Principal Component Analysis

Statistical Machine Learning (ENGG*6600*08)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Fall 2023

If there exist n training data points, i.e., {x_i}ⁿ_{i=1}, the projection of a training data point x is:

$$\mathbb{R}^{p} \ni \widetilde{\mathbf{x}} := \mathbf{U}^{\top} \mathbf{\check{x}},\tag{1}$$

where:

$$\mathbb{R}^d \ni \breve{\mathbf{x}} := \mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}},\tag{2}$$

is the centered data point and:

$$\mathbb{R}^d \ni \boldsymbol{\mu}_{\boldsymbol{x}} := \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}_i, \tag{3}$$

is the mean of training data points.

• The reconstruction of a training data point x after projection onto the PCA subspace is:

$$\mathbb{R}^d \ni \widehat{\mathbf{x}} := \mathbf{U}\mathbf{U}^\top \mathbf{\breve{x}} + \boldsymbol{\mu}_{\mathbf{x}} = \mathbf{U}\widetilde{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{x}}, \tag{4}$$

where the mean is added back because it was removed before projection.

• In PCA, all the data points should be centered, i.e., the mean should be removed first. The reason is shown in this figure.



 In some applications, centering the data does not make sense. For example, in natural language processing, the data are text and centering the data makes some negative measures which is non-sense for text. Therefore, data is not sometimes centered and PCA is applied on the non-centered data. This method is called Latent Semantic Indexing (LSI) or Latent Semantic Analysis (LSA) [1].

If we stack the n data points column-wise in a matrix X = [x₁,...,x_n] ∈ ℝ^{d×n}, we first center them:

$$\mathbb{R}^{d \times n} \ni \breve{X} := XH = X - \mu_{x}, \tag{5}$$

where $\mathbf{\check{X}} = [\mathbf{\check{x}}_1, \dots, \mathbf{\check{x}}_n] = [\mathbf{x}_1 - \mathbf{\mu}_x, \dots, \mathbf{x}_n - \mathbf{\mu}_x]$ is the centered data and:

$$\mathbb{R}^{n \times n} \ni \boldsymbol{H} := \boldsymbol{I} - (1/n) \boldsymbol{1} \boldsymbol{1}^{\top}, \tag{6}$$

is the centering matrix.

The projection and reconstruction, Eqs. (1) and (4), for the whole training data are:

$$\mathbb{R}^{p \times n} \ni \widetilde{\boldsymbol{X}} := \boldsymbol{U}^{\top} \check{\boldsymbol{X}}, \tag{7}$$

$$\mathbb{R}^{d \times n} \ni \widehat{\mathbf{X}} := \mathbf{U}\mathbf{U}^{\top} \mathbf{\check{X}} + \boldsymbol{\mu}_{x} = \mathbf{U}\widetilde{\mathbf{X}} + \boldsymbol{\mu}_{x}, \tag{8}$$

where $\widetilde{\mathbf{X}} = [\widetilde{\mathbf{x}}_1, \dots, \widetilde{\mathbf{x}}_n]$ and $\widehat{\mathbf{X}} = [\widehat{\mathbf{x}}_1, \dots, \widehat{\mathbf{x}}_n]$ are the projected data onto PCA subspace and the reconstructed data, respectively.

$$\mathbb{R}^{p} \ni \widetilde{\mathbf{x}}_{t} = \mathbf{U}^{\top} \breve{\mathbf{x}}_{t}, \tag{9}$$

$$\mathbb{R}^{d} \ni \widehat{\mathbf{x}}_{t} = \mathbf{U}\mathbf{U}^{\top} \breve{\mathbf{x}}_{t} + \boldsymbol{\mu}_{x} = \mathbf{U}\widetilde{\mathbf{x}}_{t} + \boldsymbol{\mu}_{x}, \tag{10}$$

where:

$$\mathbb{R}^d \ni \breve{\mathbf{x}}_t := \mathbf{x}_t - \boldsymbol{\mu}_x,\tag{11}$$

is the centered out-of-sample data point which is centered using the mean of training data. Note that for centering the out-of-sample data point(s), we should use the mean of the training data and not the out-of-sample data.

• If we consider the n_t out-of-sample data points, $\mathbb{R}^{d \times n_t} \ni \mathbf{X}_t = [\mathbf{x}_{t,1}, \dots, \mathbf{x}_{t,n_t}]$, all together, the projection and reconstruction of them are:

$$\mathbb{R}^{p \times n_t} \ni \widetilde{\boldsymbol{X}}_t = \boldsymbol{U}^\top \check{\boldsymbol{X}}_t, \tag{12}$$

$$\mathbb{R}^{d \times n_t} \ni \widehat{\boldsymbol{X}}_t = \boldsymbol{U} \boldsymbol{U}^\top \check{\boldsymbol{X}}_t + \boldsymbol{\mu}_x = \boldsymbol{U} \widetilde{\boldsymbol{X}}_t + \boldsymbol{\mu}_x,$$
(13)

respectively, where:

$$\mathbb{R}^{d \times n_t} \ni \breve{\boldsymbol{X}}_t := \boldsymbol{X}_t - \boldsymbol{\mu}_{\boldsymbol{X}}.$$
(14)

PCA Using Eigen-Decomposition

Projection Onto One Direction

In Eq. (4), if p = 1, we are projecting x onto only one vector u and reconstruct it. If we ignore adding the mean back, we have:

$$\widehat{\mathbf{x}} = \mathbf{u}\mathbf{u}^{\top} \widecheck{\mathbf{x}}$$

• The squared length (squared ℓ_2 -norm) of this reconstructed vector is:

$$||\widehat{\mathbf{x}}||_{2}^{2} = ||\mathbf{u}\mathbf{u}^{\top} \mathbf{\check{x}}||_{2}^{2} = (\mathbf{u}\mathbf{u}^{\top} \mathbf{\check{x}})^{\top} (\mathbf{u}\mathbf{u}^{\top} \mathbf{\check{x}})$$
$$= \mathbf{\check{x}}^{\top} \mathbf{u} \underbrace{\mathbf{u}}^{\top} \mathbf{u}^{\top} \mathbf{\check{x}} \stackrel{(a)}{=} \mathbf{\check{x}}^{\top} \mathbf{u} \mathbf{u}^{\top} \mathbf{\check{x}} \stackrel{(b)}{=} \mathbf{u}^{\top} \mathbf{\check{x}} \mathbf{\check{x}}^{\top} \mathbf{u},$$
(15)

where (a) is because \boldsymbol{u} is a unit (normal) vector, i.e., $\boldsymbol{u}^{\top}\boldsymbol{u} = ||\boldsymbol{u}||_2^2 = 1$, and (b) is because $\boldsymbol{\check{x}}^{\top}\boldsymbol{u} = \boldsymbol{u}^{\top}\boldsymbol{\check{x}} \in \mathbb{R}$.

Suppose we have n data points {x_i}ⁿ_{i=1} where {x_i}ⁿ_{i=1} are the centered data. The summation of the squared lengths of their projections {x̂_i}ⁿ_{i=1} is:

$$\sum_{i=1}^{n} ||\widehat{\mathbf{x}}_{i}||_{2}^{2} \stackrel{\text{(15)}}{=} \sum_{i=1}^{n} \boldsymbol{u}^{\top} \check{\mathbf{x}}_{i} \check{\mathbf{x}}_{i}^{\top} \boldsymbol{u} = \boldsymbol{u}^{\top} \Big(\sum_{i=1}^{n} \check{\mathbf{x}}_{i} \check{\mathbf{x}}_{i}^{\top} \Big) \boldsymbol{u}.$$
(16)

Projection Onto One Direction

• Considering
$$\mathbf{\breve{X}} = [\breve{x}_1, \dots, \breve{x}_n] \in \mathbb{R}^{d \times n}$$
, we have:

$$\mathbb{R}^{d \times d} \ni \boldsymbol{S} := \sum_{i=1}^{n} \check{\boldsymbol{x}}_{i} \check{\boldsymbol{x}}_{i}^{\top} = \check{\boldsymbol{X}} \check{\boldsymbol{X}}^{\top} \stackrel{(5)}{=} \boldsymbol{X} \boldsymbol{H} \boldsymbol{H}^{\top} \boldsymbol{X}^{\top} = \boldsymbol{X} \boldsymbol{H} \boldsymbol{H} \boldsymbol{X}^{\top} = \boldsymbol{X} \boldsymbol{H} \boldsymbol{X}^{\top}, \quad (17)$$

where **S** is called the "covariance matrix". If the data were already centered, we would have $S = XX^{\top}$.

