

# Unit 11

## Principal Component Analysis

Prof. Phil Schniter



THE OHIO STATE UNIVERSITY

ECE 4300: Introduction to Machine Learning, Sp20

# Learning objectives

- Recognize need for feature dimensionality reduction
- Understand PCA as RSS-minimizing linear approximation
  - Understand orthogonal projection
  - Recognize PCA as subspace fitting
  - Understand the role of the data-covariance eigenvectors in PCA
  - Know how to measure PCA performance using PoV
- Understand how to compute PCA using the SVD
- Understand how PCA can be used for data visualization
- Understand how the PCA coefficients can be used in supervised learning tasks

# Outline

- Dimensionality Reduction
- Principal Component Analysis (PCA)
- PCA for Data Visualization
- Computing PCA via the SVD
- Python Example: Eigenfaces and PCA-based Classification

# Dimensionality reduction

- Many modern datasets have very high dimension  $d$
- We would like to reduce the dimension (if possible) ...
  - to simplify classification/regression tasks
  - to save memory/storage space
  - to visualize data
- In this unit, we will describe dimensionality reduction via PCA
  - RSS-optimal linear dimensionality reduction

# Data representation

- Dataset:  $\{\mathbf{x}_i\}_{i=1}^n$ 
  - Each sample has  $d$  features:  $\mathbf{x}_i = [x_{i1}, \dots, x_{id}]^T \in \mathbb{R}^d$
  - Can represent dataset using the matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$
  - Will assume data is centered (i.e., mean was removed, so  $\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$ )
- Note: there are *no targets/labels* here!
  - Either they don't exist, or we are ignoring them
  - This is known as “**unsupervised learning**”
  - Previously, we considered “**supervised learning**”: classification, regression
- What if data dimension  $d$  is very large?
  - Can we reduce the dimension?

## Example: Face data

- Face images can be high dimensional
  - We will use  $d = 50 \times 37 = 1850$  pixels
  - Images from the “Labeled Faces in the Wild (LFW)” project
- As we will see, face images can be well approximated using a few coefficients
  - Can be “compressed”!
- How exactly do we do this?



## Loading the data

- The LFW face dataset is built into `sklearn`
  - The full collection contains  $n = 13000$  images (from news stories in 2000s)
  - By requiring  $\geq 70$  faces per person, we extract a subset of 1288 images

```
from sklearn.datasets import fetch_lfw_people  
lfw_people = fetch_lfw_people(min_faces_per_person=70, resize=0.4)
```

Image size =  $50 \times 37 = 1850$  pixels

Number of samples = 1288

Number of classes = 7

- Some example faces:

Donald Rumsfeld



Colin Powell



George W Bush



Gerhard Schroeder



# Outline

- Dimensionality Reduction
- Principal Component Analysis (PCA)
- PCA for Data Visualization
- Computing PCA via the SVD
- Python Example: Eigenfaces and PCA-based Classification



# PCA — Main ideas

- Main idea 1: **Linearly approximate** each feature vector  $\mathbf{x}_i \in \mathbb{R}^d$  as follows:

$$\mathbf{x}_i \approx \mathbf{B}\mathbf{z}_i \text{ with } \mathbf{z}_i \in \mathbb{R}^R, \quad i = 1 \dots n$$

- $\mathbf{B} \in \mathbb{R}^{d \times R}$  is a “dictionary” with  $R$  elements
  - $\mathbf{z}_i$  contains the coefficients
  - $R$  is the “rank” of the approximation, where  $1 \leq R \leq d$  (and ideally  $R \ll d$ )
  - linear approximation is used for simplicity
- Main idea 2: Design the approximation to **minimize RSS**:

$$(\hat{\mathbf{B}}, \{\hat{\mathbf{z}}_i\}_{i=1}^n) = \arg \min_{\mathbf{B}, \{\mathbf{z}_i\}} \left\{ \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{B}\mathbf{z}_i\|^2 \right\}$$

