

Setup

We have $X = \{x_1, \dots, x_n\}$ where X is $p \times n$ with each col corresponding to an obs.
 We want to decompose this to M dimensions where $M < p$ by projecting X onto the subspace spanned by the col vectors of $B = \{b_1, b_2, \dots, b_M\}$. Assume B is orthonormal and all col vectors of X have mean 0 and standardised variance of 1.

We want to minimize info loss (variance loss) / squared errors by selecting appropriate basis vectors & coefficients

Deriving via calculus - Part 1

Consider x_i WLOG. $x_i = B_{11}b_1 + B_{12}b_2 + \dots + B_{1M}b_M + \dots + B_{1p}b_p$
 $x_i^* = B_{11}b_1 + B_{12}b_2 + \dots + B_{1M}b_M$

Thus, we want to minimize $J = \frac{1}{n} \sum_{i=1}^n \|x_i - x_i^*\|^2$

We do this by orthogonally projecting each x_i onto our basis vectors.
 Thus $B_{ij} = x_i^T b_j = b_j^T x_i$. Order is arbitrary as B_j is a scalar

$$\text{Proof: } \frac{\partial J}{\partial B_{ij}} = \frac{\partial J}{\partial x_i^*} \cdot \frac{\partial x_i^*}{\partial B_{ij}} = \frac{-2}{n} (x_i - x_i^*)^T b_j = 0$$

$$\Rightarrow x_i^T b_j = \sum_{j=1}^M B_{ij} b_j \cdot b_j = B_{ij} //$$

So either via linear algebra ($r.s = \|r\| \|s\| \cos \theta$) or calculus, we know $B = x^T b$

$$\text{Now, } J = \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j=M+1}^p B_{ij} b_j \cdot b_j \right\|^2 = \frac{1}{n} \sum \left\| \sum b_j^T B_{ij} b_j \right\|^2 = \frac{1}{n} \sum \left\| \sum B_{ij} \right\|^2$$

$$= \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j=M+1}^p x_i^T b_j \cdot b_j \right\|^2 \stackrel{\text{ONB}}{=} \frac{1}{n} \sum \sum \|x_i^T b_j\|^2 = \frac{1}{n} \sum \sum b_j^T x_i x_i^T b_j$$

$$= \sum_{j=M+1}^p b_j^T \underbrace{\frac{1}{n} \sum_{i=1}^n x_i x_i^T}_{\text{Cov}(X)} b_j = \sum_{j=M+1}^p b_j^T \Sigma b_j$$

The covariance matrix of X ! Found by taking $\text{Var}(X)$!

Deriving via calculus - Part 2

$$\mathcal{J} = \sum_{j=M+1}^p b_j^\top \Sigma b_j, \text{ subject to } b_j^\top b_j = 1 \quad \& \quad b_i^\top b_j = 0$$

$$L = \sum_{j=M+1}^p b_j^\top \Sigma b_j - \lambda_{M+1} (b_{M+1}^\top b_{M+1} - 1) - \dots$$

$$\frac{dL}{db_j} = 2b_j^\top \Sigma - 2\lambda_{M+1} b_{M+1}^\top - \lambda_{j+1} b_{j+1}^\top - \dots = 0$$

j is arbitrary & j+1, j+2, ... represent diff basis vectors

Multiply by all basis vectors where $b_i \neq b_j$ individually to show under minimal conditions, their corresponding factor $\lambda = 0$. Since $b_i^\top b_j = 0 \Leftrightarrow b_j^\top b_i = 0$, if $i \neq j$

(left with $b_j^\top \Sigma - \lambda_j b_j^\top = 0 \Leftrightarrow \Sigma b_j = \lambda_j b_j$ (since Σ is symmetric))

Repeat this for $j = M+1$ till p .

$$\text{So } \mathcal{J} = \sum_{j=M+1}^p \lambda_j b_j^\top b_j = \sum_{j=M+1}^p \lambda_j$$

Thus, the λ 's should be the smallest e-values of Σ , and b_j the e-vectors and $B = \{b_1, \dots, b_p\}$ should be the e-vectors corresponding to the largest e-values! We have found B !

Deriving via Variance

You could also start by attempting to find B by maximising the variance of the projection of X onto the subspace spanned by B . Let us define C to be $\{c_1, \dots, c_p\}$

maximise $\text{Var}(B^\top X)$ subject to $B^\top B = BB^\top = I_p$

$$\text{Let } L = \frac{1}{n} B^\top X X^\top B - \Lambda(B^\top B - I_p)$$

$$\frac{\partial L}{\partial B} = B^\top \Sigma - \Lambda B^\top = 0 \Leftrightarrow \Sigma B = B\Lambda \Leftrightarrow \Lambda^\top = B^\top \Sigma B$$

We know the covariance matrix is always diagonalizable i.e. $\Sigma = CVC^T$ where C is the square matrix of the e.vectors of Σ and V the diagonal matrix whose elements are the corresponding e.values. Since $V = C^T \Sigma C$ and we have above $\Lambda^T = B^T \Sigma B$, $\Lambda = \Lambda^T$ must be the diagonal matrix consisting of the M largest e.values of Σ and B their corresponding e.vectors. However, we cannot write $\Sigma = B\Lambda B^T$ as B is not a square matrix and $BB^T \neq I$

Result

$B = [b_1 \ b_2 \ \dots \ b_M]$ is the matrix whose col vectors are the eigenvectors of the covariance matrix of X corresponding to the M -largest e.values. Given X is $p \times n$, define X^* as $X^* = BB^T X$, where the i th col corresponds to $x_i^* = \sum_{j=1}^M B_{ij} b_j = \sum_{j=1}^M b_j b_j^T x_i$

$B^T X = \begin{bmatrix} b_1^T x_1 & \dots & b_1^T x_n \\ b_M^T x_1 & \dots & b_M^T x_n \end{bmatrix}$, a $M \times n$ matrix where each element represents the coordinates of the projection of x onto the space (1D) spanned by the basis vector, wrt the basis vector.

Our solution $X^* = BB^T X$ is a $p \times n$ matrix which is a p-dim representation of our projected M-dim data. We have successfully orthogonally projected X onto the subspace spanned by B in a way that maximises the remaining info (variance) and minimizes the average squared reconstruction error.