# Lecture 2

# Working with High-Dimensional Data

Stan Z. Li
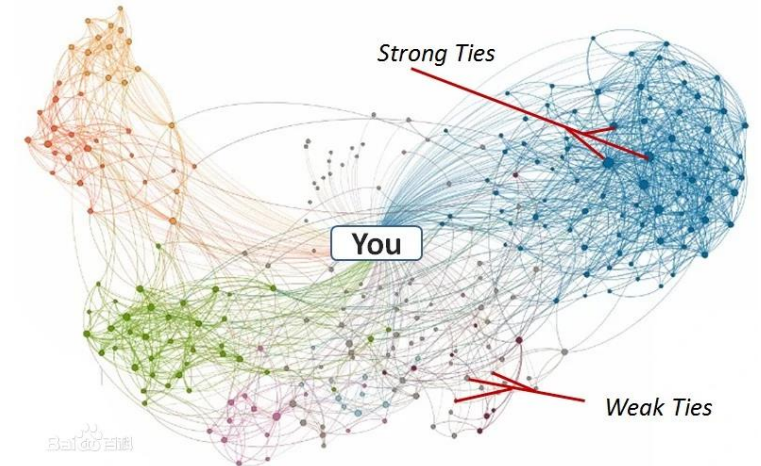
西 湖 大 學
**WESTLAKE UNIVERSITY**

# Outline

1. **High-dimensional data**

2. **Lower-dimensional patterns/manifolds**

3. **Representational learning/dimension reduction**

   - **Linear projection**

   - **Nonlinear projection/neural networks transformation**
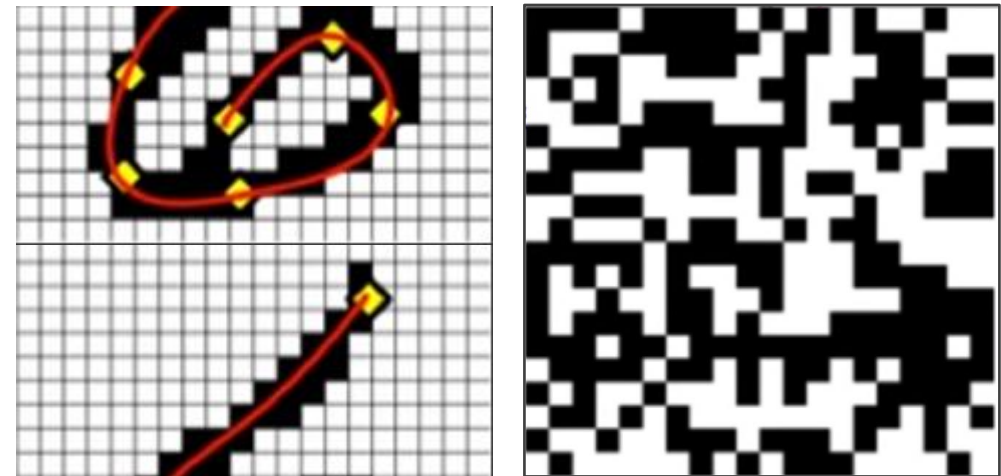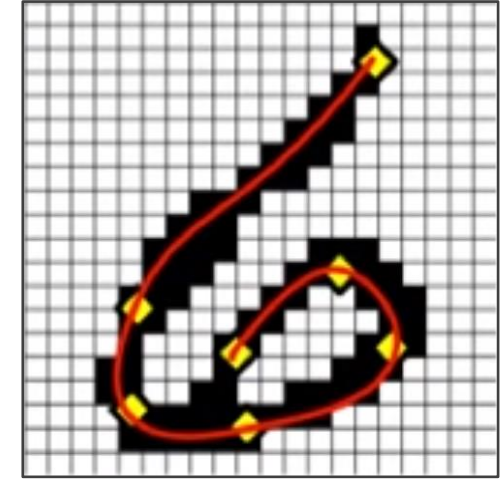
# High-Dimensional Data

