

Lecture 11

Principal Component Analysis

EE-UY 4563 / EL-GY 9123: INTRODUCTION TO MACHINE LEARNING
PROF. SUNDEEP RANGAN

Outline



Dimensionality reduction

- Principal components and directions of variance
- Approximation with PCs
- Computing PCs via the SVD
- Face example in python

Dimensionality Reduction

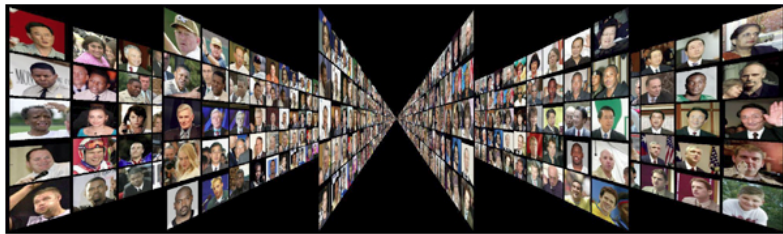
- ❑ Many modern data sets have very high dimension
- ❑ Want to reduce dimension:
 - Simplify classification / regression tasks on the data set
 - Visualize data
 - Find underlying commonalities in data

Data Definitions

- ❑ Given data: $\mathbf{x}_i, i = 1, \dots, N$
- ❑ Each sample has p features: $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$
- ❑ Represent as an $N \times p$ matrix
- ❑ Unsupervised learning
 - Samples do not have a label
 - Or, we choose to ignore the label for now
- ❑ Dimension p is large
- ❑ How do we reduce the dimension?

Example: Faces

Labeled Faces in the Wild Home



- ❑ Face images can be high-dimensional
 - We will use $50 \times 37 = 1850$ pixels
- ❑ But, there may be few degrees of freedom
- ❑ Can we reduce the dimensionality of this?
- ❑ Data Labelled Faces in the Wild project
 - <http://vis-www.cs.umass.edu/lfw>
 - Large collection of faces (13000 images)
 - Taken from web articles about 10 years ago

Loading the Data

- ❑ Lect10_PCA_Faces.ipynb
- ❑ Built-in routines to load data is scikit-learn
- ❑ Can take several minutes the first time (Be patient)

Image size = 50 x 37 = 1850 pixels
Number faces = 1288
Number classes = 7

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

```
2016-11-14 14:15:30,862 Downloading LFW metadata: http://vis-www.cs.umass.edu/lfw/pairsDevTrain.txt  
2016-11-14 14:15:30,958 Downloading LFW metadata: http://vis-www.cs.umass.edu/lfw/pairsDevTest.txt  
2016-11-14 14:15:31,028 Downloading LFW metadata: http://vis-www.cs.umass.edu/lfw/pairs.txt  
2016-11-14 14:15:31,294 Downloading LFW data (~200MB): http://vis-www.cs.umass.edu/lfw/lfw-funneled.tgz  
2016-11-14 14:20:10,056 Decompressing the data archive to C:\Users\Sundeeep\scikit_learn_data\lfw_home\lfw_funneled  
2016-11-14 14:22:08,605 Loading LFW people faces from C:\Users\Sundeeep\scikit_learn_data\lfw_home  
2016-11-14 14:22:09,735 Loading face #00001 / 01288  
2016-11-14 14:22:13,640 Loading face #01001 / 01288
```

Plotting the Data

- ❑ Some example faces
- ❑ You may be too young to remember them all

Colin Powell



Colin Powell



Hugo Chavez




George W Bush



```
def plt_face(x):  
    h = 50  
    w = 37  
    plt.imshow(x.reshape((h, w)), cmap=plt.cm.gray)  
    plt.xticks([])  
    plt.yticks([])  
  
I = np.random.permutation(n_samples)  
plt.figure(figsize=(10,20))  
nplt = 4;  
for i in range(nplt):  
    ind = I[i]  
    plt.subplot(1,nplt,i+1)  
    plt_face(X[ind])  
    plt.title(target_names[y[ind]])
```

Outline

- ☐ Dimensionality reduction
-  ☐ Principal components and directions of variance
- ☐ Approximation with PCs
- ☐ Computing PCs via the SVD
- ☐ Face example in python

