

Lecture 11

Principal Component Analysis

EE-UY 4563 / EL-GY 9123: INTRODUCTION TO MACHINE LEARNING
PROF. SUNDEEP RANGAN, WITH MODIFICATION BY YAO WANG AND
YUAN WANG

Goal: find a “good” representation



- ❑ Sufficient (for the task): $I(\text{data}; \text{task}) = I(\text{representation}; \text{task})$
- ❑ Minimal (information): $I(\text{data}; \text{representation})$ is minimal
- ❑ Invariant (to nuisances): $I(\text{nuisances}; \text{task}) = 0$

Goal: find a “good” representation

- ❑ Example task: identify Jon Snow in images.

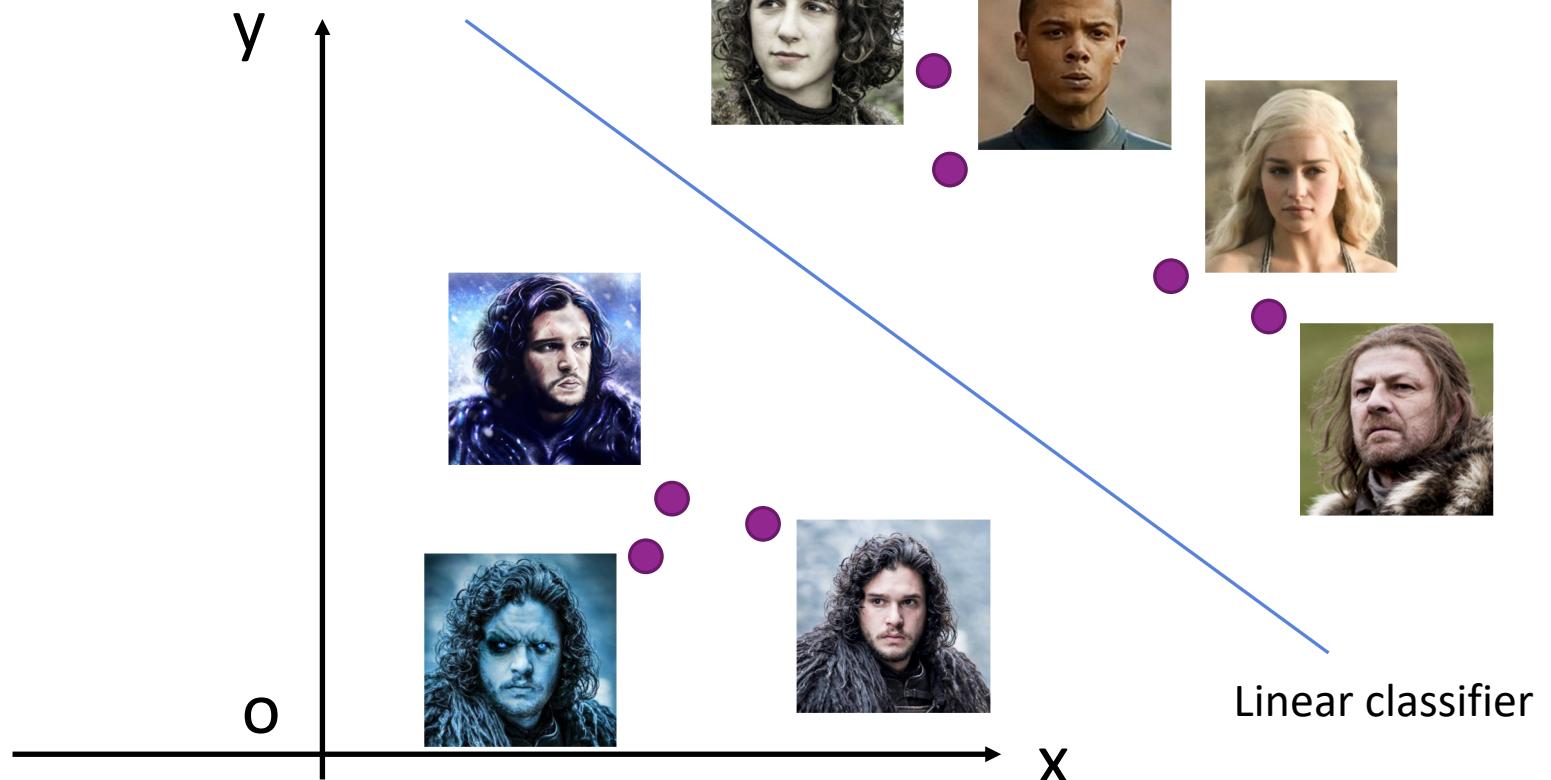


Images from the Internet

Goal: find a “good” representation

□ Example task: identify Jon Snow in images.

- Sufficient
- Minimal
- Invariant



Goal: find a “good” representation

- ## Today we focus on the simplest **linear** case:

$$\tilde{x} = \alpha x V$$

sample coefficient Dictionary / bases

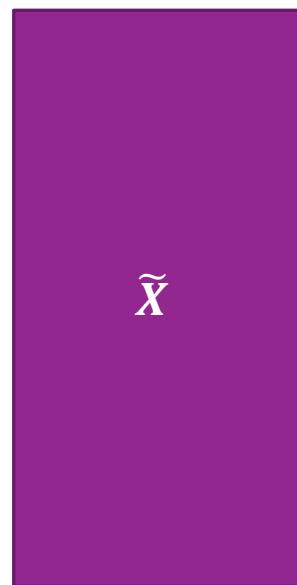
Goal: find a “good” representation

- ❑ Today we focus on the simplest linear case:

Goal: find a “good” representation

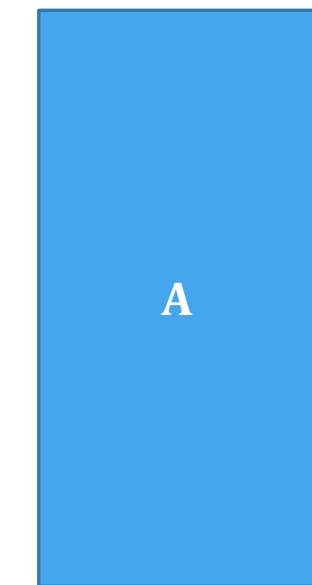
- Today we focus on the simplest **linear** case:

$$\tilde{X} = AV$$



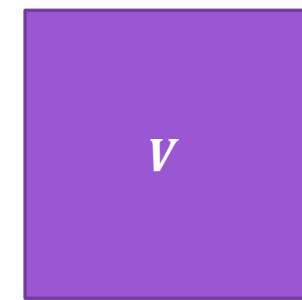
Data(many samples in rows)

=



A

x



V

Dictionary / bases

Goal: find a “good” representation

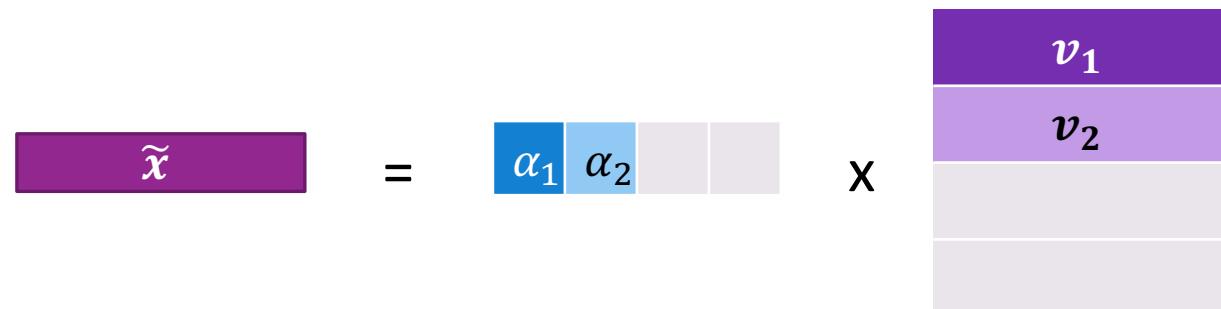
□ Today we focus on the simplest linear case:

□ Wanted properties:

- $\tilde{x} = \alpha V$
- α to be short
 - v to be typical patterns of \tilde{x}
- α is easy to compute
 - $\alpha = x^* V^{-1}$
 - $\alpha = x V^T \Leftrightarrow VV^T = I$

$$\underset{V}{\operatorname{argmin}} \underset{\alpha}{\min} \mathbb{E}(\|\tilde{x} - \alpha V\|_2)$$

$$s.t. VV^T = I, |\alpha| = d \leq |\tilde{x}| = p$$



Goal: find a “good” representation

□ Today we focus on the simplest **linear** case:

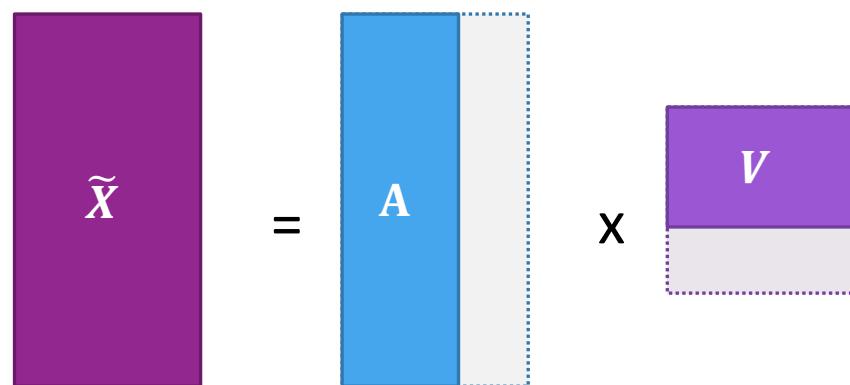
□ Wanted properties:

