

Principal Component Analysis

$$X = (N \times D)$$

examples features.

Each $x_i \in \mathbb{R}^D$

$$X \rightarrow Z_{N \times M}$$

$$M \ll D$$

Why DR?

① Reduces data size

→ Storage

→ Compute

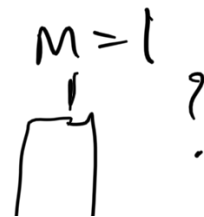
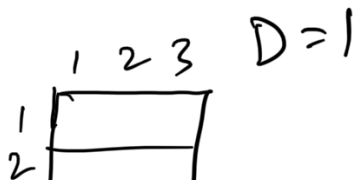
→ Transfer

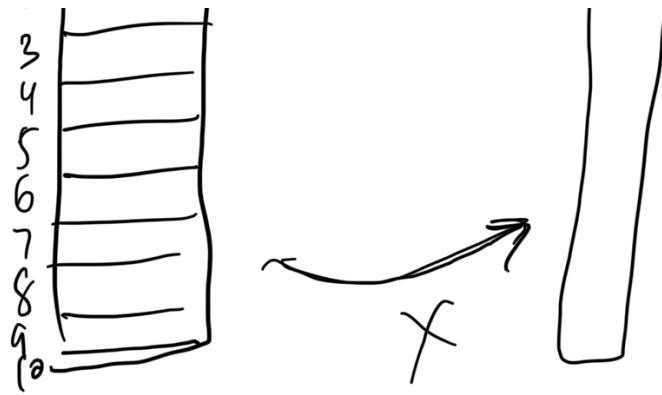
$$O(f(D))$$

② Visualization

③ Improves data

How?





DR \Rightarrow loss in information

Property to be preserved



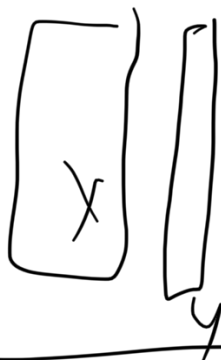
$N \times D$



$N \times M$

$M \ll D$

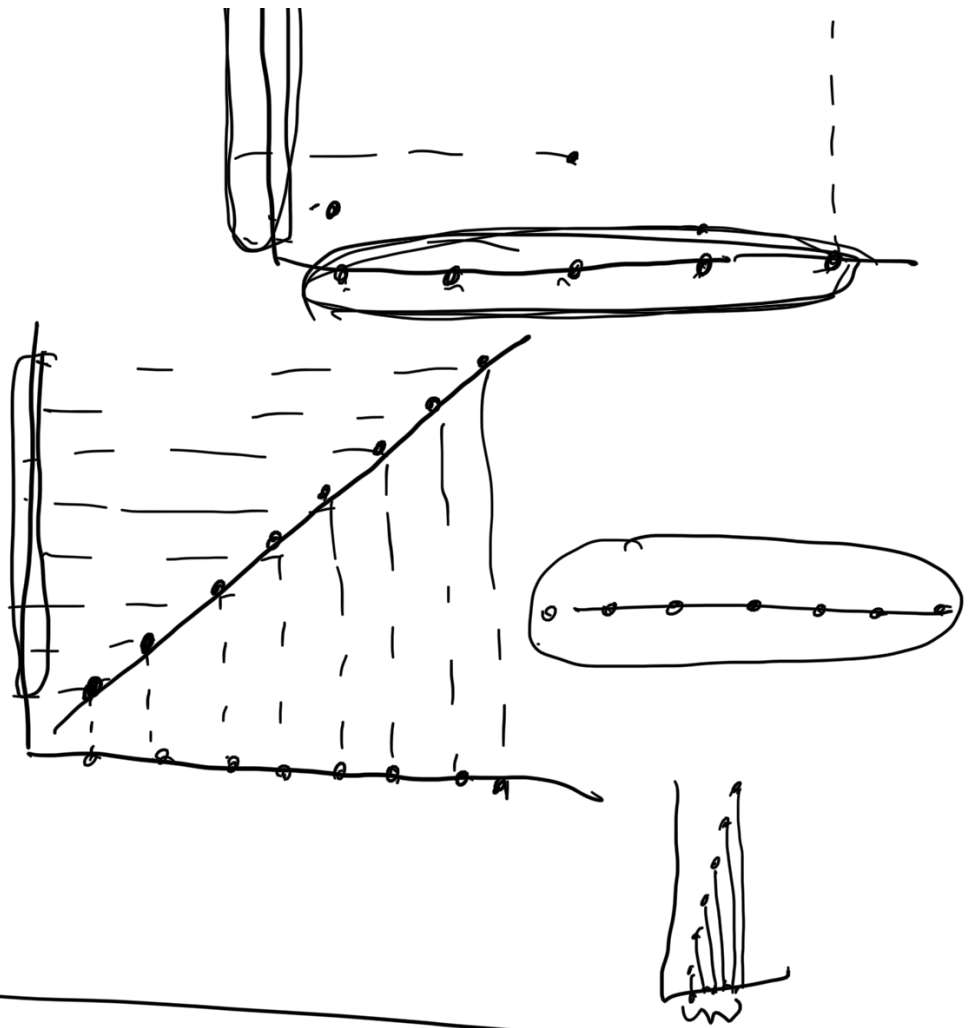
$$\|x_i - x_j\|_2 \approx \|z_i - z_j\|_2$$



