

Lecture 6: Linear Dimensionality Reduction

Shukai Gong

1 Curse of Dimensionality

We have mentioned in Section ?? that high-dimensional data has stronger data representation ability and bring about more potentials to our model. However, high-dimensional data also has some drawbacks.

Combinatorial Explosion (Discrete Example)

The number of different d -dimensional binary vectors is 2^d . Each additional dimension doubles the effort needed to try all combinations.

(The search space of chess)

High Dimensional Sampling (Continuous Example)

Sampling N samples randomly from a d -dimensional sample space, based on a distribution with identity covariance matrix $\Sigma = \mathbf{I}_d$, it can be proved that the smallest and largest Euclidean distance $d_{\min}(D)$ and $d_{\max}(D)$ between any two samples satisfy

$$\frac{d_{\min}(D)}{d_{\max}(D)} \approx \left(\frac{1}{2}\right)^{\frac{1}{d}} \rightarrow 1, d \rightarrow \infty$$

meaning that maximum distance becomes indiscernible compared to the minimum distance. **Euclidean distance functions are losing their usefulness in high dimensions.**

One method of dealing with the curse of dimensionality is **Dimensionality Reduction**, which is to map the high-dimensional data to a lower-dimensional space while preserving the essential structure of the data. Mathematically, we would like to find a **linear projection**

$$\begin{aligned} f : \mathbb{R}^D &\mapsto \mathbb{R}^L \\ \mathbf{x} &\mapsto \mathbf{z} = f(\mathbf{x}) \end{aligned}$$

where $L \ll D$. Generally, $\mathbf{z} = f(\mathbf{x}) = \mathbf{U}^\top \mathbf{x}$, where $\mathbf{U} \in \mathbb{R}^{D \times L}$. A desired projection matrix can be obtained by **minimizing reconstruction error**. We want to see a small error when remapping \mathbf{z} back to \mathbf{x} , i.e.

$$\begin{aligned} \exists g : \mathbb{R}^L &\mapsto \mathbb{R}^D \\ \mathbf{z} &\mapsto \mathbf{x} \approx g(f(\mathbf{x})) \end{aligned}$$

Another method is to consider **isometry**, which is to preserve the pairwise distance between data points in \mathcal{X} and \mathcal{Z} , i.e.

$$d_{\mathcal{Z}}(f(\mathbf{x}_i), f(\mathbf{x}_j)) \approx d_{\mathcal{X}}(\mathbf{x}_i, \mathbf{x}_j)$$

2 Linear Dimensionality Reduction

2.1 Principal Component Analysis (PCA)

Recall that in LASSO-based **supervised** feature selection, we solve the optimization problem

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

For $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D]$, the column \mathbf{x}_d contributes to the estimation of \mathbf{y} iff $\mathbf{w}_d \neq 0$. That's to say, the useful features are those with non-zero weights.

$$\hat{\mathbf{X}} = [\mathbf{x}_d]_{d:\mathbf{w}_d \neq 0} \in \mathbb{R}^{N \times L}, \quad \hat{\mathbf{w}} = [\mathbf{w}_d]_{d:\mathbf{w}_d \neq 0} \in \mathbb{R}^L, L < D$$

However, labeled training data \mathbf{y}, \mathbf{X} might be hard or expensive to get, but unlabeled training data \mathbf{X} might be more easily available. **PCA is able to extract meaningful directions from such unlabeled data.**

PCA

Principal Component Analysis (PCA) is an **unsupervised dimensionality reduction** technique. Given a matrix of data points $\mathbf{X} \in \mathbb{R}^{N \times D}$, it finds one or more **orthogonal directions** \mathbf{v}_i that capture **the largest amount of variance** in the data.

Intuitively, the directions with less variance contain less information and may be discarded without introducing too much error.

The Principal of PCA is to sequentially find the projections maximizing the preserved energy of data (Minimizing the residual that cannot be captured by the projections). Therefore, we start from the first component \mathbf{v}_1 .