- $\tilde{x} = \alpha V$
- α to be short
 - v to be typical patterns of \tilde{x}
- α is easy to compute
 - $\alpha = x''V^{-1}$
 - $\alpha = xV^T \Leftrightarrow VV^T = I$

$$\operatorname{argmin}_V \min_A \|\tilde{X} - AV\|_F$$

$$s.t. VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, d \leq p$$



Outline

- 
- Dimensionality reduction
 - ❑ Principal components and directions of variance (a geometric perspective)
 - ❑ Computing PCs via the SVD
 - ❑ Approximation with PCs
 - ❑ Face example in python
 - ❑ Extensions

Dimensionality Reduction

- ❑ Many modern data sets have very high dimension
- ❑ Want to reduce dimension:
 - The Curse of Dimensionality (task “difficulty” $\sim \exp(d)$)
 - Visualize data
 - Find underlying commonalities in data

Outline

- ❑ Dimensionality reduction
-  **Principal components and directions of variance (a geometric perspective)**
- ❑ Computing PCs via the SVD
- ❑ Approximation with PCs
- ❑ Face example in python
- ❑ Extensions

Data Definitions

- ❑ x is a random variable $x = [x_1, x_2, \dots, x_p]$
- ❑ Data points $\{x_1, x_2, \dots, x_N\}$ are independent & random samples drawn from a distribution
- ❑ Data can be represented as an $N \times p$ matrix X
- ❑ No label: unsupervised learning

- ❑ \tilde{x} is a demeaned version of x

Simple properties of $\tilde{\mathbf{x}}$

- ◻ $\tilde{\mathbf{x}}$ is a demeaned version of \mathbf{x}
- ◻ Expectation μ :
 - $\mu = \mathbb{E}(\tilde{\mathbf{x}}) = \mathbb{E}(\mathbf{x} - \bar{\mathbf{x}}) = \mathbb{E}(\mathbf{x}) - \mathbb{E}(\bar{\mathbf{x}}) = 0$

- ◻ Covariance Σ :
 - $\Sigma = cov(\tilde{\mathbf{x}}) = \mathbb{E}((\tilde{\mathbf{x}} - \mu)^T(\tilde{\mathbf{x}} - \mu)) = \mathbb{E}(\tilde{\mathbf{x}}^T \tilde{\mathbf{x}})$
 - $\Sigma \approx \frac{1}{N} \sum_{n=1}^N \tilde{\mathbf{x}}_n^T \tilde{\mathbf{x}}_n = \frac{1}{N} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \equiv Q$

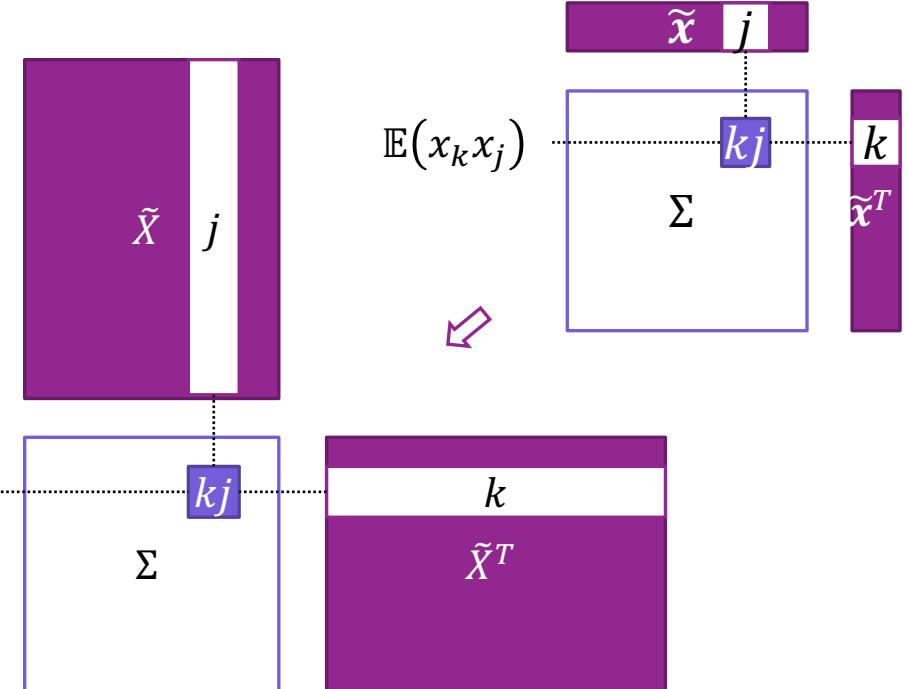
Sampling approximation:
$$\mathbb{E}(f(\mathbf{x})) = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} \approx \frac{1}{N} \sum_n f(\mathbf{x}_n)$$

Here $f(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}}^T \tilde{\mathbf{x}}$

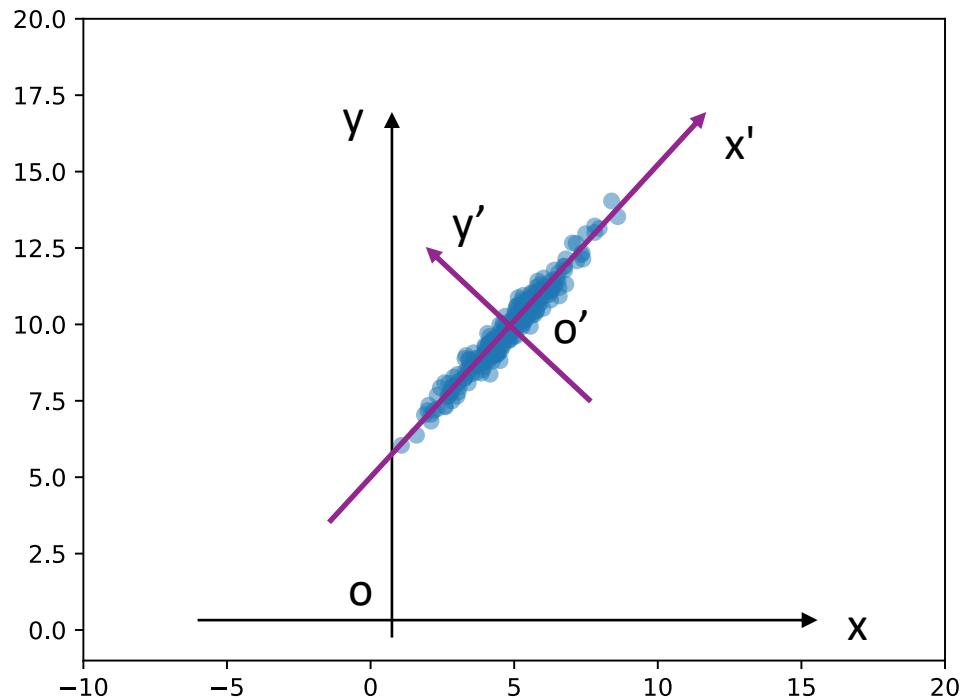
Simple manipulation

Sample covariance

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^N x_{nk} x_{nj} \\ = \frac{1}{N} \tilde{\mathbf{X}}_k \cdot \tilde{\mathbf{X}}_j \end{aligned}$$

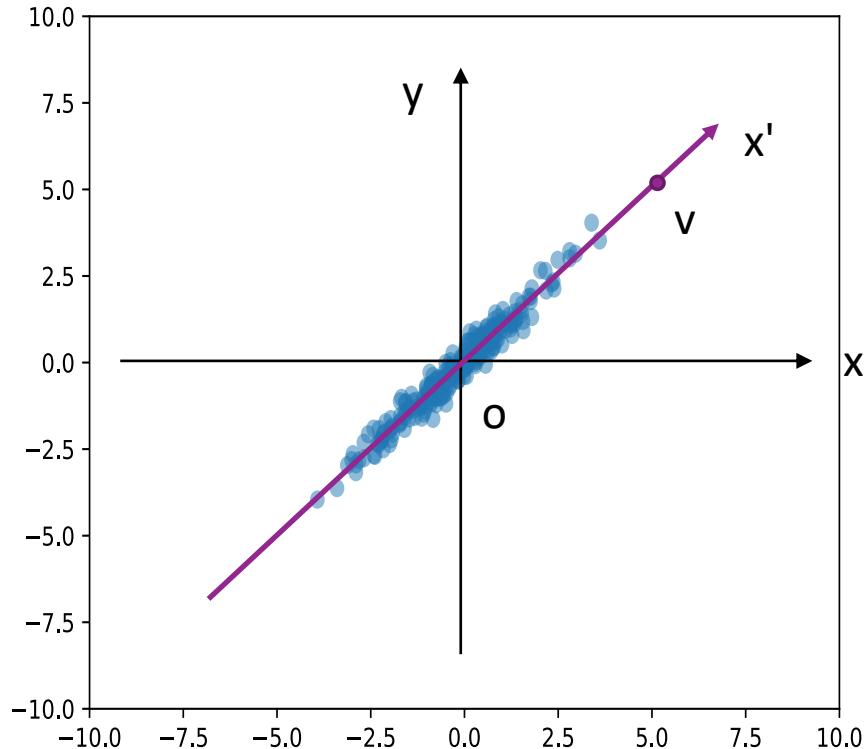


Geometric intuition



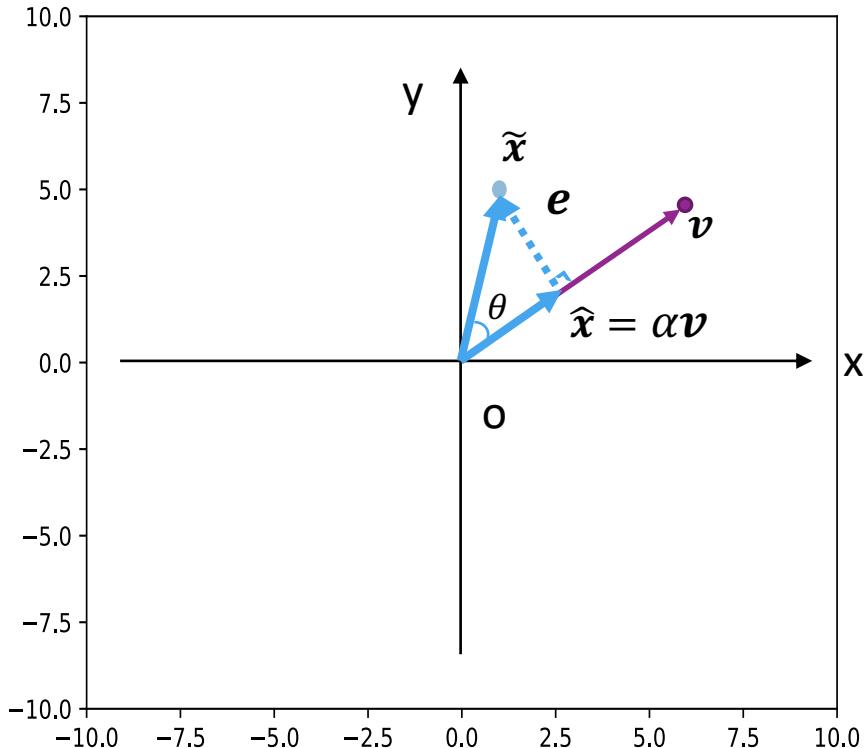
- ❑ How to find the new coordinate system?
- ❑ How specify a axis?
 - 2 points in the space in general
 - 1 point with common start point
 - Data mean serves a good candidate

