1 Principal Component Analysis (PCA)

Problem:

* Principal Component Analysis (PCA) is a feature transformation method that converts the original features f into a new set of transformed features p, ensuring their linear independence:

$$oldsymbol{f} = egin{pmatrix} f_1 \ dots \ f_k \end{pmatrix} & o & oldsymbol{p} = egin{pmatrix} p_1 \ dots \ p_m \end{pmatrix},$$

If the original features are linearly dependent, the data resides in a lower-dimensional space, meaning m < k. For clarity, we will assume m < k explicitly.

The new representation $p_1,...,p_m$ is constructed as a linear combination of the original features $f_1,...,f_k$:

$$p_s = \sum_{j=1}^k \alpha_{s,j} \cdot f_j,$$

the coefficients $\alpha_{s,j}$ form the matrix A, which defines the linear transformation from f to p.

The new, usually lower-dimensional, representation p must still be informative. This is achieved by ensuring that p can approximately restore the original features f linearly and with minimal error:

$$\hat{f}_j = \sum_{s=1}^m \beta_{j,s} \cdot p_s \approx f_j, \tag{3}$$

the coefficients $\beta_{j,s}$ form the matrix B, which defines the linear transformation from p back to f.

* The objective of PCA is to minimize the reconstruction error $\hat{f} - f$ by finding the optimal linear transformations $A : f \to p$ and $B : p \to f$:

$$R = \sum_{\boldsymbol{x} \in X^{\ell}} \left\| \hat{\boldsymbol{f}} - \boldsymbol{f} \right\|^2 = \sum_{\boldsymbol{x} \in X^{\ell}} \|BA\boldsymbol{f} - \boldsymbol{f}\|^2 \to \min_{A,B}.$$
 (4)

Linear Maps: Matrices A (dimension reducer) and B (dimension adder) are linear maps that work oppositely: A reduces the dimension of the original features f to the dimension of the principal components p, and B restores, as closely as possible, the original features from the principal components.

$$m{f} = egin{pmatrix} f_1 \ dots \ f_k \end{pmatrix} & \stackrel{A}{
ightarrow} m{p} = egin{pmatrix} p_1 \ dots \ p_m \end{pmatrix} & \stackrel{B}{
ightarrow} m{\hat{f}} = egin{pmatrix} \hat{f}_1 \ dots \ \hat{f}_k \end{pmatrix}$$

This can be written as:

$$p = Af, \qquad \hat{f} = Bp.$$

Matrix Formulation: The feature matrix F and the principal component matrix P are formed by stacking the row vectors $\mathbf{f}^{\mathsf{T}} = (f_1, ..., f_k)$ and $\mathbf{p}^{\mathsf{T}} = (p_1, ..., p_m)$:

$$F := \begin{pmatrix} \boldsymbol{f}_1^{\mathsf{T}} \\ \vdots \\ \boldsymbol{f}_\ell^{\mathsf{T}} \end{pmatrix}, \quad P := \begin{pmatrix} \boldsymbol{p}_1^{\mathsf{T}} \\ \vdots \\ \boldsymbol{p}_\ell^{\mathsf{T}} \end{pmatrix}$$

In matrix form, the linear maps A and B are applied as follows:

$$P^{\mathsf{T}} = AF^{\mathsf{T}}, \qquad \hat{F}^{\mathsf{T}} = BP^{\mathsf{T}},$$

or equivalently, by transposing:

Crumbs on the floor: Each data point is represented by three coordinates x, y, z, but z is always 0. Therefore, the data can be represented by just two coordinates:

(1)
$$f = \begin{pmatrix} x \\ y \\ 0 \end{pmatrix} \xrightarrow{A} p = \begin{pmatrix} x \\ y \end{pmatrix} \xrightarrow{B} \hat{f} = \begin{pmatrix} x \\ y \\ 0 \end{pmatrix}$$

It is straightforward to find the linear transformations A and B:

$$\underbrace{\begin{pmatrix} 1 & 0 & ? \\ 0 & 1 & ? \end{pmatrix}}_{A} \begin{pmatrix} x \\ y \\ 0 \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix}, \quad \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix}}_{B} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x \\ y \\ 0 \end{pmatrix}$$

NB: In the example above:

(2)

- * The last column of A is arbitrary, so the choice of transformations is not unique.
- * A and B are related: $\hat{f} = BAf$, thus BA = I.
- * Since A and B are non-square, they are non-invertible, so $A = B^{-1}$ does not hold.

