

The covariance matrix of the transformed representation is diagonal. This can be shown as follows.

The data matrix X ^{old basis} of the original representation will get transformed to X' in the new representation

$$\underline{X}' = \underline{X} \underline{V}$$

$$\begin{bmatrix} \equiv \\ \equiv \\ \equiv \end{bmatrix} \begin{bmatrix} | & | & | & | & | \\ \lambda_1 & \lambda_2 & \lambda_3 & & \end{bmatrix} = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix} \quad \text{X'}$$

Here V is a matrix of eigenvectors as columns arranged from left to right such that their corresponding eigenvalues would form a decreasing order. $\lambda_1 \geq \lambda_2 \geq \dots$

Let C', u' be the covariance matrix and the mean vector in the new representation. u' is a $d \times 1$ vector

original
dxd

$$\underline{C} = \frac{\underline{X}^T \underline{X}}{m} - \underline{\mu} \underline{\mu}^T$$

$$\underline{\mu}' = \underline{V}^T \underline{\mu}$$

$$\underline{X}' = \underline{X} \underline{V}$$

transformed
dxd

$$\underline{C}' = \frac{\underline{X}'^T \underline{X}'}{m} - \underline{\mu}' \underline{\mu}'^T$$

$$= \frac{(\underline{X} \underline{V})^T (\underline{X} \underline{V})}{m} - (\underline{V}^T \underline{\mu}) (\underline{V}^T \underline{\mu})^T$$

$$= \frac{\underline{V}^T \underline{X}^T \underline{X} \underline{V}}{m} - \underline{V}^T \underline{\mu} \underline{\mu}^T \underline{V}$$

Substituting $\underline{X}^T \underline{X} = m \underline{C} + m \underline{\mu} \underline{\mu}^T$ gives

$$\underline{C}' = \frac{\underline{V}^T (m \underline{C} + m \underline{\mu} \underline{\mu}^T) \underline{V}}{m} - \underline{V}^T \underline{\mu} \underline{\mu}^T \underline{V}$$

$$C' = \underline{V^T C V} + \cancel{\underline{V^T \mu \mu^T V}} - \cancel{\underline{V^T \mu \mu^T V}}$$

$$C'_{d \times d} = \underline{\Lambda}$$

$$\begin{aligned} \underline{C} &= \underline{V \Lambda V^T} \quad V^T C V \\ &\therefore \underline{V^T C V} = \underline{\Lambda} \quad V^T V \Lambda V^T = I \Lambda I = \Lambda \end{aligned}$$

Thus, the transformed representation has a diagonal covariance matrix.

\therefore the directions (i.e. the new basis vectors) that maximize the variance are also the directions that remove correlations.

In the new representation X' , only the first $k \leq d$ columns will show a significant variation in values.

data matrix $X'_{m \times d}$

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

The remaining $d-k$ features will show less variation
 \Rightarrow they take values close to the mean of the feature.

\therefore representing X' using just k features would preserve the maximum variance.

\therefore only the first k columns of the transformed

{ data matrix need to be retained.
 The remaining columns can be discarded.

Dimensionality Reduction

Thus, PCA can be used to obtain a new representation of data in which the dimensionality is reduced.

So far we have discussed an interpretation of PCA in which the new low dimensional representation could preserve the maximum variance in the data. ① variance preserving interpretation

② In another interpretation of PCA, we seek a low dimensional representation such that when the data is projected back from the low dimensional subspace to the original space, the total squared distance between the original and the recovered vectors is minimal.

That is, the reconstruction error is minimal.

Find $\underline{W}_{n \times d}$ and $\underline{U}_{d \times n}$ that $\min_{\underline{U}, \underline{W}} \sum_{i=1}^m \left\| \underset{\substack{\uparrow \\ \text{original}}}{\underline{x}_i} - \underbrace{\underline{U} \underline{W} \underline{x}_i}_{\substack{\text{recovered} \\ \text{vector}}} \right\|^2$

The proof is given in the text book.

It first shows that $\underline{W} = \underline{U}^T$

Then it shows that minimizing the reconstruction error is equivalent to maximizing $\text{trace}(\underline{U}^T \underbrace{\sum_i \underline{x}_i \underline{x}_i^T}_A \underline{U})$

Then it shows that $\text{trace}(\underline{U}^T \sum_i \underline{x}_i \underline{x}_i^T \underline{U}) \leq \sum_{i=1}^m \lambda_i$

where λ_i are the largest n eigenvalues of the eigenvalue decomposition of matrix $\underline{A} = \sum_i \underline{x}_i \underline{x}_i^T$

The proof shows that setting the columns of matrix U as the n leading eigenvectors of A will ensure that

$$\text{trace}(\underline{U}^T \underline{A} \underline{U}) = \sum_{i=1}^n \lambda_i$$

It can be showed that the reconstruction loss is the sum of discarded eigenvalues $\sum_{i=n+1}^d \lambda_i$ distortion

An efficient way to apply PCA when $d \gg n$

Consider vectors \underline{x}_i which have been mean-centered.
The covariance matrix $\underline{C} = \underline{X}^T \underline{X}$ where $\underline{X} \equiv$ data matrix $m \times d$.

Clearly, $\underline{C}_{d \times d}$ is very large.

Instead, we formulate another matrix $\underline{B} = \underline{X} \underline{X}^T$

\underline{B} is $m \times m$ matrix.

The $(i, j)^{\text{th}}$ element of \underline{B} is given as $\langle \underline{x}_i, \underline{x}_j \rangle$.

\underline{B} is also called as the Gram matrix.

We obtain the eigen decomposition of \underline{B} .

If \underline{u} is an eigenvector of \underline{B} then $\underline{B} \underline{u} = \lambda \underline{u}$

$$\underline{X} \underline{X}^T \underline{u} = \lambda \underline{u}$$

$$\underline{X}^T (\underline{X} \underline{X}^T) \underline{u} = \lambda \underline{X}^T \underline{u}$$

since $\underline{C} \equiv \underline{X}^T \underline{X}$, we have