- Images, Videos, Text, Audio,

- Web pages, Social Networks

- Molecular Structures

- DNA Sequences
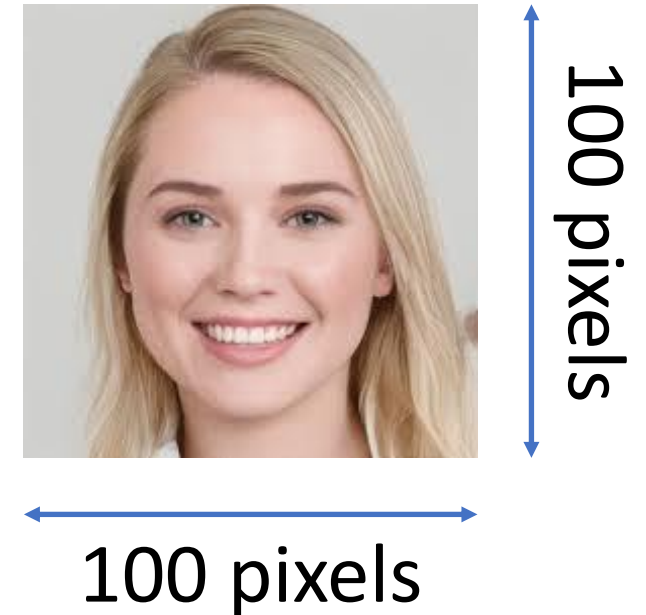
- Protein Sequence-Structures

# Handwritten Digit images

- **Image size 20x20 = 400**

- **Pixel values in {0,1}**

- **Image Space $\mathbf{S}$ = {0,1} $^{400}$**

- **#$\mathbf{S}$ = $2.58 \times 10^{120}$**

- **Only a tiny portion of $\mathbf{S}$ is of digits**

- **The digit pattern lives in a**

   **low dim subspace (manifold)**

# Face Image Data

- Image size 100x100 = $10^4$ pixels
- RGB image size $3 \times 10^4$ pixels
- **Dimensionality = $3 \times 10^4$**
- Pixel values in {0,...,255}
- #Possibility = $256^{30,000} \cong$ infinity
- Only a tiny portion is of faces
- **Face pattern lives in low dim subspace**
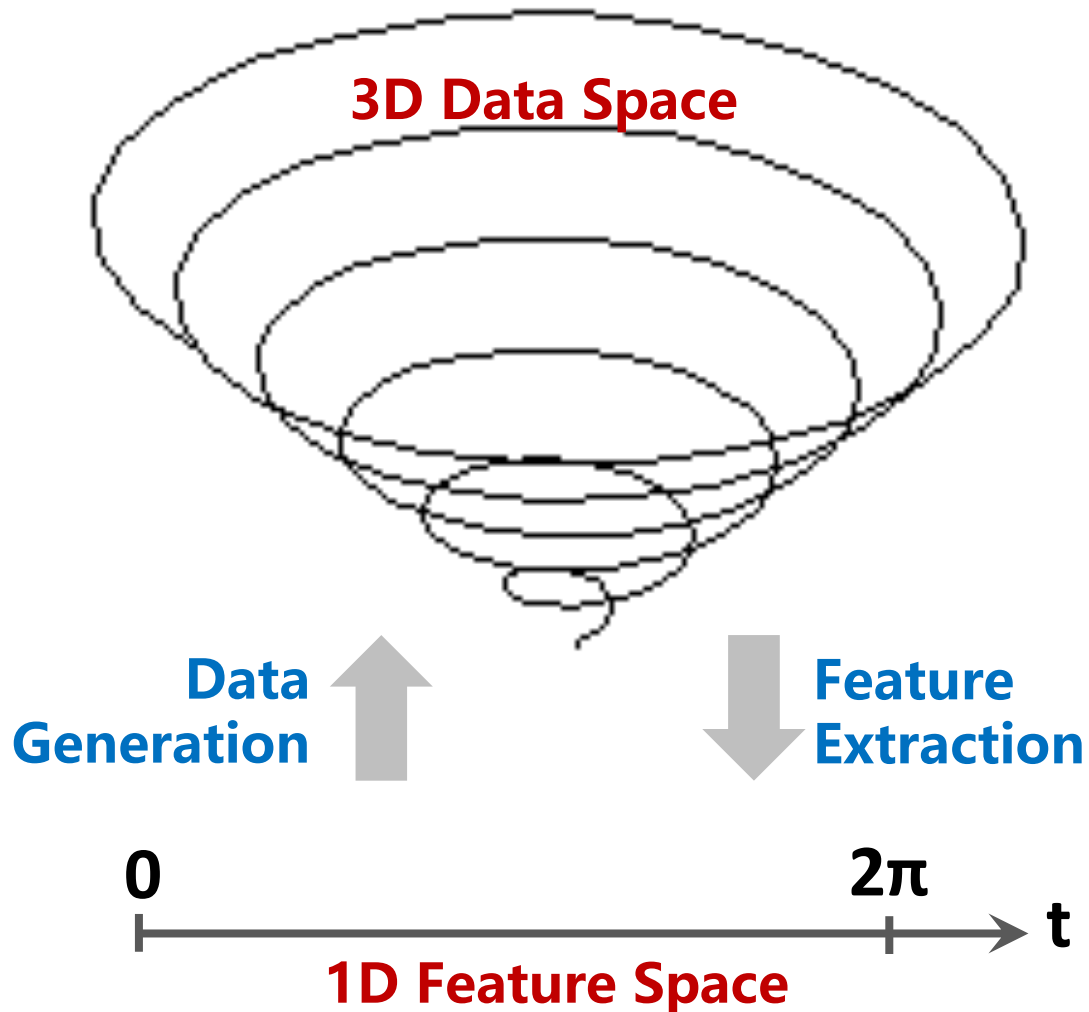