- RSS is used for simplicity
  - RSS is *not* optimal for many tasks (e.g., feature reduction prior to classification)
- Known as “**principal component analysis**” (PCA)

# PCA — Solution

- The optimal  $R$ -element approximation dictionary is

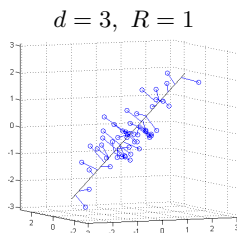
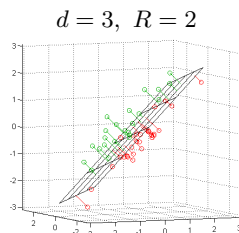
$$\widehat{\mathbf{B}} = \mathbf{V}_R \triangleq [\mathbf{v}_1, \dots, \mathbf{v}_R]$$

where  $\{\mathbf{v}_r\}_{r=1}^R$  are the eigenvectors corresponding to the  $R$  largest eigenvalues  $\{\lambda_r\}_{r=1}^R$  of the sample covariance matrix

$$\mathbf{Q} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \sum_{j=1}^d \lambda_j \mathbf{v}_j \mathbf{v}_j^T$$

- The optimal approximation coefficients are  $\widehat{\mathbf{z}}_i = \mathbf{V}_R^T \mathbf{x}_i$

- The PCA approximation projects  $\mathbf{x}_i \in \mathbb{R}^d$  onto the subspace spanned by  $\{\mathbf{v}_r\}_{r=1}^R$ , which are the  $R$  “principal components” of  $\mathbf{Q}$



## PCA — Derivation ...

- Our PCA derivation uses two steps:
  - 1 Optimize the coefficients  $\{z_i\}$  for an arbitrary fixed dictionary  $B$
  - 2 Optimize the dictionary  $B$
- When optimizing  $z_i$ , the problem reduces to the familiar LS problem:

$$\hat{z}_i = \arg \min_{z_i} \|x_i - Bz_i\|^2 = (B^T B)^{-1} B^T x_i$$

- Plugging  $\{\hat{z}_i\}_{i=1}^n$  back into the original problem yields

$$\begin{aligned} \hat{B} &= \arg \min_B \left\{ \sum_{i=1}^n \|x_i - B(B^T B)^{-1} B^T x_i\|^2 \right\} \\ &= \arg \min_B \left\{ \sum_{i=1}^n \underbrace{\| (I - B(B^T B)^{-1} B^T) x_i \|^2}_{\triangleq P_B^\perp} \right\} \end{aligned}$$

- How do we interpret  $P_B^\perp$ ?

# Orthogonal projection

- Consider the **subspace**  $\text{colsp}(B) \triangleq \{Bz \text{ s.t. } z \in \mathbb{R}^R\} \subset \mathbb{R}^d$ 
  - the set of all linear combinations of the columns of  $B$
- The **orthogonal projection** of  $x \in \mathbb{R}^d$  onto  $\text{colsp}(B)$  ...
  - is the vector  $\hat{x}$  such that

$$x = \hat{x} + e, \text{ where } \hat{x} \in \text{colsp}(B), e \perp \text{colsp}(B)$$

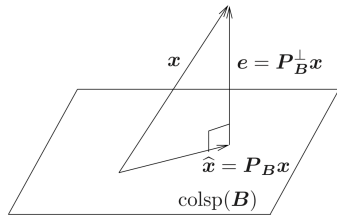
- which can be computed via

$$\hat{x} = P_B x \text{ and } e = P_B^\perp x$$

- using the **orthogonal projection** operators

$$P_B \triangleq B(B^T B)^{-1} B^T \text{ and } P_B^\perp \triangleq I - P_B$$

- These operators are **symmetric** and **idempotent**, i.e.,  $P_B = P_B^T$  &  $P_B = P_B^2$
- Also,  $P_B$  has  $R$  eigenvalues = 1, and all other eigenvalues = 0, while  $P_B^\perp$  has  $R$  eigenvalues = 0, and all other eigenvalues = 1