Derivation of Principle Components

The first L principal components of a collection of data points $\mathbf{X} \in \mathbb{R}^{N \times D}$ are **the eigenvectors corresponding to the largest L orthonormal eigenvalues of the $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{D \times D}$.**

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix with N rows of D -dimensional data. \mathbf{x}_i^\top are considered to be i.i.d. samples from some random vector \mathbf{x} .

First, we process the data to be columnwise zero-mean, i.e. $\mathbf{X}^\top \mathbf{1}_N = \mathbf{0}_D$. This is achieved by subtracting the average of all the rows, i.e. $\bar{\mathbf{x}}^\top = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^\top$ from each row.

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} \xrightarrow{\text{zero-meaned}} \mathbf{X}' = \begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1D} - \bar{x}_D \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots & x_{2D} - \bar{x}_D \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} - \bar{x}_1 & x_{N2} - \bar{x}_2 & \cdots & x_{ND} - \bar{x}_D \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^\top - \bar{\mathbf{x}}^\top \\ \mathbf{x}_2^\top - \bar{\mathbf{x}}^\top \\ \vdots \\ \mathbf{x}_N^\top - \bar{\mathbf{x}}^\top \end{bmatrix} = \mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^\top$$

Still we denote \mathbf{X} as the zero-meaned data matrix for convenience.

Second, since \mathbf{X} is zero-meaned, the sample variance of the **datapoints' projections** onto a **unit vector** \mathbf{v} is given by

$$\begin{aligned} \text{Var}(\mathbf{x}^\top \mathbf{v}) &= \mathbb{E}[(\mathbf{x}^\top \mathbf{v})^2] - \mathbb{E}[\mathbf{x}^\top \mathbf{v}]^2 = \mathbb{E}[(\mathbf{x}^\top \mathbf{v})^2] \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{v})^2 = \frac{1}{N} \|\mathbf{X}\mathbf{v}\|^2 = \frac{1}{N} \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} \end{aligned}$$

With the motivation of maximizing the variance, we have the optimization problem

$$\max_{\mathbf{v}} \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}, \quad \text{s.t. } \mathbf{v}^\top \mathbf{v} = 1$$

The Lagrangian of the optimization problem is

$$\begin{aligned} \mathcal{L}(\mathbf{v}, \lambda) &= \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} + \lambda(\mathbf{v}^\top \mathbf{v} - 1) \\ \Rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{v}} &= 2\mathbf{X}^\top \mathbf{X} \mathbf{v} - 2\lambda \mathbf{v} = 0 \Rightarrow \mathbf{X}^\top \mathbf{X} \mathbf{v} = \lambda \mathbf{v} \end{aligned}$$

that the first component \mathbf{v}_1 satisfies $\mathbf{X}^\top \mathbf{X} \mathbf{v}_1 = \lambda \mathbf{v}_1$, i.e. \mathbf{v}_1 is an eigenvector of $(\mathbf{X}^\top \mathbf{X})_{D \times D}$. Since our constraint is $\mathbf{v}^\top \mathbf{v} = 1$, we have

$$\begin{aligned} \mathbf{v}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}_1 &= \lambda \mathbf{v}_1^\top \mathbf{v}_1 = \lambda \\ \Rightarrow \max_{\mathbf{v}} \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} &= \lambda_{\max}(\mathbf{X}^\top \mathbf{X}) \end{aligned}$$

So the first component \mathbf{v}_1 is the orthonormal eigenvector corresponding to the largest eigenvalue of $\mathbf{X}^\top \mathbf{X}$. Furthermore, it can be proved that the i -th component \mathbf{v}_i is the eigenvector corresponding to the i -th largest eigenvalue of $\mathbf{X}^\top \mathbf{X}$.

Therefore, **The first L principal components of a collection of data points are the eigenvectors corresponding to the largest L orthonormal eigenvalues of the $\mathbf{X}^\top \mathbf{X}$.** It can also be observed that The first L principal components are the top- L columns of \mathbf{V} , the right-singular matrix of \mathbf{X} . Therefore, we can obtain PCs by performing SVD on \mathbf{X} .

