# Unit 11 Principal Component Analysis

EE-UY 4563 / EL-GY 6143: INTRODUCTION TO MACHINE LEARNING PROF. SUNDEEP RANGAN





# Learning Objectives

- □ Identify cases to use dimensionality reduction
- ☐ Mathematically describe principal components representations of data
- ☐ Compute principal components via SVDs
- □ Compute PC components in python
- □Add PCA transforms as a pre-processing step to classification and regression
- ☐ Implement low-rank transforms for recommender systems



#### Outline

Why dimensionality reduction?

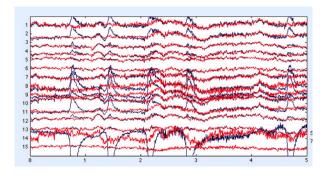
- ☐ Principal components and directions of variance
- □ Approximation with PCs
- ☐ Computing PCs via the SVD
- ☐ Face example in python
- ☐ Training models from PCs
- □ Low rank approximations and recommender systems





# High-Dimensional Data

- ☐ Many data sets have very high dimension
- ☐ Training can be difficult
  - Especially when number of samples is small
  - Classifier needs many parameters



**EEG** 

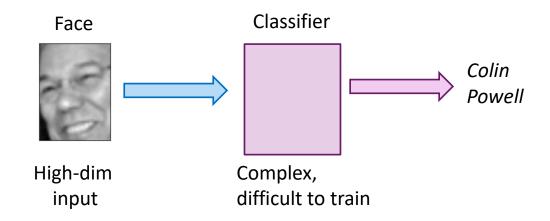
Ex: 32 channels x 1 kHz x 10s



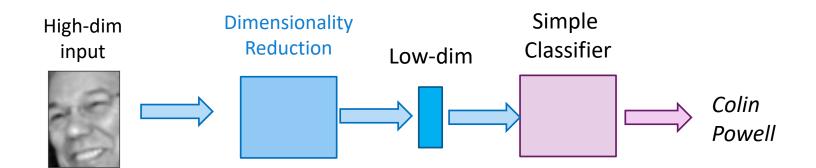
Face recognition with high-resolution images

# Problems with High-Dimensions

- □ Consider face recognition
- □Input is high-dimensional
  - Esp. for high resolution image
- ☐ Resulting classifier:
  - Requires many parameters
  - Difficult to train
  - Needs many samples
  - Computationally complex



#### **Dimensionality Reduction**



- □ Dimensionality reduction:
  - Reduce the input dimension to lower dimensional representation
- ☐ Can build simpler classifier
- □ Low-dimensional representational also good for:
  - Visualizing data
  - Clustering and other unsupervised tasks
  - Finding underlying structure of the data





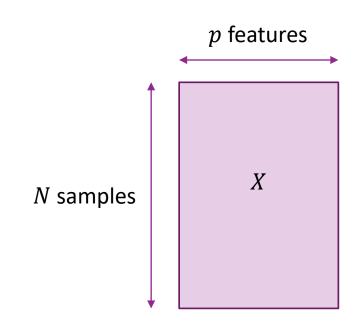
#### Outline

- □ Dimensionality reduction
- Principal components and directions of variance
- □ Approximation with PCs
- ☐ Computing PCs via the SVD
- ☐ Face example in python
- ☐ Training models from PCs
- □ Low rank approximations and recommender systems



#### **Data Definitions**

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



# Projections

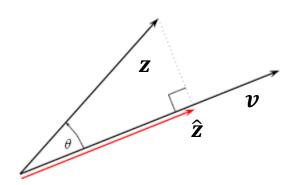
- □PCA reduces dimensionality by "projecting" data to a lower dim subspace
- $\square$  Projection: Given vectors z and v, the projection of z onto v is:

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

- $\circ \ \alpha = \text{coefficient of the projection}$
- $\square$ Theorem:  $Proj_{v}(z)$  is closest point in V to z:

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

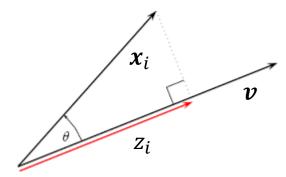
•  $V = {\alpha v | \alpha \in R}$  = vectors on the line spanned by v



#### Maximal Directional Variance

- $\square$  Given data:  $x_i$ , i=1,...,N and direction v with ||v||=1
- $\square$  Let  $z_i = v^T x_i$  = coefficient of the projection of  $x_i$  onto v
- $\square$  Sample mean and variance in direction  $oldsymbol{v}$  is :
  - Sample mean  $\bar{z} = \frac{1}{N} \sum_{i=1}^{N} z_i$
  - Sample variance  $s_z^2 = \frac{1}{N} \sum_{i=1}^{N} (z_i \bar{z})^2$

Problem: Find the direction v that maximizes the variance  $s_z^2$ 



- **□**Why?
  - Captures the most variation of the data
  - Provides the best vector for dimensionality reduction

# Sample Covariance Matrix

- Sample mean of the data:  $\overline{x} = \frac{1}{N} \sum_{i=1}^{N} x_i$
- $\square$ Sample covariance matrix: Matrix 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  $\ell$  in the dataset
- Matrix is  $p \times p$
- ☐ Sample covariance is given by

$$\mathbf{Q} = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i - \overline{\mathbf{x}}) (\mathbf{x}_i - \overline{\mathbf{x}})^T = \frac{1}{N} \widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}}$$