# The role of projection in PCA

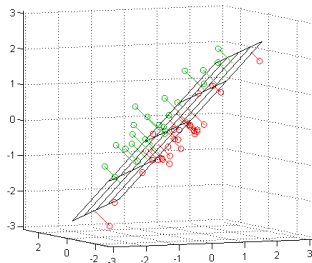
- Back to the PCA problem:

$$\hat{B} = \arg \min_B \left\{ \sum_{i=1}^n \|P_B^\perp x_i\|^2 \right\}$$

- We now recognize  $P_B^\perp x_i$  as projection error
- Thus, PCA chooses  $B$  to minimize the sum-squared projection error
- To go further, let's reformulate the RSS cost:

$$\begin{aligned} J(B) &= \sum_i \|P_B^\perp x_i\|^2 = \sum_i x_i^\top (P_B^\perp)^\top P_B^\perp x_i = \sum_i x_i^\top P_B^\perp x_i = \sum_i \text{tr}(x_i^\top P_B^\perp x_i) \\ &= \sum_i \text{tr}(P_B^\perp x_i x_i^\top) = \text{tr}\left(P_B^\perp \sum_i x_i x_i^\top\right) = n \text{tr}\left(P_B^\perp \underbrace{\frac{1}{n} \sum_{i=1}^n x_i x_i^\top}_{\text{sample covariance mtx } Q}\right) \end{aligned}$$

where  $\text{tr}(A) = \sum_j [A]_{jj}$  and  $\text{tr}(AC) = \text{tr}(CA)$ .



sample covariance mtx  $Q$

# Eigen-decomposition of sample covariance matrices

- The sample covariance mtx  $Q$  is **positive semi-definite**, i.e.,  $x^T Q x \geq 0 \forall x$ 
  - Proof:  $x^T Q x = x^T (\frac{1}{n} \sum_i x_i x_i^T) x = \frac{1}{n} \sum_i (x^T x_i)(x_i^T x) = \frac{1}{n} \sum_i (x^T x_i)^2 \geq 0$
- All positive semi-definite matrices have an **eigen-decomposition** of the form

$$Q = V \Lambda V^T \quad \text{where} \quad \begin{cases} V \text{ is orthogonal (i.e., } V V^T = I_d = V^T V) \\ \Lambda = \text{Diag}(\lambda_1, \dots, \lambda_d) \text{ with } \lambda_j \geq 0 \forall j \\ \text{Will assume w.l.o.g. that } \{\lambda_j\} \text{ sorted from high to low} \end{cases}$$

## Theorem (Eckart-Young, 1936)

The optimal  $B$  can be constructed from the  $R$  **principal eigenvectors** of  $Q$ :

$$\hat{B} = \arg \min_B n \operatorname{tr}(P_B^\perp Q) = [v_1, \dots, v_R] \triangleq V_R,$$

More precisely, the optimal  $\hat{B}$  is *any*  $B$  for which  $\operatorname{colsp}(B) = \operatorname{colsp}(V_R)$

# Simple proof of Eckart-Young

- Recall that we want to minimize the RSS cost

$$J(\mathbf{B}) = n \operatorname{tr}(\mathbf{P}_B^\perp \mathbf{Q}) = n \operatorname{tr}((\mathbf{I}_d - \mathbf{P}_B) \mathbf{Q}) = n \operatorname{tr}(\mathbf{Q}) - n \operatorname{tr}(\mathbf{P}_B \mathbf{Q})$$

- Equivalently, we can *maximize* the utility

$$U(\mathbf{B}) \triangleq \operatorname{tr}(\mathbf{P}_B \mathbf{Q}) = \operatorname{tr}(\mathbf{P}_B \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top) = \operatorname{tr}(\mathbf{V}^\top \mathbf{P}_B \mathbf{V} \mathbf{\Lambda}) = \sum_{j=1}^d \alpha_j \lambda_j$$

$$\text{for } \alpha_j \triangleq [\mathbf{V}^\top \mathbf{P}_B \mathbf{V}]_{jj} = \mathbf{v}_j^\top \mathbf{P}_B \mathbf{v}_j \in [0, 1]$$

