

UNSUPERVISED LEARNING ((h 10))

3/5/21

So far : supervised learning (regression , classification)

$$x_1, \dots, x_p \rightarrow y$$

Unsupervised : x_1, \dots, x_p

- Finding patterns / structure
- Visualizing the data
- exploratory data analysis

① Principal Component Analysis (PCA)

→ reducing the dimension of the data

② Clustering :

→ Discover unknown subgroups of data

PCA

x_1, \dots, x_p features

n observations

We would like to do some visualization:

- We could do it pairwise (Junction pairs in R)
- $\binom{p}{2} = \frac{p(p-1)}{2}$ possible pairs of features
- $p = 10 \rightarrow 45$ plots
- informative up to a certain point

PCA tries to find a low-dimensional representation of the data

Two components / ingredients of PCA:

- ① Mathematical component: We want to rewrite the data in a different basis
- ② statistical: We want to choose a basis so that most of the information in the data is contained in the first "components"

MATHEMATICAL VIEW OF PCA

It consists of :

- Eigen decomposition of Covariance Matrix
- Singular Value Decomposition of Data Matrix

The two methods are very much related but they lead to different software implementations.

COVARIANCE MATRIX

$(X_1, \dots, X_p) = \underline{x}$ p random variables

$$\Sigma = (\sigma_{ij})_{\substack{i=1, \dots, p \\ j=1, \dots, p}} = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1p} \\ \vdots & \ddots & \vdots \\ \sigma_{p1} & \dots & \sigma_{pp} \end{pmatrix}, \quad \sigma_{ij} = \text{Cov}(X_i, X_j) = E[(X_i - E(X_i))(X_j - E(X_j))]$$

The matrix is symmetric and positive definite :

① $\sigma_{ij} = \text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i) = \sigma_{ji}$ (symmetric)

② For any $\underline{v} \in \mathbb{R}^p$,

$$\underline{v}^t \Sigma \underline{v} = \text{Var}(\underbrace{\underline{v}^t \underline{x}}_{v_1 X_1 + \dots + v_p X_p}) > 0$$

(positive definite, unless there is perfect multicollinearity - $\underline{v}^t \underline{x}$ constant)

(Spectral Theorem)

Σ symmetric \rightarrow p real eigenvalues \rightarrow

$$\Sigma = U D U^t$$

EIGENDECOMPOSITION
OF Σ

$$= (\underline{u}_1, \dots, \underline{u}_p) \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & \ddots & \lambda_p \end{pmatrix} U^t$$

$\underline{u}_1, \dots, \underline{u}_p$ are eigenvectors of length 1
and orthogonal

$\lambda_1, \dots, \lambda_p$ eigenvalues

Some remarks :

① U is orthogonal : $U^t U = U U^t = I_p$

Orthogonal matrices have special properties in the context of linear transformations:

① Preserve scalar product : $(U \underline{w}) \cdot (U \underline{v}) = \underline{w} \cdot (U^t U \underline{v}) = \underline{w} \cdot \underline{v}$

② $\det(U^t U) = 1 \Rightarrow \det(U) = \pm 1$

$\cancel{\frac{\underline{w}}{\det(U)}}, \quad U = (\underline{w}, \underline{v})$
or $\cancel{\underline{v}}$ reflections

So U is associated to rotations

② Σ positive definite $\Rightarrow \lambda_1, \dots, \lambda_p > 0$

③ In practice Σ unknown and needs to be estimated from data

Data

$$X_{n \times p} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \parallel & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix}$$

n observations of p variables

SAMPLE COVARIANCE

$$S = (s_{kj})_{k,j}$$

$$s_{kj} = \frac{\sum_{i=1}^n (x_{ik} - \bar{x}_k)(x_{ij} - \bar{x}_j)}{n} \quad \text{with } \bar{x}_k = \frac{\sum_{i=1}^n x_{ik}}{n} \text{ sample means}$$

Consider centered data : $\bar{x}_k = 0 \forall j$

$$\Rightarrow s_{kj} = \frac{\sum_{i=1}^n x_{ik} x_{ij}}{n} \quad (*) \quad k, j = 1, \dots, p$$

$$\Rightarrow S = \frac{\sum_{i=1}^n \underline{x_i} \cdot \underline{x_i}^T}{n} \quad \text{with } \underline{x_i} = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} \quad i^{\text{th}} \text{ "row" of } X$$

$\underline{x_i} \cdot \underline{x_i}^T$ has rank 1
so S does not have rank
larger than n

$$\begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} (x_{i1} \mid \cdots \mid x_{ip})$$

So if $p > n$, S is not full rank \rightarrow ill-conditioned

$$\Rightarrow S = \frac{\underline{X}^T \underline{X}}{n}$$

$$s_{kj} = \frac{1}{n} (\underline{x}_{1k} \mid \cdots \mid \underline{x}_{nk}) \begin{pmatrix} x_{1j} \\ \vdots \\ x_{nj} \end{pmatrix} = \frac{\sum_{i=1}^n x_{ik} x_{ij}}{n}$$