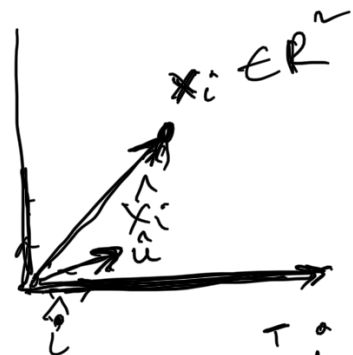
reduce

$$\mathbb{R}^D \rightarrow \mathbb{R}^M$$



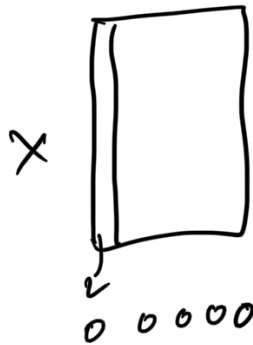


$u \rightarrow$ unit vector
 $\|u\|_2 = 1$



u

We will assume that our data is mean-centered



$$\hat{u}^T \left(\frac{1}{N} \sum_{i=1}^N (x_i^T \hat{u})^2 \right)$$

variance of the new data along \hat{u}

$$\text{Variance along } \hat{u} = \hat{u}^T \left(\frac{1}{N} \sum x_i x_i^T \right) \hat{u}$$

Sample covariance matrix of x

$$\frac{1}{N} \sum (x_i - \mu)(x_i - \mu)^T$$

$$\mu = 0$$

$$S = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$$

$$\begin{aligned} \max_{\hat{u}} & \hat{u}^T S \hat{u} \\ \text{s.t.} & \hat{u}^T \hat{u} = 1 \end{aligned}$$

$$\begin{aligned} \|\hat{u}\|_2 &= 1 \\ (\hat{u}^T \hat{u})^2 &= 1 \end{aligned}$$

$$\max_{\hat{u}, \lambda} \frac{d}{d\hat{u}} \hat{u}^T S \hat{u} - \lambda (\hat{u}^T \hat{u} - 1)$$

$$2S\hat{u} - 2\lambda\hat{u} = 0$$

$$\frac{d}{d\hat{u}} \hat{u}^T S \hat{u} = 2S\hat{u}$$

$$\begin{aligned} S\hat{u} &= \lambda\hat{u} \\ \hat{u}^T \hat{u} &= 1 \end{aligned}$$

$$\frac{d}{d\lambda}$$

Eigen-vector analysis

\hat{u} will be the eigen vector of S .

$$\hat{u}_{\max} : \hat{u}^T S \hat{u}$$

$$= \hat{u}^T S \hat{u}$$

$$= \hat{u}^T \lambda \hat{u} = \lambda \hat{u}^T \hat{u} = \lambda$$

Variance along the eigen vector \hat{u} is λ .