Geometric intuition



- ❑ Demean, so that the origin is 0.
- ❑ What is the key property of the new axis?

Geometric intuition

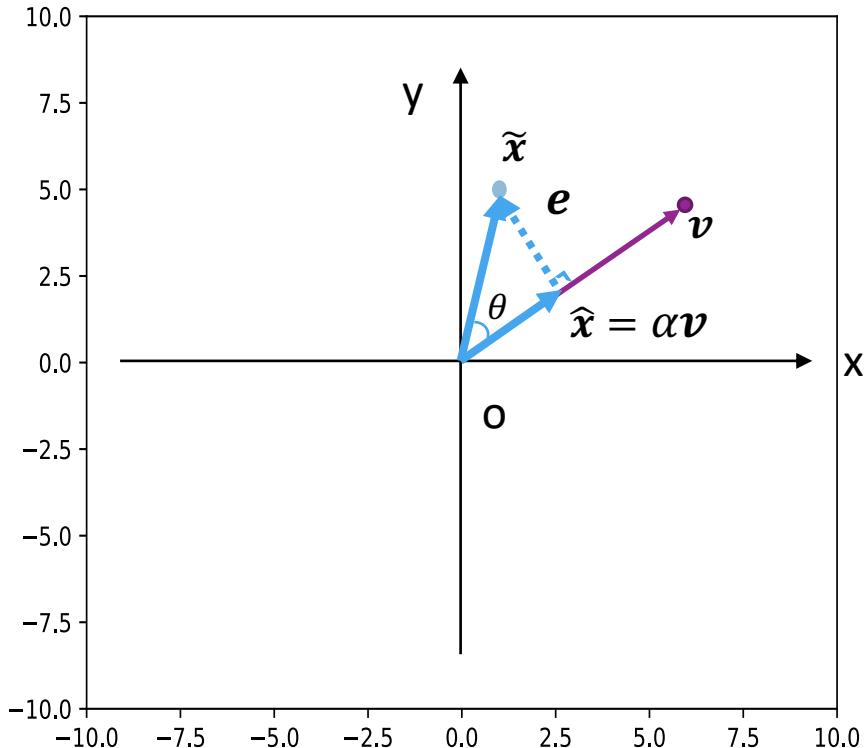


- ❑ Demean, so that the origin is 0.
- ❑ What is the key property of the new axis?
 - For a fixed v , chose α to approximate \tilde{x} with $\hat{x} = \alpha v$
$$\tilde{x} = \hat{x} + e$$
 - The approximation error $\|e\|^2$ is minimized when $e \perp v$:
$$\|e\|^2 = (\|\tilde{x}\| \sin \theta)^2$$

$$\|\hat{x}\|^2 = (\|\tilde{x}\| \cos \theta)^2$$

$$\|e\|^2 + \|\hat{x}\|^2 = \|\tilde{x}\|^2$$
 - Minimizing $\|e\|^2$ is equivalent to maximizing $\|\hat{x}\|^2$
$$\max_v \mathbb{E}(\|\hat{x}\|^2)$$

Geometric intuition

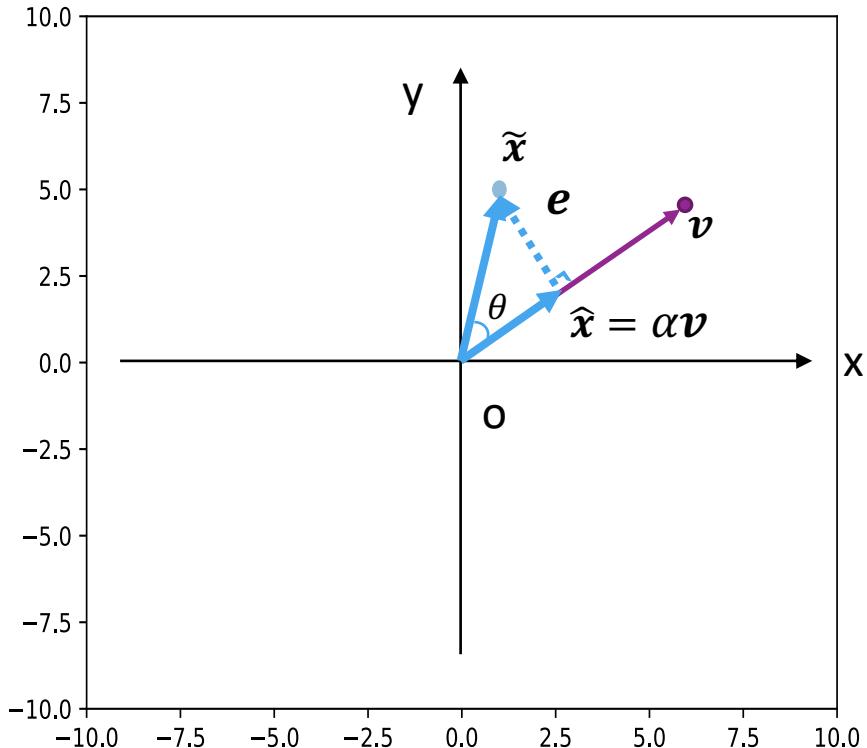


- ❑ Demean, so that the origin is 0.
- ❑ What is the key property of the new axis?

$$\begin{aligned} & \max_{\mathbf{v}} \mathbb{E}(\|\hat{\mathbf{x}}\|^2) \\ &= \max_{\mathbf{v}} \mathbb{E}((\|\tilde{\mathbf{x}}\| \cos(\theta))^2) \\ &= \max_{\mathbf{v}} \mathbb{E}\left(\left(\frac{\tilde{\mathbf{x}} \cdot \mathbf{v}}{\|\mathbf{v}\|}\right)^2\right) \\ &= \max_{\mathbf{v}} \mathbb{E}((\tilde{\mathbf{x}} \cdot \mathbf{v})^2) \\ &\text{s.t. } \|\mathbf{v}\| = 1 \end{aligned}$$

- $\hat{\mathbf{x}} = \text{Proj}_v(\tilde{\mathbf{x}}) = \|\tilde{\mathbf{x}}\| \cos(\theta) \times \frac{\mathbf{v}}{\|\mathbf{v}\|}$
- $\tilde{\mathbf{x}} \cdot \mathbf{v} = \|\mathbf{v}\| \|\tilde{\mathbf{x}}\| \cos(\theta) \rightarrow \|\tilde{\mathbf{x}}\| \cos(\theta) = \frac{\tilde{\mathbf{x}} \cdot \mathbf{v}}{\|\mathbf{v}\|}$
- If $\|\mathbf{v}\| = 1$, $\|\tilde{\mathbf{x}}\| \cos(\theta) = \tilde{\mathbf{x}} \cdot \mathbf{v}$

Geometric intuition



- ❑ Demean, so that the origin is 0.
- ❑ What is the key property of the new axis?

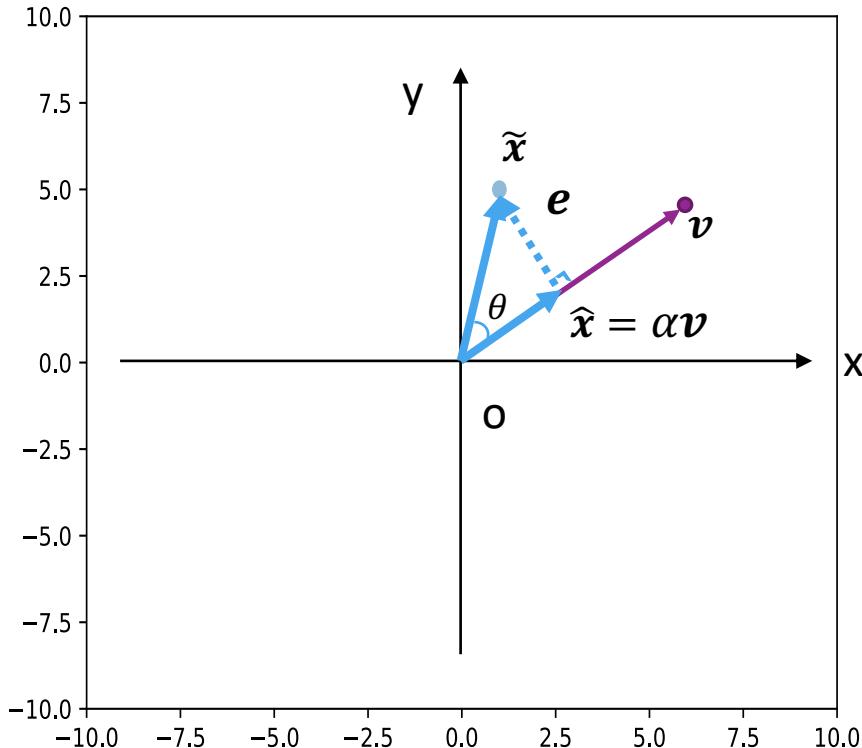
$$\max_{\mathbf{v}} \mathbb{E}((\tilde{\mathbf{x}} \cdot \mathbf{v})^2) \text{ s.t. } \|\mathbf{v}\| = 1$$

$$\begin{aligned} & var(\tilde{\mathbf{x}} \cdot \mathbf{v}) \\ &= \mathbb{E}((\tilde{\mathbf{x}} \cdot \mathbf{v} - \mathbb{E}(\tilde{\mathbf{x}} \cdot \mathbf{v}))^2) \\ &= \mathbb{E}((\tilde{\mathbf{x}} \cdot \mathbf{v})^2) \end{aligned}$$

Minimizing approximation error
is equivalent to
maximizing directional variance
(coordinate variance long a direction).