Projections

□ Given a vectors \mathbf{z} and \mathbf{v}

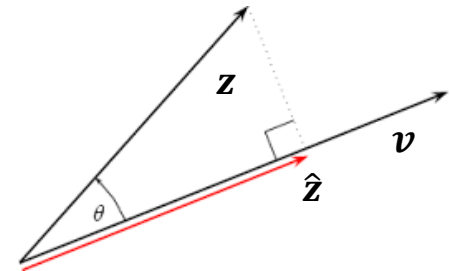
□ **Projection** of \mathbf{z} onto \mathbf{v} is:

$$\hat{\mathbf{z}} = \text{Proj}_{\mathbf{v}}(\mathbf{z}) = \alpha \mathbf{v}, \quad \alpha = \frac{\mathbf{v}^T \mathbf{z}}{\mathbf{v}^T \mathbf{v}} = \frac{\|\mathbf{z}\|}{\|\mathbf{v}\|} \cos \theta$$

□ Let $V = \{\alpha \mathbf{v} | \alpha \in \mathbb{R}\}$ = vectors on the line spanned by \mathbf{v}

□ **Theorem**: $\text{Proj}_{\mathbf{v}}(\mathbf{z})$ is closest point in V to \mathbf{z} :

$$\hat{\mathbf{z}} = \arg \min_{\mathbf{w} \in V} \|\mathbf{z} - \mathbf{w}\|^2$$



Sample Covariance Matrix

□ Let $\tilde{\mathbf{X}}$ = data matrix with sample mean removed.

- Rows: $\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$

□ Sample covariance matrix: Matrix \mathbf{Q} with components:

$$Q_{k\ell} = \frac{1}{N} \sum_{i=1}^N (x_{ik} - \bar{x}_k)(x_{i\ell} - \bar{x}_\ell)$$

- Covariance between feature k and ℓ in the dataset
- Matrix is $p \times p$

□ **Theorem:** Sample covariance is given by

$$\mathbf{Q} = \frac{1}{N} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$$

- Proof on board
- Compute sample covariance via a matrix product

Directional Variance

- Given data: $\mathbf{x}_i, i = 1, \dots, N$ and direction \mathbf{v} with $\|\mathbf{v}\| = 1$
- How much does \mathbf{x}_i vary in the direction \mathbf{v} ?
- Let $z_i = \mathbf{v}^T \mathbf{x}_i$ = projection of \mathbf{x}_i onto \mathbf{v}
- Sample mean and variance in direction \mathbf{v} is (proof on board):
 - Sample mean $\bar{z} = \mathbf{v}^T \bar{\mathbf{x}}$
 - Sample variance $s_z^2 = \mathbf{v}^T \mathbf{S} \mathbf{v}$

Maximizing Directional Variance

□ What directions \mathbf{v} maximize the variance?

□ Formulate as an optimization problem:

$$\max_{\mathbf{v}} \mathbf{v}^T \mathbf{Q} \mathbf{v} \quad \text{s.t.} \quad \|\mathbf{v}\| = 1$$

□ Let $\mathbf{v}_1, \dots, \mathbf{v}_p$ be the eigenvectors of \mathbf{Q} : $\mathbf{Q} \mathbf{v}_j = \lambda_j \mathbf{v}_j$

□ Sort eigenvalues in descending order: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$

- Can show that eigenvalues are real and non-negative

□ **Theorem:** Any local maxima of the variance directional is an eigenvector

- $\mathbf{v} = \mathbf{v}_j$ for some j and $\mathbf{v}^T \mathbf{Q} \mathbf{v} = \lambda_j$
- Proof on board

Principal Components

□ **Principal components:** The eigenvectors of Q , v_1, \dots, v_p

- Always normalized $\|v_j\| = 1$

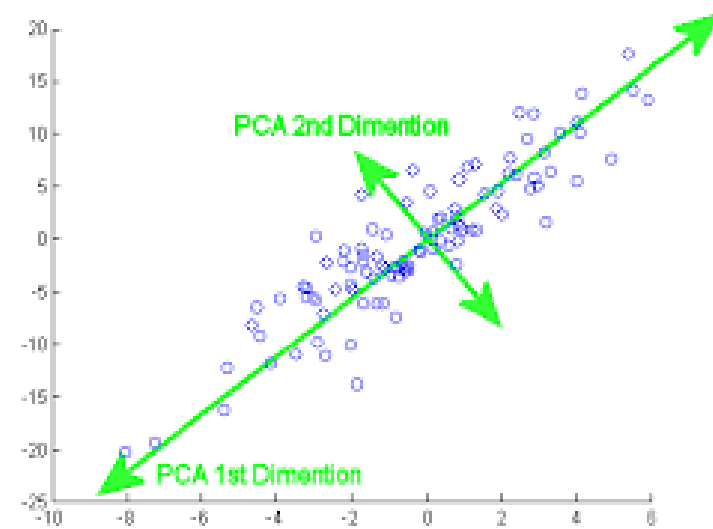
□ Sorted by eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$