S is a $D \times D$ matrix.

It has D eigen vectors

and D eigen values.

$$\begin{bmatrix} \underline{\hat{u}}_1 & \underline{\hat{u}}_2 & \dots & \hat{u}_D \\ \underline{\lambda_1} & \underline{\lambda_2} & \dots & \lambda_D \end{bmatrix} \begin{cases} \hat{u}_i^T \hat{u}_i = 1 \\ \hat{u}_i^T \hat{u}_j = 0 \end{cases}$$

The eigen vector corresponding to the largest eigen-value gives the direction of maximal variance.

1st Principal Components

$$\begin{aligned} \hat{u} &= D \times 1 \\ X \hat{u} &= N \times 1 \end{aligned}$$

$N \times D$ $N \times 1$

$$\begin{aligned} \lambda_1 & \text{ variance} \\ \lambda_2 & \text{ variance} \\ & \vdots \\ \lambda_D & \text{ variance} \end{aligned}$$

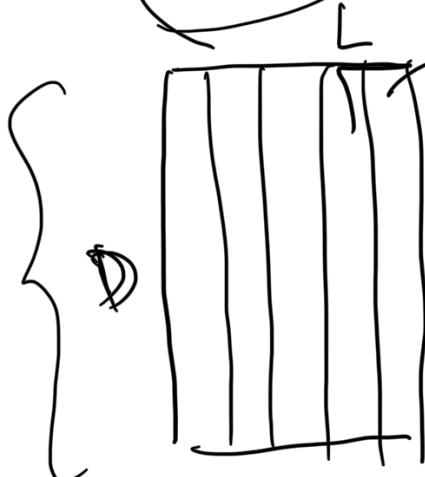
$$\sum_{i=1}^D \lambda_i \equiv \text{Total variance}$$