Geometric intuition



- ❑ Demean, so that the origin is 0.
- ❑ What is the key property of the new axis?

$$\begin{aligned} & \max_{\mathbf{v}} \mathbb{E}((\tilde{\mathbf{x}} \cdot \mathbf{v})^2) \text{ s.t. } \|\mathbf{v}\| = 1 \\ &= \max_{\mathbf{v}} \mathbb{E}((\tilde{\mathbf{x}} \mathbf{v}^T)^2) \\ &= \max_{\mathbf{v}} \mathbb{E}((\tilde{\mathbf{x}} \mathbf{v}^T)^T \tilde{\mathbf{x}} \mathbf{v}^T) \\ &= \max_{\mathbf{v}} \mathbb{E}(\mathbf{v} \tilde{\mathbf{x}}^T \tilde{\mathbf{x}} \mathbf{v}^T) \\ &= \max_{\mathbf{v}} \mathbf{v} \mathbb{E}(\tilde{\mathbf{x}}^T \tilde{\mathbf{x}}) \mathbf{v}^T \\ &\approx \max_{\mathbf{v}} \mathbf{v} Q \mathbf{v}^T \text{ s.t. } \|\mathbf{v}\| = 1 \end{aligned}$$

Q is the sample covariance

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

Directional variance can be calculated from the sample covariance.



Optimization

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

❑ **Idea:** shrink the solution set of \mathbf{v} to a few candidates and then exhaustive search

❑ **Theorem:** any local maxima of the directional variance is an eigenvector of Q

- Let $\mathbf{v}_1, \dots, \mathbf{v}_p$ (column vectors) be the eigenvectors of $Q : Q\mathbf{v}_j = \lambda_j \mathbf{v}_j$
- Sort eigenvalues in descending order: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$
- Can show that eigenvalues are real and non-negative

Optimization

- ❑ Lagrange multiplier
- ❑ For an constrained problem: $\max_{\nu} f(\nu) \text{ s.t. } g(\nu) = 0$
- ❑ Build the Lagrange function: $L(\nu, \lambda) = f(\nu) + \lambda g(\nu)$
- ❑ For any optimal ν^* for the original problem, there exists a λ^* such that (ν^*, λ^*) is a stationary point for the Lagrange function
- ❑ Thus, the method of Lagrange multipliers yields a necessary condition for optimal point for the original constrained problem.

Optimization

❑ For our constrained problem: $\max_{\boldsymbol{v}} \boldsymbol{v} Q \boldsymbol{v}^T$ s.t. $\|\boldsymbol{v}\| = 1$

❑ Lagrange function: $L(\boldsymbol{v}, \lambda) = \boldsymbol{v} Q \boldsymbol{v}^T + \lambda(\boldsymbol{v} \boldsymbol{v}^T - 1)$

❑ Find the stationary point of the Lagrange function:

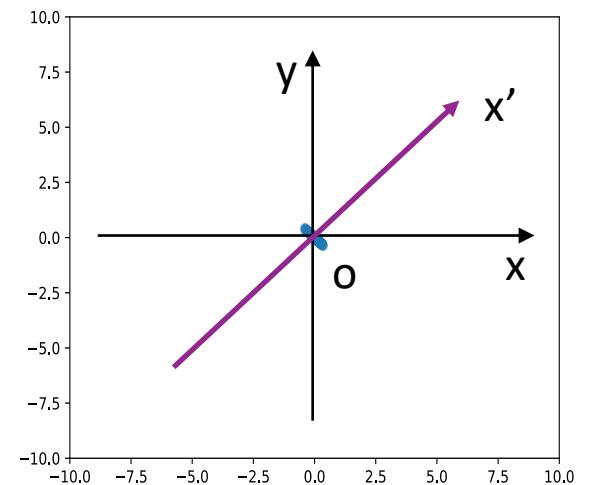
$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{v}} &= 2\boldsymbol{v} Q + 2\lambda\boldsymbol{v} = 0 \\ \rightarrow \boldsymbol{v} Q &= -\lambda\boldsymbol{v} \\ \rightarrow Q^T \boldsymbol{v}^T &= -\lambda\boldsymbol{v}^T \\ \rightarrow Q \boldsymbol{v}^T &= -\lambda\boldsymbol{v}^T\end{aligned}$$

❑ Any local maxima of the variance directional is an eigenvector

Optimization

- ❑ The candidate set of \mathbf{v} is the eigenvectors of Q : $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$
- ❑ The corresponding cost value $\mathbf{v}^T Q \mathbf{v}$: $\{\lambda_1, \dots, \lambda_p\}$
- ❑ Note from here onwards \mathbf{v} represent column vectors
- ❑ Since $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$
- ❑ \mathbf{v}_1 is the solution for the problem $\max_{\mathbf{v}} \mathbf{v}^T Q \mathbf{v} \quad s.t. \|\mathbf{v}\| = 1$
- ❑ After removing all variance explained by \mathbf{v}_1 , find the next max variance direction.

$$\mathbf{v}^T Q \mathbf{v} = \mathbf{v}^T \lambda \mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda$$



Principal Components

- ❑ Principal components: The eigenvectors of Q , $\mathbf{v}_1, \dots, \mathbf{v}_p$
- ❑ Sorted by eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$
- ❑ Each vector is of dimension p
- ❑ Represents directions of maximal variance

Principal Components

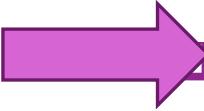
- ❑ Key properties of Q

- Symmetric: $Q_{ij} = Q_{ji}$
 - Eigenvalues are real and non-negative

- ❑ Key properties of PCs (here \mathbf{v}_j are row vectors)

- $\mathbf{v}_j \mathbf{v}_k^T = 1$ if $j = k$
 - $\mathbf{v}_j \mathbf{v}_k^T = 0$ if $j \neq k$
 - $VV^T = I$

Outline

- ❑ Dimensionality reduction
- ❑ Principal components and directions of variance (a geometric perspective)
-  Computing PCs via the SVD
- ❑ Approximation with PCs
- ❑ Face example in python
- ❑ Extensions

Computing PCs via the SVD

- SVD provides a numerically more stable way to calculate PCs

$$\begin{aligned} Q &= \frac{1}{N} \tilde{X}^T \tilde{X} \\ &= \frac{1}{N} V S^T U^T U S V^T \\ &= V \frac{S^T S}{N} V^T \\ \rightarrow QV &= V \frac{S^T S}{N} \end{aligned}$$

- Singular-Value Decomposition (SVD)
 $\forall M \in \mathbb{C}^{m \times n}, \exists r, U, S, V \text{ so that } M = USV^T$
 - $U \in \mathbb{C}^{m \times r}, U^T U = I$
 - $S \in \mathbb{C}^{r \times r}$ is a diagonal matrix
 - $V \in \mathbb{C}^{n \times r}, V^T V = I$
 - $r \leq \min(m, n)$ (Rank of the matrix)
- $\tilde{X} = USV^T$

Computing PCs via the SVD

- ❑ SVD provides a numerically more stable way to calculate PCs

$$QV = V \frac{S^T S}{N} \quad \text{or} \quad Q\boldsymbol{v}_i = \frac{S_{ii}^T S_{ii}}{N} \boldsymbol{v}_i$$

- ❑ Rows of V^T are eigenvectors of Q , e.g. PCs

- ❑ $\frac{S^T S}{N}$ are the eigenvalues (real and nonnegative)

➤ Singular-Value Decomposition (SVD)
 $\forall M \in \mathbb{C}^{m \times n}, \exists r, U, S, V \text{ so that}$

$$M = USV^T$$

- $U \in \mathbb{C}^{m \times r}, U^T U = I$
- $S \in \mathbb{C}^{r \times r}$ is a diagonal matrix
- $V \in \mathbb{C}^{n \times r}, V^T V = I$
- $r \leq \min(m, n)$ (Rank of the matrix)

➤ $\tilde{X} = USV^T$

Outline

- ❑ Dimensionality reduction
- ❑ Principal components and directions of variance (a geometric perspective)
- ❑ Computing PCs via the SVD
-  Approximation with PCs
- ❑ Face example in python
- ❑ Extensions