SINGULAR VALUE DECOMPOSITION (of $X_{n \times p}$)

Any matrix X can be factorized as

$$X = V \cdot D \cdot U^t$$

The diagram illustrates the SVD factorization $X = V \cdot D \cdot U^t$. The matrices are represented as follows:

- V : A rectangle labeled $n \times p$, representing an orthogonal matrix.
- D : A rectangle labeled $n \times n$, representing a rectangular diagonal matrix. It has a vertical line down the middle with $p+1 \dots n$ written above it. The main diagonal contains elements d_{11}, \dots, d_{pp} , and the off-diagonal elements are marked with zeros (0).
- U^t : A square labeled $p \times p$, representing an orthogonal matrix.

$d_{11}, \dots, d_{pp} \geq 0$
singular values

In order to show that any matrix X admits at least one SVD, we can construct one in this way:

① Set V as the matrix of eigenvectors of $\underbrace{XX^t}_{n \times n}$ and U the matrix of unit eigenvectors of $\underbrace{X^t X}_{p \times p}$.

② The singular values in D are the square root of the eigenvalues of $X^t X$.

Proof Firstly, $X^t X$ is symmetric and positive semi-definite, since

$$\textcircled{1} \quad (X^t X)^t = X^t (X^t)^t = X^t X$$

$$\textcircled{2} \quad \underline{v}^t (X^t X) \underline{v} = (\underline{v}^t X^t) (X \underline{v}) = (X \underline{v})^t (X \underline{v}) = \|X \underline{v}\|^2 \geq 0 \quad \text{for any } \underline{v}$$

Therefore, I can consider the eigen decomposition of $X^t X$:

$$X^t X = U \Lambda U^t \quad \text{with} \quad \Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & 0 \\ 0 & & \lambda_p \end{pmatrix}, \quad \lambda_1, \dots, \lambda_p \geq 0$$

$U = (U_1, \dots, U_p)$ unit eigenvectors of $X^t X$

$$(X^t X) U_i = \lambda_i U_i$$

$\frac{X^t X}{n}$
sample covariance

Let $d_{ii} = \sqrt{\lambda_i}$ and suppose that $X^t X$ is full rank (proof can be generalized), i.e.
 $X^t X$ positive definite $\Rightarrow \lambda_1, \dots, \lambda_p > 0$

Then

* $\underline{v}_i = \frac{X \underline{u}_i}{d_{ii}}, \quad i = 1, \dots, p$

and create a matrix $V = (\underline{v}_1, \dots, \underline{v}_p, \underline{v}_{p+1}, \dots, \underline{v}_n)$ with $\underline{v}_{p+1}, \dots, \underline{v}_n$ $n-p$ orthogonal vectors (not so important as they are multiplied with the "zero" part of D)

By construction, $\underline{v}_1, \dots, \underline{v}_p$ are unit eigenvectors of XX^t :

$$\begin{aligned} \textcircled{1} \quad (XX^t) \underline{v}_i &= (XX^t) \frac{X \underline{u}_i}{d_{ii}} = \frac{X(X^t X \underline{u}_i)}{d_{ii}} = \frac{X \lambda_i \underline{u}_i}{d_{ii}} = d_{ii} X \underline{u}_i \\ &= (d_{ii}^2) \underline{v}_i \quad \underline{v}_i \text{ is an eigenvector of } XX^t \text{ with eigenvalue } d_{ii}^2 = \lambda_i \end{aligned}$$

$$\textcircled{2} \quad \|\underline{v}_i\|^2 = \left(\frac{X \underline{u}_i}{d_{ii}} \right)^t \left(\frac{X \underline{u}_i}{d_{ii}} \right) = \frac{\underline{u}_i^t (X^t X \underline{u}_i)}{d_{ii}^2} = \frac{\underline{u}_i^t d_{ii}^2 \underline{u}_i}{d_{ii}^2} = \|\underline{u}_i\|^2 = 1$$

$$\text{Let } D = \begin{pmatrix} d_{11} & 0 & \dots & 0 \\ 0 & \ddots & & 0 \\ \vdots & & d_{pp} & 0 \\ 0 & 0 & \dots & 0 \end{pmatrix} ; \quad v_i := \frac{x_{ii}}{d_{ii}}$$

In matrix form:

$$V = X U D^{-1} \quad \text{with} \quad D^{-1} := \begin{pmatrix} \frac{1}{d_{11}} & \dots & 0 \\ 0 & \ddots & \frac{1}{d_{pp}} \\ \vdots & & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}$$

Equivalently:

$$VD = XU \Rightarrow VDU^T = X \underbrace{UU^T}_{I} = X$$

$$\Rightarrow X = V \cdot D \cdot U^T$$

Interpretation of SVD : $X = V D U^t$

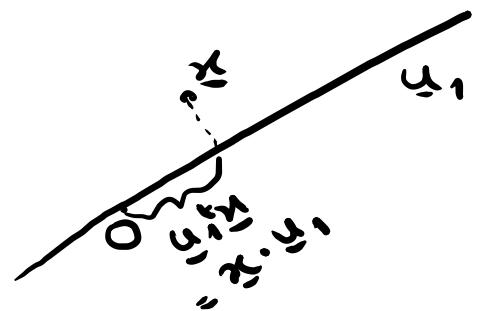
① Since $XU = VD$, the matrix VD provides the representation of the data X in the new basis U .