□ Each vector is of dimension p


□ Key property: Vectors are orthogonal

- $v_j^T v_k = 0$ if $j \neq k$
- Proof on board

□ Represents directions of maximal variance



Outline

- ☐ Dimensionality reduction
- ☐ Principal components and directions of variance
-  ☐ Approximation with PCs
- ☐ Computing PCs via the SVD
- ☐ Face example in python

Low-Dimensional Representations

□ Given data $x_i, i = 1, \dots, N$

□ **Problem:** Find **basis vectors** $v_j, j = 1, \dots, d$ such that:

$$x_i \approx \bar{x} + \sum_{j=1}^d \alpha_{ij} v_j$$

- Sample mean + linear combination of basis vectors
- $\alpha_i = (\alpha_{i1}, \dots, \alpha_{id})$ is an approximate **coordinates** of x_i in basis (v_1, \dots, v_d)

□ Dimensionality reduction:

- If $d \ll p$ we have represented v_i with a smaller number of coefficients.

Orthonormal Sets and Bases

□ **Definition:** A set of vectors $\mathbf{v}_1, \dots, \mathbf{v}_d$ are an **orthonormal set** if:

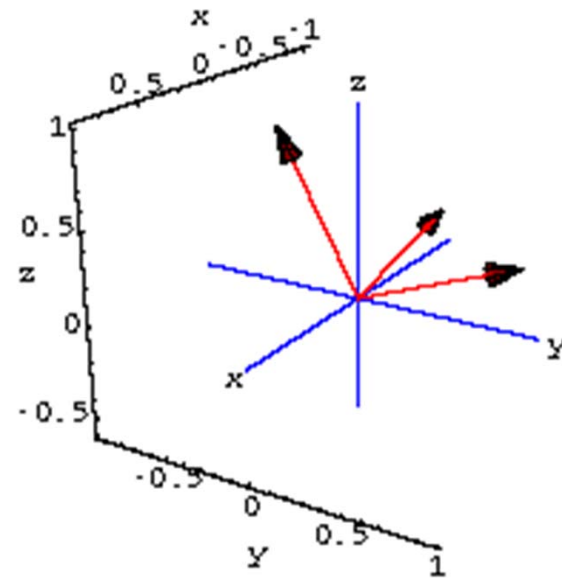
- $\|\mathbf{v}_j\| = 1$ for all j (unit length)
- $\mathbf{v}_j^T \mathbf{v}_k = 0$ if $j \neq k$ (perpendicular to one another)

□ If $d = p$ then $\mathbf{v}_1, \dots, \mathbf{v}_p$ is called an **orthonormal basis**

□ **Matrix form:** If $\mathbf{V} = [\mathbf{v}_1 \dots \mathbf{v}_d]$, then $\mathbf{V}^T \mathbf{V} = \mathbf{I}_d$

□ If $d = p$, then \mathbf{V} is an **orthogonal matrix**

□ **Key property:** the PCs form an orthonormal basis



Coefficients in an Orthonormal Basis

□ Suppose $\mathbf{v}_1, \dots, \mathbf{v}_p$ is an orthonormal basis

□ Given a vector \mathbf{z} , can write

$$\mathbf{z} = \sum_{j=1}^p \alpha_j \mathbf{v}_j, \quad \alpha_j = \mathbf{v}_j^T \mathbf{z}$$

- Simple expression for computing coefficients in an orthonormal basis

□ Matrix form:

$$\boldsymbol{\alpha} = \mathbf{V}^T \mathbf{z}, \quad \mathbf{z} = \mathbf{V} \boldsymbol{\alpha}$$

Approximating the Data Matrix

- Given data $x_i, i = 1, \dots, N$
- Let v_1, \dots, v_p be the PCs
- Find coefficient expansion of each data sample:

$$x_i = \bar{x} + \sum_{j=1}^p \alpha_{ij} v_j, \quad \alpha_{ij} = v_j^T (x_i - \bar{x})$$