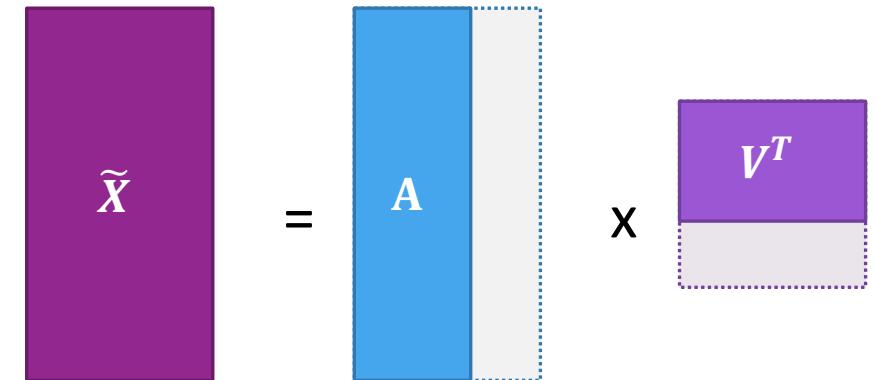
Average Approximation Error

- Given data $\tilde{\mathbf{x}}_i, i = 1, \dots, N$
- Let $\mathbf{v}_1, \dots, \mathbf{v}_p$ be the PCs
- Find coefficient expansion of each data sample:

$$\tilde{\mathbf{x}}_i = \sum_{j=1}^p \alpha_{ij} \mathbf{v}_j, \quad \alpha_{ij} = \mathbf{v}_j^T \tilde{\mathbf{x}}_i$$

- Now consider approximation with d coefficients:

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



Average Approximation Error

□ Error in sample i :

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

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

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

- Sum of the smallest $p - d$ eigenvalues

Proportion of Variance (PoV)

- ❑ Total variance of data set: $\frac{1}{N} \sum_{i=1}^N \|\tilde{\mathbf{x}}_i\|^2 = \sum_{j=1}^p \lambda_j$
- ❑ Average approximation error: $\frac{1}{N} \sum_{i=1}^N \|\tilde{\mathbf{x}}_i - \hat{\mathbf{x}}_i\|^2 = \sum_{j=d+1}^p \lambda_j$

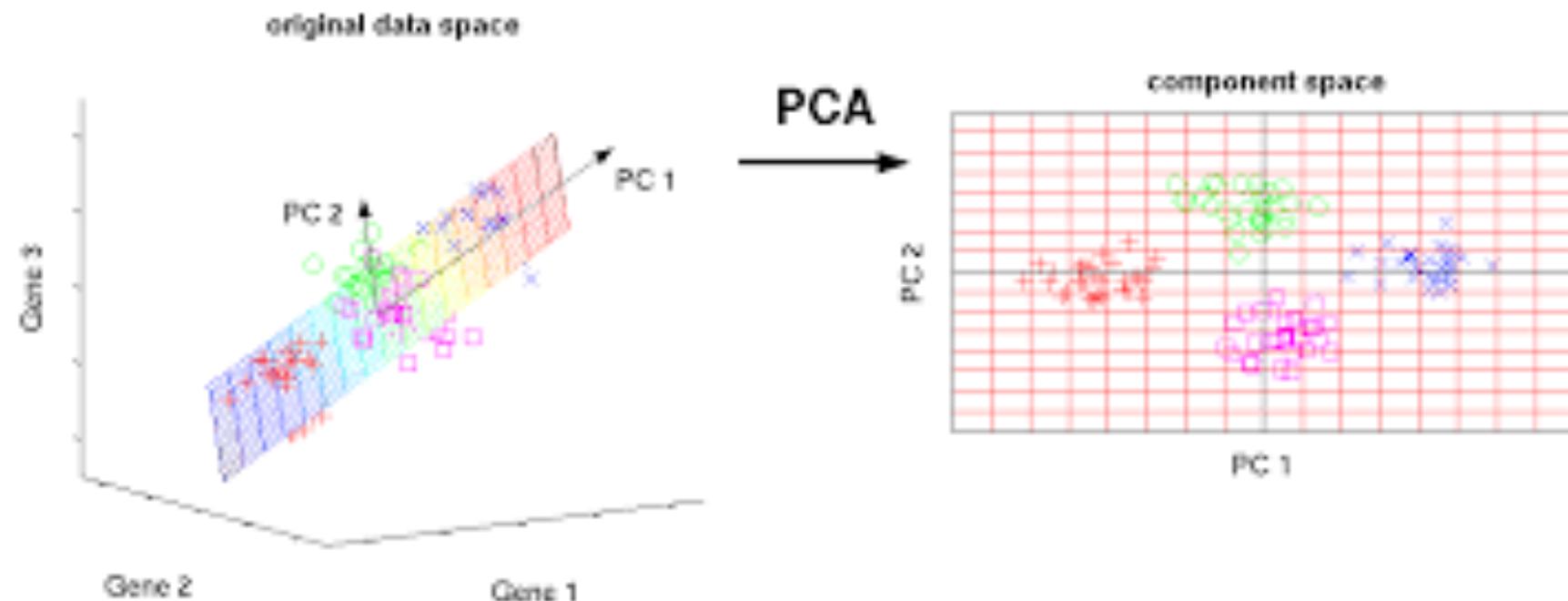
- ❑ The proportion of variance explained by d PCs is:

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

- 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

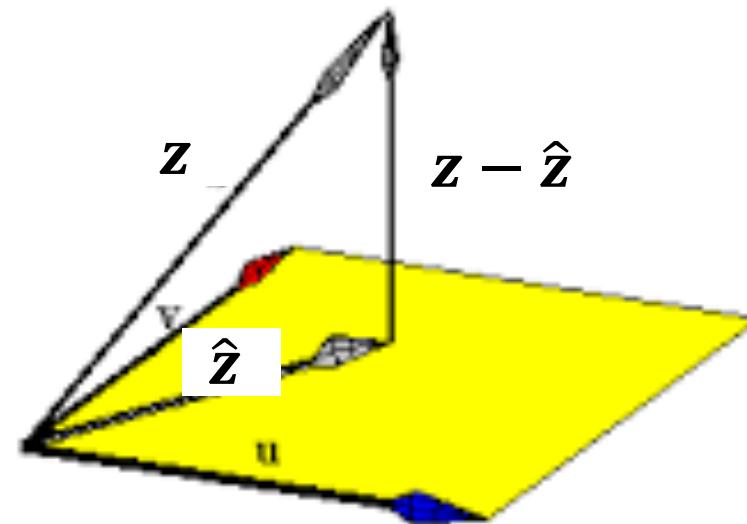
□ Let V be set of all linear combinations

$$\sum_{j=1}^d \alpha_j v_j$$

- V is a vector space
- Called the span of v_1, \dots, 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 v_1, \dots, v_d

Outline

- ❑ Dimensionality reduction
- ❑ Principal components and directions of variance (a geometric perspective)
- ❑ Computing PCs via the SVD
- ❑ Approximation with PCs
-  Face example in python
- ❑ Extensions

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

- ❑ 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\Sundeep\scikit_learn_data\lfw_home\lfw_funneled
2016-11-14 14:22:08,605 Loading LFW people faces from C:\Users\Sundeep\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



```
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]])
```

Computing the SVD

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

❑ Note efficient use of python broadcasting

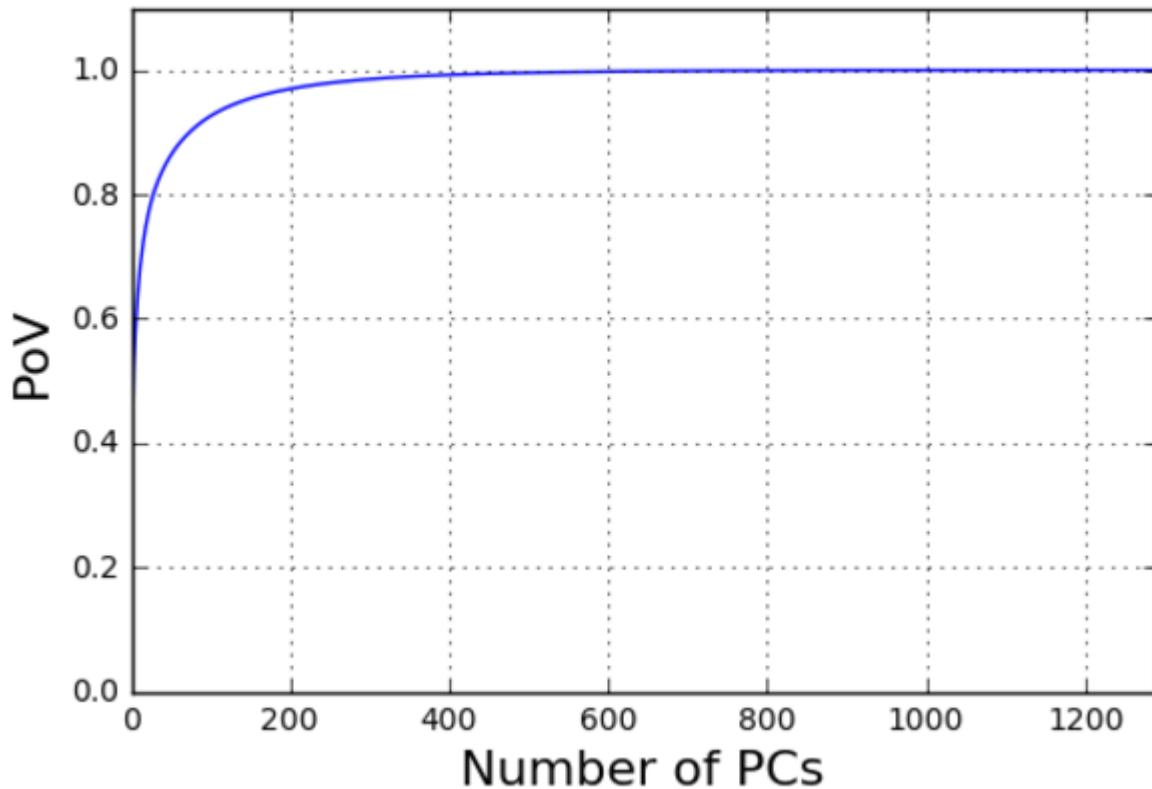
Then, we compute an SVD. Note that in python the SVD re

