

Functional Data Analysis (Lecture 4) - PCA

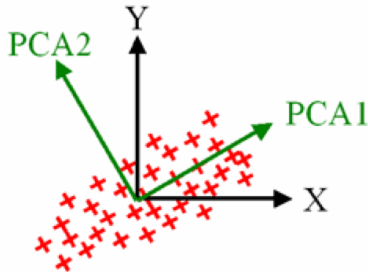
Zoltán Szabó

November 8, 2016

One-page summary

- Covered topics:
 - PEN_L -regularized least squares,
 - smoothing with constraints,
 - curve registration.

- Today:
 - ① dimensionality reduction,
 - ② principal component analysis (PCA; in \mathbb{R}^d first):
 - see continuous registration.



PCA

PCA – intuition

- Given: a set of observations $X = \{\mathbf{x}_i\}_{i=1}^N \subset \mathbb{R}^D$.
- Goal: find the best d -dimensional subspace approximating X .
- $d \ll D$: compression (images, music, ...).



- We are looking for the best **one-dimensional projection**.



- \mathbb{E} := empirical expectation (population: similarly).
- In other words: $\mathbb{E}\mathbf{x} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$.
- Assumption: $\mathbb{E}\mathbf{x} = \mathbf{0}$

- We are looking for the best **one-dimensional projection**.



- \mathbb{E} := empirical expectation (population: similarly).
- In other words: $\mathbb{E}\mathbf{x} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$.
- Assumption: $\mathbb{E}\mathbf{x} = \mathbf{0}$
 - centering: $\mathbf{x} \rightarrow \mathbf{x} - \mathbb{E}\mathbf{x}$.

- One-dimensional projection:
 - \mathbf{w} with $\|\mathbf{w}\|_2 = 1$,
 - $\hat{\mathbf{x}} = \underbrace{\langle \mathbf{w}, \mathbf{x} \rangle}_{\text{=score}} \mathbf{w}$.
- Mean of the projection is zero:

$$\mathbf{0} \stackrel{?}{=} \mathbb{E} \hat{\mathbf{x}} = \mathbb{E} [\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}]$$

- One-dimensional projection:

- \mathbf{w} with $\|\mathbf{w}\|_2 = 1$,

- $\hat{\mathbf{x}} = \underbrace{\langle \mathbf{w}, \mathbf{x} \rangle}_{\text{=score}} \mathbf{w}$.

- Mean of the projection is zero:

$$\mathbf{0} \stackrel{?}{=} \mathbb{E} \hat{\mathbf{x}} = \mathbb{E} [\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}] = \langle \mathbf{w}, \underbrace{\mathbb{E} \mathbf{x}}_{=\mathbf{0}} \rangle \mathbf{w}.$$

PCA: MSE of the residual

Residual:

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2$$

Residual:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &= (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})^T (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})\end{aligned}$$

PCA: MSE of the residual

Residual:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &= (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})^T (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}) \\ &= \|\mathbf{x}\|_2^2 - 2 \langle \mathbf{w}, \mathbf{x} \rangle^2 + \underbrace{\langle \mathbf{w}, \mathbf{x} \rangle^2}_{=1} \|\mathbf{w}\|_2^2\end{aligned}$$

Residual:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &= (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})^T (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}) \\ &= \|\mathbf{x}\|_2^2 - 2 \langle \mathbf{w}, \mathbf{x} \rangle^2 + \underbrace{\langle \mathbf{w}, \mathbf{x} \rangle^2}_{=1} \|\mathbf{w}\|_2^2 \\ &= \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2.\end{aligned}$$

PCA: MSE of the residual

Residual:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &= (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})^T (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}) \\ &= \|\mathbf{x}\|_2^2 - 2 \langle \mathbf{w}, \mathbf{x} \rangle^2 + \underbrace{\langle \mathbf{w}, \mathbf{x} \rangle^2 \|\mathbf{w}\|_2^2}_{=1} \\ &= \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2.\end{aligned}$$

MSE:

$$\min_{\mathbf{w}} \leftarrow \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \mathbb{E} \left[\|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \right]$$

PCA: MSE of the residual

Residual:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &= (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})^T (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}) \\ &= \|\mathbf{x}\|_2^2 - 2 \langle \mathbf{w}, \mathbf{x} \rangle^2 + \underbrace{\langle \mathbf{w}, \mathbf{x} \rangle^2 \|\mathbf{w}\|_2^2}_{=1} \\ &= \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2.\end{aligned}$$