Plugging Eq. (17) in Eq. (16) gives us:

$$\sum_{i=1}^{n} ||\widehat{\boldsymbol{x}}_{i}||_{2}^{2} = \boldsymbol{u}^{\top} \boldsymbol{S} \boldsymbol{u}.$$
(18)

- Note that we can also say that $u^{\top} Su$ is the variance of the projected data onto PCA subspace. In other words, $u^{\top} Su = \mathbb{V}ar(u^{\top} \check{X})$. This makes sense because when some non-random thing (here u) is multiplied to the random data (here \check{X}), it will have squared (quadratic) effect on variance, and $u^{\top} Su$ is quadratic in u.
- Therefore, u[⊤]Su can be interpreted in two ways: (I) the squared length of reconstruction and (II) the variance of projection.

Projection Onto One Direction

• We want to find a projection direction *u* which maximizes the squared length of reconstruction (or variance of projection):

$$\begin{array}{ll} \underset{\boldsymbol{u}}{\text{maximize}} & \boldsymbol{u}^{\top} \boldsymbol{S} \boldsymbol{u}, \\ \\ \text{subject to} & \boldsymbol{u}^{\top} \boldsymbol{u} = 1, \end{array}$$
 (19)

where the constraint ensures that the \boldsymbol{u} is a unit (normal) vector as we assumed beforehand.

Using Lagrange multiplier [2], we have:

$$\mathcal{L} = \boldsymbol{u}^{\top} \boldsymbol{S} \boldsymbol{u} - \lambda (\boldsymbol{u}^{\top} \boldsymbol{u} - 1),$$

Taking derivative of the Lagrangian and setting it to zero gives:

$$\mathbb{R}^{p} \ni \frac{\partial \mathcal{L}}{\partial \boldsymbol{u}} = 2\boldsymbol{S}\boldsymbol{u} - 2\lambda\boldsymbol{u} \stackrel{\text{set}}{=} 0 \implies \boldsymbol{S}\boldsymbol{u} = \lambda\boldsymbol{u}.$$
(20)

- The Eq. (20) is the eigen-decomposition of S where u and λ are the leading eigenvector and eigenvalue of S, respectively [3].
- Note that the leading eigenvalue is the largest one. The reason of being leading is that we are maximizing in the optimization problem.
- As a conclusion, if projecting onto one PCA direction, the PCA direction **u** is the leading eigenvector of the covariance matrix.
- Note that the "PCA direction" is also called "principal direction" or "principal axis" in the literature. The dimensions (features) of the projected data onto PCA subspace are called "principal components".

Projection Onto Span of Several Directions

In Eq. (4) or (8), if p > 1, we are projecting x or X onto PCA subspace with dimensionality more than one and then reconstruct back. If we ignore adding the mean back, we have:

$$\widehat{\boldsymbol{X}} = \boldsymbol{U}\boldsymbol{U}^{\top}\breve{\boldsymbol{X}}.$$

- It means that we project every column of \check{X} , i.e., \check{x} , onto a space spanned by the *p* vectors $\{u_1, \ldots, u_p\}$ each of which is *d*-dimensional. Therefore, the projected data are *p*-dimensional and the reconstructed data are *d*-dimensional.
- The squared length (squared Frobenius Norm) of this reconstructed matrix is:

$$\begin{aligned} \|\widehat{\boldsymbol{X}}\|_{F}^{2} &= \|\boldsymbol{U}\boldsymbol{U}^{\top} \, \check{\boldsymbol{X}}\|_{F}^{2} = \operatorname{tr}((\boldsymbol{U}\boldsymbol{U}^{\top} \, \check{\boldsymbol{X}})^{\top}(\boldsymbol{U}\boldsymbol{U}^{\top} \, \check{\boldsymbol{X}})) \\ &= \operatorname{tr}(\check{\boldsymbol{X}}^{\top} \, \boldsymbol{U} \, \underbrace{\boldsymbol{U}}^{\top} \, \underbrace{\boldsymbol{U}}^{\top} \, \boldsymbol{U} \, \boldsymbol{U}^{\top} \, \check{\boldsymbol{X}}) \stackrel{(a)}{=} \operatorname{tr}(\check{\boldsymbol{X}}^{\top} \, \boldsymbol{U} \, \boldsymbol{U}^{\top} \, \check{\boldsymbol{X}}) \stackrel{(b)}{=} \operatorname{tr}(\boldsymbol{U}^{\top} \, \check{\boldsymbol{X}} \, \check{\boldsymbol{X}}^{\top} \, \boldsymbol{U}), \end{aligned}$$

where tr(.) denotes the trace of matrix, (a) is because \boldsymbol{U} is an orthogonal matrix (its columns are orthonormal), and (b) is because tr($\boldsymbol{\check{X}}^{\top} \boldsymbol{U} \boldsymbol{U}^{\top} \boldsymbol{\check{X}}$) = tr($\boldsymbol{\check{X}} \boldsymbol{\check{X}}^{\top} \boldsymbol{U} \boldsymbol{U}^{\top}$) = tr($\boldsymbol{U}^{\top} \boldsymbol{\check{X}} \boldsymbol{\check{X}}^{\top} \boldsymbol{U}$). According to Eq. (17), the $\boldsymbol{S} = \boldsymbol{\check{X}} \boldsymbol{\check{X}}^{\top}$ is the covariance matrix; therefore:

$$||\widehat{\boldsymbol{X}}||_{F}^{2} = \operatorname{tr}(\boldsymbol{U}^{\top}\boldsymbol{S}\boldsymbol{U}).$$
(21)

Projection Onto Span of Several Directions

We want to find several projection directions {u₁,..., u_p}, as columns of U ∈ ℝ^{d×p}, which maximize the squared length of reconstruction (or variance of projection):

$$\begin{array}{l} \underset{\boldsymbol{U}}{\text{maximize}} \quad \mbox{tr}(\boldsymbol{U}^{\top}\boldsymbol{S}\,\boldsymbol{U}), \\ \\ \text{subject to} \quad \boldsymbol{U}^{\top}\boldsymbol{U} = \boldsymbol{I}, \end{array}$$

where the constraint ensures that the \boldsymbol{U} is an orthogonal matrix as we assumed beforehand.

• Using Lagrange multiplier [2], we have:

$$\mathcal{L} = \operatorname{tr}(\boldsymbol{U}^{\top}\boldsymbol{S}\,\boldsymbol{U}) - \operatorname{tr}(\boldsymbol{\Lambda}^{\top}(\boldsymbol{U}^{\top}\boldsymbol{U} - \boldsymbol{I})),$$

where $\mathbf{\Lambda} \in \mathbb{R}^{p \times p}$ is a diagonal matrix $\operatorname{diag}([\lambda_1, \dots, \lambda_p]^\top)$ including the Lagrange multipliers.

$$\mathbb{R}^{d \times p} \ni \frac{\partial \mathcal{L}}{\partial \boldsymbol{U}} = 2\boldsymbol{S}\boldsymbol{U} - 2\boldsymbol{U}\boldsymbol{\Lambda} \stackrel{\text{set}}{=} 0 \implies \boldsymbol{S}\boldsymbol{U} = \boldsymbol{U}\boldsymbol{\Lambda}.$$
 (23)

• The Eq. (23) is the eigen-decomposition of S where the columns of U and the diagonal of Λ are the eigenvectors and eigenvalues of S, respectively [3]. The eigenvectors and eigenvalues are sorted from the leading (largest eigenvalue) to the trailing (smallest eigenvalue) because we are maximizing in the optimization problem. As a conclusion, if projecting onto the PCA subspace or $span\{u_1, \ldots, u_p\}$, the PCA directions $\{u_1, \ldots, u_p\}$ are the sorted eigenvectors of the covariance matrix of data X.