- Now consider approximation with d coefficients:

$$\hat{x}_i = \bar{x} + \sum_{j=1}^d \alpha_{ij} v_j$$

Average Approximation Error

□ Let $\hat{\mathbf{x}}_i$ = approximation with d PCs

□ Error in sample i :

$$\mathbf{x}_i - \hat{\mathbf{x}}_i = \sum_{j=d+1}^p \alpha_{ij} \mathbf{v}_j$$

□ **Theorem:** Average error with a d PC approximation is:

$$\frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \sum_{j=d+1}^p \lambda_j$$

- Sum of the smallest $p - d$ eigenvalues

Proportion of Variance (PoV)

□ Total variance of data set: $\frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 = \sum_{j=1}^p \lambda_j$

□ Average approximation error: $\frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \sum_{j=d+1}^p \lambda_j$

□ The **proportion of variance** explained by d PCs is:

$$PoV(d) = \frac{\sum_{j=1}^d \lambda_j}{\sum_{j=1}^p \lambda_j}$$

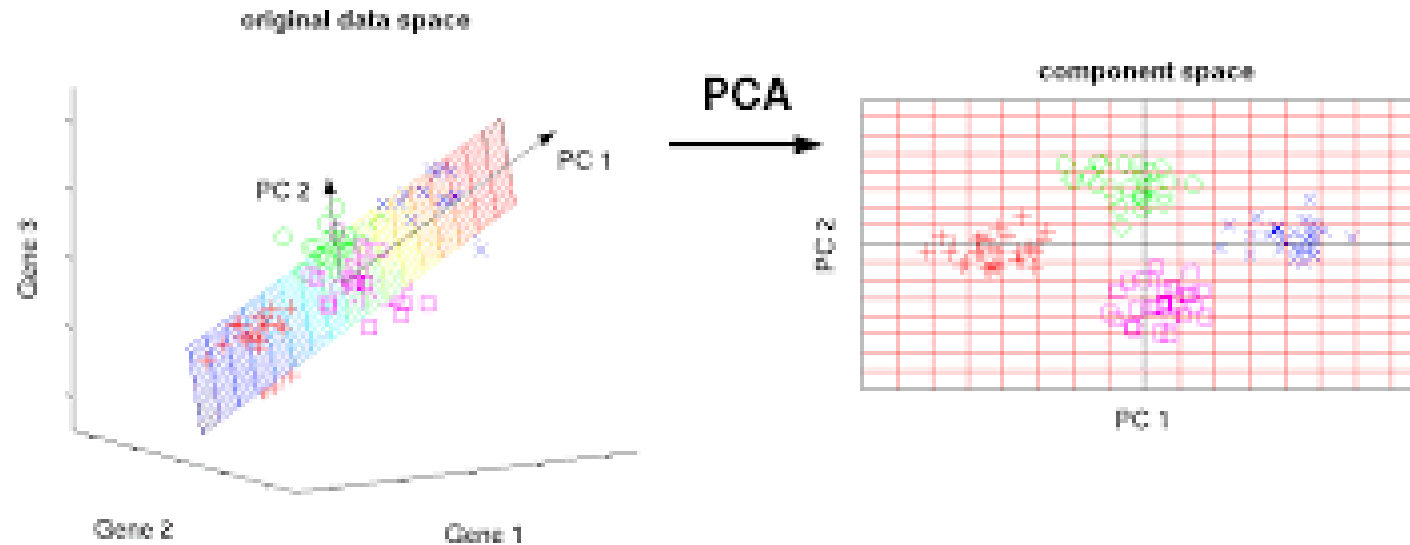
- Measure of approximation error in using d PCs

□ Example: Suppose eigenvalues of sample covariance matrix are 10, 4, 0.2, 0.1, 0, 0, ...

- What is the PoV for $d = 1, 2, 3, \dots$

Visualizing the Representation

- Finds a low-dimensional representation



Geometry of Approximations

□ Approximation can be interpreted geometrically

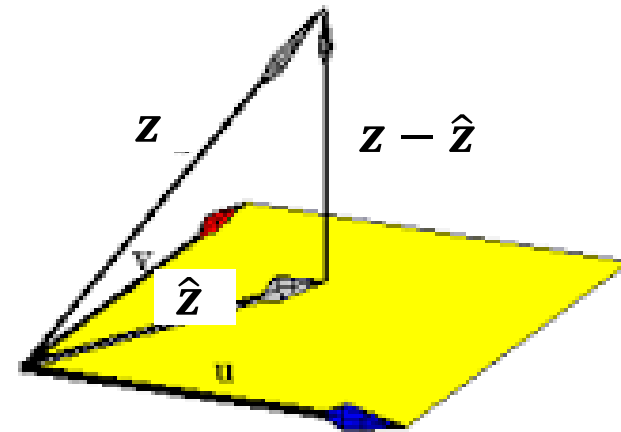
□ Let V be set of all linear combinations