- Notice also that

$$\begin{aligned} \sum_{j=1}^d \alpha_j &= \operatorname{tr}(\mathbf{V}^\top \mathbf{P}_B \mathbf{V}) = \operatorname{tr}(\mathbf{V} \mathbf{V}^\top \mathbf{P}_B) = \operatorname{tr}(\mathbf{P}_B) = \operatorname{tr}(\mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top) \\ &= \operatorname{tr}(\mathbf{B}^\top \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1}) = \operatorname{tr}(\mathbf{I}_R) = R. \end{aligned}$$

- Thus we can consider the simplified optimization problem:

$$\text{Find } \{\alpha_j\} \text{ with } \alpha_j \in [0, 1] \text{ and } \sum_{j=1}^d \alpha_j = R \text{ that maximizes } \sum_{j=1}^d \alpha_j \lambda_j$$

# Simple proof of Eckart-Young (cont.)

- For the optimization problem ...

Find  $\{\alpha_j\}$  with  $\alpha_j \in [0, 1]$  and  $\sum_{j=1}^d \alpha_j = R$  that maximizes  $\sum_{j=1}^d \alpha_j \lambda_j$

- Think of  $\lambda_j$  as the reward for the  $j$ th item, and  $\alpha_j$  a purchasing variable
  - You must buy  $R$  units total, and between 0 and 1 units of each item
  - Question: What purchase is most rewarding?
  - Answer: One unit each of the  $R$  best items! i.e.,  $\alpha_j = \begin{cases} 1 & \text{if } j = 1 \dots R \\ 0 & \text{if } j = R+1 \dots d \end{cases}$
  - Recall that  $\{\lambda_j\}$  are ordered from high to low
- If  $\mathbf{v}_j$  denotes the  $j$ th eigenvector of  $\mathbf{Q}$ , these optimal  $\{\alpha_j\}$  are attained when

$$\mathbf{B} = [\mathbf{v}_1, \dots, \mathbf{v}_R] \triangleq \mathbf{V}_R, \quad \text{since } \alpha_j = \mathbf{v}_j^\top \mathbf{P}_B \mathbf{v}_j = \begin{cases} 1 & \text{if } j = 1 \dots R \\ 0 & \text{if } j = R+1 \dots d \end{cases}$$



# Summary of PCA

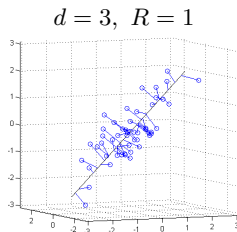
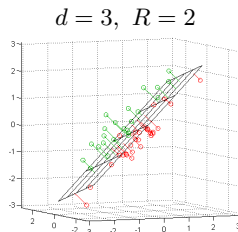
- Summary: Given centered data  $\{\mathbf{x}_i\}_{i=1}^n$ , PCA approximates  $\mathbf{x}_i$  as

$$\hat{\mathbf{x}}_i = \mathbf{V}_R \mathbf{z}_i \quad \text{with} \quad \mathbf{z}_i = \mathbf{V}_R^\top \mathbf{x}_i$$

where  $\mathbf{V}_R$  contains the  $R$  principal eigenvectors of the sample covariance mtx

$$\mathbf{Q} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{n} \mathbf{X}^\top \mathbf{X} \quad \text{with} \quad \mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$$

- These eigenvectors are called the “**principal components**”
- The PCA approximation projects  $\mathbf{x}_i \in \mathbb{R}^d$  onto the subspace spanned by the  $R$  principal components of  $\mathbf{Q}$



# Performance of PCA

- How do we quantify the performance of PCA for a given rank  $R$ ?
  - This could help in choosing  $R$

- The total variance of our (centered) data  $\{\mathbf{x}_i\}$  is

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 = \frac{1}{n} \sum_{i=1}^n \text{tr}(\mathbf{x}_i \mathbf{x}_i^T) = \text{tr}(\mathbf{Q}) = \text{tr}(\mathbf{V} \mathbf{\Lambda} \mathbf{V}^T) = \sum_{j=1}^d \lambda_j$$