Truncating U

Truncating **U**

If rank(S) = d: we have p = d (we have d non-zero eigenvalues of S), so that U ∈ ℝ^{d×d}. It means that the dimensionality of the PCA subspace is d, equal to the dimensionality of the original space. Why does this happen? That is because rank(S) = d means that the data are spread wide enough in all dimensions of the original space up to a possible rotation. Therefore, the dimensionality of PCA subspace is equal to the original dimensionality; however, PCA might merely rotate the coordinate axes. In this case, U ∈ ℝ^{d×d} is a square orthogonal matrix so that ℝ^{d×d} ∋ UU^T = UU⁻¹ = I and ℝ^{d×d} ∋ U^T U = U⁻¹U = I because rank(U) = d, rank(UU^T) = d, and rank(U^TU) = d. That is why in the literature, PCA is also referred to as coordinate rotation.



If rank(S) < d and n > d: it means that we have enough data points but the data points exist on a subspace and do not fill the original space wide enough in every direction. In this case, U ∈ ℝ^{d×p} is not square and rank(U) = p < d (we have p non-zero eigenvalues of S). Therefore, ℝ^{d×d} ∋ UU^T ≠ I and ℝ^{p×p} ∋ U^TU = I because rank(U) = p, rank(UU^T) = p < d, and rank(U^TU) = p.

Truncating **U**

- If rank(S) ≤ n 1 < d: it means that we do not have enough data points to properly represent the original space and the points have an "intrinsic dimensionality". For example, we have two three-dimensional points which are one a two-dimensional line (subspace). So, similar to previous case, the data points exist on a subspace and do not fill the original space wide enough in every direction. The discussions about U, UU^T, and U^TU are similar to previous case.
- Note that we might have rank(S) = d and thus U ∈ ℝ^{d×d} but want to "truncate" the matrix U to have U ∈ ℝ^{d×p}. Truncating U means that we take a subset of best (leading) eigenvectors rather than the whole d eigenvectors with non-zero eigenvalues. In this case, again we have UU^T ≠ I and U^TU = I. The intuition of truncating is this: the variance of data might be noticeably smaller than another direction; in this case, we can only keep the p < d top eigenvectors (PCA directions) and "ignore" the PCA directions with smaller eigenvalues to have U ∈ ℝ^{d×p}.



• From all the above analyses, we conclude that as long as the columns of the matrix $U \in \mathbb{R}^{d \times p}$ are orthonormal, we always have $U^{\top}U = I$ regardless of the value p. If the orthogonal matrix U is not truncated and thus is a square matrix, we also have $UU^{\top} = I$.

If we center the data, the residual (reconstruction error) becomes r = x̃ − x̂ because the reconstructed data will also be centered according to Eq. (4). According to Eqs. (2), and (4), we have:

$$\mathbf{r} = \mathbf{x} - \hat{\mathbf{x}} = \check{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{x}} - \boldsymbol{U}\boldsymbol{U}^{\top}\check{\mathbf{x}} - \boldsymbol{\mu}_{\mathbf{x}} = \check{\mathbf{x}} - \boldsymbol{U}\boldsymbol{U}^{\top}\check{\mathbf{x}}.$$
 (24)

• This figure shows the projection of a two-dimensional point (after the data being centered) onto the first principal direction, its reconstruction, and its reconstruction error. As can be seen in this figure, the reconstruction error is different from least square error in linear regression.



• For *n* data points, we have:

$$\boldsymbol{R} := \boldsymbol{X} - \widehat{\boldsymbol{X}} = \check{\boldsymbol{X}} + \boldsymbol{\mu}_{x} - \boldsymbol{U}\boldsymbol{U}^{\top}\check{\boldsymbol{X}} - \boldsymbol{\mu}_{x} = \check{\boldsymbol{X}} - \boldsymbol{U}\boldsymbol{U}^{\top}\check{\boldsymbol{X}},$$
(25)

where $\mathbb{R}^{d \times n} \ni \boldsymbol{R} = [\boldsymbol{r}_1, \dots, \boldsymbol{r}_n]$ is the matrix of residuals.

• If we want to minimize the reconstruction error subject to the orthogonality of the projection matrix **U**, we have:

minimize
$$\| \boldsymbol{\check{X}} - \boldsymbol{U} \boldsymbol{U}^{\top} \boldsymbol{\check{X}} \|_{F}^{2},$$

subject to $\boldsymbol{U}^{\top} \boldsymbol{U} = \boldsymbol{I}.$ (26)

• The objective function can be simplified:

$$\begin{aligned} ||\breve{X} - UU^{\top}\breve{X}||_{F}^{2} \\ &= \operatorname{tr}((\breve{X} - UU^{\top}\breve{X})^{\top}(\breve{X} - UU^{\top}\breve{X})) = \operatorname{tr}((\breve{X}^{\top} - \breve{X}^{\top}UU^{\top})(\breve{X} - UU^{\top}\breve{X})) \\ &= \operatorname{tr}(\breve{X}^{\top}\breve{X} - 2\breve{X}^{\top}UU^{\top}\breve{X} + \breve{X}^{\top}U\underbrace{U^{\top}U}_{I}U^{\top}\breve{X}) \\ &= \operatorname{tr}(\breve{X}^{\top}\breve{X} - \breve{X}^{\top}UU^{\top}\breve{X}) = \operatorname{tr}(\breve{X}^{\top}\breve{X}) - \operatorname{tr}(\breve{X}^{\top}UU^{\top}\breve{X}) = \operatorname{tr}(\breve{X}^{\top}\breve{X}) - \operatorname{tr}(\breve{X}\breve{X}^{\top}UU^{\top}). \end{aligned}$$

$$\bullet \text{ Using Lagrange multiplier [2], we have:}$$

$$\mathcal{L} = \operatorname{tr}(\breve{\boldsymbol{X}}^{\top}\breve{\boldsymbol{X}}) - \operatorname{tr}(\breve{\boldsymbol{X}}\breve{\boldsymbol{X}}^{\top}\boldsymbol{\boldsymbol{U}}\boldsymbol{\boldsymbol{U}}^{\top}) + \operatorname{tr}(\boldsymbol{\Lambda}^{\top}(\boldsymbol{\boldsymbol{U}}^{\top}\boldsymbol{\boldsymbol{U}}-\boldsymbol{\boldsymbol{I}})),$$

where $\mathbf{\Lambda} \in \mathbb{R}^{\rho \times \rho}$ is a diagonal matrix $\operatorname{diag}([\lambda_1, \dots, \lambda_\rho]^\top)$ containing the Lagrange multipliers.

• Equating the derivative of Lagrangian to zero gives:

$$\mathbb{R}^{d \times p} \ni \frac{\partial \mathcal{L}}{\partial \boldsymbol{U}} = -2\check{\boldsymbol{X}}\check{\boldsymbol{X}}^{\top}\boldsymbol{U} + 2\boldsymbol{U}\boldsymbol{\Lambda} \stackrel{\text{set}}{=} \boldsymbol{0}$$

$$\implies \check{\boldsymbol{X}}\check{\boldsymbol{X}}^{\top}\boldsymbol{U} = \boldsymbol{U}\boldsymbol{\Lambda},$$

$$\stackrel{(17)}{\Longrightarrow} \boldsymbol{S}\boldsymbol{U} = \boldsymbol{U}\boldsymbol{\Lambda},$$
(27)

which is again the eigenvalue problem [3] for the covariance matrix S.

• We had the same eigenvalue problem in PCA. Therefore, **PCA subspace is the best linear** projection in terms of reconstruction error. In other words, **PCA** has the least squared error in reconstruction.