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

❑ 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

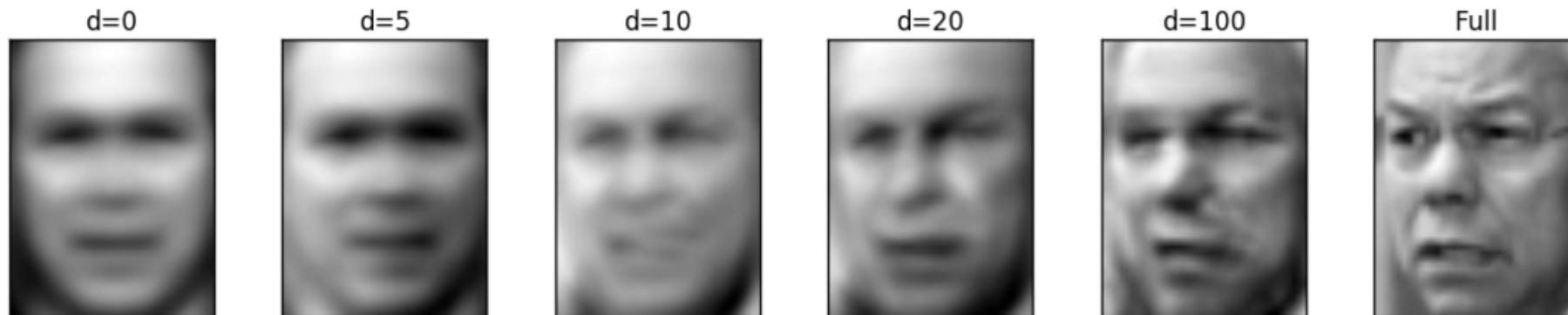
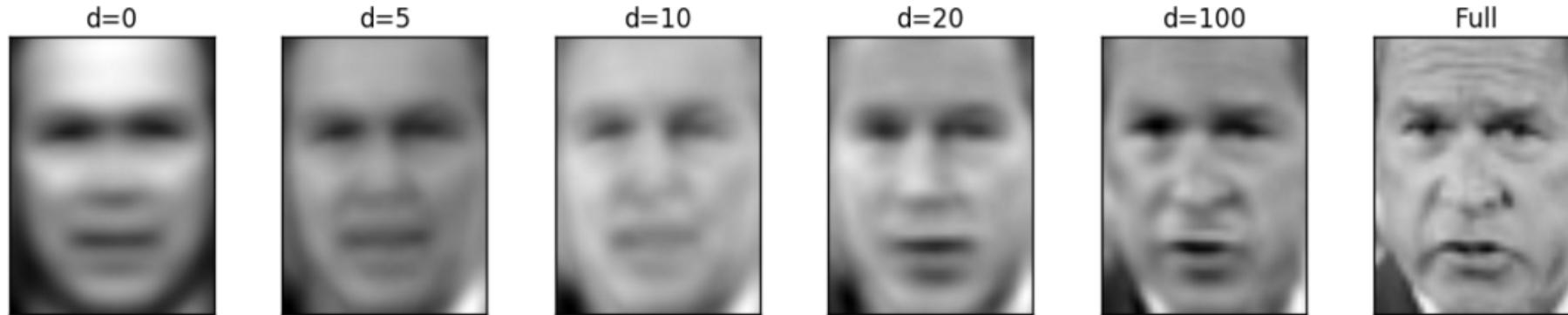
    # Plot the true face
    plt.subplot(nplt,nd+1,iplt+1)
    plt_face(X[ind,:])
    plt.title('Full')
    iplt += 1
```

❑ Selection of figure sizes for subplot

❑ Note: Efficient computing of approx

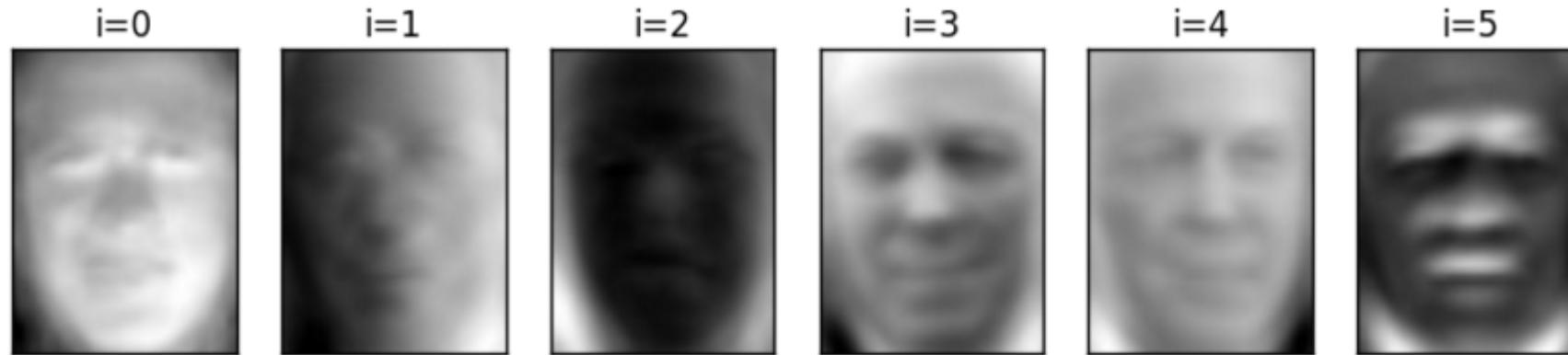


Plotting the Approximations



Plotting the PCs

- ❑ The PCs can be plotted as well
 - ❑ These are known as Eigenfaces!
-



How do we use Principle Components?

- ❑ PCs tell us “latent” factors of the data
 - Each sample is a weighted average of these PCs (eg. eigen-faces)
- ❑ Can use the coefficients as low-dimension features for classification
 - Early success with face recognition (using projection onto eigenfaces)
- ❑ Can use the coefficients for compression
 - JPEG image coder: uses DCT basis at block level
 - JPEG wavelet coder: uses Wavelet basis over entire image
- ❑ Denoising
 - Project noisy data onto the basis vectors
 - Reconstruct from low freq. basis vectors
- ❑ Data normalization: scaled coefficients will be independent and have unit variance.

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

Outline

- ❑ Dimensionality reduction
- ❑ Principal components and directions of variance (a geometric perspective)
- ❑ Computing PCs via the SVD
- ❑ Approximation with PCs
- ❑ Face example in python



❑ Extensions

Orthonormal Sets and Bases

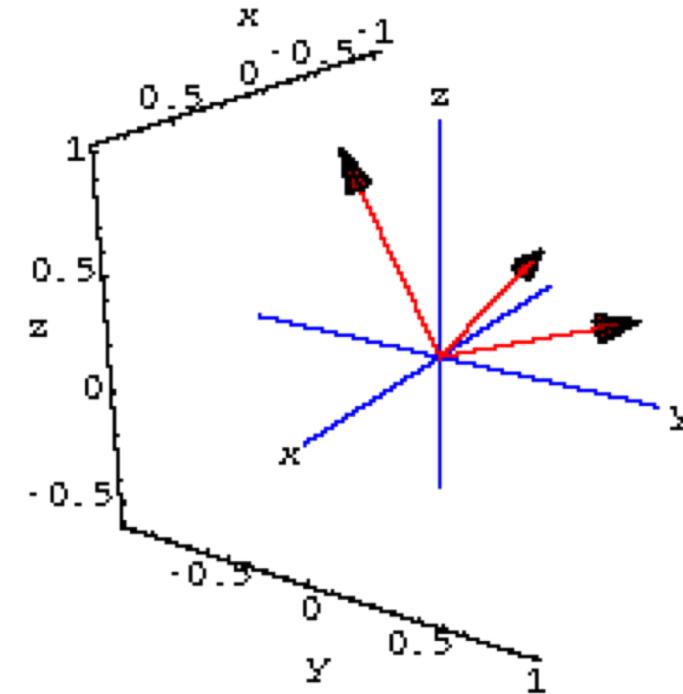
□ **Definition:** A set of vectors $\mathbf{v}_1, \dots, \mathbf{v}_d$ are an **orthonormal set** if:

- $\|\mathbf{v}_j\| = 1$ 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 $\mathbf{V}^T \mathbf{V} = I_d$

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



Coefficients in an Orthonormal Basis

- ❑ For an orthonormal set V

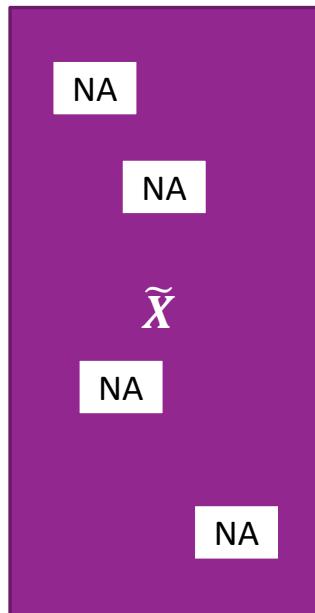
$$\alpha = V^T x, \quad x = V\alpha$$

- ❑ Examples of Orthonormal Basis:

- PCs, a data driven orthonormal basis
- Discrete Fourier Transform (normalized)
- Discrete Cosine Transform
- Wavelet transform

Wiberg's Method

- PCA with missing data



$$\operatorname{argmin}_V \min_A \|\tilde{X} - AV\|_F$$