first L PCs

$$\sum_{i=1}^L \lambda_i$$

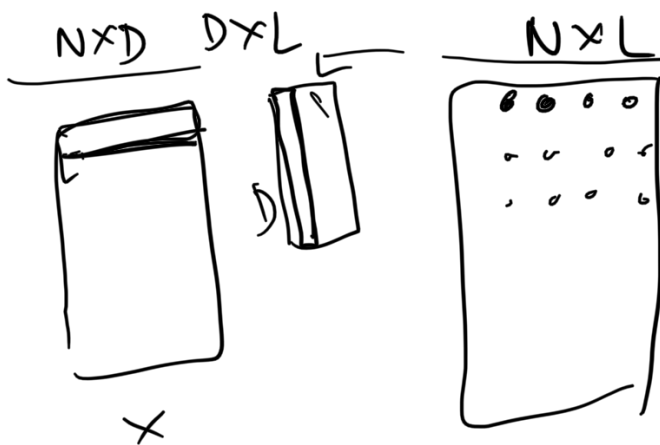
$\times 100$

$$\sum_{i=1}^D \lambda_i$$



each column is a PC

$$XW = Z \leftarrow \text{reduce data set.}$$

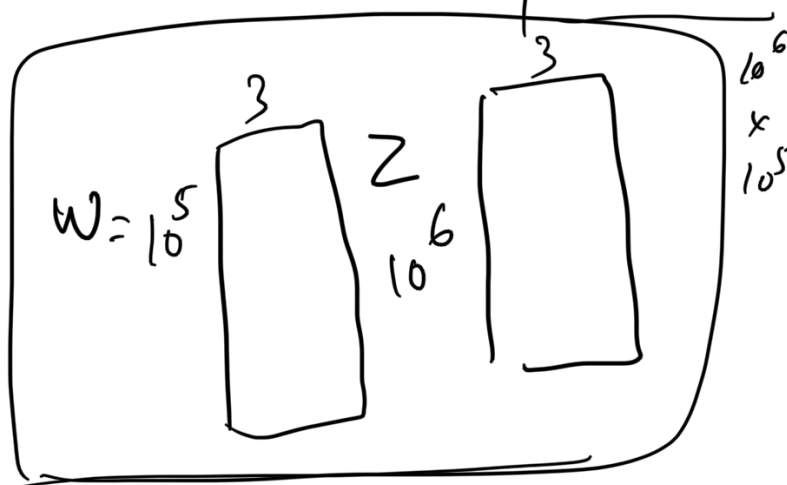


PCA is a linear algorithm.

$$Z_i = W^T x_i$$

$$\begin{array}{c}
 \left(\begin{array}{cc} L \times 1 & L \times D \end{array} \right) \begin{array}{c} \uparrow \\ D \times 1 \end{array} \\
 X \quad \begin{array}{c} \underline{\underline{X_i}} \\ \xrightarrow{\quad} \end{array} \begin{array}{c} \underline{\underline{Z_i}} \\ \wedge \\ \underline{\underline{X_i}} \end{array} \\
 \quad \quad \quad \begin{array}{c} \underline{\underline{Z_i}} \\ \xrightarrow{\quad} \end{array} \begin{array}{c} \underline{\underline{X_i}} \\ \underline{\underline{\quad}} \end{array}
 \end{array}$$

$$\begin{array}{c}
 \hat{X}_i = W Z_i \\
 \begin{array}{cc} D \times 1 & D \times L \quad L \times 1 \end{array} \\
 \hookrightarrow \text{reconstruction}
 \end{array}$$



$$10^5 \times 3 + 10^6 \times 3$$

$$\hat{X} = Z W^T$$

Reconstruction loss from PCA

$$\hat{X}_i = W Z_i$$

$$J = \frac{1}{N} \sum_{i=1}^N \|x_i - \hat{x}_i\|_2^2$$

Non-linear PCA

$$\Phi = \{ \phi_1, \phi_2, \phi_3, \dots \}$$

$$\phi_1(x_i) \rightarrow \mathbb{R}$$

$$\phi_2(x_i) \rightarrow \mathbb{R}$$

$$\phi_3(x_i) \rightarrow \mathbb{R}$$

$$\Phi(x_i) \rightarrow \mathbb{R}^M$$

$$k(x_i, x_j) \equiv \Phi(x_i)^T \Phi(x_j)$$

$$X \begin{matrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{matrix} \rightarrow \begin{matrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_N) \end{matrix}$$

$$C = \frac{1}{N} \sum_{i=1}^N \phi(x_i) \phi(x_i)^T$$

$(M \times M)$

$$C \underline{\underline{v_k}} = \lambda_k \underline{\underline{v_k}}$$

$(M \times 1)$

$$\frac{1}{N} \sum \phi(x_i) \phi(x_i)^T \underline{\underline{v_k}} = \lambda_k \underline{\underline{v_k}}$$

$$\sum \phi(x_i) \left(\frac{\phi(x_i)^T \underline{\underline{v_k}}}{N} \right) = \lambda_k \underline{\underline{v_k}}$$

$$\text{let } \underline{\underline{a_{ik}}} = \frac{\phi(x_i)^T \underline{\underline{v_k}}}{N \lambda_k}$$

$$\underline{\underline{v_k}} = \sum a_{ki} \phi(x_i)$$

$$\sum_{i=1}^N k(x_l, x_i) \sum_{j=1}^N a_{kj} k(x_i, x_j)$$

$$= N \lambda_k \sum_{i=1}^N k(x_l, x_i)$$

$l=1$

$$K^2 a_k = \lambda_k N K a_k$$

$$\underline{\underline{K a_k}} = \lambda_k \underbrace{N}_{\text{scalar}} \underline{\underline{a_k}}$$

K - kernel matrix: $N \times N$

We Directly calculate K

How can we ensure that implicitly mapped dataset is mean-centered?