• We saw that PCA is the best in reconstruction error for *linear* projection. If we have m > 1 successive linear projections, the reconstruction is:

$$\widehat{\mathbf{X}} = \underbrace{\mathbf{U}_{1}\cdots\mathbf{U}_{m}}_{\text{reconstruct}}\underbrace{\mathbf{U}_{m}^{\top}\cdots\mathbf{U}_{1}^{\top}}_{\text{project}} \check{\mathbf{X}} + \boldsymbol{\mu}_{x}, \qquad (28)$$

which can be seen as an undercomplete *autoencoder* [4] with 2m layers without activation function (or with identity activation functions $f(\mathbf{x}) = \mathbf{x}$). The $\mu_{\mathbf{x}}$ is modeled by the intercepts included as input to the neurons of autoencoder layers.

• As we do not have any non-linearity between the projections, we can define:

•

$$\ddot{\boldsymbol{U}} := \boldsymbol{U}_1 \cdots \boldsymbol{U}_m \implies \ddot{\boldsymbol{U}}^\top = \boldsymbol{U}_m^\top \cdots \boldsymbol{U}_1^\top, \qquad (29)$$

$$\widehat{\boldsymbol{X}} = \ddot{\boldsymbol{U}}\ddot{\boldsymbol{U}}^{\top}\boldsymbol{X} + \boldsymbol{\mu}_{\boldsymbol{x}}.$$
(30)

The Eq. (30) shows that the whole autoencoder can be reduced to an undercomplete autoencoder with one hidden layer where the weight matrix is \ddot{U} . In other words, in autoencoder neural network, every layer excluding the activation function behaves as a linear projection.

 Comparing the Eqs. (8) and (30) shows that the whole autoencoder is reduced to PCA. Therefore, PCA is equivalent to an undercomplete autoencoder with one hidden layer without activation function. Therefore, if we trained weights of such an autoencoder by back-propagation [5], they will be roughly equal to the PCA directions.





PCA Using Singular Value Decomposition

PCA Using Singular Value Decomposition

• The PCA can be done using Singular Value Decomposition (SVD) of \check{X} , rather than eigen-decomposition of *S*. Consider the complete SVD of \check{X} :

$$\mathbb{R}^{d \times n} \ni \check{\mathbf{X}} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top}, \tag{31}$$

where the columns of $\boldsymbol{U} \in \mathbb{R}^{d \times d}$ (called left singular vectors) are the eigenvectors of $\boldsymbol{\check{X}} \boldsymbol{\check{X}}^{\top}$, the columns of $\boldsymbol{V} \in \mathbb{R}^{n \times n}$ (called right singular vectors) are the eigenvectors of $\boldsymbol{\check{X}}^{\top} \boldsymbol{\check{X}}$, and the $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times n}$ is a rectangular diagonal matrix whose diagonal entries (called singular values) are the square root of eigenvalues of $\boldsymbol{\check{X}} \boldsymbol{\check{X}}^{\top}$ and/or $\boldsymbol{\check{X}}^{\top} \boldsymbol{\check{X}}$. See Preliminaries slides for proof of this claim.

- According to Eq. (17), the XX^T is the covariance matrix S. In Eq. (23), we saw that the eigenvectors of S are the principal directions. On the other hand, here, we saw that the columns of U are the eigenvectors of XX^T. Hence, we can apply SVD on X and take the left singular vectors (columns of U) as the principal directions.
- An interesting thing is that in SVD of \check{X} , the columns of U are automatically sorted from largest to smallest singular values (eigenvalues) and we do not need to sort as we did in using eigenvalue decomposition for the covariance matrix.

Determining the Number of Principal Directions

Determining the Number of Principal Directions

- Usually in PCA, the components with smallest eigenvalues are cut off to reduce the data. There are different methods for estimating the best number of components to keep (denoted by *p*), such as using Bayesian model selection [6], scree plot [7], and comparing the ratio λ_j / Σ^d_{k=1} λ_k with a threshold [8] where λ_i denotes the eigenvalue related to the *j*-th principal component.
- Here, we explain the two methods of scree plot and the ratio.
- The scree plot [7] is a plot of the eigenvalues versus sorted components from the leading (having largest eigenvalue) to trailing (having smallest eigenvalue). A threshold for the vertical (eigenvalue) axis chooses the components with the large enough eigenvalues and removes the rest of the components. A good threshold is where the eigenvalue drops significantly. In most of the datasets, a significant drop of eigenvalue occurs.
- Another way to choose the best components is the ratio [8]:

$$\frac{\lambda_j}{\sum_{k=1}^d \lambda_k},\tag{32}$$

for the j-th component. Then, we sort the features from the largest to smallest ratio and select the p best components or up to the component where a significant drop of the ratio happens.

Dual Principal Component Analysis

Dual Principal Component Analysis

 Assume the case where the dimensionality of data is high and much greater than the sample size, i.e., d ≫ n. In this case, consider the incomplete SVD of X
:

$$\check{\boldsymbol{X}} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\top}, \tag{33}$$

where here, $\boldsymbol{U} \in \mathbb{R}^{d \times p}$ and $\boldsymbol{V} \in \mathbb{R}^{n \times p}$ contain the *p* leading left and right singular vectors of \boldsymbol{X} , respectively, where *p* is the number of "non-zero" singular values of \boldsymbol{X} and usually $p \ll d$.

- Here, the $\Sigma \in \mathbb{R}^{p \times p}$ is a square matrix having the *p* largest non-zero singular values of \check{X} .
- As the Σ is a square diagonal matrix and its diagonal includes non-zero entries (is full-rank), it is invertible [9]. Therefore, $\Sigma^{-1} = \text{diag}([\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_p}]^{\top})$ if we have $\Sigma = \text{diag}([\sigma_1, \dots, \sigma_p]^{\top})$.

Dual Principal Component Analysis: Projection

We had:

$$\breve{X} = U\Sigma V^{\top},$$

Recall Eq. (7) for projection onto PCA subspace: X̃ = U^TX̃. On the other hand, according to Eq. (33), we have:

$$\ddot{\mathbf{X}} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top} \implies \mathbf{U}^{\top} \ddot{\mathbf{X}} = \underbrace{\mathbf{U}^{\top} \mathbf{U}}_{\mathbf{I}} \mathbf{\Sigma} \mathbf{V}^{\top} = \mathbf{\Sigma} \mathbf{V}^{\top}.$$
 (34)

• According to Eqs. (7) and (34), we have:

$$\tilde{\boldsymbol{X}} = \boldsymbol{\Sigma} \boldsymbol{V}^{\top}$$
(35)

The Eq. (35) can be used for projecting data onto PCA subspace instead of Eq. (7). This is projection of training data in dual PCA.

Dual Principal Component Analysis: Reconstruction

• We had:

$$\breve{X} = U \Sigma V^{\top},$$

• According to Eq. (33), we have:

$$\begin{split}
\check{\mathbf{X}} &= \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \implies \check{\mathbf{X}} \mathbf{V} = \mathbf{U} \mathbf{\Sigma} \underbrace{\mathbf{V}^\top \mathbf{V}}_{I} = \mathbf{U} \mathbf{\Sigma} \\
\implies \mathbf{U} &= \check{\mathbf{X}} \mathbf{V} \mathbf{\Sigma}^{-1}.
\end{split}$$
(36)

Plugging Eq. (36) in Eq. (8) gives us:

$$\widehat{\mathbf{X}} = \mathbf{U}\widetilde{\mathbf{X}} + \mu_{x} \stackrel{(36)}{=} \widetilde{\mathbf{X}} \mathbf{V} \mathbf{\Sigma}^{-1} \widetilde{\mathbf{X}} + \mu_{x}$$

$$\stackrel{(35)}{=} \widetilde{\mathbf{X}} \mathbf{V} \underbrace{\mathbf{\Sigma}^{-1} \mathbf{\Sigma}}_{I} \mathbf{V}^{\top} + \mu_{x}$$

$$\implies \widehat{\mathbf{X}} = \widecheck{\mathbf{X}} \mathbf{V} \mathbf{V}^{\top} + \mu_{x}.$$
(37)

The Eq. (37) can be used for reconstruction of data instead of Eq. (8). This is reconstruction of training data in dual PCA.