$$s.t. VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, d \leq p$$



$$\operatorname{argmin}_V \min_A \|W \circ (\tilde{X} - AV)\|_F$$

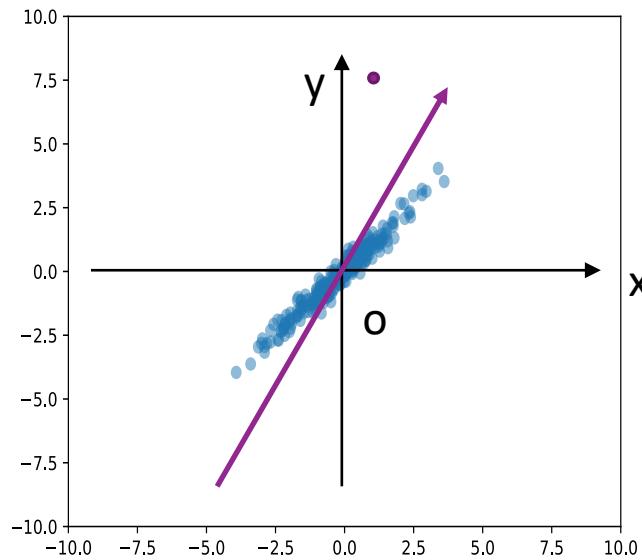
$$s.t. VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, d \leq p$$

$$W_{ij} = 1 \text{ if valid data, otherwise } 0$$

Robust PCA

- PCA is known sensitive to outlier



$$\operatorname{argmin}_V \min_A \|\tilde{X} - AV\|_F$$

$$s.t. VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, d \leq p$$



$$\operatorname{argmin}_V \min_A \|L - AV\|_F + \lambda \|S\|_1$$

$$s.t. \tilde{X} = L + S$$

$$VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, d \leq p$$

Sparse orthonormal transforms

□ From low rank to sparsity

$$\tilde{X} = \begin{array}{c|c|c} A & & \\ \hline & \ddots & \\ & & V \end{array} X$$

$$\tilde{X} = \begin{array}{c|c} A & \\ \hline & V \end{array} X$$

$$\operatorname{argmin}_{V} \min_{A} \|\tilde{X} - AV\|_F$$

$$s.t. VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, d \leq p$$



$$\operatorname{argmin}_{V} \min_{A} \|\tilde{X} - AV\|_F + \lambda \|A\|_1$$

$$s.t. VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times p}, V \in \mathbb{R}^{p \times p}$$

KSVD

- More parsimonious representation with redundant dictionary

$$\tilde{X} = \begin{matrix} A \\ \vdots \end{matrix} X \begin{matrix} V \\ \vdots \end{matrix}$$

$$\tilde{X} = \begin{matrix} A \\ \vdots \end{matrix} X \begin{matrix} V \\ \vdots \end{matrix}$$

$$\operatorname{argmin}_V \min_A \|\tilde{X} - AV\|_F$$

$$s.t. VV^T = I$$

$$\tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, d \leq p$$



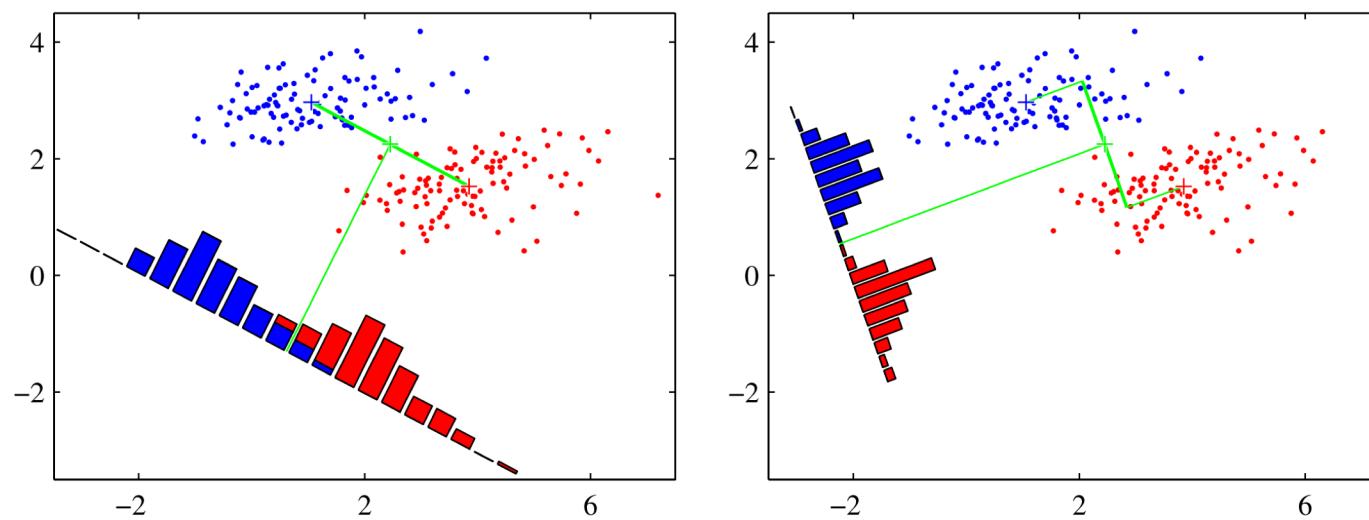
$$\operatorname{argmin}_V \min_A \|\tilde{X} - AV\|_F$$

$$s.t. \tilde{X} \in \mathbb{R}^{N \times p}, A \in \mathbb{R}^{N \times d}, V \in \mathbb{R}^{d \times p}, (d \geq p)$$



Linear discriminative analysis

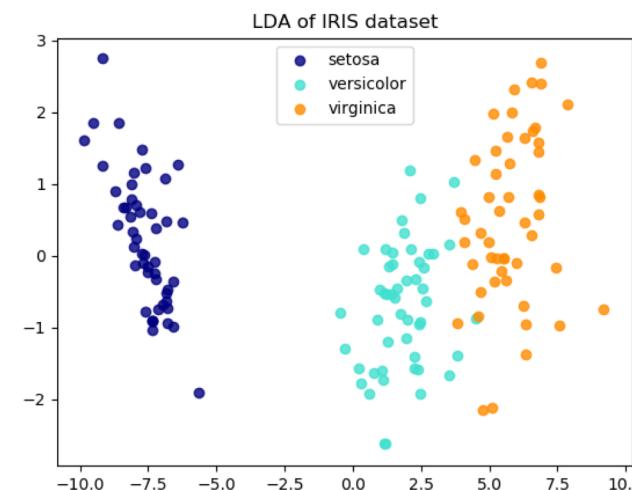
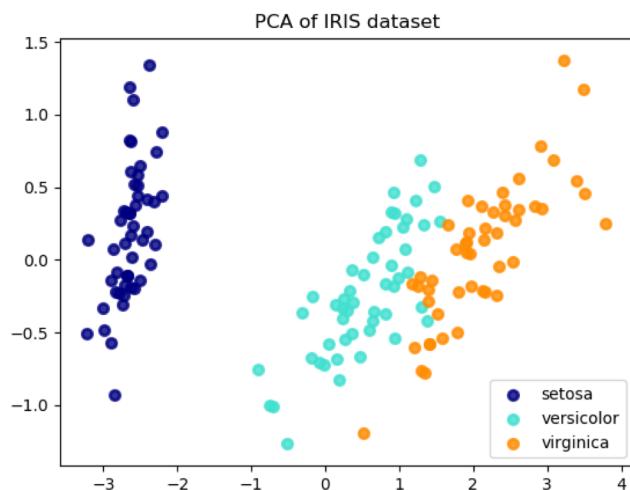
- ❑ Projects to the most discriminative direction
- ❑ Supervised learning with labels
- ❑ Simultaneously maximizing between classes variance and minimizing within class variance



*figure from Bishop, Pattern Recognition and Machine Learning

Linear discriminative analysis

- ❑ Projects to the most discriminative direction
- ❑ Supervised learning with labels