$$\sum_{j=1}^d \alpha_j \mathbf{v}_j$$

- V is a vector space
- Called the span of $\mathbf{v}_1, \dots, \mathbf{v}_d$

□ Given \mathbf{z} , $\hat{\mathbf{z}}$ is the closest vector in V to \mathbf{z}

□ Called the projection of \mathbf{z} onto V

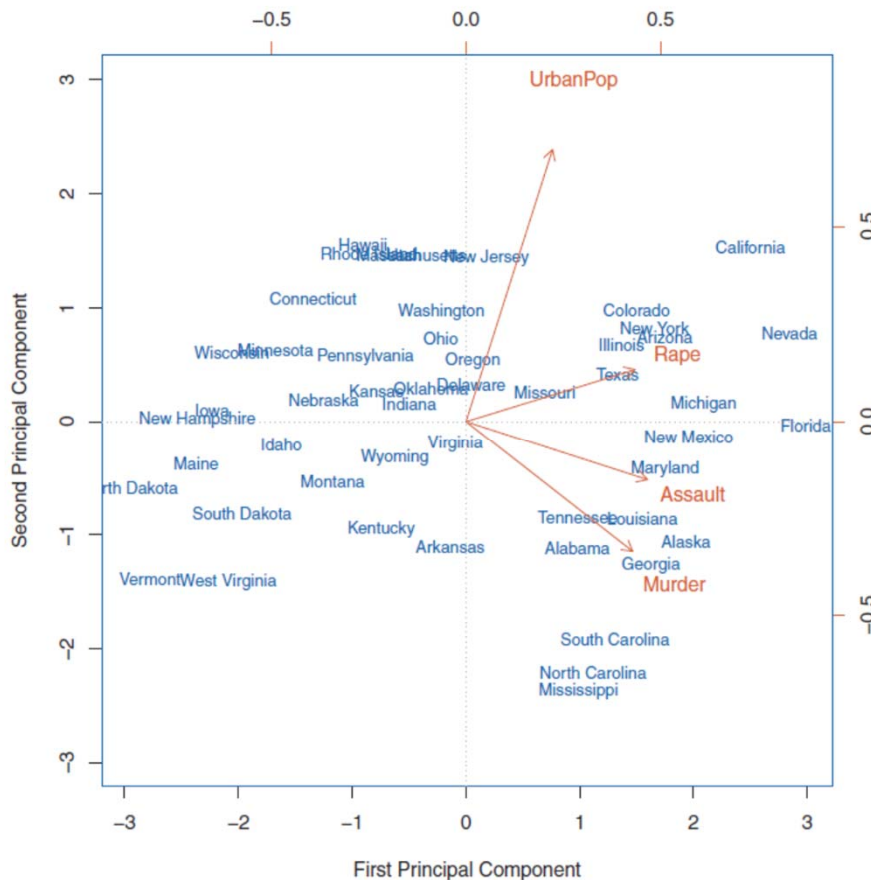


Space spanned by $\mathbf{v}_1, \dots, \mathbf{v}_d$

Latent Representations

- ❑ Each record is of the form: $\mathbf{x}_i \approx \bar{\mathbf{x}} + \sum_{j=1}^d \alpha_{ij} \mathbf{v}_j$
- ❑ Variance in \mathbf{x}_i explained by small number of “latent components”
 - Coefficients α_{ij} are the latent representations of \mathbf{x}_i
- ❑ Example:
 - \mathbf{x}_i = list of movie preferences for customer i
 - Movie preferences are highly correlated.
 - Could be explained by small number of components (action, romance, presence of stars, ...)
 - PCA can be used to find these out

Example: USArrests



Arrests per capita in four categories

- One record per US state


Visualize PCA in a **biplot**

- See the scores (i.e. coefficients of each state)
- Loading (PC vectors)

Fig from ISL 10.1



Outline

- Dimensionality reduction
- Principal components and directions of variance
- Approximation with PCs
-  Computing PCs via the SVD
- Face example in python

Singular Value Decomposition

- Given matrix $A \in R^{M \times N}$
- SVD is $A = USV^T$, where
 - $U \in R^{M \times r}$, $U^T U = I_r$
 - $V \in R^{N \times r}$, $V^T V = I_r$
 - $S = \text{diag}(s_1, \dots, s_r)$, sorted $s_1 \geq s_2 \geq \dots \geq s_r \geq 0$. Called the singular values
- All matrices have an SVD
- Number of singular values $r \leq \min(M, N)$

Computing the PCA via SVD

□ Let $\tilde{\mathbf{X}}$ = data matrix with sample mean removed.