- The total variance of our PCA approximation  $\{\hat{\mathbf{x}}_i\}$  is

$$\frac{1}{n} \sum_{i=1}^n \|\hat{\mathbf{x}}_i\|^2 = \frac{1}{n} \sum_{i=1}^n \text{tr}(\mathbf{V}_R \mathbf{V}_R^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{V}_R \mathbf{V}_R^T) = \text{tr}(\mathbf{V}_R^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{V}_R) = \sum_{j=1}^R \lambda_j$$

- Thus the **proportion of variance** explained by the  $R$  principal components is

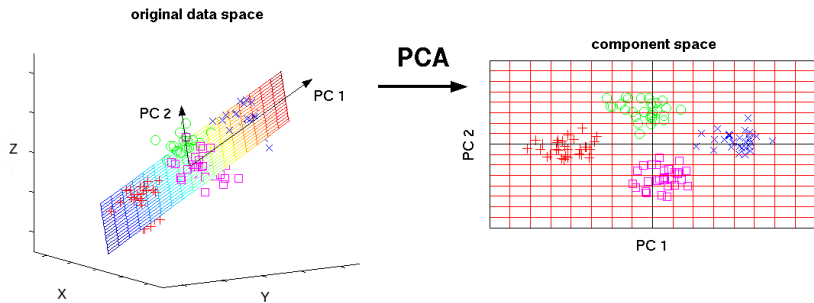
$$\text{PoV}(R) \triangleq \frac{\sum_{j=1}^R \lambda_j}{\sum_{j=1}^d \lambda_j} \quad \dots \text{want this to be close to 1}$$

# Outline

- Dimensionality Reduction
- Principal Component Analysis (PCA)
- PCA for Data Visualization
- Computing PCA via the SVD
- Python Example: Eigenfaces and PCA-based Classification

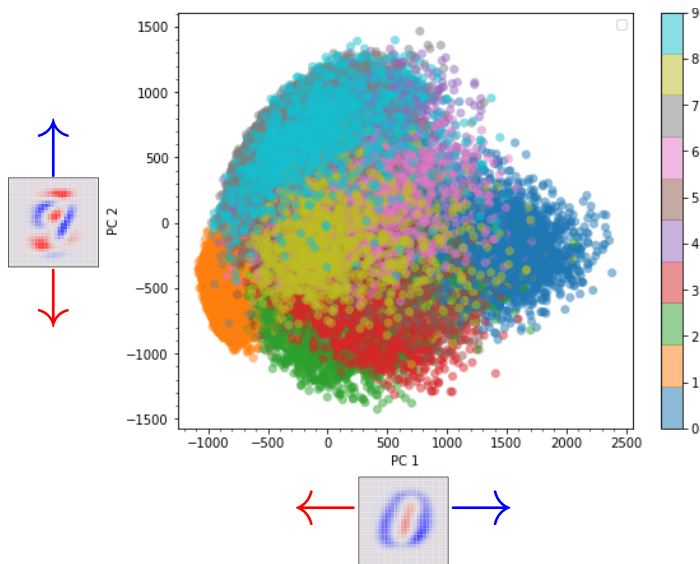
# PCA for data visualization

- When  $d \geq 3$ , the original dataset  $\{\mathbf{x}_i\}_{i=1}^n$  can be difficult to visualize
- But when  $R \leq 2$ , the PCA coordinates  $\{\mathbf{z}_i\}_{i=1}^n$  can be easily visualized:



- May also be possible to visualize the principal components  $\{\mathbf{v}_j\}_{j=1}^R$  themselves

# Example: PCA visualization of MNIST



# Outline

- Dimensionality Reduction
- Principal Component Analysis (PCA)
- PCA for Data Visualization
- **Computing PCA via the SVD**
- Python Example: Eigenfaces and PCA-based Classification