- ❑ Simultaneously maximizing between classes variance and minimizing within class variance



*figure from <http://scikit-learn.org>

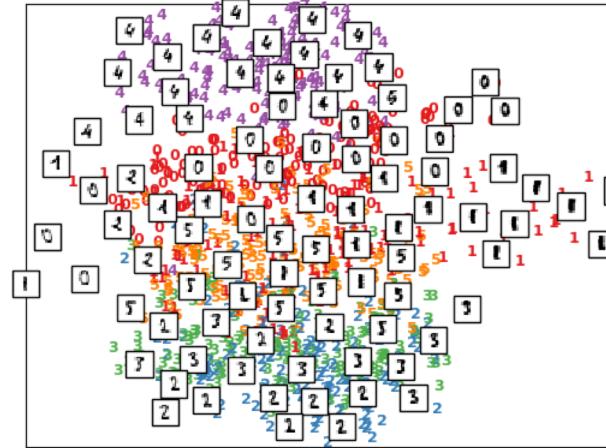
Linear discriminative analysis

- ❑ Projects to the most discriminative direction
- ❑ Supervised learning with labels
- ❑ Simultaneously maximizing between classes variance and minimizing within class variance

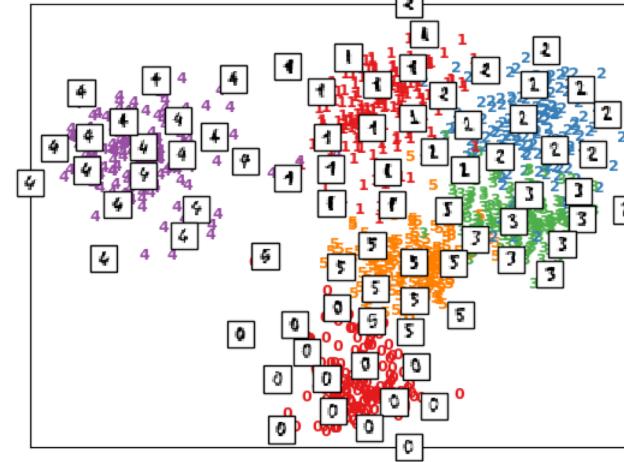
A selection from the 64-dimensional digits dataset

0	1	2	3	4	5	0	1	3	3	0	1	2	3	4	5	0	5
5	5	0	4	1	3	5	1	0	0	2	2	2	0	1	2	3	3
4	4	1	5	0	5	4	4	0	0	1	3	2	1	1	3	1	4
3	4	4	0	5	3	4	5	4	2	2	2	5	5	4	0	0	1
2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	5
0	4	4	3	5	4	1	3	0	2	0	1	3	3	3	4	4	4
1	5	0	5	2	2	0	0	1	3	1	4	3	1	4	4	3	4
0	5	7	4	4	4	1	1	2	5	5	4	4	0	0	1	2	3
5	0	1	2	3	4	5	0	2	3	4	5	0	5	5	5	0	4
3	5	4	0	0	2	2	0	1	4	2	1	3	3	3	4	4	5
5	2	2	0	0	4	3	2	4	3	4	3	1	4	3	1	4	5
3	8	5	4	4	2	2	2	5	5	4	4	0	3	0	1	2	3
0	1	2	3	4	5	0	1	2	3	4	5	0	3	0	1	2	3
5	1	0	0	1	2	1	2	1	0	1	3	3	3	4	4	4	5
1	2	0	0	1	3	2	1	4	3	1	3	4	3	1	4	5	0
1	2	0	0	1	3	2	1	4	3	1	3	4	3	1	4	5	0
1	5	4	4	2	1	4	3	1	3	4	3	1	4	3	1	4	5
2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	4
0	0	2	1	0	1	3	3	3	4	4	4	4	4	4	4	4	5
0	0	1	1	0	1	3	3	3	4	4	4	4	4	4	4	4	5
0	4	2	2	1	5	5	4	0	0	1	2	3	4	5	0	1	2

Principal Components projection of the digits (time 0.01s)



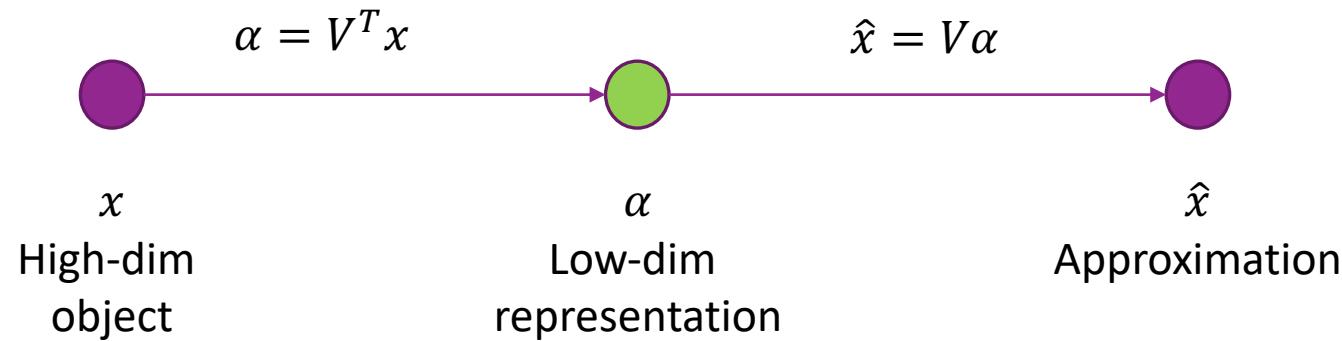
Linear Discriminant projection of the digits (time 0.01s)



*figure from <http://scikit-learn.org>

State-of-the-Art: Auto-Encoders

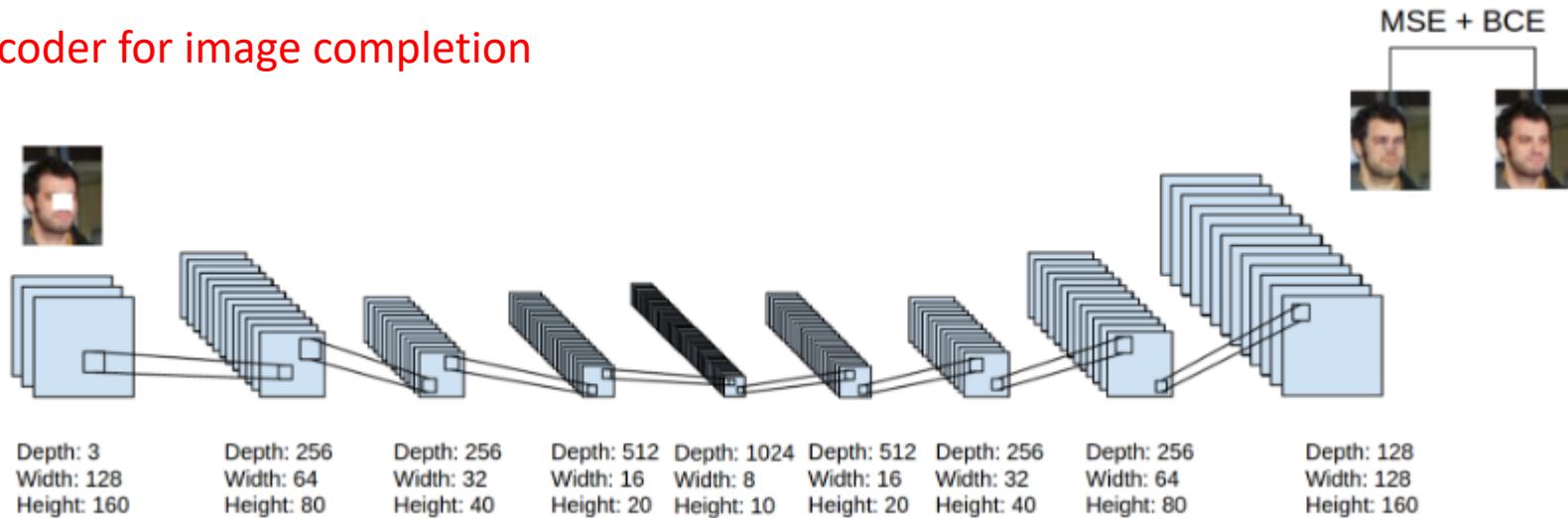
- ❑ PCA is a simple example of an **autoencoder**
- ❑ Tries to find low-dim representation
- ❑ Restricted to linear transforms
- ❑ Not very good for images and complex data



Deep Auto-Encoders

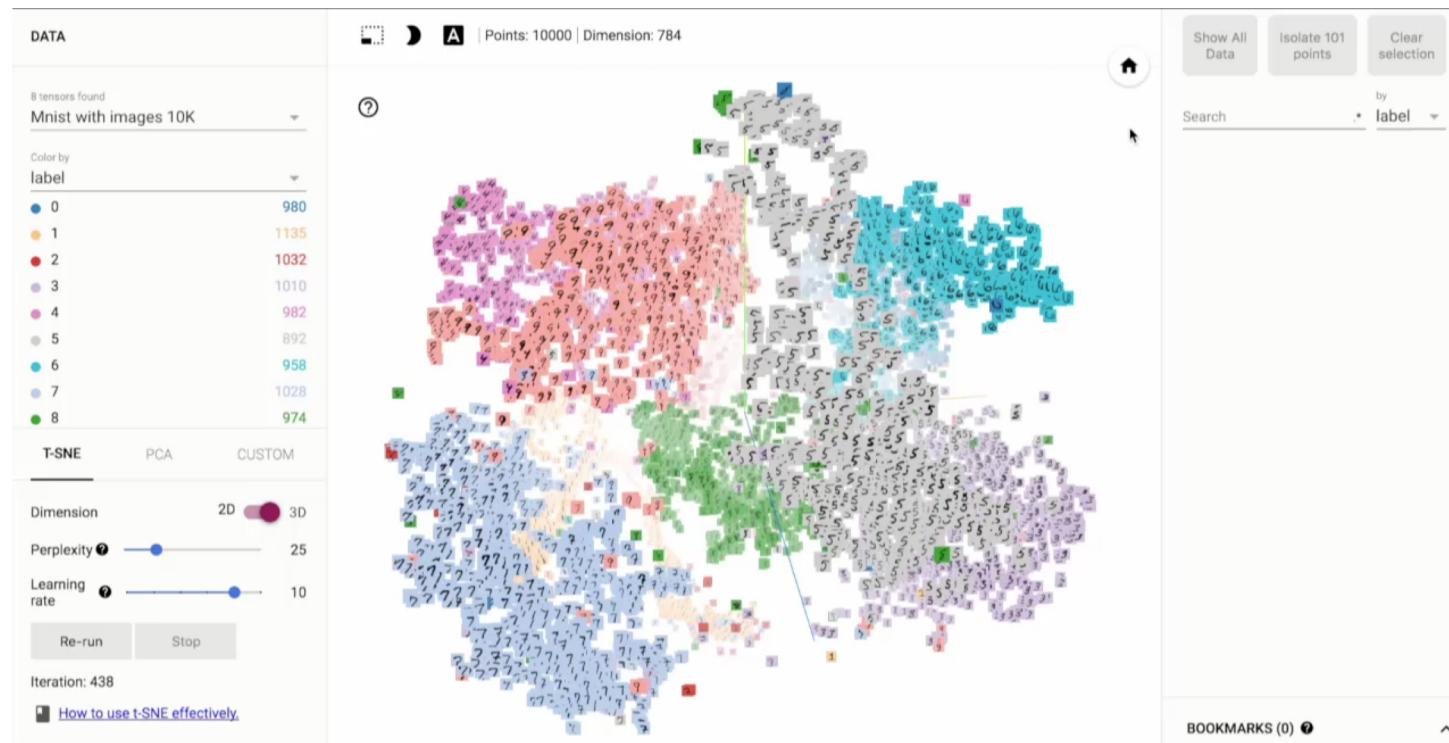
- Can use deep networks for learning complex latent representations and their inverses
 - http://www.cc.gatech.edu/~hays/7476/projects/Avery_Wenchens/
 - <https://swarbrickjones.wordpress.com/2016/01/13/enhancing-images-using-deep-convolutional-generative-adversarial-networks-dcgans/> (Code in Theano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

Autoencoder for image completion



Representation with deep learning

□ MNIST in 3D



*figure from <https://tensorflow.org>