100 pixels

100 pixels

# Manifold Assumption

**High-Dimensional Data**: Images, Web pages, Gene sequences, ….

**Dimension Reduction into Coordinate System of a Lower Dim**

- **For representation learning（feature extraction）**

- **For data visualization – in 2D or 3D**

**Manifold Assumption: an interesting pattern in high dimensional data resides on a low dimensional manifold**
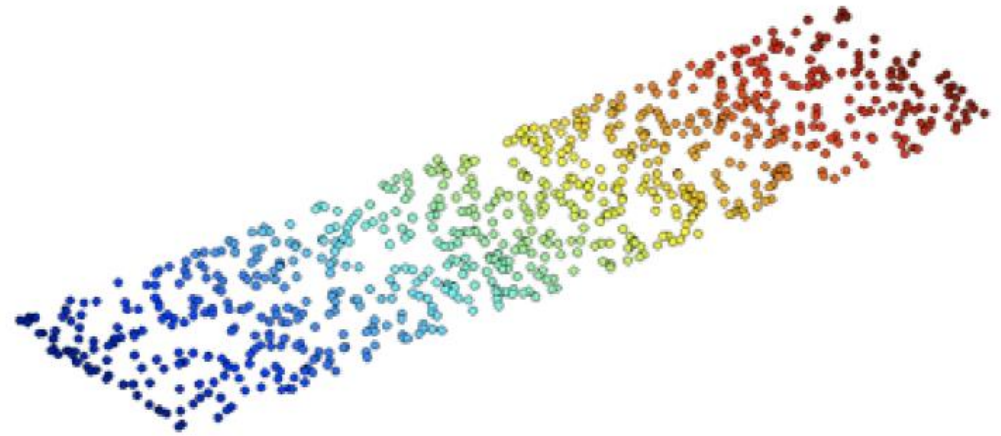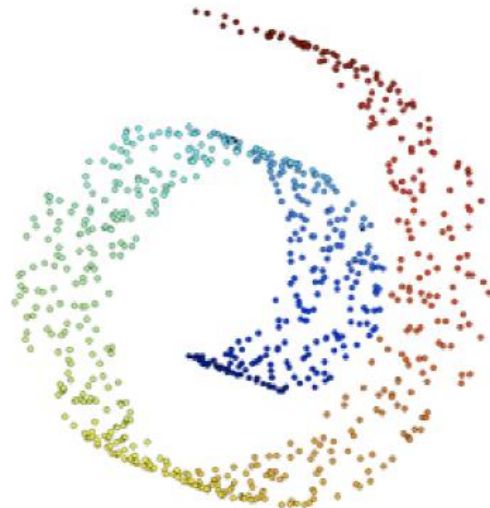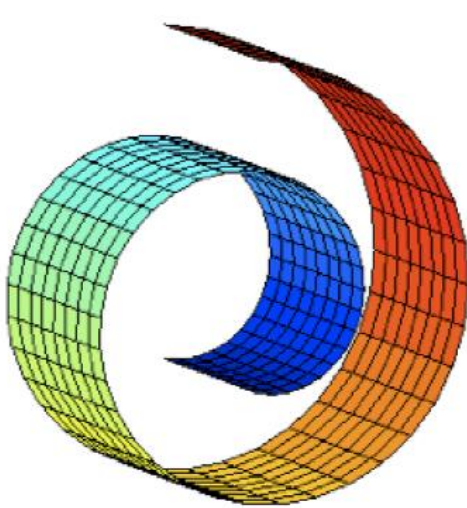
# Manifold in Hi-D Data Space: 1D Curve in 3D Space



