

Dimensionality Reduction/Feature Extraction

Sunday 28 April 2024 5:47 AM

Curse Of Dimensionality

Sunday 28 April 2024 6:05 AM

- When developing ML models, the aim is to maximise the amount of info provided to the algorithm in order to produce accurate & informative predictions.
- The temptation is to blindly include as many predictors as possible into the model. However this does not always improve performance & in some circumstances can reduce accuracy.
- COD refers to the situation where providing additional predictors to ML models (optimal # of features) can lead to reduced performance.
- When allocating new variables to a ML model, if the number of training data points remain constant then issues can develop. It is typically necessary to collect more training samples which is necessary to cover enough of the problem space that a model needs to properly learn the dataset.
- Data becomes sparse in higher dimensions.

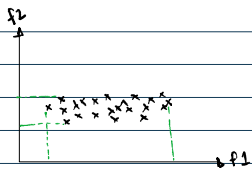
→ Unsupervised machine problem

→ PCA is a technique which can transform a higher dimension data into lower dimension data while keeping the essence of the data.

→ Benefits

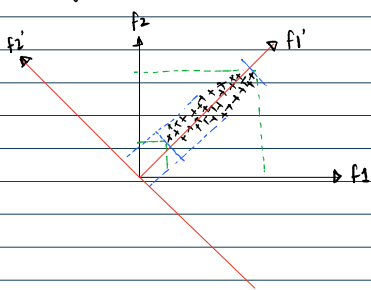
- faster algorithm execution
- visualization

→ choosing feature/axis on which the dataset displays more variance i.e. losing less information



$$X = \begin{bmatrix} f_1 & f_2 \end{bmatrix}_{n \times 2} \rightarrow X' = \begin{bmatrix} f_1' \end{bmatrix}_{n \times 1}$$

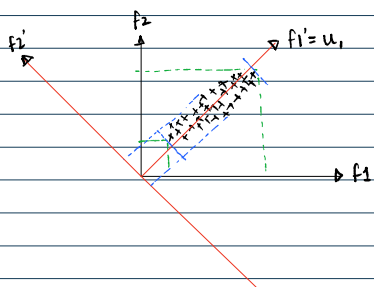
→ What happens if the variances are similar?



① $f_1' \perp f_2'$
find a direction f_1' such that the variance of x_i 's projected onto f_1' is maximal.
spread on $f_2' \ll$ spread on f_1'

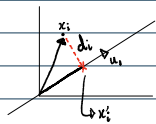
$$X = \begin{bmatrix} f_1 & f_2 \end{bmatrix}_{n \times 2} \rightarrow X' = \begin{bmatrix} f_1' \end{bmatrix}_{n \times 1}$$

→ Mathematical formulation



find u_1 s.t. $\text{var} \{ \text{proj}_{u_1} x_i \}_{i=1}^n$ is maximum

u_1 = unit-vector
 $\|u_1\| = 1$



$$x_i' = \text{projection of } x_i \text{ onto } u_1 \rightarrow \frac{u_1 \cdot x_i}{\|u_1\|^2} = u_1^T x_i$$

$$D = \{x_i\}_{i=1}^n \rightarrow D' = \{x_i'\}_{i=1}^n$$

$$x_i' = u_1^T x_i$$

$$\bar{x}_i' = u_1^T \bar{x}_i$$

→ mean $\{x_i\}_{i=1}^n$

$$\max \left(\text{var} \{ u_1^T x_i \}_{i=1}^n \right) = \max \left(\frac{1}{n} \sum_{i=1}^n (u_1^T x_i - u_1^T \bar{x})^2 \right)$$

→ 0 (as we mean centered)

$$\rightarrow \max_{u_1} \left| \frac{1}{n} \sum_{i=1}^n (u_1^T x_i) \right| \text{ such that } u_1^T u_1 = 1 = \|u_1\|^2$$

→ variance maximization

$$\text{OR } \min_{u_1} \sum_{i=1}^n \overbrace{x_i'^2 - (u_1^T x_i)^2}^{d_i^2} \text{ such that } u_1^T u_1 = 1$$