Once we have computed the principal component \mathbf{v}_i , we can use them as a **new coordinate system**. The k -th principal component of a D -dimensional datapoint $\mathbf{x}_i \in \mathbb{R}^D$ is $\mathbf{x}_i^\top \mathbf{v}_k$, the scalar projection of \mathbf{x}_i onto the k -th principal component \mathbf{v}_k . Denote $\mathbf{V}_{D \times L} = [\mathbf{v}_1, \dots, \mathbf{v}_L]$, we can compute all the PCs of all the datapoints by

$$\mathbf{Z}_{N \times L} = \mathbf{X}_{N \times D} \mathbf{V}_{D \times L}, \quad L < D$$

We have successfully reduced the dimensionality of the data from D to L ! **A reconstruct of \mathbf{X} from \mathbf{Z} will be**

$$\hat{\mathbf{X}}_{N \times D} = \mathbf{Z}_{N \times L} \mathbf{V}_{L \times D}^\top = \mathbf{X}_{N \times D} \mathbf{V}_{D \times L} \mathbf{V}_{L \times D}^\top$$

2.2 Whitening

Let's revisit Whitening. Given $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_D]$ with D features, we similarly process the data to be column-wise zero-meaned:

$$\mathbf{X}' = \mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^\top$$

where $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$. Then the estimate covariance matrix measuring the covariance between feature i and j is given by

$$\hat{\Sigma} = \frac{1}{N-1} (\mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^\top)^\top (\mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^\top) = \frac{1}{N-1} \mathbf{X}'^\top \mathbf{X}'$$

Then we can **whiten the data** by

$$\tilde{\mathbf{X}} = (\mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^\top) \hat{\Sigma}^{-\frac{1}{2}} = \mathbf{X}' \hat{\Sigma}^{-\frac{1}{2}}$$

This is very similar to PCA. By EVD we know

$$\mathbf{X}'^\top \mathbf{X}' = \mathbf{V} \Lambda \mathbf{V}^\top = (N-1) \hat{\Sigma}$$

Therefore

$$(N-1) \tilde{\mathbf{X}} = (N-1) \mathbf{X}' \hat{\Sigma}^{-\frac{1}{2}} = \underbrace{\mathbf{X}'}_{\text{Shifting}} \underbrace{\mathbf{V}}_{\text{PCA Matrix}} \underbrace{\Lambda^{-\frac{1}{2}}}_{\text{Scaling}} \mathbf{V}^\top$$

2.3 Data Denoising by PCA

Data Denoising by PCA

We claim that PCA is least-square data denoising in statistical ML. Suppose an i.i.d Gaussian noise model for observed data \mathbf{X}

$$\mathbf{X}_{\text{noisy}} = \mathbf{X}_{\text{clean}} + \mathbf{E}, \mathbf{E} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

The least square data denoising problem is denoted as

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X}} \|\mathbf{X}_{\text{noisy}} - \mathbf{X}\|_F^2$$

When the feasible domain corresponds to a low-rank constraint

$$\Omega := \{\mathbf{X} \in \mathbb{R}^{N \times D} : \text{rank}(\mathbf{X}) \leq L\}$$

We have the **closed-form solution**

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \Omega} \|\mathbf{X}_{\text{noisy}} - \mathbf{X}\|_F^2 = \mathbf{U}_L \Sigma_L \mathbf{V}_L^\top, \text{ where } \mathbf{X}_{\text{noisy}} = \mathbf{U} \Sigma \mathbf{V}^\top$$

where Σ_L is the top- L singular values matrix and $\mathbf{U}_L, \mathbf{V}_L$ contains only the top- L singular vectors.

Intuitively, noises and unimportant information are usually captured by the last few principal components. Therefore, we only keep the first L principal components to denoise the data.