Crumbs on the table.: Now, the third coordinate equals the table height h = 1:

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \xrightarrow{A} \begin{pmatrix} x \\ y \end{pmatrix}, \quad \begin{pmatrix} x \\ y \end{pmatrix} \xrightarrow{B} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}$$

Here, A is the same as before, but no B can restore the original vector exactly.

(5) Formally, if *B* exists, we could write the system of equations:

(6)
$$\begin{pmatrix} \beta_{1,1} & \beta_{1,2} \\ \beta_{2,1} & \beta_{2,2} \\ \beta_{3,1} & \beta_{3,2} \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \Rightarrow \begin{cases} 1x + 0y = x \\ 0x + 1y = y \\ \beta_{3,1}x + \beta_{3,2}y = 1 \end{cases}$$

- * The coefficients in the first two equations are determined by the identities x = x and y = y.
- * The third equation cannot yield 1 for all x, y since it lacks a bias term.
- Approximate solution.: In the example above, we could find *B* as the pseudoinverse $B = A^+ = (A^T A)^{-1} A^T$, but:
 - * The original vector will only be restored approximately, so $AB \approx I$.

(8)

* Since the choice of A is arbitrary, the choice of B is also arbitrary. This freedom allows us to impose additional constraints on the transformations.

$$P = FA^{\mathsf{T}}, \quad \hat{F} = PB^{\mathsf{T}}.$$

Substituting P into \hat{F} yields the following equation:

$$\hat{F} = FA^{\mathsf{T}} B^{\mathsf{T}} = F(AB)^{\mathsf{T}}, \qquad (1$$

The approximation \hat{F} equals F exactly if AB = I. Ideally, A would equal B^{-1} , but in general, A and B are non-square and therefore non-invertible.

Pseudoinverse matrix: AB = I holds if B is the pseudoinverse of A:

$$B = A^{+} = (A^{\mathsf{T}} A)^{-1} A^{\mathsf{T}} .$$
 (1)

 A^+ is exact if A has full rank, but in general, it does not, so the solution is only approximate:

$$AB \approx I.$$
 (12)

Geometric Interpretation: Matrices A and B resemble transition matrices between bases:

- * A transforms vectors from the original basis of features $f_1, ..., f_k$ into a new space with the basis of principal components $p_1, ..., p_m$. However, since these bases are in different dimensional spaces, this is only an analogy.
- * *B* performs the reverse transformation, converting from the principal component basis back to the original basis (approximately).

Since A and B are related by the pseudoinverse operation and perform inverse transformations, we can focus on one of the matrices. Let it be B.

The basis transition matrix stores the vectors of the new basis in the coordinates of the old basis. As the linear map B transforms principal components into the original features (approximately):

$$f \approx Bp,$$
 (13)

it acts similarly to a basis transition matrix from f to p, storing the orthogonal basis of principal axes in the coordinates of the original space.

Any basis consists of linearly independent, or orthogonal, vectors, meaning that B stores orthogonal vectors, and $B^{\mathsf{T}} B = \Lambda$ is diagonal.

Since the choice of B is not unique, we can use this freedom to demand that $B^{\mathsf{T}} B$ be not just diagonal Λ , but the identity matrix I:

$$\exists B : B^{\mathsf{T}} B = I, \tag{14}$$

This implies that B stores not just orthogonal vectors but an **orthonormal** basis of principal components.

Risk Minimization: The objective of PCA is to minimize the restoration error. In this notation, the empirical risk depends on A and B:

$$R := \left\| \hat{F} - F \right\|^{2}$$
$$= \left\| FA^{\mathsf{T}} B^{\mathsf{T}} - F \right\|^{2} \to \min_{A,B}.$$
 (15)

We can reformulate the objective in terms of the new coordinates P and the transition matrix B by substituting $P = FA^{\mathsf{T}}$, which at least reduces one matrix multiplication:

$$R = \left\| PB^{\mathsf{T}} - F \right\|^2 \to \min_{P,B}.$$
 (16)