→ Distance Minimization

→ Covariance matrix [square & symmetrical]

$$\begin{matrix} & \begin{matrix} x & y & z \end{matrix} \\ \begin{matrix} x \\ y \\ z \end{matrix} & \begin{bmatrix} v_x & c_{x,y} & c_{x,z} \\ c_{y,x} & v_y & c_{y,z} \\ c_{z,x} & c_{z,y} & v_z \end{bmatrix} \end{matrix}$$

defines both the spread & orientation of our data

$$\text{cov}(X) = \frac{X^T X}{n-1}$$

→ * make sure features are mean centered / column standardized

→ Eigenvectors are vectors which maintain their direction after linear transformations

$$A\vec{v} = \lambda\vec{v}$$

→ Eigenvalues are the scaling factor of the eigenvectors after linear transformations

→ The largest eigenvector of the covariance matrix always points into the direction of the largest variance of the data & the magnitude of this vector equals the corresponding eigenvalue. The second largest eigenvector is always \perp to the largest eigenvector, & points into the direction of the second largest spread of the data.

1) mean centering the data

2) find covariance matrix

2) find the eigen values/eigen vectors of cov matrix

4) Transform data $U^T \cdot X$

from sklearn.decomposition import PCA

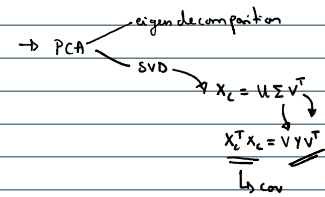


$$\left(\frac{\lambda_i}{\sum_{i=1}^n \lambda_i} \right) \times 100$$
, optimal principal components = total explained variance is 90%
→ % of variance explained

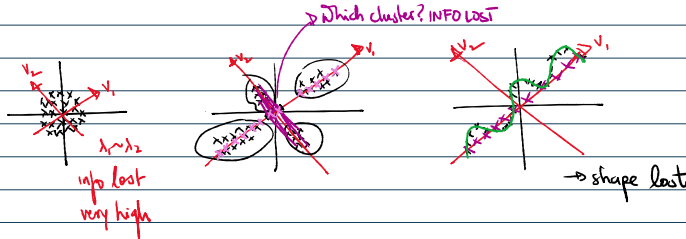
→ When does PCA fail/limitations?

→ Kernel PCA [Research]

→ Incremental & Randomized PCA [Research]



→ Explainability is lost/Difficult esp. Pharma & Finance
→ effects Decision Making



Let the real values data matrix X be of $n \times p$ size, where n is the number of samples and p is the number of variables. Let us assume that it is centered, i.e. column means have been subtracted and are now equal to zero.

Then the $p \times p$ covariance matrix C is given by $C = X^T X / (n - 1)$. It is a symmetric matrix and so it can be diagonalized:

$$C = V \Lambda V^T,$$

where V is a matrix of eigenvectors (each column is an eigenvector) and Λ is a diagonal matrix with eigenvalues λ_i in the decreasing order on the diagonal. The eigenvectors are called *principal axes* or *principal directions* of the data. Projections of the data on the principal axes are called *principal components*, also known as *PC scores*; these can be seen as new, transformed, variables. The j -th principal component is given by j -th column of XV . The coordinates of the i -th data point in the new PC space are given by the i -th row of XV .

If we now perform singular value decomposition of X , we obtain a decomposition

$$X = USV^T,$$

where U is a unitary matrix (with columns called left singular vectors), S is the diagonal matrix of singular values s_i and V columns are called right singular vectors. From here one can easily see that

$$C = VSU^TUSV^T/(n-1) = V \frac{S^2}{n-1} V^T,$$

meaning that right singular vectors V are principal directions (eigenvectors) and that singular values are related to the eigenvalues of covariance matrix via $\lambda_i = s_i^2/(n-1)$. Principal components are given by $XV = USV^T V = US$.