2.4 Other projection and factorization models

Robust PCA

Consider sparse noise

$$\mathbf{X}_{\text{noisy}} = \mathbf{X}_{\text{clean}} + \mathbf{E}, \mathbf{E} \sim \mathcal{L}(0, \sigma^2 \mathbf{I})$$

Still we have the feasible domain corresponds to a low-rank constraint

$$\Omega := \{\mathbf{X} \in \mathbb{R}^{N \times D} : \text{rank}(\mathbf{X}) \leq L\}$$

The denoised data is given by

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \Omega} \|\mathbf{X}_{\text{noisy}} - \mathbf{X}\|_1$$

Non-negative Matrix Factorization (NMF)

The denoised data is given by

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \Omega} \|\mathbf{X}_{\text{noisy}} - \mathbf{X}\|_F^2$$

where the feasible domain is

$$\Omega := \{\mathbf{X} = \mathbf{U}\mathbf{V}^\top : \text{rank}(\mathbf{X}) = L, \mathbf{U}, \mathbf{V} \geq 0\}$$

Subspace Clustering

Subspace clustering is similar to LASSO, but the supervised signal \mathbf{y} is the data itself:

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{D \times D}} \|\mathbf{X} - \mathbf{X}\mathbf{W}\|_F^2 + \lambda \|\mathbf{W}\|_1 + \lambda_2 \|\mathbf{W}\|_*$$

where $\|\mathbf{W}\|_*$ is the nuclear norm of \mathbf{W} , i.e. the sum of its singular values.

Before introducing **Compressive Sensing**, let's first look at an example of image recovering.

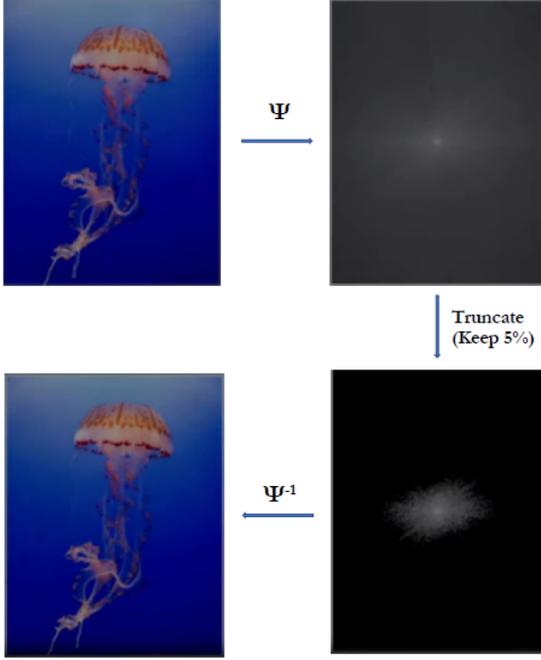


Image 1: Showing traditional Image compressing and decompressing using a basis function Ψ over a High Resolution-Image.

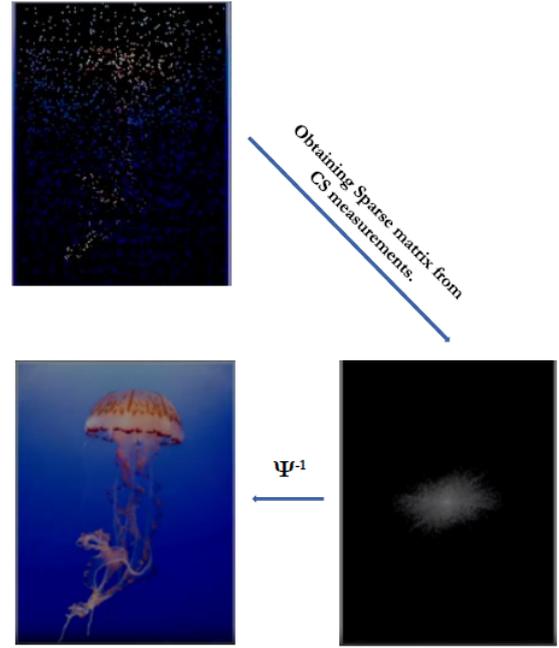


Fig 2: Recovering Original Image \mathbf{x} by sampling few pixels ($\Phi \mathbf{x}$) while sensing and solving for the sparse representation over a basis Ψ .