MSE:

$$\begin{aligned}\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \mathbb{E} \left[\|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \right] \\ &= \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{w}} - \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2.\end{aligned}$$

PCA: MSE of the residual

Residual:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ &= (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})^T (\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}) \\ &= \|\mathbf{x}\|_2^2 - 2 \langle \mathbf{w}, \mathbf{x} \rangle^2 + \underbrace{\langle \mathbf{w}, \mathbf{x} \rangle^2 \|\mathbf{w}\|_2^2}_{=1} \\ &= \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2.\end{aligned}$$

MSE:

$$\begin{aligned}\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \mathbb{E} \left[\|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \right] \\ &= \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{w}} - \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2.\end{aligned}$$

⇔ Maximize the mean squared projection.

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual : $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}} \Leftrightarrow$

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual : $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}} \Leftrightarrow$

Maximize mean squared projection : $\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \rightarrow \max_{\mathbf{w}} \Leftrightarrow$

By using $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$:

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{(\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual : $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}} \Leftrightarrow$

Maximize mean squared projection : $\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \rightarrow \max_{\mathbf{w}} \Leftrightarrow$

Maximize variance of the projection : $\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) \rightarrow \max_{\mathbf{w}}.$

PCA: max. variance \rightarrow optimization

Using the bilinearity of covariance:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\langle \mathbf{w}, \mathbf{x} \rangle, \langle \mathbf{w}, \mathbf{x} \rangle)$$

PCA: max. variance \rightarrow optimization

Using the bilinearity of covariance:

$$\begin{aligned} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) &= \text{cov}(\langle \mathbf{w}, \mathbf{x} \rangle, \langle \mathbf{w}, \mathbf{x} \rangle) \\ &= \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) \end{aligned}$$

PCA: max. variance \rightarrow optimization

Using the bilinearity of covariance:

$$\begin{aligned} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) &= \text{cov}(\langle \mathbf{w}, \mathbf{x} \rangle, \langle \mathbf{w}, \mathbf{x} \rangle) \\ &= \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) \\ &= \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} . \end{aligned}$$

PCA: max. variance \rightarrow optimization

Using the bilinearity of covariance:

$$\begin{aligned} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) &= \text{cov}(\langle \mathbf{w}, \mathbf{x} \rangle, \langle \mathbf{w}, \mathbf{x} \rangle) \\ &= \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) \\ &= \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} . \end{aligned}$$

Lagrange function, solving for 'derivatives = 0':

PCA: max. variance \rightarrow optimization

Using the bilinearity of covariance:

$$\begin{aligned} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) &= \text{cov}(\langle \mathbf{w}, \mathbf{x} \rangle, \langle \mathbf{w}, \mathbf{x} \rangle) \\ &= \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) \\ &= \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} . \end{aligned}$$

Lagrange function, solving for 'derivatives = 0':

$$L(\mathbf{w}, \lambda) = \underbrace{\mathbf{w}^T \Sigma \mathbf{w}}_{=\text{objective}} - \lambda \underbrace{(\mathbf{w}^T \mathbf{w} - 1)}_{=\text{condition}} \Rightarrow$$

PCA: max. variance \rightarrow optimization

Using the bilinearity of covariance:

$$\begin{aligned} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) &= \text{cov}(\langle \mathbf{w}, \mathbf{x} \rangle, \langle \mathbf{w}, \mathbf{x} \rangle) \\ &= \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) \\ &= \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} . \end{aligned}$$

Lagrange function, solving for 'derivatives = 0':

$$\begin{aligned} L(\mathbf{w}, \lambda) &= \underbrace{\mathbf{w}^T \Sigma \mathbf{w}}_{=\text{objective}} - \lambda \underbrace{(\mathbf{w}^T \mathbf{w} - 1)}_{=\text{condition}} \Rightarrow \\ 0 &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1), \end{aligned}$$

PCA: max. variance \rightarrow optimization

Using the bilinearity of covariance:

$$\begin{aligned} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) &= \text{cov}(\langle \mathbf{w}, \mathbf{x} \rangle, \langle \mathbf{w}, \mathbf{x} \rangle) \\ &= \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) \\ &= \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} . \end{aligned}$$