□ Take SVD: $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{S}\mathbf{V}^T$

□ Properties:

- Sample covariance matrix is $\mathbf{Q} = \frac{1}{N}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}} = \frac{1}{N}\mathbf{V}\mathbf{S}^2\mathbf{V}^T$
- Eigenvalues are $\lambda_j = s_j^2/N$
- PCs are \mathbf{v}_j , columns of \mathbf{V}
- Coefficients are $\alpha = \tilde{\mathbf{X}}\mathbf{V} = \mathbf{U}\mathbf{S}$

Finding the Basis Vectors

- ❑ Consider problem of finding one basis vector, \mathbf{v}
- ❑ Given basis vector \mathbf{u} , the minimum approximation error for the i -th sample is:

$$\min_{\alpha_i} \|\mathbf{x}_i - \bar{\mathbf{x}} - \alpha_i \mathbf{v}\|^2$$

- Represents how well \mathbf{x}_i can be represented by the vector \mathbf{v}

- ❑ Define the **average approximation error**:

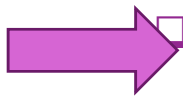
$$J(\mathbf{v}) := \frac{1}{N} \sum_{i=1}^N \min_{\alpha_i} \|\mathbf{x}_i - \bar{\mathbf{x}} - \alpha_i \mathbf{v}\|^2$$

- ❑ Select \mathbf{v} to minimize approximation error

$$\min_{\mathbf{v}} J(\mathbf{v})$$

Outline

- Dimensionality reduction
- Principal components and directions of variance
- Approximation with PCs
- Computing PCs via the SVD
- Face example in python



Computing the SVD

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

Then, we compute an SVD. Note that in python the SVD re

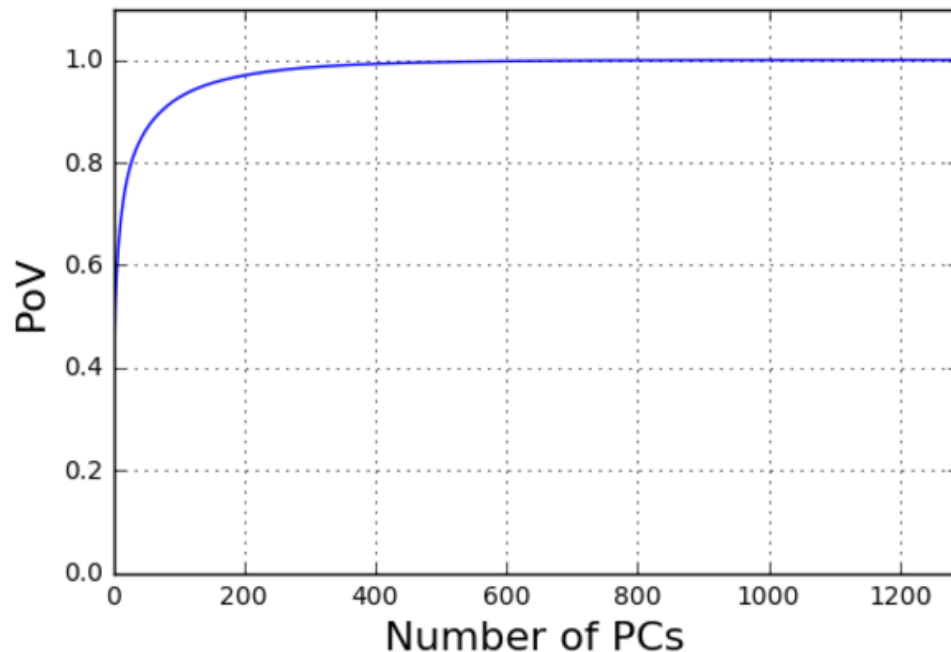
```
U,S,V = np.linalg.svd(Xs, full_matrices=False)
```