# The singular value decomposition (SVD)

- Given *any* matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  ...
- The standard **SVD** decomposes  $\mathbf{X}$  using square  $\mathbf{U}$  &  $\mathbf{V}$  as follows:

$$\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T \quad \text{where} \quad \begin{cases} \mathbf{U} \in \mathbb{R}^{n \times n} & \text{obeys } \mathbf{U} \mathbf{U}^T = \mathbf{I}_n = \mathbf{U}^T \mathbf{U} \\ \mathbf{S} \in \mathbb{R}^{n \times d} & \text{obeys } \mathbf{S} = \text{Diag}(s_1, \dots, s_m) \\ & \text{where } m = \min\{n, d\} \\ & \text{and } s_1 \geq s_2 \geq s_3 \geq \dots \geq 0 \\ \mathbf{V} \in \mathbb{R}^{d \times d} & \text{obeys } \mathbf{V} \mathbf{V}^T = \mathbf{I}_d = \mathbf{V}^T \mathbf{V} \end{cases}$$

- The “**economy SVD**” uses square  $\mathbf{S}_r$  and possibly tall  $\mathbf{U}_r$  &  $\mathbf{V}_r$ :

$$\mathbf{X} = \mathbf{U}_r \mathbf{S}_r \mathbf{V}_r^T \quad \text{where} \quad \begin{cases} \mathbf{U}_r \in \mathbb{R}^{n \times r} & \text{obeys } \mathbf{U}_r^T \mathbf{U}_r = \mathbf{I}_r \\ \mathbf{S}_r \in \mathbb{R}^{r \times r} & \text{obeys } \mathbf{S}_r = \text{Diag}(s_1, \dots, s_r) \\ & \text{where } r \triangleq \text{rank}(\mathbf{X}) \leq \min\{n, d\} \\ & \text{and } s_1 \geq s_2 \geq \dots \geq s_r > 0 \\ \mathbf{V}_r \in \mathbb{R}^{d \times r} & \text{obeys } \mathbf{V}_r^T \mathbf{V}_r = \mathbf{I}_r \end{cases}$$

# Computing PCA via the standard SVD

- As before, assume that  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$  is centered
  - That is, the sample mean of every column equals zero
- Then the sample covariance matrix has the (sorted) eigen-decomposition

$$\mathbf{Q} = \frac{1}{n} \mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top \quad \text{where} \quad \begin{cases} \mathbf{V} \mathbf{V}^\top = \mathbf{I}_d = \mathbf{V}^\top \mathbf{V} \\ \mathbf{\Lambda} = \text{Diag}(\lambda_1, \dots, \lambda_d) \\ \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq 0 \end{cases}$$

- Plugging in the standard SVD  $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$ , we find

$$\mathbf{Q} = \frac{1}{n} \mathbf{V} \mathbf{S}^\top \underbrace{\mathbf{U}^\top \mathbf{U}}_{\mathbf{I}_n} \mathbf{S} \mathbf{V}^\top = \mathbf{V} \left( \frac{1}{n} \mathbf{S}^\top \mathbf{S} \right) \mathbf{V}^\top \quad \text{where} \quad \begin{cases} \mathbf{V} \mathbf{V}^\top = \mathbf{I}_d = \mathbf{V}^\top \mathbf{V} \\ \frac{1}{n} \mathbf{S}^\top \mathbf{S} = \text{Diag}\left(\frac{s_1^2}{n}, \dots, \frac{s_d^2}{n}\right) \\ \frac{s_1^2}{n} \geq \frac{s_2^2}{n} \geq \frac{s_3^2}{n} \geq \dots \geq 0 \end{cases}$$

- So, we can use the standard SVD to compute the eigen-decomposition of  $\mathbf{Q}$ 
  - The  $\mathbf{V}$  matrices are the same (if evals distinct and sorted) and  $\lambda_j = s_j^2/n$



