

Appendix A for EigenGAN

In Sec. 1-8, we derive the analytical result of maximum likelihood estimation (MLE) for the linear case of the proposed EigenGAN. Sec. 8 discusses the MLE result and the relation among the linear EigenGAN, the Principal Component Analysis (PCA) [1], and the Probabilistic PCA [2].

1. The Likelihood

The linear EigenGAN relates a d -dimension observation vector \mathbf{x} to a corresponding q -dimension ($q \leq d$) latent variables \mathbf{z} by an affine transform \mathbf{UL} and a translation $\boldsymbol{\mu}$, which is formulated as

$$\mathbf{x} = \mathbf{ULz} + \boldsymbol{\mu} + \sigma\boldsymbol{\epsilon}, \quad (1)$$

with constraints:

$$\mathbf{z} \sim \mathcal{N}_q(\mathbf{0}, \mathbf{I}), \quad (2)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}_d(\mathbf{0}, \mathbf{I}), \text{ independent of } \mathbf{z}, \quad (3)$$

$$\mathbf{U}^T\mathbf{U} = \mathbf{I}, \mathbf{U} \text{ is of size } d \times q, \quad (4)$$

$$\mathbf{L} \text{ is a } q \times q \text{ diagonal matrix.} \quad (5)$$

The noise vector $\boldsymbol{\epsilon}$ in Eq. (1) is introduced to compensate the missing energy (missing rank) since the rank of the latent variables is no more than the rank of the observation ($q \leq d$). According to Eq. (1)-(3), the probability density function of \mathbf{x} is

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\mathbf{C}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad (6)$$

where

$$\mathbf{C} = \mathbf{UL}^2\mathbf{U}^T + \sigma^2\mathbf{I}, \quad (7)$$

$$\mathbf{C}^{-1} = \mathbf{UMU}^T + \sigma^{-2}\mathbf{I}, \quad (8)$$

$$\mathbf{M} = (\mathbf{L}^2 + \sigma^2\mathbf{I})^{-1} - \sigma^{-2}\mathbf{I}. \quad (9)$$

Then for n observations $\{\mathbf{x}\}_{i=1}^n$, the log-likelihood is

$$\mathcal{L}_1 = -\frac{n}{2} \left\{ d \log 2\pi + \log |\mathbf{C}| + \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\}. \quad (10)$$

According to Eq. (7), only the square value of \mathbf{L} affects the probability density function, therefore we can assume the elements of \mathbf{L} to be non-negative. Further, for convenience of the following analysis, without loss of generality, we organize \mathbf{L} by grouping and sorting it by the value of the diagonal elements:

$$\mathbf{L} = \text{diag} (l_1 \mathbf{I}_{d_1}, l_2 \mathbf{I}_{d_2}, \dots, l_p \mathbf{I}_{d_p}), \quad (11)$$

where $l_1 > l_2 > \dots > l_p \geq 0$; \mathbf{I}_{d_j} denotes a $d_j \times d_j$ identity matrix, $d_j \neq 0$, and $d_1 + d_2 + \dots + d_p = q$. According to Eq. (9) and (11), \mathbf{M} also has a grouped form:

$$\mathbf{M} = \text{diag} \left(((l_1^2 + \sigma^2)^{-1} - \sigma^{-2}) \mathbf{I}_{d_1}, \dots, ((l_p^2 + \sigma^2)^{-1} - \sigma^{-2}) \mathbf{I}_{d_p} \right). \quad (12)$$

And we can also define a block form of \mathbf{U} accordingly:

$$\mathbf{U} = [\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_p], \quad (13)$$

where \mathbf{U}_i is of size $d \times d_i$.

2. Determination of $\boldsymbol{\mu}$

The partial derivative of the log-likelihood \mathcal{L}_1 (10) with respect to $\boldsymbol{\mu}$ is

$$\frac{\partial \mathcal{L}_1}{\partial \boldsymbol{\mu}} = \sum_{i=1}^n \mathbf{C}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}). \quad (14)$$

Then the stationary point is

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \bar{\mathbf{x}}. \quad (15)$$

Since \mathcal{L}_1 is a concave function of $\boldsymbol{\mu}$, the above stationary point is also the global maximum point.

