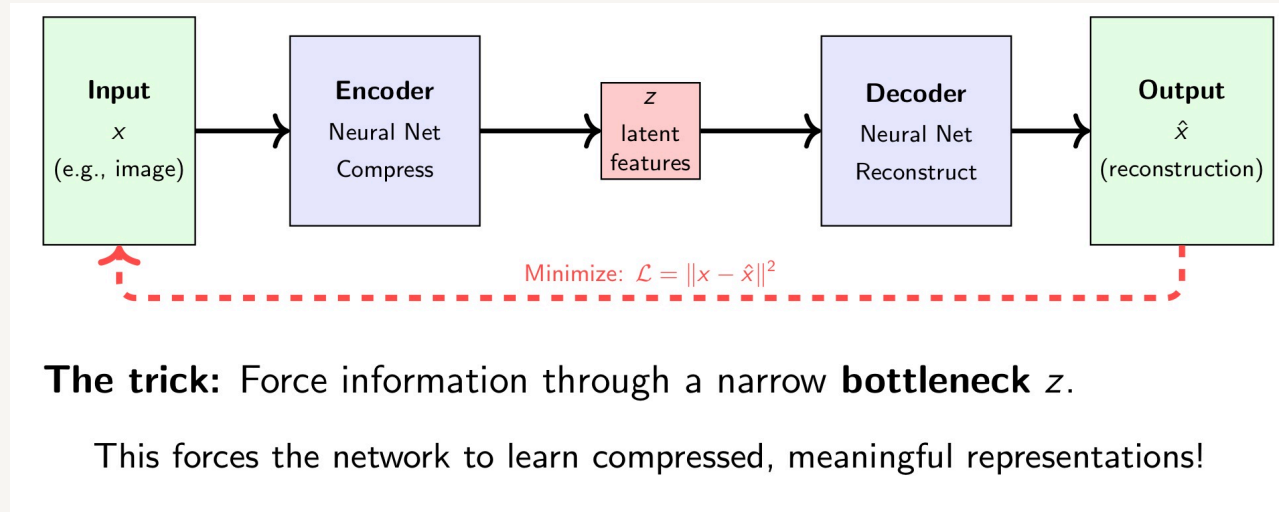
Solution (**t-sne**): tries to build 2d map where neighbors in high d remain neighbors in 2d.

3- Autoencoders

- **Autoencoders** is a type of neural network that learns to:

↳ **Encode(encoders)**: Compress data into a compact form and then

↳ **Decode(decoder)**: Reconstruct it to match the original input



The trick: Force information through a narrow **bottleneck** z.

This forces the network to learn compressed, meaningful representations!

Autoencoders Applications

Application.

- Dimensionality Reduction
- Image Generation
- Denoising
- Anomaly Detection

How it works

Train full network, then discard the decoder, use encoder to map input x to z (low dimension)

Train full network, discard the encoder, sample random vector from z then use this vector as input of decoder to generate new images

Take clean image as ground truth add to it noise and feed it to the full Autoencoders to map the noisy image to the clean ground truth

The Network was trained to map normal data with low error, so when an outlier come, the Autoencoder will output a high reconstruction error indicating this is outlier