Dual Principal Component Analysis: Out-of-sample Projection

• Recall Eq. (9) for projection of an out-of-sample point x_t onto PCA subspace:

$$\widetilde{\boldsymbol{x}}_t = \boldsymbol{U}^\top \breve{\boldsymbol{x}}_t.$$

According to Eq. (36):

$$\boldsymbol{U}=\boldsymbol{\breve{X}}\boldsymbol{V}\boldsymbol{\Sigma}^{-1},$$

we have:

$$\boldsymbol{U}^{\top} \stackrel{(36)}{=} \boldsymbol{\Sigma}^{-\top} \boldsymbol{V}^{\top} \boldsymbol{\check{X}}^{\top} \stackrel{(a)}{=} \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^{\top} \boldsymbol{\check{X}}^{\top}$$
(38)

$$\stackrel{(9)}{\Longrightarrow} \widetilde{\mathbf{x}}_t = \mathbf{\Sigma}^{-1} \mathbf{V}^\top \breve{\mathbf{X}}^\top \breve{\mathbf{x}}_t, \tag{39}$$

where (a) is because Σ^{-1} is diagonal and thus symmetric. The Eq. (39) can be used for projecting out-of-sample data point onto PCA subspace instead of Eq. (9). This is out-of-sample projection in dual PCA.

Considering all the nt out-of-sample data points, the projection is:

$$\widetilde{\boldsymbol{X}}_{t} = \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^{\top} \breve{\boldsymbol{X}}^{\top} \breve{\boldsymbol{X}}_{t}.$$
(40)

Dual Principal Component Analysis: Out-of-sample Reconstruction

• Recall Eq. (10) for reconstruction of an out-of-sample point x_t :

$$\widehat{\mathbf{x}}_t = \mathbf{U}\mathbf{U}^{\top} \widecheck{\mathbf{x}}_t + \mathbf{\mu}_x = \mathbf{U}\widetilde{\mathbf{x}}_t + \mathbf{\mu}_x.$$

According to Eqs. (36) and (38), we have:

$$UU^{\top} = \check{X}V\Sigma^{-1}\Sigma^{-1}V^{\top}\check{X}^{\top}$$
$$\stackrel{(10)}{\Longrightarrow} \hat{x}_{t} = \check{X}V\Sigma^{-2}V^{\top}\check{X}^{\top}\check{x}_{t} + \mu_{x}.$$
(41)

The Eq. (41) can be used for reconstruction of an out-of-sample data point instead of Eq. (10). This is out-of-sample reconstruction in dual PCA.

• Considering all the *n_t* out-of-sample data points, the reconstruction is:

$$\widehat{\boldsymbol{X}}_{t} = \widecheck{\boldsymbol{X}} \boldsymbol{V} \boldsymbol{\Sigma}^{-2} \boldsymbol{V}^{\top} \widecheck{\boldsymbol{X}}^{\top} \widecheck{\boldsymbol{X}}_{t} + \boldsymbol{\mu}_{x}.$$
(42)

Why is Dual PCA Useful?

• The dual PCA can be useful for two reasons:

As can be seen in Eqs. (35), (37), (39), and (41):

$$\begin{split} \widetilde{\boldsymbol{X}} &= \boldsymbol{\Sigma} \boldsymbol{V}^{\top}, \\ \widehat{\boldsymbol{X}} &= \breve{\boldsymbol{X}} \boldsymbol{V} \boldsymbol{V}^{\top} + \boldsymbol{\mu}_{x}, \\ \widetilde{\boldsymbol{x}}_{t} &= \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^{\top} \breve{\boldsymbol{X}}^{\top} \breve{\boldsymbol{x}}_{t}, \\ \widehat{\boldsymbol{x}}_{t} &= \breve{\boldsymbol{X}} \boldsymbol{V} \boldsymbol{\Sigma}^{-2} \boldsymbol{V}^{\top} \breve{\boldsymbol{X}}^{\top} \breve{\boldsymbol{x}}_{t} + \boldsymbol{\mu}_{x}, \end{split}$$

the formulae for dual PCA only include V and not U. The columns of V are the eigenvectors of $\check{X}^{\top}\check{X} \in \mathbb{R}^{n \times n}$ and the columns of U are the eigenvectors of $\check{X}\check{X}^{\top} \in \mathbb{R}^{d \times d}$. In case the dimensionality of data is much high and greater than the sample size, i.e., $n \ll d$, computation of eigenvectors of $\check{X}^{\top}\check{X}$ is easier and faster than $\check{X}\check{X}^{\top}$ and also requires less storage. Therefore, dual PCA is more efficient than direct PCA in this case in terms of both speed and storage. Note that the results of PCA and dual PCA are exactly the same.

Some inner product forms, such as $\check{X}^{\top}\check{x}_t$, have appeared in the formulae of dual PCA. This provides opportunity for kernelizing the PCA to have kernel PCA using the so-called kernel trick. As will be seen in the next section, we use dual PCA in formulation of kernel PCA.

Kernel Principal Component Analysis

Kernel Principal Component Analysis

The PCA is a linear method because the projection is linear. In case the data points exist
on a non-linear sub-manifold, the linear subspace learning might not be completely
effective.



- In order to handle this problem of PCA, we have two options. We should either change PCA to become a nonlinear method or we can leave the PCA to be linear but change the data hoping to fall on a linear or close to linear manifold.
- Here, we do the latter so we change the data. We increase the dimensionality of data by mapping the data to feature space with higher dimensionality hoping that in the feature space, it falls on a linear manifold. This is referred to as "blessing of dimensionality" in the literature [10] which is pursued using kernels [11]. This PCA method which uses the kernel of data is named "kernel PCA" [12].

Kernel and Centered Kernel

Kernel:

$$k(\mathbf{x}_1, \mathbf{x}_2) := \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_2), \tag{43}$$

$$\boldsymbol{K}(\boldsymbol{X}_1, \boldsymbol{X}_2) := \boldsymbol{\Phi}(\boldsymbol{X}_1)^\top \boldsymbol{\Phi}(\boldsymbol{X}_2). \tag{44}$$

If we want the pulled data Φ(X) to be centered, i.e.:

$$\check{\Phi}(\boldsymbol{X}) := \Phi(\boldsymbol{X})\boldsymbol{H},\tag{45}$$

we should double center the kernel matrix because if we use centered pulled data, we have:

$$\check{\Phi}(\boldsymbol{X})^{\top}\check{\Phi}(\boldsymbol{X}) = \left(\Phi(\boldsymbol{X})\boldsymbol{H}\right)^{\top}\left(\Phi(\boldsymbol{X})\boldsymbol{H}\right) = \boldsymbol{H}\Phi(\boldsymbol{X})^{\top}\Phi(\boldsymbol{X})\boldsymbol{H} = \boldsymbol{H}\boldsymbol{K}_{x}\boldsymbol{H},$$

which is the double-centered kernel matrix. Thus:

$$\breve{\boldsymbol{K}}_{\boldsymbol{X}} := \boldsymbol{H}\boldsymbol{K}_{\boldsymbol{X}}\boldsymbol{H} = \breve{\boldsymbol{\Phi}}(\boldsymbol{X})^{\top}\breve{\boldsymbol{\Phi}}(\boldsymbol{X}), \tag{46}$$

where \breve{K}_x denotes the double-centered kernel matrix.

Kernel Principal Component Analysis: Projection

• We apply incomplete SVD on the centered pulled (mapped) data $\check{\Phi}(X)$:

$$\mathbb{R}^{t \times n} \ni \breve{\Phi}(\boldsymbol{X}) = \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^{\top}, \tag{47}$$

where $\boldsymbol{U} \in \mathbb{R}^{t \times p}$ and $\boldsymbol{V} \in \mathbb{R}^{n \times p}$ contain the *p* leading left and right singular vectors of $\check{\boldsymbol{\Phi}}(\boldsymbol{X})$, respectively, where *p* is the number of "non-zero" singular values of $\check{\boldsymbol{\Phi}}(\boldsymbol{X})$ and usually $p \ll t$. Here, the $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ is a square matrix having the *p* largest non-zero singular values of $\check{\boldsymbol{\Phi}}(\boldsymbol{X})$.