3. Determination of \mathbf{U} : Part (1)

Substituting Eq. (15) into the log-likelihood \mathcal{L}_1 (10), we obtain a new objective:

$$\mathcal{L}_2 = -\frac{n}{2} \{d \log 2\pi + \log |\mathbf{C}| + \text{tr}(\mathbf{S}\mathbf{C}^{-1})\}, \quad (16)$$

where

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^{\text{T}}, \quad (17)$$

i.e., the covariance matrix of the data. According to Eq. (4), the maximization of \mathcal{L}_2 (16) with respect to \mathbf{U} is a constrained optimization as below:

$$\begin{aligned} & \max_{\mathbf{U}} \quad \mathcal{L}_2 \\ & \text{subject to} \quad \mathbf{U}^{\text{T}}\mathbf{U} = \mathbf{I} \end{aligned}$$

Introducing the Lagrange multiplier \mathbf{H} , the Lagrangian function is

$$\begin{aligned} \mathcal{L}_{\mathbf{U}} &= \mathcal{L}_2 + \text{tr}(\mathbf{H}^{\text{T}}(\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{I})) \\ &= -\frac{n}{2} \{d \log 2\pi + \log |\mathbf{C}| + \text{tr}(\mathbf{S}\mathbf{C}^{-1})\} + \text{tr}(\mathbf{H}^{\text{T}}(\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{I})). \end{aligned} \quad (18)$$

Then the partial derivative of $\mathcal{L}_{\mathbf{U}}$ with respect to \mathbf{U} is

$$\frac{\partial \mathcal{L}_{\mathbf{U}}}{\partial \mathbf{U}} = -n \left\{ \mathbf{U} \left((\mathbf{L}^2 + \sigma^2 \mathbf{I})^{-1} \mathbf{L}^2 + \frac{\mathbf{H} + \mathbf{H}^{\text{T}}}{2} \right) + \mathbf{SUM} \right\}. \quad (19)$$

At the stationary point,

$$\mathbf{SUM} = -\mathbf{U} \left((\mathbf{L}^2 + \sigma^2 \mathbf{I})^{-1} \mathbf{L}^2 + \frac{\mathbf{H} + \mathbf{H}^{\text{T}}}{2} \right). \quad (20)$$

Left multiplying the above equation by \mathbf{U}^{T} and using $\mathbf{U}^{\text{T}}\mathbf{U} = \mathbf{I}$, we obtain

$$\mathbf{U}^{\text{T}}\mathbf{SUM} = -(\mathbf{L}^2 + \sigma^2 \mathbf{I})^{-1} \mathbf{L}^2 - \frac{\mathbf{H} + \mathbf{H}^{\text{T}}}{2}. \quad (21)$$

The right-hand side of above equation is a symmetric matrix, therefore the left-hand side $\mathbf{U}^{\text{T}}\mathbf{SUM}$ is also symmetric. Furthermore, since both $\mathbf{U}^{\text{T}}\mathbf{S}\mathbf{U}$ and \mathbf{M} are also symmetric, to satisfy the symmetry of $\mathbf{U}^{\text{T}}\mathbf{SUM}$, according to the form of \mathbf{M} in Eq. (12), $\mathbf{U}^{\text{T}}\mathbf{S}\mathbf{U}$ must have a similar block diagonal form:

$$\mathbf{U}^{\text{T}}\mathbf{S}\mathbf{U} = \text{diag}(\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_p) \quad (22)$$

$$= \text{diag}(\mathbf{Q}_1^{\text{T}}\mathbf{\Lambda}_1\mathbf{Q}_1, \mathbf{Q}_2^{\text{T}}\mathbf{\Lambda}_2\mathbf{Q}_2, \dots, \mathbf{Q}_p^{\text{T}}\mathbf{\Lambda}_p\mathbf{Q}_p), \quad (23)$$

where \mathbf{A}_j is a $d_j \times d_j$ symmetric matrix, and $\mathbf{Q}_j^\top \mathbf{\Lambda}_j \mathbf{Q}_j$ is the eigendecomposition of \mathbf{A}_j . Using Eq. (20), (21), and (23), we can derive

$$\text{SUM} = \mathbf{U} \cdot \text{diag} \left(\mathbf{Q}_1^\top \mathbf{\Lambda}_1 \mathbf{Q}_1, \mathbf{Q}_2^\top \mathbf{\Lambda}_2 \mathbf{Q}_2, \dots, \mathbf{Q}_p^\top \mathbf{\Lambda}_p \mathbf{Q}_p \right) \cdot \mathbf{M}. \quad (24)$$