□ Note efficient use of python broadcasting

□ SVD:

- Use full_matrices (avoids computing zero SVs)
- Matrix V is what we call V^T (Different from MATLAB)

Finding the PoV



- Most variance explained in about 400 components
- Some reduction

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

plt.plot(PoV)
plt.grid()
plt.axis([1,n_samples,0, 1.1])
plt.xlabel('Number of PCs', fontsize=16)
plt.ylabel('PoV', fontsize=16)
```

Plotting Approximations

```
nplt = 2          # number of faces to plot
ds = [0,5,10,20,100] # number of SVD approximations

# Select random faces
inds = np.random.permutation(n_samples)
inds = inds[:nplt]
nd = len(ds)

# Set figure size
plt.figure(figsize=(1.8 * (nd+1), 2.4 * nplt))
plt.subplots_adjust(bottom=0, left=.01, right=.99, top=.90, hspace=.35)

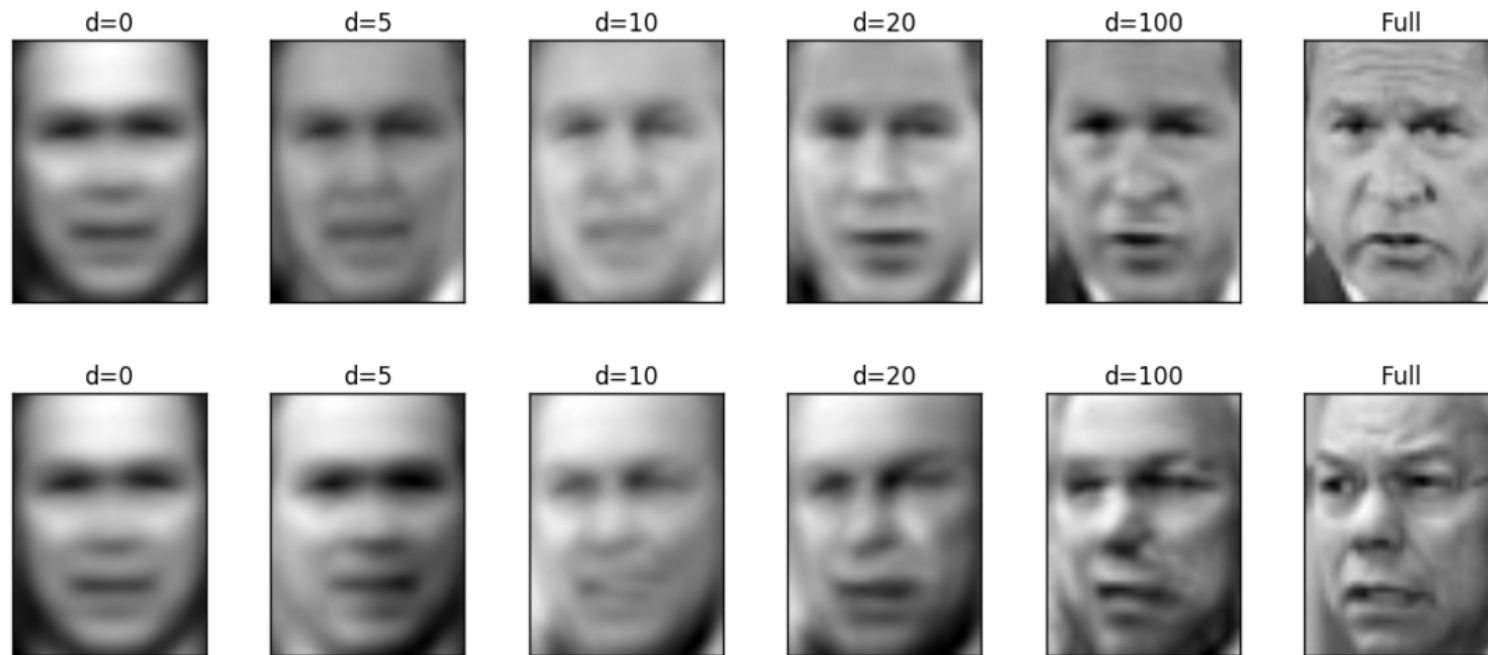
# Loop over figures
iplt = 0
for ind in inds:
    for d in ds:
        plt.subplot(nplt,nd+1,iplt+1)
        Xhati = (U[ind,:d]*S[None,:d]).dot(V[:,d,:]) + Xmean
        plt_face(Xhati)
        plt.title('d={0:d}'.format(d))
        iplt += 1