• However, as mentioned before, the pulled data are not necessarily available so Eq. (47) cannot be done. The kernel, however, is available. Therefore, we apply eigen-decomposition [3] to the double-centered kernel:

$$\breve{K}_{\times}V = V\Lambda, \tag{48}$$

where the columns of V and the diagonal of Λ are the eigenvectors and eigenvalues of \breve{K}_x , respectively.

- The columns of \mathbf{V} in Eq. (47) are the right singular vectors of $\mathbf{\Phi}(\mathbf{X})$ which are equivalent to the eigenvectors of $\mathbf{\Phi}(\mathbf{X})^{\top}\mathbf{\Phi}(\mathbf{X}) = \mathbf{K}_{\times}$, according to the Preliminaries slides. Also, according to those slides, the diagonal of $\mathbf{\Sigma}$ in Eq. (47) is equivalent to the square root of eigenvalues of \mathbf{K}_{\times} .
- Therefore, in practice where the pulling function is not necessarily available, we use Eq. (48) in order to find the V and Σ in Eq. (47).

Kernel Principal Component Analysis: Projection

We had:

$$\breve{K}_{X}V=V\Lambda.$$

• The Eq. (48) can be restated as:

$$\breve{\boldsymbol{K}}_{\times}\boldsymbol{V}=\boldsymbol{V}\boldsymbol{\Sigma}^{2},\tag{49}$$

to be compatible to Eq. (47):

$$\mathbb{R}^{t\times n}\ni \breve{\Phi}(\boldsymbol{X})=\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\top}.$$

- It is noteworthy that because of using Eq. (49) instead of Eq. (47), the projection directions *U* are not available in kernel PCA to be observed or plotted.
- Similar to what we did for Eq. (35):

$$\begin{split} \check{\Phi}(\boldsymbol{X}) &= \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\top} \\ \implies \boldsymbol{U}^{\top}\check{\Phi}(\boldsymbol{X}) = \underbrace{\boldsymbol{U}^{\top}\boldsymbol{U}}_{\boldsymbol{I}}\boldsymbol{\Sigma}\boldsymbol{V}^{\top} = \boldsymbol{\Sigma}\boldsymbol{V}^{\top} \\ \therefore \qquad \Phi(\widetilde{\boldsymbol{X}}) = \boldsymbol{U}^{\top}\check{\Phi}(\boldsymbol{X}) = \boldsymbol{\Sigma}\boldsymbol{V}^{\top}, \end{split}$$
(50)

where $\boldsymbol{\Sigma}$ and \boldsymbol{V} are obtained from Eq. (49). The Eq. (50) is projection of the training data in kernel PCA.

Kernel Principal Component Analysis: Reconstruction

Similar to what we did for Eq. (37):

$$\check{\Phi}(X) = U\Sigma V^{\top} \implies \check{\Phi}(X)V = U\Sigma \underbrace{V^{\top}V}_{I} = U\Sigma$$

$$\implies U = \check{\Phi}(X)V\Sigma^{-1}.$$
(51)

• Therefore, the reconstruction is:

$$\Phi(\widehat{X}) = U\Phi(\widetilde{X}) + \mu_{x} \stackrel{(51)}{=} \check{\Phi}(X)V\Sigma^{-1}\Phi(\widetilde{X}) + \mu_{x}$$

$$\stackrel{(50)}{=} \check{\Phi}(X)V\sum_{I} V^{\top} + \mu_{x}$$

$$\Longrightarrow \Phi(\widehat{X}) = \check{\Phi}(X)VV^{\top} + \mu_{x}.$$
(52)

However, the $\check{\Phi}(X)$ is not available necessarily; therefore, we cannot reconstruct the training data in kernel PCA.

Kernel Principal Component Analysis: Out-of-sample Projection

• Similar to what we did for Eq. (39):

$$\boldsymbol{U}^{\top} \stackrel{(51)}{=} \boldsymbol{\Sigma}^{-\top} \boldsymbol{V}^{\top} \check{\boldsymbol{\Phi}}(\boldsymbol{X})^{\top} \stackrel{(a)}{=} \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^{\top} \check{\boldsymbol{\Phi}}(\boldsymbol{X})^{\top}
\Longrightarrow \phi(\tilde{\boldsymbol{x}}_{t}) = \boldsymbol{U}^{\top} \check{\boldsymbol{\phi}}(\boldsymbol{x}_{t}) = \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^{\top} \check{\boldsymbol{\Phi}}(\boldsymbol{X})^{\top} \check{\boldsymbol{\phi}}(\boldsymbol{x}_{t}),
\Longrightarrow \phi(\tilde{\boldsymbol{x}}_{t}) = \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^{\top} \check{\boldsymbol{k}}_{t},$$
(53)

where (a) is because Σ^{-1} is diagonal and thus symmetric and the $\check{k}_t \in \mathbb{R}^n$ is calculated by (see Appendix C in our tutorial [13] for proof):

$$\mathbb{R}^n \ni \check{\boldsymbol{k}}_t = \boldsymbol{k}_t - \frac{1}{n} \boldsymbol{1}_{n \times n} \boldsymbol{k}_t - \frac{1}{n} \boldsymbol{K} \boldsymbol{1}_n + \frac{1}{n^2} \boldsymbol{1}_{n \times n} \boldsymbol{K} \boldsymbol{1}_n.$$
(54)

The Eq. (53) is the projection of out-of-sample data in kernel PCA.
Considering all the nt out-of-sample data points, Xt, the projection is:

$$\phi(\widetilde{\boldsymbol{X}}_t) = \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^\top \breve{\boldsymbol{K}}_t, \tag{55}$$

where \breve{K}_t is calculated by (see Appendix C in our tutorial [13] for proof):

$$\mathbb{R}^{n \times n_t} \ni \breve{\boldsymbol{K}}_t = \boldsymbol{K}_t - \frac{1}{n} \boldsymbol{1}_{n \times n} \boldsymbol{K}_t - \frac{1}{n} \boldsymbol{K} \boldsymbol{1}_{n \times n_t} + \frac{1}{n^2} \boldsymbol{1}_{n \times n} \boldsymbol{K} \boldsymbol{1}_{n \times n_t},$$
(56)

where $\mathbb{R}^{n \times n} \ni \mathbf{1}_{n \times n} := \mathbf{1}_n \mathbf{1}_n^\top$, $\mathbb{R}^{n \times n_t} \ni \mathbf{1}_{n \times n_t} := \mathbf{1}_n \mathbf{1}_{n_t}^\top$, $\mathbb{R}^n \ni \mathbf{1}_n := [1, \dots, 1]^\top$, and $\mathbb{R}^{n_t} \ni \mathbf{1}_{n_t} := [1, \dots, 1]^\top$.

Kernel Principal Component Analysis: Out-of-sample Reconstruction

• Similar to what we did for Eq. (41):

$$\implies \boldsymbol{U}\boldsymbol{U}^{\top} \stackrel{(51)}{=} \boldsymbol{\check{\Phi}}(\boldsymbol{X})\boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}^{-1}\boldsymbol{V}^{\top}\boldsymbol{\check{\Phi}}(\boldsymbol{X})^{\top}$$
$$\implies \boldsymbol{\phi}(\hat{\boldsymbol{x}}_{t}) = \boldsymbol{\check{\Phi}}(\boldsymbol{X})\boldsymbol{V}\boldsymbol{\Sigma}^{-2}\boldsymbol{V}^{\top}\boldsymbol{\check{\Phi}}(\boldsymbol{X})^{\top}\boldsymbol{\check{\phi}}(\boldsymbol{x}_{t}) + \boldsymbol{\mu}_{x}$$
$$\implies \boldsymbol{\phi}(\hat{\boldsymbol{x}}_{t}) = \boldsymbol{\check{\Phi}}(\boldsymbol{X})\boldsymbol{V}\boldsymbol{\Sigma}^{-2}\boldsymbol{V}^{\top}\boldsymbol{\check{k}}_{t} + \boldsymbol{\mu}_{x}, \qquad (57)$$

where the $\check{k}_t \in \mathbb{R}^n$ is calculated by Eq. (54).