Suppose figure \mathbf{x} is Fourier transformed Φ to frequency domain \mathbf{s} , and after truncating the signal in \mathbf{s} , we still can expect to reconstruct a high quality image after the removal of some high frequency signals of lower energy.

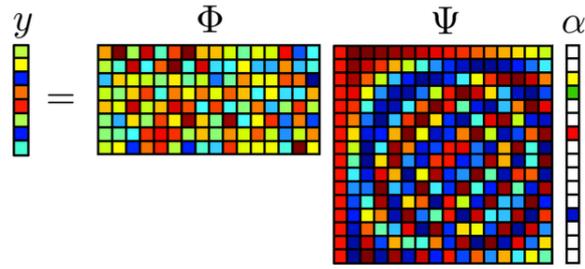
However, in practical, we want to reconstruct a high-quality original image \mathbf{x}' using a low-quality observation \mathbf{y} , i.e., an undersampling of the true image \mathbf{x} , by transforming it to get a sparse space \mathbf{s} of \mathbf{x} . Mathematically

Compressive Sensing

For an unknown signal $\mathbf{x}_{N \times 1}$, suppose it can be sparsely represented in sparse basis $\Psi_{N \times K}$, i.e. $\mathbf{x} = \Psi \mathbf{s}$, where coefficient vector $\mathbf{s}_{K \times 1}$ is sparse. We expect to recover \mathbf{s} from observation $\mathbf{y}_{M \times 1}$. The optimization problem is

$$\begin{aligned} \hat{\mathbf{s}} &= \arg \min_{\mathbf{s}} \|\mathbf{y}_{M \times 1} - \Phi_{M \times N} \Psi_{N \times K} \mathbf{s}_{K \times 1}\|_2^2 + \lambda \|\mathbf{s}\|_1 \\ &= \arg \min_{\mathbf{s}} \|\mathbf{y}_{M \times 1} - \Theta_{M \times K} \mathbf{s}_{K \times 1}\|_2^2 + \lambda \|\mathbf{s}\|_1 \end{aligned}$$

where Φ is a **random** measurement matrix, Ψ is the sparse basis matrix, and \mathbf{s} is the sparse coefficient vector. Specifically, $\Theta_{M \times K} = \Phi_{M \times N} \Psi_{N \times K}$ is named as sensing matrix.



If $M = N$ here, we can directly solve $\mathbf{y} = \Theta \mathbf{s}$. But normally $M < N$, there are infinitely many solutions to $\mathbf{y} = \Theta \mathbf{s}$. Therefore, we add a L_1 regularization term to the optimization problem to find the sparsest solution.

Furthermore, when Φ satisfies the **Restricted Isometry property (RIP)**,

Restricted Isometry Property (RIP)

A matrix Φ satisfies the RIP with restricted isometry constant $\delta < 1$ if for all K -sparse vectors $\mathbf{x}_1, \mathbf{x}_2$, i.e. $\forall \mathbf{x}_1, \mathbf{x}_2 \in \{\mathbf{x} \in \mathbb{R}^N : \|\mathbf{x}\|_0 \leq S\}$, we have

$$(1 - \delta) \leq \frac{\|\Phi \mathbf{x}_1 - \Phi \mathbf{x}_2\|_2^2}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2} \leq (1 + \delta)$$

we have

- **Stable Recovery:** When the rows of sensing matrix $M \geq \frac{1}{C} \cdot S \cdot \log \frac{N}{S}$ leads to an exact reconstruction with probability $1 - \mathcal{O}(N^{-M})$.
- Solving $\min_{\mathbf{s}} \|\mathbf{y} - \Theta \mathbf{s}\|_2^2 + \lambda \|\mathbf{s}\|_1$ is equivalent to solving $\min_{\mathbf{s}} \|\mathbf{y} - \Theta \mathbf{s}\|_2^2 + \lambda \|\mathbf{s}\|_0$

Commonly used sensing matrices Θ include **Sub-Gaussian matrices (Gaussian, Bernoulli)** and **Fourier matrices**.