**3D Data Space**

**Data Generation**

**Feature Extraction**

**0**   **2π**   **t**

**1D Feature Space**

**Conical Helix:**

x=t*cos(6t), y=t*sin(6t), z=t

0 ≤ t ≤ 2π

**1D line segment**

Latent variable t

# 2D Manifold in 3D Space



**Swiss Roll:**

x=φcos(φ), y=φsin(φ), z=ψ
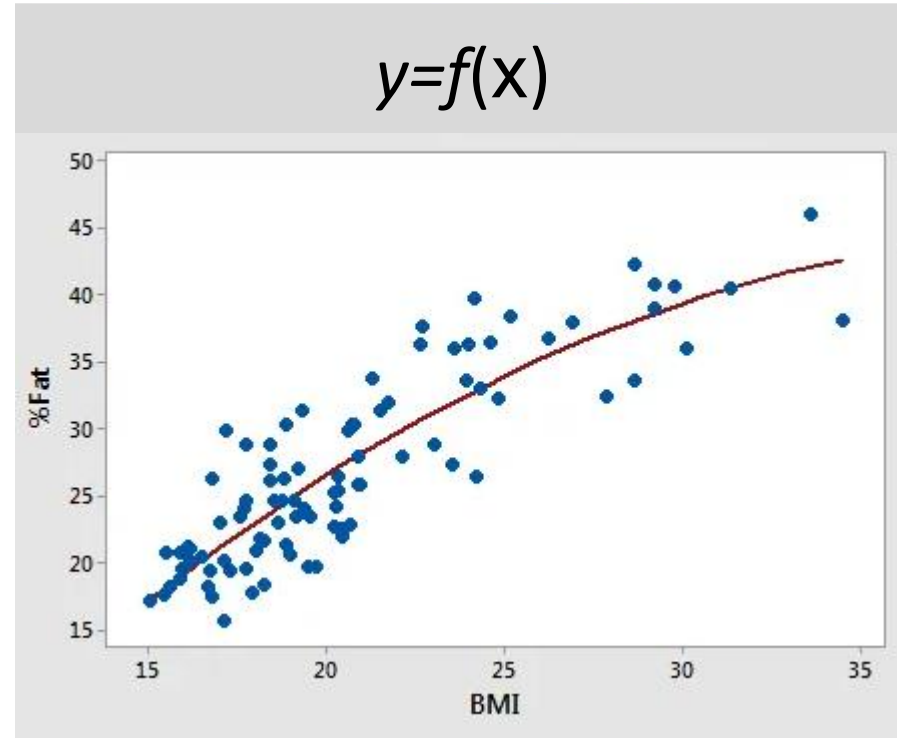
1.5π ≤ φ ≤ 4.5π,  0 ≤ ψ ≤ 10