Substituting Eq. (12) and (13) into Eq. (24), we obtain

$$\left((l_j^2 + \sigma^2)^{-1} - \sigma^{-2} \right) \mathbf{S} \mathbf{U}_j \mathbf{Q}_j^\top = \left((l_j^2 + \sigma^2)^{-1} - \sigma^{-2} \right) \mathbf{U}_j \mathbf{Q}_j^\top \mathbf{\Lambda}_j \quad (25)$$

$$\implies \mathbf{S} \mathbf{U}_j \mathbf{Q}_j^\top = \mathbf{U}_j \mathbf{Q}_j^\top \mathbf{\Lambda}_j, \quad j = 1, 2, \dots, p', \quad (26)$$

where

$$p' = \begin{cases} p, & l_p > 0, \\ p - 1, & l_p = 0. \end{cases} \quad (27)$$

Eq. (26) tells us that, the columns of $\mathbf{U}_j \mathbf{Q}_j^\top$ are eigenvectors of \mathbf{S} , and the diagonal elements of $\mathbf{\Lambda}_j$ are the corresponding eigenvalues. Further, since $(\mathbf{U}_j \mathbf{Q}_j^\top)^\top \mathbf{U}_j \mathbf{Q}_j^\top = \mathbf{I}$, these eigenvectors are orthonormal. Let $\mathbf{V}_j = \mathbf{U}_j \mathbf{Q}_j^\top$, we obtain the stationary point:

$$\mathbf{U}_j = \mathbf{V}_j \mathbf{Q}_j, \quad j = 1, 2, \dots, p', \quad (28)$$

where the columns of \mathbf{V}_j are orthonormal eigenvectors of \mathbf{S} with corresponding eigenvalues as $\mathbf{\Lambda}_j = \text{diag} \left(\lambda_{j1}, \lambda_{j2}, \dots, \lambda_{jd_j} \right)$, and \mathbf{Q}_j is an arbitrary orthogonal matrix. Note that if $p' = p - 1$, i.e., $l_p = 0$, \mathbf{U}_p is an arbitrary matrix.

4. Determination of $\mathbf{L} = \text{diag} \left(l_1 \mathbf{I}_{d_1}, l_2 \mathbf{I}_{d_2}, \dots, l_p \mathbf{I}_{d_p} \right)$

Substituting Eq. (26), Eq. (7)-(8), and Eq. (11)-(12) into \mathcal{L}_2 (16) and after some manipulation, a new objective is derived:

$$\begin{aligned} \mathcal{L}_3 = & -\frac{n}{2} \left\{ d \log 2\pi + \sum_{j=1}^{p'} \left[d_j \log (l_j^2 + \sigma^2) + (l_j^2 + \sigma^2)^{-1} \text{tr} (\mathbf{\Lambda}_j) \right] \right. \\ & \left. + (d - q') \log \sigma^2 + \sigma^{-2} \left(\text{tr} (\mathbf{S}) - \sum_{j=1}^{p'} \text{tr} (\mathbf{\Lambda}_j) \right) \right\}, \quad (29) \end{aligned}$$

where

$$q' = \sum_{j=1}^{p'} d_j. \quad (30)$$

Then the partial derivative of \mathcal{L}_3 with respect to l_j is

$$\frac{\partial \mathcal{L}_3}{\partial l_j} = -n \left\{ \frac{d_j l_j}{l_j^2 + \sigma^2} - \frac{l_j \text{tr}(\Lambda_j)}{(l_j^2 + \sigma^2)^2} \right\}, \quad j = 1, 2, \dots, p', \quad (31)$$

and the stationary point is

$$l_j^2 = \frac{\text{tr}(\Lambda_j)}{d_j} - \sigma^2, \quad j = 1, 2, \dots, p'. \quad (32)$$

5. Determination of σ

Substituting Eq. (32) into \mathcal{L}_3 (29), we obtain a new objective:

$$\begin{aligned} \mathcal{L}_4 = & -\frac{n}{2} \left\{ d \log 2\pi + \sum_{j=1}^{p'} \left[d_j \log \frac{\text{tr}(\Lambda_j)}{d_j} + d_j \right] \right. \\ & \left. + (d - q') \log \sigma^2 + \sigma^{-2} \left(\text{tr}(\mathbf{S}) - \sum_{j=1}^{p'} \text{tr}(\Lambda_j) \right) \right\}. \end{aligned} \quad (33)$$