Lagrange function, solving for 'derivatives = 0':

$$\begin{aligned} L(\mathbf{w}, \lambda) &= \underbrace{\mathbf{w}^T \Sigma \mathbf{w}}_{=\text{objective}} - \lambda \underbrace{(\mathbf{w}^T \mathbf{w} - 1)}_{=\text{condition}} \Rightarrow \\ 0 &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1), \\ \mathbf{0} &= \frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = 2\Sigma \mathbf{w} - 2\lambda \mathbf{w}. \end{aligned}$$

PCA: max. variance \rightarrow optimization

To sum up: $\Sigma \mathbf{w} = \lambda \mathbf{w}$, $\|\mathbf{w}\|_2 = 1$. $\Rightarrow \mathbf{w}^*$: eigenvector associated to $\lambda_{\max}(\Sigma)$.

PCA: $d \geq 1$

PCA ($d \geq 1$): basis, approximation

- Goal: approximate with a d -dimensional subspace.
- ONB in the subspace: $\{\mathbf{w}_i\}_{i=1}^d$, $\langle \mathbf{w}_i, \mathbf{w}_j \rangle = \delta_{ij}$.
- Approximation ($\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_d]$):

$$\hat{\mathbf{x}} = \sum_{i=1}^d \underbrace{\langle \mathbf{w}_i, \mathbf{x} \rangle}_{\text{=scores}} \mathbf{w}_i = \mathbf{W}\mathbf{W}^T \mathbf{x}.$$

PCA ($d \geq 1$): residual

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \mathbf{W}\mathbf{W}^T\mathbf{x}\|_2^2$$

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \mathbf{W}\mathbf{W}^T\mathbf{x}\|_2^2 = (\mathbf{x} - \mathbf{W}\mathbf{W}^T\mathbf{x})^T (\mathbf{x} - \mathbf{W}\mathbf{W}^T\mathbf{x})$$

Using $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}\|_2^2 = (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x})^T (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}) \\ &= \mathbf{x}^T \underbrace{(\mathbf{I} - \mathbf{W}\mathbf{W}^T)(\mathbf{I} - \mathbf{W}\mathbf{W}^T)}_{=\mathbf{I} - 2\mathbf{W}\mathbf{W}^T + \mathbf{W}\mathbf{W}^T \mathbf{W}\mathbf{W}^T = \mathbf{I} - \mathbf{W}\mathbf{W}^T} \mathbf{x}\end{aligned}$$

Using $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}\|_2^2 = (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x})^T (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}) \\ &= \mathbf{x}^T \underbrace{(\mathbf{I} - \mathbf{W}\mathbf{W}^T)(\mathbf{I} - \mathbf{W}\mathbf{W}^T)}_{=\mathbf{I} - 2\mathbf{W}\mathbf{W}^T + \mathbf{W}\mathbf{W}^T \mathbf{W}\mathbf{W}^T = \mathbf{I} - \mathbf{W}\mathbf{W}^T} \mathbf{x} \\ &= \|\mathbf{x}\|_2^2 - \|\mathbf{W}^T \mathbf{x}\|_2^2,\end{aligned}$$

Using $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}\|_2^2 = (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x})^T (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}) \\ &= \mathbf{x}^T \underbrace{(\mathbf{I} - \mathbf{W}\mathbf{W}^T)(\mathbf{I} - \mathbf{W}\mathbf{W}^T)}_{=\mathbf{I} - 2\mathbf{W}\mathbf{W}^T + \mathbf{W}\mathbf{W}^T \mathbf{W}\mathbf{W}^T = \mathbf{I} - \mathbf{W}\mathbf{W}^T} \mathbf{x} \\ &= \|\mathbf{x}\|_2^2 - \|\mathbf{W}^T \mathbf{x}\|_2^2,\end{aligned}$$

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \mathbb{E} \|\mathbf{x}\|_2^2 - \mathbb{E} \|\mathbf{W}^T \mathbf{x}\|_2^2.$$