**Manifold: 2D rectangle**

generated by two latent

variables φ, ψ

# Geodesic Distance on Manifolds
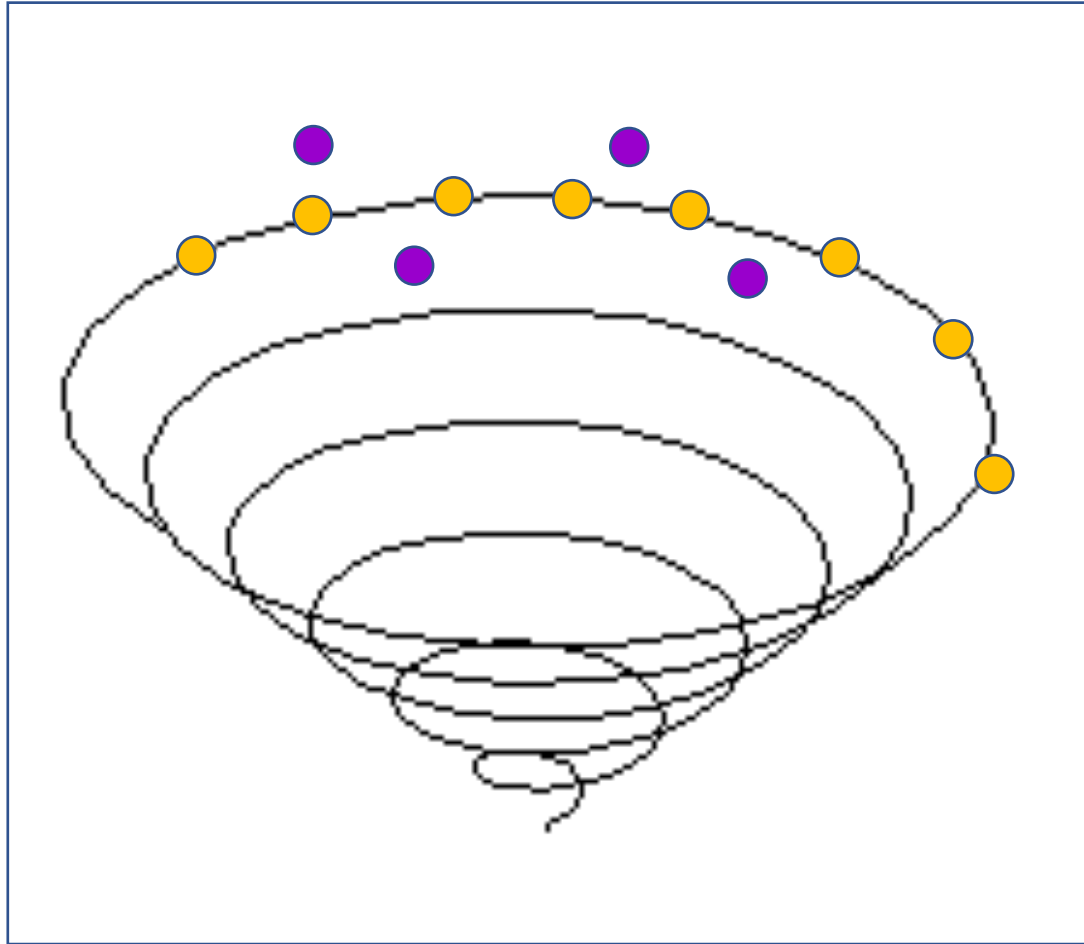
# Data Samples on Manifold
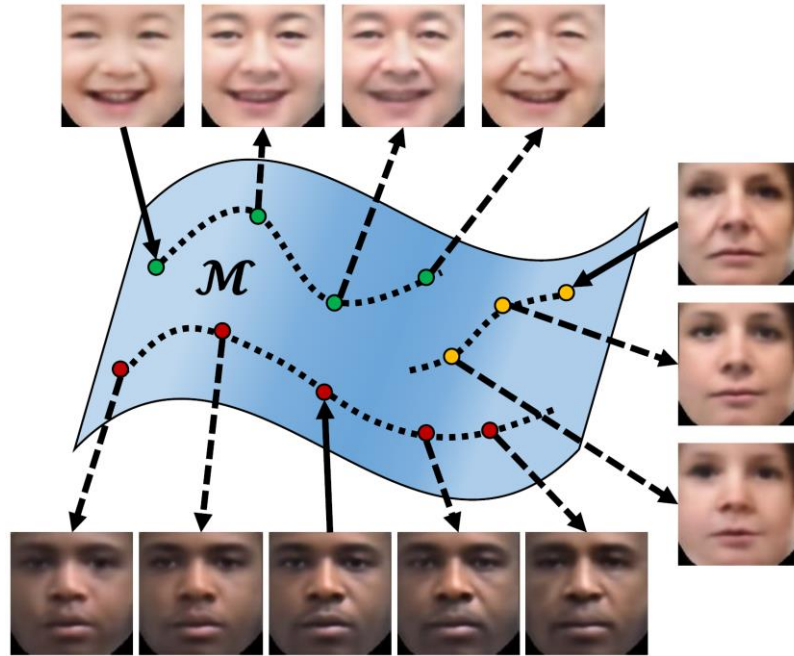
$y=f(\mathbf{X})$ sampled to $\{(\mathbf{X}_i , y_i) \mid i = 1,...,n\}$

# Samples on Face Manifold in Data Space

# Samples Close to the Face Manifold
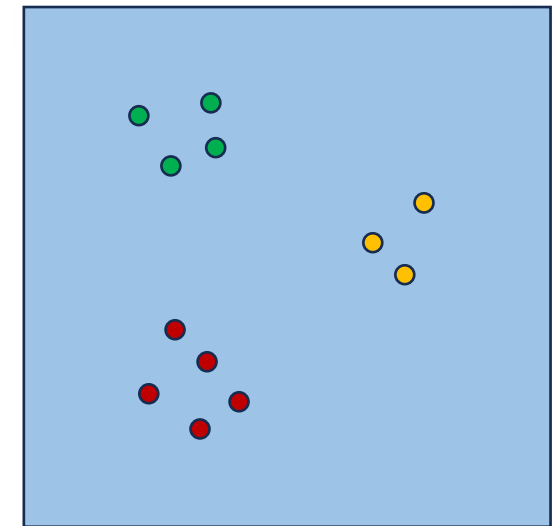
# Low-Dim Manifold/Surface in High-Dim Space



**Feature Extraction By Encoder Neural Network**

**Data Generation by Decoder Neural Network**

Samples on low-dim but complex manifold in high-dim data space

Features in lower-dim Euclidean embedding space

# Scientific Modeling