Then the partial derivative of \mathcal{L}_4 with respect to σ is

$$\frac{\partial \mathcal{L}_4}{\partial \sigma} = -n \left\{ \frac{d - q'}{\sigma} - \frac{1}{\sigma^3} \left(\text{tr}(\mathbf{S}) - \sum_{j=1}^{p'} \text{tr}(\Lambda_j) \right) \right\}, \quad (34)$$

and the stationary point is

$$\sigma^2 = \frac{\text{tr}(\mathbf{S}) - \sum_{j=1}^{p'} \text{tr}(\Lambda_j)}{d - q'}. \quad (35)$$

6. Determination of Λ_j

Substituting Eq. (35) into \mathcal{L}_4 (33), we obtain a new objective:

$$\mathcal{L}_5 = -\frac{n}{2} \left\{ d \log 2\pi + \sum_{j=1}^{p'} d_j \log \frac{\text{tr}(\Lambda_j)}{d_j} + (d - q') \log \frac{\text{tr}(\mathbf{S}) - \sum_{j=1}^{p'} \text{tr}(\Lambda_j)}{d - q'} + d \right\}. \quad (36)$$

According to Sec. 3, the diagonal elements of $\mathbf{\Lambda}_j = \text{diag}(\lambda_{j1}, \lambda_{j2}, \dots, \lambda_{jd_j})$, $j = 1, \dots, p'$ are the eigenvalues of \mathbf{S} , therefore the problem here is to select the suitable eigenvalues from \mathbf{S} and separate them into different $\mathbf{\Lambda}_j$ s to maximize \mathcal{L}_5 (36). Using Jensen's inequality:

$$\log \frac{\text{tr}(\mathbf{\Lambda}_j)}{d_j} = \log \frac{\lambda_{j1} + \dots + \lambda_{jd_j}}{d_j} \geq \frac{\log \lambda_{j1} + \dots + \log \lambda_{jd_j}}{d_j}, \quad (37)$$

and the equality holds if and only if $\lambda_{j1} = \dots = \lambda_{jd_j}$. That means, no matter how we select the eigenvalues, only grouping them by the same values can maximize \mathcal{L}_5 (36). Therefore the optimal grouping is

$$\begin{aligned} \mathbf{\Lambda}_j &= \text{diag}(\lambda_{j1}, \lambda_{j2}, \dots, \lambda_{jd_j}) \\ &= \text{diag}(\lambda_j, \lambda_j, \dots, \lambda_j) \\ &= \lambda_j \mathbf{I}_{d_j}, \end{aligned} \quad (38)$$

where λ_j is an eigenvalue of \mathbf{S} whose algebraic multiplicity $\geq d_j$.

Now, the left problem is to select the eigenvalues $\lambda_j, j = 1, \dots, p'$. Substituting Eq. (38) into \mathcal{L}_5 (36), we obtain

$$\begin{aligned} \mathcal{L}_6 &= -\frac{n}{2} \left\{ d \log 2\pi + \sum_{j=1}^{p'} d_j \log \lambda_j + (d - q') \log \frac{\text{tr}(\mathbf{S}) - \sum_{j=1}^{p'} \text{tr}(\mathbf{\Lambda}_j)}{d - q'} + d \right\} \\ &= -\frac{n}{2} \left\{ d \log 2\pi + \text{tr}(\log \mathbf{S}) - \sum_{i=q'+1}^d \log \gamma_i + (d - q') \log \frac{\sum_{i=q'+1}^d \gamma_i}{d - q'} + d \right\} \\ &= -\frac{n}{2} \left\{ d \log 2\pi + \text{tr}(\log \mathbf{S}) - (d - q') \left(\frac{\sum_{i=q'+1}^d \log \gamma_i}{d - q'} - \log \frac{\sum_{i=q'+1}^d \gamma_i}{d - q'} \right) + d \right\}, \end{aligned} \quad (39)$$

where $q' = d_1 + d_2 + \dots + d_{p'}$ and $\gamma_i, i = q' + 1, \dots, d$ are the rest eigenvalues not