• Considering all the n_t out-of-sample data points, X_t , the reconstruction is:

$$\boldsymbol{\Phi}(\widehat{\boldsymbol{X}}_{t}) = \check{\boldsymbol{\Phi}}(\boldsymbol{X}) \boldsymbol{V} \boldsymbol{\Sigma}^{-2} \boldsymbol{V}^{\top} \check{\boldsymbol{K}}_{t} + \boldsymbol{\mu}_{x},$$
(58)

where \breve{K}_t is calculated by Eq. (56).

- In Eq. (57), the Φ(X) appeared at the left of expression, is not available necessarily; therefore, we cannot reconstruct an out-of-sample point in kernel PCA.
- According to Eqs. (52) and (57), we conclude that kernel PCA is not able to reconstruct any data, whether training or out-of-sample.

Why is Kernel PCA Useful?

- Finally, it is noteworthy that as the choice of the best kernel might be hard, the kernel PCA is not "always" effective in practice [9].
- However, it provides us some useful theoretical insights for explaining the PCA, Multi-Dimensional Scaling (MDS) [14], Isomap [15], Locally Linear Embedding (LLE) [16], and Laplacian Eigenmap (LE) [17] as special cases of kernel PCA with their own kernels (see [18] and chapter 2 in [19]).

Supervised Principal Component Analysis Using HSIC

Hilbert-Schmidt Independence Criterion

- Suppose we want to measure the **dependence** of two random variables.
- Measuring the correlation between them is easier because correlation is just "linear" dependence.
- According to [20], two random variables are independent if and only if any bounded continuous functions of them are uncorrelated. Therefore, if we map the two random variables x and y to two different ("separable") Reproducing Kernel Hilbert Spaces (RKHSs) and have φ(x) and φ(y), we can measure the correlation of φ(x) and φ(y) in Hilbert space to have an estimation of dependence of x and y in the original space.
- An empirical estimation of the Hilbert-Schmidt Independence Criterion (HSIC) is [21]:

$$HSIC := \frac{1}{(n-1)^2} \operatorname{tr}(\ddot{\mathbf{K}}_x \mathbf{H} \mathbf{K}_y \mathbf{H}),$$
(59)

where $\ddot{\mathbf{K}}_x$ and \mathbf{K}_y are the kernels over \mathbf{x} and \mathbf{y} , respectively. In other words, $\ddot{\mathbf{K}}_x = \phi(\mathbf{x})^{\top} \phi(\mathbf{x})$ and $\mathbf{K}_y = \phi(\mathbf{y})^{\top} \phi(\mathbf{y})$.

• The term $1/(n-1)^2$ is used for normalization. The **H** is the centering matrix:

$$\mathbb{R}^{n \times n} \ni \boldsymbol{H} = \boldsymbol{I} - (1/n) \boldsymbol{1} \boldsymbol{1}^{\top}.$$
 (60)

- The $HK_{y}H$ double centers the K_{y} in HSIC.
- The HSIC (Eq. (59)) measures the dependence of two random variable vectors x and y. Note that HSIC = 0 and HSIC > 0 mean that x and y are independent and dependent, respectively. The greater the HSIC, the greater dependence they have.

- Supervised PCA (SPCA) [22] uses the HSIC. We have the data X = [x₁,..., x_n] ∈ ℝ^{d×n} and the labels Y = [y₁,..., y_n] ∈ ℝ^{ℓ×n}, where ℓ is the dimensionality of the labels and we usually have ℓ = 1. However, in case the labels are encoded (e.g., one-hot-encoded) or SPCA is used for regression (e.g., see [23]), we have ℓ > 1.
- SPCA tries to maximize the dependence of the projected data points U^TX and the labels
 Y. It uses a linear kernel for the projected data points:

$$\ddot{\boldsymbol{K}}_{\boldsymbol{X}} = (\boldsymbol{U}^{\top}\boldsymbol{X})^{\top}(\boldsymbol{U}^{\top}\boldsymbol{X}) = \boldsymbol{X}^{\top}\boldsymbol{U}\boldsymbol{U}^{\top}\boldsymbol{X}, \tag{61}$$

and an arbitrary kernel K_{y} over Y.

• For classification task, one of the best choices for the K_y is delta kernel [22] where the (i, j)-th element of kernel is:

$$\boldsymbol{K}_{\boldsymbol{y}} = \delta_{\boldsymbol{y}_i, \boldsymbol{y}_j} := \begin{cases} 1 & \text{if } \boldsymbol{y}_i = \boldsymbol{y}_j, \\ 0 & \text{if } \boldsymbol{y}_i \neq \boldsymbol{y}_j, \end{cases}$$
(62)

where δ_{y_i,y_j} is the Kronecker delta which is one if the x_i and x_j belong to the same class.
 Another good choice for kernel in classification task in SPCA is an arbitrary kernel (e.g., linear kernel K_y = Y^TY) over Y where the columns of Y are one-hot encoded. This is a good choice because the distances of classes will be equal; otherwise, some classes will fall closer than the others for no reason and fairness between classes goes away.

● The SPCA can also be used for regression (e.g., see [23]) and that is one of the advantages of SPCA. In that case, a good choice for K_y is an arbitrary kernel (e.g., linear kernel K_y = Y^TY) over Y where the columns of the Y, i.e., labels, are the observations in regression. Here, the distances of observations have meaning and should not be manipulated.

HSIC was:

$$HSIC := \frac{1}{(n-1)^2} \operatorname{tr}(\ddot{\boldsymbol{K}}_{\boldsymbol{X}} \boldsymbol{H} \boldsymbol{K}_{\boldsymbol{Y}} \boldsymbol{H}).$$

The HSIC in SPCA case becomes:

$$HSIC = \frac{1}{(n-1)^2} \operatorname{tr}(\boldsymbol{X}^{\top} \boldsymbol{U} \boldsymbol{U}^{\top} \boldsymbol{X} \boldsymbol{H} \boldsymbol{K}_{y} \boldsymbol{H}).$$
(63)

where $\boldsymbol{U} \in \mathbb{R}^{d \times p}$ is the unknown projection matrix for projection onto the SPCA subspace and should be found. The desired dimensionality of the subspace is p and usually $p \ll d$.

• We should maximize the HSIC in order to maximize the dependence of $U^{\top}X$ and Y. Hence:

$$\begin{array}{ll} \underset{\boldsymbol{U}}{\text{maximize}} & \operatorname{tr}(\boldsymbol{X}^{\top}\boldsymbol{U}\boldsymbol{U}^{\top}\boldsymbol{X}\boldsymbol{H}\boldsymbol{K}_{y}\boldsymbol{H}), \\ \text{subject to} & \boldsymbol{U}^{\top}\boldsymbol{U} = \boldsymbol{I}, \end{array}$$

$$(64)$$

where the constraint ensures that the ${\pmb U}$ is an orthogonal matrix, i.e., the SPCA directions are orthonormal.

We had:

 $\begin{array}{ll} \underset{\boldsymbol{U}}{\text{maximize}} & \operatorname{tr}(\boldsymbol{X}^{\top}\boldsymbol{U}\boldsymbol{U}^{\top}\boldsymbol{X}\boldsymbol{H}\boldsymbol{K}_{y}\boldsymbol{H}),\\ \\ \underset{\boldsymbol{U}}{\text{subject to}} & \boldsymbol{U}^{\top}\boldsymbol{U}=\boldsymbol{I}. \end{array}$

• Using Lagrangian [2], we have:

$$\begin{aligned} \mathcal{L} &= \operatorname{tr}(\boldsymbol{X}^{\top}\boldsymbol{U}\boldsymbol{U}^{\top}\boldsymbol{X}\boldsymbol{H}\boldsymbol{K}_{y}\boldsymbol{H}) - \operatorname{tr}(\boldsymbol{\Lambda}^{\top}(\boldsymbol{U}^{\top}\boldsymbol{U}-\boldsymbol{I})) \\ &\stackrel{(a)}{=} \operatorname{tr}(\boldsymbol{U}\boldsymbol{U}^{\top}\boldsymbol{X}\boldsymbol{H}\boldsymbol{K}_{y}\boldsymbol{H}\boldsymbol{X}^{\top}) - \operatorname{tr}(\boldsymbol{\Lambda}^{\top}(\boldsymbol{U}^{\top}\boldsymbol{U}-\boldsymbol{I})), \end{aligned}$$

where (a) is because of the cyclic property of trace and $\mathbf{\Lambda} \in \mathbb{R}^{p \times p}$ is a diagonal matrix $\operatorname{diag}([\lambda_1, \dots, \lambda_p]^\top)$ including the Lagrange multipliers.