Using $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}\|_2^2 = (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x})^T (\mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x}) \\ &= \mathbf{x}^T \underbrace{(\mathbf{I} - \mathbf{W}\mathbf{W}^T)(\mathbf{I} - \mathbf{W}\mathbf{W}^T)}_{=\mathbf{I} - 2\mathbf{W}\mathbf{W}^T + \mathbf{W}\mathbf{W}^T \mathbf{W}\mathbf{W}^T = \mathbf{I} - \mathbf{W}\mathbf{W}^T} \mathbf{x} \\ &= \|\mathbf{x}\|_2^2 - \|\mathbf{W}^T \mathbf{x}\|_2^2,\end{aligned}$$

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \mathbb{E} \|\mathbf{x}\|_2^2 - \mathbb{E} \|\mathbf{W}^T \mathbf{x}\|_2^2.$$

Thus $\min_{\mathbf{W}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow \max_{\mathbf{W}} \mathbb{E} \|\mathbf{W}^T \mathbf{x}\|_2^2.$

PCA ($d \geq 1$): mean squared projection \rightarrow variance

Let $\mathbf{y} = \mathbf{W}^T \mathbf{x}$:

$$\mathbb{E} \|\mathbf{y}\|_2^2 - \|\mathbb{E}\mathbf{y}\|_2^2 = \text{var}(\mathbf{y})?$$

Let $\mathbf{y} = \mathbf{W}^T \mathbf{x}$:

$$\begin{aligned} \mathbb{E} \|\mathbf{y}\|_2^2 - \|\mathbb{E}\mathbf{y}\|_2^2 &= \text{var}(\mathbf{y})? \\ &= \mathbb{E} \left[\sum_i y_i^2 \right] - \sum_i (\mathbb{E} y_i)^2 = \sum_i \text{var}(y_i) \Rightarrow \end{aligned}$$

PCA ($d \geq 1$): mean squared projection \rightarrow variance

Let $\mathbf{y} = \mathbf{W}^T \mathbf{x}$:

$$\begin{aligned} \mathbb{E} \|\mathbf{y}\|_2^2 - \|\mathbb{E}\mathbf{y}\|_2^2 &= \text{var}(\mathbf{y})? \\ &= \mathbb{E} \left[\sum_i y_i^2 \right] - \sum_i (\mathbb{E} y_i)^2 = \sum_i \text{var}(y_i) \Rightarrow \end{aligned}$$

$$\mathbb{E} \|\mathbf{W}^T \mathbf{x}\|_2^2 - \underbrace{\left\| \mathbb{E}[\mathbf{W}^T \mathbf{x}] \right\|_2^2}_{=\mathbf{W}^T \mathbb{E}\mathbf{x}=\mathbf{0}} = \sum_i \text{var} \left((\mathbf{W}^T \mathbf{x})_i \right) \rightarrow \max_{\mathbf{W}}.$$

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \Sigma = \text{cov}(\mathbf{x}).$$

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \Sigma = \text{cov}(\mathbf{x}).$$

- Σ : symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS, $\lambda_i \geq 0$.

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \Sigma = \text{cov}(\mathbf{x}).$$

- Σ : symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS, $\lambda_i \geq 0$.
- Variance decomposition: $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$.

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \Sigma = \text{cov}(\mathbf{x}).$$

- Σ : symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS, $\lambda_i \geq 0$.
- Variance decomposition: $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$.
- **Energy** preserved using d components: $\sum_{i=1}^d \lambda_i \Rightarrow$

$$R^2 = R^2(d) := \frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^D \lambda_i} \in [0, 1].$$

- The d principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \Sigma = \text{cov}(\mathbf{x}).$$

- Σ : symmetric, positive semi-definite $\Rightarrow \{\mathbf{w}_i\}$: ONS, $\lambda_i \geq 0$.
- Variance decomposition: $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$.
- **Energy** preserved using d components: $\sum_{i=1}^d \lambda_i \Rightarrow$

$$R^2 = R^2(d) := \frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^D \lambda_i} \in [0, 1].$$

- In practice: choose d such that $R^2 \approx 0.8 - 0.9$.

Recursive formulation (deflation approach):

- Assume: we have \mathbf{w}_0 .
- Ask for the first PC of the residual (\mathbf{w}_1).
- Iterate.

This leads to an equivalent definition.

- Dimensionality reduction.
- Approximation with a 'small' dimensional subspace.
- Top eigenvectors/-values of the observation covariance.

We covered the ' \mathbb{R}^d part' of Chapter 8 in [1], and [3].