been selected. Maximizing \mathcal{L}_6 (39) requires maximizing

$$\mathcal{F} = \frac{\sum_{i=q'+1}^d \log \gamma_i}{d - q'} - \log \frac{\sum_{i=q'+1}^d \gamma_i}{d - q'}, \quad (40)$$

which only requires $\gamma_i, i = q' + 1, \dots, d$ to be adjacent in ordered eigenvalues. However according to Eq. (32), we need $\lambda_j > \sigma^2, j = 1, \dots, p'$, and then from Eq. (35), the only choice to maximize \mathcal{F} is to let $\gamma_i, i = q' + 1, \dots, d$ be the $d - q'$ smallest eigenvalues. Meanwhile, larger q' leads to larger \mathcal{F} , therefore,

$$p' = p \quad (41)$$

$$q' = q = \sum_{j=1}^p d_j \quad (42)$$

7. Determination of \mathbf{U} : Part (2)

According to Eq. (28) and Eq. (38), the columns of \mathbf{V}_j are orthonormal eigenvectors of \mathbf{S} corresponding to a same eigenvalue λ_j . Since \mathbf{Q}_j is an arbitrary orthogonal matrix, the column of $\mathbf{U}_j = \mathbf{V}_j \mathbf{Q}_j$ are still orthonormal eigenvectors corresponding to the eigenvalue λ_j .

8. Summary and Discussion

Summarizing the above analysis (Eq. (15), (32), (35), and Sec. 7), the global maximum of the likelihood with respect to the model parameters is

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad (43)$$

$$\sigma^2 = \frac{\text{tr}(\mathbf{S}) - \text{tr}(\boldsymbol{\Lambda})}{d - q}, \quad (44)$$

$$\mathbf{L}^2 = \boldsymbol{\Lambda} - \sigma^2 \mathbf{I}, \quad (45)$$

$$\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_q], \quad (46)$$

where the elements of the diagonal matrix $\boldsymbol{\Lambda}$ is the q largest eigenvalues of the data covariance \mathbf{S} , and $\mathbf{u}_1, \dots, \mathbf{u}_q$ are the principal q eigenvectors corresponding to $\boldsymbol{\Lambda}$. As

can be seen, under maximum likelihood estimation, the basis vectors \mathbf{U} of our linear model are exactly the same as that learned by PCA [1]. Moreover, diagonal elements of \mathbf{L} represent the “importance” or “energy” of the corresponding basis vectors, and from Eq. (45), when $\sigma \rightarrow 0$, the elements of \mathbf{L}^2 approach the q largest eigenvalues. Besides, as shown in Eq. (44), energy (σ^2) of the noise is the average of the discard eigenvalues, which exactly compensates the energy missed by the subspace model.

Our model can be viewed as constrained case of Probabilistic PCA (PPCA) [2]:

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \sigma\boldsymbol{\epsilon}, \quad (47)$$

whose maximum likelihood estimation is

$$\mathbf{W} = \mathbf{V} (\boldsymbol{\Lambda} - \sigma^2\mathbf{I})^{\frac{1}{2}} \mathbf{Q}, \quad (48)$$

where the columns of \mathbf{V} are the principal eigenvectors of the data covariance, $\boldsymbol{\Lambda}$ is a diagonal matrix whose elements are the corresponding eigenvalues, \mathbf{Q} is an arbitrary orthogonal matrix. Therefore, MLE result of PPCA is nondeterministic due to the arbitrary \mathbf{Q} . Although \mathbf{W} contains information of the principal eigenvectors, the columns of \mathbf{W} itself do not show explicit property of the orthogonality. Our model (1) restricts \mathbf{W} of PPCA (47) to the special form of \mathbf{UL} where \mathbf{U} has orthonormal columns and \mathbf{L} is diagonal matrix. In consequence, MLE result of our model is deterministic (Eq. (43)-(46)). Moreover, our model can build a linear subspace with the principal eigenvectors as the basis vectors explicitly, which is very suitable for extension to the nonlinear case to learn layer-wise interpretable dimensions, as introduced in the main text.

References

- [1] Ian T Jolliffe. *Principal component analysis*. 1986. 1, 8
- [2] Michael E Tipping and Christopher M Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622, 1999. 1, 8