# Computing PCA via the economy SVD

- Remember: the economy SVD computes only the top  $r$  singular vectors, i.e.,  $U_r$  and  $V_r$ , where  $r = \text{rank}(\mathbf{X}) \leq \min(n, d)$ 
  - Thus, it runs faster than the standard SVD
- For PCA, need to compute only the top  $R$  singular vectors  $V_R$ , where  $R \leq r$ 
  - $R$  is a design choice, and typically  $R \ll r$
- Thus it's more efficient to use the economy SVD when computing PCA:

$$(U_r, S_r, V_r^T) = \text{economy-svd}(\mathbf{X})$$

$$V_R = \text{first } R \text{ columns of } V_r$$

$$z_i = V_R^T x_i \quad \forall i = 1 \dots n$$

# Standardizing the PCA coefficients

- After computing the PCA coefficients  $\{z_i\}$ , we might use them for classification or regression (assuming we also have some targets  $\{y_i\}$ )
- In that case, it would be good to **standardize**  $\{z_i\}$ 
  - Could be accomplished by computing  $\tilde{z}_i \triangleq Q_z^{-\frac{1}{2}}(z_i - \bar{z}) \forall i$ 
    - where  $\bar{z} \triangleq \frac{1}{n} \sum_{i=1}^n z_i$  is the sample mean
    - and  $Q_z \triangleq \frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(z_i - \bar{z})^\top$  is the sample covariance
  - Why? Because  $\frac{1}{n} \sum_{i=1}^n \tilde{z}_i = \mathbf{0}$  and  $\frac{1}{n} \sum_{i=1}^n \tilde{z}_i \tilde{z}_i^\top = Q_z^{-\frac{1}{2}} Q_z Q_z^{-\frac{1}{2}} = I_R$
- But there is a shortcut!
  - $\bar{z} = \frac{1}{n} \sum_{i=1}^n V_R^\top x_i = V_R^\top (\frac{1}{n} \sum_{i=1}^n x_i) = V_R^\top \mathbf{0} = \mathbf{0}$ , since  $\{x_i\}$  were centered
  - $Q_z = \frac{1}{n} \sum_{i=1}^n z_i z_i^\top = V_R^\top (\frac{1}{n} \sum_{i=1}^n x_i x_i^\top) V_R = V_R^\top Q V_R = \frac{1}{n} V_R^\top V_r S_r^2 V_r^\top V_R = \text{Diag}(\frac{s_1^2}{n}, \dots, \frac{s_R^2}{n})$  since  $V_R^\top V_r = [I_R \ \mathbf{0}_{R \times (r-R)}]$
  - Thus  $\tilde{z}_i \triangleq \text{Diag}(\frac{\sqrt{n}}{s_1}, \dots, \frac{\sqrt{n}}{s_R}) z_i$

# Outline

- Dimensionality Reduction
- Principal Component Analysis (PCA)
- PCA for Data Visualization
- Computing PCA via the SVD
- Python Example: Eigenfaces and PCA-based Classification

# PCA on the LFW face dataset

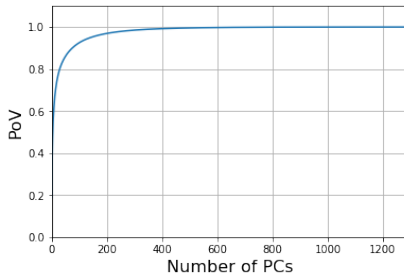
- First we center the data  $\{x_i\}$
- Then we compute the economy SVD
- Note that  $V_r^T$  is wide, as we expect
- Then we compute the eigenvalues  $\{\lambda_j\}_{j=1}^r$
- And finally we compute the Proportion-of-Variance
- The PoV plot suggests that  $R = 400$  principle components capture nearly all the variance of our data

```
Xmean = np.mean(X,0)
Xs = X - Xmean[None,:]
```

```
U,S,VT = np.linalg.svd(Xs, full_matrices=False)
VT.shape

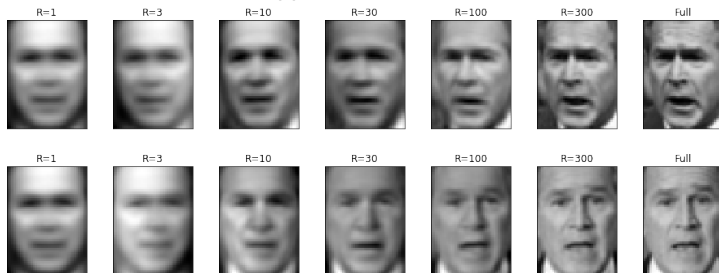
(1288, 1850)
```

```
lam = S**2 / n_samples
PoV = np.cumsum(lam)/np.sum(lam)
```