# Plot the true face
plt.subplot(nplt,nd+1,iplt+1)
plt_face(X[ind,:])
plt.title('Full')
iplt += 1
```

□ Selection of figure sizes for subplot

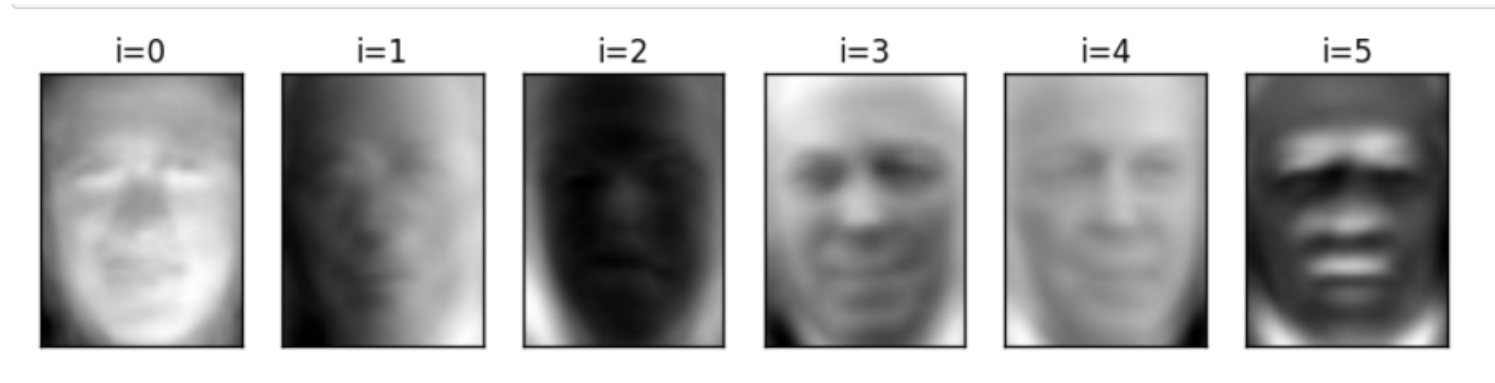
□ Note: Efficient computing of approx

Plotting the Approximations



Plotting the PCs

- The PCs can be plotted as well

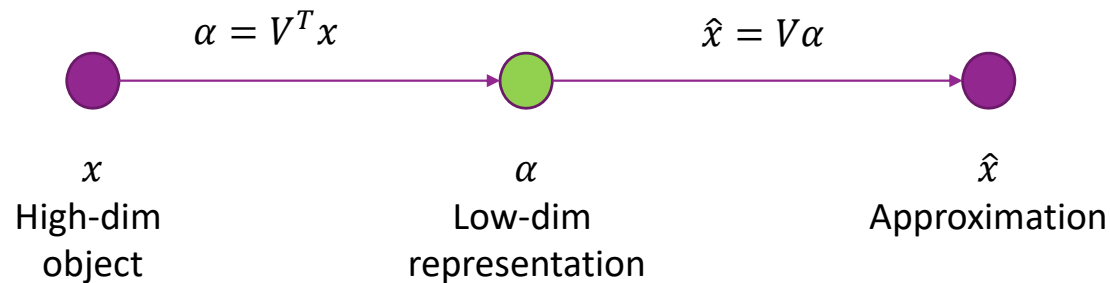


Other Considerations

- ❑ Normalization: For most data, it is essential to standardize before computing PC
 - Otherwise, components with large values dominate small ones
- ❑ Sklearn has built in PCA routine (will explore in lab)
- ❑ Some texts do not subtract mean
 - Will be picked up as one of the PCs

State-of-the-Art: Auto-Encoders

- ❑ PCA is a simple example of an **autoencoder**
- ❑ Tries to find low-dim representation
- ❑ Restricted to linear transforms
- ❑ Not very good for images and complex data



Deep Auto-Encoders

- Can use deep networks for learning complex latent representations and their inverses
 - http://www.cc.gatech.edu/~hays/7476/projects/Avery_Wenchen/
 - <https://swarbrickjones.wordpress.com/2016/01/13/enhancing-images-using-deep-convolutional-generative-adversarial-networks-dcgans/> (Code in Theano not tensorflow)