• Setting the derivative of Lagrangian to zero gives:

$$\mathbb{R}^{d \times p} \ni \frac{\partial \mathcal{L}}{\partial \boldsymbol{U}} = 2\boldsymbol{X} \boldsymbol{H} \boldsymbol{K}_{\boldsymbol{y}} \boldsymbol{H} \boldsymbol{X}^{\top} \boldsymbol{U} - 2\boldsymbol{U} \boldsymbol{\Lambda} \stackrel{\text{set}}{=} 0$$
$$\implies \boldsymbol{X} \boldsymbol{H} \boldsymbol{K}_{\boldsymbol{y}} \boldsymbol{H} \boldsymbol{X}^{\top} \boldsymbol{U} = \boldsymbol{U} \boldsymbol{\Lambda}, \tag{65}$$

which is the eigen-decomposition of XHK_yHX^{\top} where the columns of U and the diagonal of Λ are the eigenvectors and eigenvalues of XHK_yHX^{\top} , respectively [3].

- The eigenvectors and eigenvalues are sorted from the leading (largest eigenvalue) to the trailing (smallest eigenvalue) because we are maximizing in the optimization problem.
- As a conclusion, if projecting onto the SPCA subspace or span{u₁,..., u_p}, the SPCA directions {u₁,..., u_p} are the sorted eigenvectors of XHK_yHX[⊤]. In other words, the columns of the projection matrix U in SPCA are the p leading eigenvectors of XHK_yHX[⊤].

• Similar to what we had in PCA, the projection, projection of out-of-sample, reconstruction, and reconstruction of out-of-sample in SPCA are:

$$\widetilde{\boldsymbol{X}} = \boldsymbol{U}^{\top} \boldsymbol{X}, \tag{66}$$

$$\widetilde{\boldsymbol{x}}_t = \boldsymbol{U}^\top \boldsymbol{x}_t, \tag{67}$$

$$\widehat{\boldsymbol{X}} = \boldsymbol{U}\boldsymbol{U}^{\top}\boldsymbol{X} = \boldsymbol{U}\widetilde{\boldsymbol{X}},\tag{68}$$

$$\widehat{\boldsymbol{x}}_t = \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{x}_t = \boldsymbol{U}\widetilde{\boldsymbol{x}}_t, \tag{69}$$

respectively.

- In SPCA, there is no need to center the data as the centering is already handled by *H* in HSIC.
- Considering all the n_t out-of-sample data points, the projection and reconstruction are:

$$\widetilde{\boldsymbol{X}}_t = \boldsymbol{U}^\top \boldsymbol{X}_t, \tag{70}$$

$$\widehat{\boldsymbol{X}}_t = \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{X}_t = \boldsymbol{U}\widetilde{\boldsymbol{X}}_t, \tag{71}$$

respectively.

PCA is a special case of SPCA!

- Not considering the similarities of the labels means that we do not care about the class labels so we are unsupervised.
- If we do not consider the similarities of labels, the kernel over the labels becomes the identity matrix, $K_{\gamma} = I$.
- According to Eq. (65), SPCA is the eigen-decomposition of XHK_yHX[⊤]. In this case, this matrix becomes:

$$XHK_{y}HX^{\top} = XHIHX^{\top} = XHIH^{\top}X^{\top}$$
$$= XHH^{\top}X^{\top} = (XH)(XH)^{\top}$$
$$\stackrel{(5)}{=} \breve{X}\breve{X}^{\top} \stackrel{(17)}{=} S,$$

which is the covariance matrix whose eigenvectors are the PCA directions.

• Thus, if we do not consider the similarities of labels, i.e., we are unsupervised, SPCA reduces to PCA as expected.

More Information

 For reading about dual SPCA, kernel SPCA, and eigenfaces, see our tutorial: "Unsupervised and supervised principal component analysis: Tutorial" [13]

Acknowledgment

- Some slides are based on our tutorial paper: "Unsupervised and supervised principal component analysis: Tutorial" [13]
- Some slides of this slide deck are inspired by teachings of Prof. Ali Ghodsi at University of Waterloo, Department of Statistics.
- The code of PCA in my GitHub page (in Python language): https://github.com/bghojogh/Principal-Component-Analysis
- PCA in sklearn: https: //scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html

References

- S. T. Dumais, "Latent semantic analysis," Annual review of information science and technology, vol. 38, no. 1, pp. 188–230, 2004.
- [2] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [3] B. Ghojogh, F. Karray, and M. Crowley, "Eigenvalue and generalized eigenvalue problems: Tutorial," arXiv preprint arXiv:1903.11240, 2019.
- [4] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [5] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [6] T. P. Minka, "Automatic choice of dimensionality for pca," in Advances in neural information processing systems, pp. 598–604, 2001.
- [7] R. B. Cattell, "The scree test for the number of factors," *Multivariate behavioral research*, vol. 1, no. 2, pp. 245–276, 1966.
- [8] H. Abdi and L. J. Williams, "Principal component analysis," Wiley interdisciplinary reviews: computational statistics, vol. 2, no. 4, pp. 433–459, 2010.
- [9] A. Ghodsi, "Dimensionality reduction: a short tutorial," Department of Statistics and Actuarial Science, Univ. of Waterloo, Ontario, Canada, vol. 37, 2006.

References (cont.)

- [10] D. L. Donoho, "High-dimensional data analysis: The curses and blessings of dimensionality," AMS math challenges lecture, 2000.
- [11] T. Hofmann, B. Schölkopf, and A. J. Smola, "Kernel methods in machine learning," The annals of statistics, pp. 1171–1220, 2008.
- [12] B. Schölkopf, A. Smola, and K.-R. Müller, "Kernel principal component analysis," in International conference on artificial neural networks, pp. 583–588, Springer, 1997.
- [13] B. Ghojogh and M. Crowley, "Unsupervised and supervised principal component analysis: Tutorial," arXiv preprint arXiv:1906.03148, 2019.
- [14] M. A. Cox and T. F. Cox, "Multidimensional scaling," in *Handbook of data visualization*, pp. 315–347, Springer, 2008.
- [15] J. B. Tenenbaum, V. De Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *science*, vol. 290, no. 5500, pp. 2319–2323, 2000.
- [16] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," science, vol. 290, no. 5500, pp. 2323–2326, 2000.
- [17] M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural computation*, vol. 15, no. 6, pp. 1373–1396, 2003.
- [18] J. H. Ham, D. D. Lee, S. Mika, and B. Schölkopf, "A kernel view of the dimensionality reduction of manifolds," in *International Conference on Machine Learning*, 2004.

References (cont.)

- [19] H. Strange and R. Zwiggelaar, Open Problems in Spectral Dimensionality Reduction. Springer, 2014.
- [20] M. Hein and O. Bousquet, "Kernels, associated structures and generalizations," Max-Planck-Institut fuer biologische Kybernetik, Technical Report, 2004.
- [21] A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf, "Measuring statistical dependence with Hilbert-Schmidt norms," in *International conference on algorithmic learning theory*, pp. 63–77, Springer, 2005.
- [22] E. Barshan, A. Ghodsi, Z. Azimifar, and M. Z. Jahromi, "Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds," *Pattern Recognition*, vol. 44, no. 7, pp. 1357–1371, 2011.
- [23] B. Ghojogh and M. Crowley, "Instance ranking and numerosity reduction using matrix decomposition and subspace learning," in Advances in Artificial Intelligence: 32nd Canadian Conference on Artificial Intelligence, Canadian AI 2019, Springer, 2019.