Take the generic row of X , \underline{x} . Then

$$U^t \underline{x} = \begin{pmatrix} \underline{u}_1^t \underline{x} \\ \vdots \\ \underline{u}_p^t \underline{x} \end{pmatrix} = \begin{matrix} t \\ \uparrow \end{matrix}$$

p "principal components", projections of the point \underline{x} on the new p directions

$$\underline{u}_1, \dots, \underline{u}_p$$



② SVD is not unique :

- D can be arranged from the largest singular value to the smallest (unique)
- V and U^t are not unique (think of the case of repeated eigenvalues)

$$\textcircled{3} \quad X = V D U^t = \underbrace{\underline{v}_1 d_{11} \underline{u}_1^t + \dots + \dots}_{\substack{\text{largest} \\ n \times 1 \quad 1 \times 1 \quad 1 \times p \\ n \times p}} + \dots + \underbrace{\underline{v}_p d_{pp} \underline{u}_p^t}_{\text{smallest}}$$

Proof :

$$DU^t = \begin{pmatrix} d_{11} & & 0 \\ & \ddots & \\ 0 & & d_{pp} \end{pmatrix} \begin{pmatrix} \underline{u}_1^t \\ \vdots \\ \underline{u}_p^t \end{pmatrix} = \begin{pmatrix} d_{11} \underline{u}_1^t \\ \vdots \\ d_{pp} \underline{u}_p^t \\ 0 \end{pmatrix}$$

So

$$(X)_{ij} = \underbrace{V_{[i,:]} \cdot (DU^t)_{[:,j]}}_{\substack{\text{i-th row} \\ \text{of } V}} = \sum_{k=1}^p v_{ik} (DU^t)_{kj} = \sum_{k=1}^p v_{ik} d_{kk} (\underline{u}_k^t)_j$$

$$= \left(\sum_{k=1}^p \underline{v}_k d_{kk} \underline{u}_k^t \right)_{ij} \Rightarrow$$

$$X = \underline{v}_1 d_{11} \underline{u}_1^t + \dots + \underline{v}_p d_{pp} \underline{u}_p^t$$

Consider $\tilde{X} = \underline{v}_1 d_{11} \underline{u}_1^t + \dots + \underline{v}_k d_{kk} \underline{u}_k^t$. Is $\tilde{X} \approx X$?
 Lower-dimensional representation of X

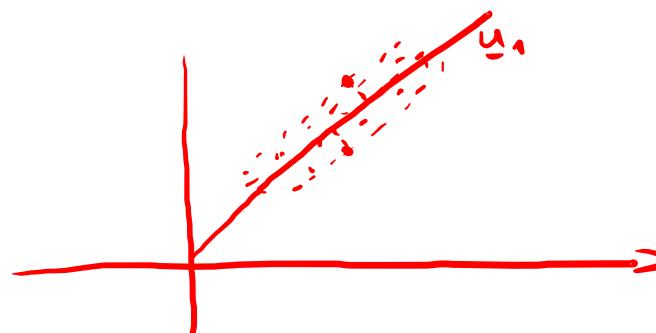
PCA - STATISTICAL VIEW

The choice of U orthonormal basis leads to P UNCORRELATED variables

$$x_1, \dots, x_p \xrightarrow[U^t x = t]{\downarrow} t_1, \dots, t_p$$

$$\text{Var}(t) = \text{Var}(U^t x) = U^t \text{Var}(x) U = U^t \Sigma U$$

$$\begin{aligned} \widehat{\text{Var}}(t) &= U^t \frac{x^t x}{n} U = \frac{U^t (V D U^t)^t (V D U^t) U}{n} = \\ &= \frac{U^t U D^t V^t V D U^t U}{n} = \frac{D^t D}{n} = \frac{1}{n} \begin{pmatrix} d_{11}^2 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & d_{pp}^2 \end{pmatrix} = \frac{1}{n} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} \end{aligned}$$



$$(ABC)^t = C^t B^t A^t$$

$\var(t_1)$
 $\var(t_i, t_j)$
 $\var(t_p)$

So the components t_1, \dots, t_p are uncorrelated and "sorted" according to their variance.