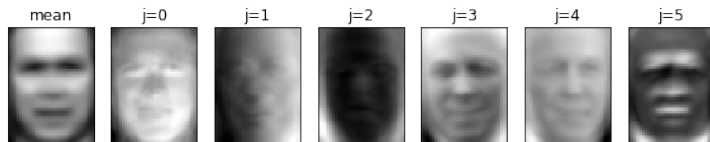


# PCA approximation and eigenfaces

- We now show the PCA approximations versus  $R$  for two faces:



- And the mean & top 5 principle components  $\{v_j\}$  (i.e., “eigenfaces”):



# Face recognition using the PCA coefficients

We now demonstrate classification (i.e., face recognition) via PCA coefficients

- Split data into training and test:

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, stratify = y, random_state=43)
```

- Center the training:

- Perform economy SVD:

```
n_samples, _ = X_train.shape
Xtr_mean = np.mean(X_train,0)
Xtr = X_train - Xtr_mean[None,:]
Utr,Str,VTtr = np.linalg.svd(Xtr, full_matrices=False)
```

- Choose  $R = 100$ :

- Compute PCA coefficients

$$\mathbf{z}_i^T = \mathbf{x}_i^T \mathbf{V}_R \quad \forall i:$$

```
npc = 100
eigenfaces = VTtr[:npc,:]
Ztr = Xtr.dot(eigenfaces.T)
```

- Standardize the PCA

$$\text{coefficients: } \tilde{\mathbf{z}}_i = \text{Diag} \left( \frac{\sqrt{n}}{s_1}, \dots, \frac{\sqrt{n}}{s_R} \right) \mathbf{z}_i:$$

```
Ztr_s = Ztr / Str[None,:npc] * np.sqrt(n_samples)
```

## Face recognition using the PCA coefficients (cont.)

- Tune an SVM classifier over regularization  $C$  & RBF kernel width  $\gamma$ :

```
param_grid = {'C': [1, 3, 10, 30, 100, 300],
              'gamma': [0.00001, 0.00003, 0.0001, 0.0003, 0.001, 0.003, 0.01], }
clf = GridSearchCV(SVC(kernel='rbf', class_weight='balanced'), param_grid, cv=5, iid=False)
clf = clf.fit(Ztr_s, y_train)
print("Best estimator found by grid search:")
print(clf.best_estimator_)
print("Cross validation accuracy with best estimator:")
print(clf.best_score_)
```

Best estimator found by grid search:

```
SVC(C=3, cache_size=200, class_weight='balanced', coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma=0.003, kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

Cross validation accuracy with best estimator:

0.8506081737123565

- After pre-processing the test data in the same way as training, classify it:

```
Xts = X_test - Xtr_mean[None,:]
Zts = Xts.dot(eigenfaces.T)
Zts_s = Zts / Str[None,:npc] * np.sqrt(n_samples)
y_hat = clf.predict(Zts_s)
acc = np.mean(y_hat==y_test)
print("The model accuracy on the test set is %f" % acc)
```

The model accuracy on the test set is 0.829193

# Learning objectives

- Recognize need for feature dimensionality reduction
- Understand PCA as RSS-minimizing linear approximation
  - Understand orthogonal projection
  - Recognize PCA as subspace fitting
  - Understand the role of the data-covariance eigenvectors in PCA
  - Know how to measure PCA performance using PoV
- Understand how to compute PCA using the SVD
- Understand how PCA can be used for data visualization
- Understand how the PCA coefficients can be used in supervised learning tasks