

TAM 598 Lecture 18 :

Unsupervised Learning - Dimensionality Reduction

Announcements:

- No class on Wednesday
- HW 5 covers lectures 17-20; due on Fri Apr 18

Dimensionality Reduction -

- we have observations $\tilde{x}_{1:n}$, each of which is a high-dimensional vector $\tilde{x}_i \in \mathbb{R}^D$, with $D \gg 1$.
- our goal is to describe the data set with a smaller number of dimensions without losing too much information, ie to **project** each \tilde{x}_i to a d -dimensional vector \tilde{z}_i where $d \ll D$.

why? for visualization

to do clustering, density estimation, ...
for supervised learning tasks

Principal Component Analysis :

- ④ we want a linear map from $\tilde{x}_i \in \mathbb{R}^D$ to $\tilde{z}_i \in \mathbb{R}^d$

- ④ use an affine projection map $\mathbb{R}^D \rightarrow \mathbb{R}^d$

$$\tilde{z} = f(\tilde{x}) = \underline{W}^T (\tilde{x} - \underline{x}_0)$$

$$\underline{W} = \begin{bmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{bmatrix}^D_d$$

where \underline{W} is a $D \times d$ matrix

\underline{x}_0 is the empirical mean of

the data $\underline{x}_0 = \frac{1}{n} \sum_{i=1}^n \tilde{x}_i$

- ④ use a reconstruction map $\mathbb{R}^d \rightarrow \mathbb{R}^D$

$$\tilde{x} = g(\tilde{z}) = \underline{V} \tilde{z} + \underline{x}_0$$

$$\underline{V} \text{ is a } D \times d \text{ matrix} \quad \begin{bmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{bmatrix}^D_d$$

To find the matrices \underline{W} , \underline{V} we minimize the reconstruction error.

$$\begin{aligned} L(\underline{W}, \underline{V}, \underline{x}_0) &= \frac{1}{n} \sum_{i=1}^n \| \underline{x}_i - g(f(\underline{x}_i)) \|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \| \underline{x}_i - \underline{V} \underline{W}^T \underline{x}_i - \underline{x}_0 \|^2 \end{aligned}$$

Taking $\frac{\partial L}{\partial \underline{W}} = 0$, $\frac{\partial L}{\partial \underline{V}} = 0$ and solving:

$$(1) \quad \underline{W} = \underline{V}$$

$$\underline{x}_i - \underline{x}_0 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$i=1: \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} [123]$$

$$\underline{x}_2 - \underline{x}_0 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$$

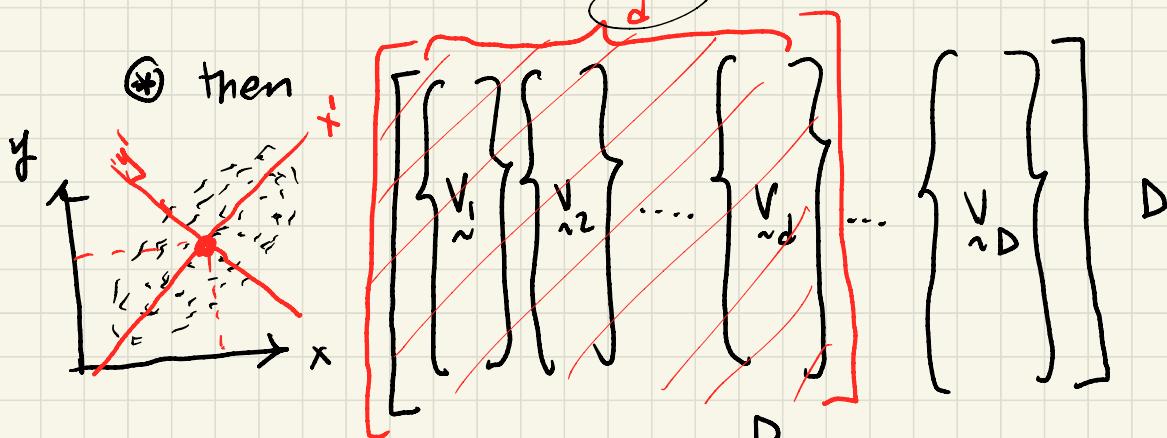
(2) \underline{V} can be constructed from eigenvectors / eigenvalues of empirical covariance matrix $\underline{C} = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \underline{x}_0)(\underline{x}_i - \underline{x}_0)^T$

$D \times D$

Specifically:

* let u_i and λ_i be the i^{th} eigenvalue of $\underline{\underline{U}}$, sorted so

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$



column i is
 $v_i = \sqrt{\lambda_i} u_i$

V is $D \times d$

* our projection map becomes

$$\underline{\underline{z}} = f(\underline{\underline{x}}) = \sum_{i=1}^d \sqrt{\lambda_i} u_i^T \underline{\underline{x}} = \underline{\underline{V}}^T \underline{\underline{x}}$$

* the reconstruction map is $\underline{\underline{x}} = g(\underline{\underline{z}}) = \underline{\underline{x}}_0 + \sum_{i=1}^d \sqrt{\lambda_i} u_i^T \underline{\underline{z}}$

④ The sum of the first d eigenvectors

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

Σ in principal coord system

$$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix}$$

tells us how much variance is explained with a model that keeps the first d PCA components

⑤ the minimum reconstruction error is

$$L(\underline{w}, \underline{v}, \underline{x}_0) = \frac{1}{n} \sum_{i=1}^n \| \underline{x}_i - g(f(\underline{x}_i)) \|^2$$

$$= \sum_{j=d+1}^D \lambda_j$$

probabilistic interpretation: assume our data points $x_{1:n}$ are generated by a linear Gaussian model

and latent variables z_i are generated by a Gaussian prior

We can maximize the marginal likelihood:

which is a gaussian. Maximizing $\log p(x_{1:n})$ gives the same result as before, now with variance