<https://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues/140579#140579>

To summarize:

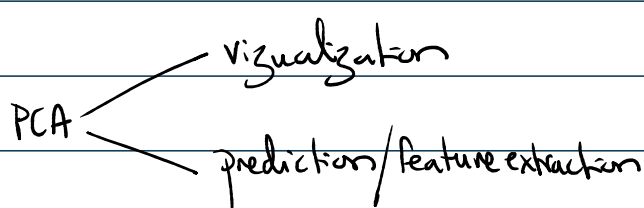
1. If $X = USV^T$, then the columns of V are principal directions/axes (eigenvectors).
2. Columns of US are principal components ("scores").
3. Singular values are related to the eigenvalues of covariance matrix via $\lambda_i = s_i^2/(n-1)$. Eigenvalues λ_i show variances of the respective PCs.
4. Standardized scores are given by columns of $\sqrt{n-1}U$ and loadings are given by columns of $VS/\sqrt{n-1}$. See e.g. [here](#) and [here](#) for why "loadings" should not be confused with principal directions.
5. The above is correct only if X is centered. Only then is covariance matrix equal to $X^T X / (n-1)$.
6. The above is correct only for X having samples in rows and variables in columns. If variables are in rows and samples in columns, then U and V exchange interpretations.
7. If one wants to perform PCA on a correlation matrix (instead of a covariance matrix), then columns of X should not only be centered, but standardized as well, i.e. divided by their standard deviations.
8. To reduce the dimensionality of the data from p to $k < p$, select k first columns of U , and $k \times k$ upper-left part of S . Their product $U_k S_k$ is the required $n \times k$ matrix containing first k PCs.
9. Further multiplying the first k PCs by the corresponding principal axes V_k^T yields $X_k = U_k S_k V_k^T$ matrix that has the original $n \times p$ size but is of lower rank (of rank k). This matrix X_k provides a reconstruction of the original data from the first k PCs. It has the lowest possible reconstruction error, see my answer [here](#).
10. Strictly speaking, U is of $n \times n$ size and V is of $p \times p$ size. However, if $n > p$ then the last $n-p$ columns of U are arbitrary (and corresponding rows of S are constant zero); one should therefore use an economy size (or thin) SVD that returns U of $n \times p$ size, dropping the useless columns. For large $n \gg p$ the matrix U would otherwise be unnecessarily huge. The same applies for an opposite situation of $n \ll p$.

t-SNE (t-distributed Stochastic Neighbor Embedding) is an unsupervised non-linear dimensionality reduction technique for data exploration and visualizing high-dimensional data.

Neighborhood preserving embedding

→ t-SNE aims to preserve local structure while PCA aims to preserve the global structure of the data. → neighborhood relationships guaranteed, relationship of points outside neighbourhood not guaranteed.

→ We can make t-sne preserve the global structure by tuning the 'perplexity' parameter



t-sne → visualization **NOT FOR** prediction/feature extraction

→ Complicated Math formulation [LATER]

→ Crowding Problem: Sometimes it's impossible to preserve distances in all the neighbourhoods.

→ t-distribution used to resolve this problem. [Try it's best]

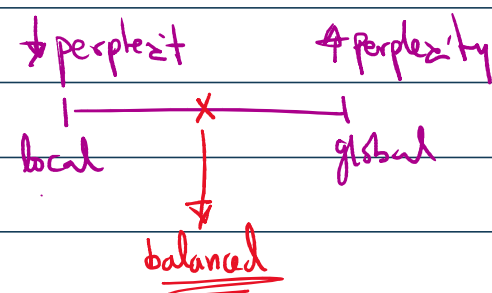
→ Role of t-distribution → complicated

→ t-sne is an iterative algorithm.

****** → Always run t-sne algo with multiple perplexity values ******

→ Run atleast for 1000 iterations to reach stable state

→ # of neighbors whose distances you wish to preserve



→ On the same dataset with the same parameter settings, if run multiple times output may differ.

t-sne is not a deterministic algorithm [stochastic/probabilistic]

→ t-sne tries to expand dense cluster & shrink sparse clusters s.t. the densities of groups of points are roughly similar.

WE CANNOT READ CLUSTER SIZES FROM t-sne

→ Distances b/w clusters might not mean anything [inter-cluster distances]

→ Don't fall in the trap of trying to make sense of random/noisy/junk data.

→ t-sne tries to spread points evenly

→ Always re-run t-sne with p as stepsize & see if shape is stable

* standardize features

→ Neighbourhoods of all points should be preserved if t-sne is working perfectly.

[Learn Later]

Saturday 11 May 2024 11:37 PM

→ Linear Discriminant Analysis *

→ UMAP *

→ MDS

→ Sammon mapping

→ Graph-based tech