By differentiating R with respect to P and B, we can find the values of P and B at the extremum:

$$BA = A^{+}A = (A^{\top} A)^{-1}(A^{\top} A) = I$$

Basis Transition Matrix.: If in vector space V, there are two bases: the old one $\mathcal{O}: \boldsymbol{\omega}_1, ..., \boldsymbol{\omega}_n$ and the new one $\mathcal{N}: \boldsymbol{\nu}_1, ..., \boldsymbol{\nu}_n$, the vectors of the new basis can be represented as linear combinations of the old basis vectors:

$$\begin{cases} \nu_1 = \alpha_{1,1}\omega_1 + \dots + \alpha_{1,n}\omega_n \\ \vdots \\ \nu_n = \alpha_{n,1}\omega_1 + \dots + \alpha_{n,n}\omega_n \end{cases}$$

The coefficients $\alpha_{s,j}$ are the coordinates of the new basis vectors in the coordinate system of the old basis. These coefficients form the basis transition matrix (by columns!):

$$A = \begin{pmatrix} \alpha_{1,1} & \dots & \alpha_{n,1} \\ \vdots & \ddots & \vdots \\ \alpha_{1,n} & \dots & \alpha_{n,n} \end{pmatrix}$$

This matrix transforms coordinates between bases:

$$\begin{aligned} \boldsymbol{\nu}_{1} \}_{\mathcal{O}} &= \begin{pmatrix} \alpha_{1,1} \\ \vdots \\ \alpha_{1,n} \end{pmatrix}_{\mathcal{O}} = A \begin{pmatrix} 1 \\ 0 \\ \vdots \end{pmatrix}_{\mathcal{N}} = A \{ \boldsymbol{\nu}_{1} \}_{\mathcal{N}} \\ \{ \boldsymbol{v} \}_{\mathcal{O}} = A \{ \boldsymbol{v} \}_{\mathcal{N}}, \quad \{ \boldsymbol{v} \}_{\mathcal{N}} = A^{-1} \{ \boldsymbol{v} \}_{\mathcal{O}} \end{aligned}$$

The choice of matrix B is flexible, allowing us to impose additional constraints. For example, we can require that $B^T B$ be diagonal or even the identity matrix:

$$B^T B = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

The objective R depends only on the product PB^T , which can result from multiplying any number of different pairs of matrices:

$$PB^{\mathsf{T}} = PIB^{\mathsf{T}} = \underbrace{(P^*R)}_{P} \underbrace{\left(R^{-1}B^{*^{\mathsf{T}}}\right)}_{B^{\mathsf{T}}}$$
(20)

We will use the freedom in choosing R and let $P^{\mathsf{T}} P$ and $B^{\mathsf{T}} B$ be diagonal:

- * P stores the principal components in their respective coordinates.
- * B stores the orthonormal "basis" of principal components in the coordinates of the original space, so $B^{\mathsf{T}} B = I$.

$$\begin{cases} P^{\mathsf{T}} \ P = \Lambda \\ B^{\mathsf{T}} \ B = I \end{cases}$$
(21)

Now, we can further simplify the expressions for P and B:

$$P = FB(B^{\mathsf{T}} B)^{-1} = FBI,$$

$$B = F^{\mathsf{T}} P(P^{\mathsf{T}} P)^{-1} = F^{\mathsf{T}} P\Lambda^{-1}.$$
(22)

Eliminate P:

Eliminate B:

$$B\Lambda = F^{\mathsf{T}} F B$$
 (23) $P\Lambda = F F^{\mathsf{T}} P$ (25)

This means that the columns of B are eigenvectors of $F^{\mathsf{T}} F$:

This means that the columns of
$$P$$
 are eigenvectors of FF^{T} :

$$\boldsymbol{b}_{j} \cdot \boldsymbol{\lambda}_{j} = (F^{\mathsf{T}} F) \boldsymbol{b}_{j}. \tag{24} \qquad \boldsymbol{p}_{j} \cdot \boldsymbol{\lambda}_{j} = (FF^{\mathsf{T}}) \boldsymbol{p}_{j}. \tag{26}$$

 $S = P^{\mathsf{T}} P$ is symmetric, *i.e.* $S^{T} = S$

Earlier, we showed that B could be chosen to store an orthonormal basis, but this wasn't strictly necessary.

It can be demonstrated analytically that it is sufficient to choose R such that $B^T B$ is diagonal, which is enough to ensure $B^T B = I$. This will determine the form of B, which can then be interpreted as a matrix storing an orthonormal basis.

As the proof involves boring linear algebra, we relied on geometric intuition instead (though formal proof is possible!).