

# Lecture Notes on ML.

付呈宇. 中山大学数学学院

2018.1.31. 广东·广州

## Section 2: Principal Components Analysis.

When faced with a large set of correlated predictor variables, PCA allow us to summarize this set with a smaller number of representative variables that collectively explain most of the variability in the original dataset. PCA is an unsupervised approach, since it involves only a set of features  $(x_1, x_2, \dots, x_p)$  and no associated response  $Y$ .

Say we have a following dataset:

$$\begin{pmatrix} X_{11} & X_{12} & X_{13} & \dots & X_{1p} \\ X_{21} & X_{22} & X_{23} & \dots & X_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & X_{n3} & \dots & X_{np} \end{pmatrix} = X.$$

where each feature's mean is normalized to zero.

we give each observation a new feature, say:

$$Z_{i1} = \sum_{j=1}^p \phi_{j1} X_{ij}, \text{ where } \sum_{j=1}^p |\phi_{j1}|^2 = 1$$

we want to maximize the variance of  $(Z_{11}, Z_{21}, \dots, Z_{n1})$

$$\text{i.e. } \max_{\phi \in \mathbb{R}^{1 \times p}} \frac{1}{n} \sum_{i=1}^n (Z_{i1} - 0)^2, \text{ where } \|\phi\|_2 = 1, \phi = (\phi_{j1})^T.$$

while :

$$\frac{1}{n} \sum_{i=1}^n (z_{i1})^2 = \frac{1}{n} \sum_{i=1}^n \left( \phi_1 \cdot X_i^T \right)^2 \quad (X_i = (X_{i1} \ X_{i2} \ \dots \ X_{ip}))$$

$$= \frac{1}{n} \sum_{i=1}^n \phi_1 \cdot (X_i^T \cdot X_i) \phi_1^T$$

$$= \phi_1 \left( \frac{1}{n} \sum_{i=1}^n X_i^T \cdot X_i \right) \phi_1^T$$

we define  $\left( \bar{\Sigma} \triangleq \frac{1}{n} \sum_{i=1}^n X_i^T X_i \right)$ .

Our problem now converted into :

$$\begin{aligned} \max_{\substack{\phi_1 \in \mathbb{R}^{1 \times p} \\ \|\phi_1\|=1}} Q(\phi_1) &= \max_{\substack{\phi_1 \in \mathbb{R}^{1 \times p} \\ \|\phi_1\|=1}} \phi_1 \bar{\Sigma} \phi_1^T \end{aligned}$$

use Lagrange Multiplier Method, it's easy to show that  $Q(\phi_1)$  is maximized if  $\phi_1$  is the principal eigen-vector of  $\bar{\Sigma}$  ~~and~~ <sup>and</sup>  $\max Q = \max \lambda_i$  where  $\lambda_i$  is the eigen-value of  $\bar{\Sigma}$ . (HW.)

After Finding the way to construct a highly variant feature for each observation, a natural question is how to construct another variant feature?

We just need to choose top  $k$  eigen vectors of  $\bar{\Sigma}$



Then, each observation is mapped into a  $k$ -dimensional feature space where we can choose  $k \ll p$ .

Hence PCA is referred as a dimension reduction algorithm.

- Application 1.

If we reduce the dimension of observation to  $k=2$  or  $3$ , then we can visualize the data.

- Application 2.

Preprocess dataset to reduce its dimension before running a supervised learning algorithm. Apart from computational benefits, reducing dimension can reduce the complexity of the supervised model and therefore ~~will~~ help avoid over-fitting.

- Application 3.

$$\begin{array}{c} x_i \xrightarrow{\text{PCA}} y_i \\ \uparrow \\ \text{Face image} \end{array} \quad \text{then similar}(x_i, x_j) \triangleq \underline{\|y_i - y_j\|_2}$$

it turns out that this is a surprisingly good face matching algorithm...