

Unsupervised Learning *

Principal component analysis

PCA seeks a small number of dimensions ~~that are~~ where the observations vary along each dimension. Each dimension in PCA is a linear combination of P features $z_i = \phi_{i1}x_1 + \dots + \phi_{iP}x_P$ where $\sum_{j=1}^P \phi_{ij}^2 = 1$.
 \Rightarrow Maximize $\left\{ \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^P \phi_{ji} x_{ij} \right)^2 \right\}$ subject to $\sum_{j=1}^P \phi_{ji}^2 = 1$ (sample variance of the n values of z_i).
 $\phi_{11}, \dots, \phi_{P1}$

The loading vector ϕ_i with elements $\phi_{i1}, \phi_{i2}, \dots, \phi_{iP}$ defines a direction in feature space along which the data vary the most. If we project the n data points x_1, \dots, x_n onto this direction, the projected values are the principal component scores z_{i1}, \dots, z_{in} , themselves.

The first M principal component score vectors and the first M principal component loading vectors provide the best M -dimensional approximation (in terms of Euclidean distance) to the i th observation x_{ij} . This representation can be written $x_{ij} \approx \sum_{m=1}^M z_{im} \phi_{jm}$.

Hence the \neq between original data and approximation is $\sum_{j=1}^P \sum_{i=1}^n \left(x_{ij} - \sum_{m=1}^M z_{im} \phi_{jm} \right)^2 = \text{approximation error}$.
 $\sum_{j=1}^P \frac{1}{n} \sum_{i=1}^n x_{ij}^2 = \sum_{m=1}^M \frac{1}{n} \sum_{i=1}^n z_{im}^2 + \frac{1}{n} \sum_{j=1}^P \sum_{i=1}^n \left(x_{ij} - \sum_{m=1}^M z_{im} \phi_{jm} \right)^2$
Var of data = Var of first M PCs (to be max)

We see when PCA can be seen as minimizing the approximation error or minimizing the variance.
MSE of M -dimensional approximation (to be min)
Proportion of variance explained of the m th principal component is

$$1 - \frac{\sum_{j=1}^P \sum_{i=1}^n \left(x_{ij} - \sum_{m=1}^M z_{im} \phi_{jm} \right)^2}{\sum_{j=1}^P \sum_{i=1}^n x_{ij}^2} = 1 - \frac{RSS}{TSS}$$

where TSS is the total sum of squared elements of x , and RSS represents the residual sum of squares of the M -dimensional approximation given by the principal components. Recalling the definition of R^2 , this means we can interpret the PVE as the R^2 of approximating $\text{Var } x$ given the first M principal components.

Before performing PCA, the variables should be centered to have zero mean. Further more they should have been individually scaled. This is in contrast to some other supervised and unsupervised techniques such as linear regression in which scaling has no effect.

Matrix completion

- create a complete data matrix of dimension $n \times P$ of which the (i,j) elements equal $\hat{x}_{ij} \begin{cases} x_{ij} & \text{if } (i,j) \in O \\ \hat{x}_{ij} & \text{if } (i,j) \notin O \end{cases}$
- Repeat steps (a) - (c) until the objective fails to decrease:
 - Solve
$$\text{minimize}_{A \in \mathbb{R}^{n \times M}, B \in \mathbb{R}^{M \times P}} \left\{ \sum_{j=1}^P \sum_{i=1}^n \left(\hat{x}_{ij} - \sum_{m=1}^M a_{im} b_{jm} \right)^2 \right\}$$
 - For each element $(i,j) \in O$ set $\hat{x}_{ij} \leftarrow \sum_{m=1}^M a_{im} b_{jm}$
 - compute the objective
$$\sum_{(i,j) \in O} \left(x_{ij} - \sum_{m=1}^M a_{im} b_{jm} \right)^2$$
- Return the estimated missing entries $\hat{x}_{ij}, (i,j) \notin O$

Clustering Methods

In general, cluster observations on the basis of the features in order to identify subgroups amongst the observations. or we can cluster features on the basis of the observations in order to discover subgroups among the features (transposing the matrix).

k-means clustering:

The idea behind k-means clustering is that a good clustering is one for which the within cluster variation is as small as possible

$$\text{minimize}_{c_1, \dots, c_k} \left\{ \sum_{h=1}^k W(c_h) \right\} \text{ with } W(c_h) = \frac{1}{|c_h|} \sum_{i: c_{ij}=c_h}^p (x_{ij} - x_{ij}^*)^2$$

1. Randomly assign a number from 1 to k to each of the observations. These serve as initial cluster assignment for the observations.
 2. Iterate until the cluster assignment stop changing:
 - (a) For each cluster, compute the cluster centroid (the median of p feature means for the observation in the h th cluster)
 - (b) Assign each observation to the cluster whose centroid is closest (Euclidean distance).
- K mean find a local optima, it should then be repeated \rightarrow then select the best solution.

Hierarchical clustering

bottom-up or agglomerative clustering (dendrogram - upside down tree). The term hierarchical refers to the fact that clusters obtained by cutting the dendrogram at a given height are necessarily nested with the clusters. However, on an arbitrary data set, this assumption of hierarchical clustering might be unrealistic. i.e. the true structure are not nested, in the sense that the best division into three groups doesn't result from taking the best division into two groups and splitting up one of those groups

1. Begin with n observations and a measure (Euclidean) of all $\binom{n}{2} = \frac{n(n-1)}{2}$ pairwise dissimilarities. treat each observation as its own cluster.
2. For $i = n-1, \dots, 2$
 - (a) Examine all pairwise inter-cluster dissimilarities among the i clusters and identify the pair of clusters that least dissimilarity. Fuse them.
 - (b) Compute the new pairwise inter-cluster dissimilarities among the $i-1$ remaining clusters.

Linkage

- Complete \rightarrow Maximal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B and record the largest of the dissimilarities.
- Single \rightarrow Minimal intercluster dissimilarity. Compute all pairwise dissimilarities between the observation in cluster A and the observation in cluster B, and record the smallest of these dissimilarities. Single linkage can result in extended, trailing clusters in which single observations are one at a time.
- Average \rightarrow Mean intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observation in cluster B, and record the average of these.
- Centroid \rightarrow Dissimilarity between the centroid for cluster A (a mean of vector of length p) and the centroid for cluster B. Centroid linkage can result in undesirable emission.
- Instead of Euclidean distance for dissimilarity we can use correlation based distance (similar if correlated).

- Validating the model \rightarrow compute p-values but not that great, no consensus on the method.
- Suppose outlier that do not belong into each cluster, how to cluster the rest without outlier \rightarrow mixture model \rightarrow soft version of k-means.