- $\circ$   $\widetilde{\textbf{\textit{X}}}$  = data matrix with sample mean removed (rows:  $\widetilde{\textbf{\textit{x}}}_i = \textbf{\textit{x}}_i \overline{\textbf{\textit{x}}}$  )
- Compute sample covariance via a matrix product

#### Sample Covariance and Directional Variance

- $\square$  Let  $z_i = v^T x_i$  = coefficient of the projection of  $x_i$  onto v
- lacksquare Can compute the sample mean and variance of  $z_i$  from  $\overline{x}$  and  $oldsymbol{Q}$
- ■Sample mean of the coefficients:

$$\bar{z} = \frac{1}{N} \sum_{i=1}^{N} z_i = \frac{1}{N} \sum_{i=1}^{N} v^T x_i = v^T \left[ \frac{1}{N} \sum_{i=1}^{N} x_i \right] = v^T \overline{x}$$

□Sample variance of the coefficients:

$$S_Z^2 = \frac{1}{N} \sum_{\substack{i=1 \ N}}^N (z_i - \overline{z})^2 = \frac{1}{N} \sum_{i=1}^N (v^T (x_i - \overline{x}))^2$$
$$= \frac{1}{N} \sum_{\substack{i=1 \ N}}^N v^T (x_i - \overline{x}) (x_i - \overline{x})^T v$$
$$= v^T Q v$$

## Maximizing Directional Variance

- $\square$  From previous slide: Directional variance  $s_z^2 = \frac{1}{N} \sum_{i=1}^{N} (z_i \bar{z})^2 = v^T Q v$
- lacktriangle Maximizing directional variance can be formulated as an optimization problem:  $\max_{m{v}} m{v}^T m{Q} m{v} \;\; ext{s.t.} \;\; \|m{v}\| = 1$
- lacksquare Let  $oldsymbol{v}_1$ , ...,  $oldsymbol{v}_p$  be the eigenvectors of  $oldsymbol{Q}:oldsymbol{Q}oldsymbol{v}_j=\lambda_joldsymbol{v}_j$
- □ Sort eigenvalues in descending order:  $\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_p$ 
  - Can show that eigenvalues are real and non-negative
- ☐ Theorem: Any local maxima of the variance directional is an eigenvector
  - $v = v_j$  for some j and  $v^T Q v = \lambda_j$
  - Proof below

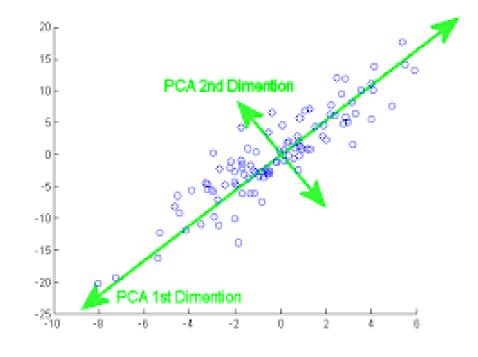
# Visualizing Principal Components

- $\square$  Principal components: The eigenvectors of  $\emph{\textbf{Q}}$ ,  $\emph{\textbf{v}}_1$ , ...,  $\emph{\textbf{v}}_p$ 
  - $\circ$  Always normalized  $\|v_i\|=1$
  - Sorted by eigenvalues  $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p$
  - Each vector is of dimension p
- ☐ Key property: Vectors are orthogonal

$$\circ \ \boldsymbol{v}_i^T \boldsymbol{v}_k = 0 \text{ if } j \neq k$$

- □ Represents directions of decreasing variance:
  - $v_1$ : PC 1 = Direction of max variance
  - $v_2$ : PC 2 = Direction of second most variance
  - $v_3$ : PC 3 = Direction of third most variance

0



# Proof PCs = Eigenvectors of Q

□PC constrained optimization problem:

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

- □ Define Lagrangian:  $L(v, \lambda) = v^T Q v \lambda [||v||^2 1]$
- ☐At any local maxima:

$$\frac{\partial L}{\partial v} = 0 \Rightarrow \mathbf{Q}v - \lambda v = \mathbf{0}$$

lacksquare This shows that  $oldsymbol{v}$  is an eigenvector of  $oldsymbol{Q}$ 

#### In-Class Exercise

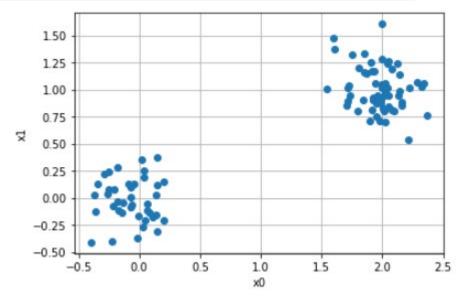
#### Exercise 1

We begin by showing how to compute and visualize PCs manually on a simple 2-dim synthentic data set. First, run the following code to generate 100 samples of synthetic data. Each data point has d=2 dimensions.

```
p = 0.5
std = 0.2
s = np.array([2,1])
d = 2
ns = 100

U = np.random.normal(0,std,(ns,d))
v = np.random.uniform(0,1,ns)
X = U + (v < p)[:,None]*s[None,:]</pre>
```

Create a scatter plot of the data of the two features of the data, X[:,0] and X[:,1].



#### Outline

- □ Dimensionality reduction
- ☐ Principal components and directions of variance
- Approximation with PCs
- ☐ Computing PCs via the SVD
- ☐ Face example in python
- ☐ Training models from PCs
- □ Low rank approximations and recommender systems

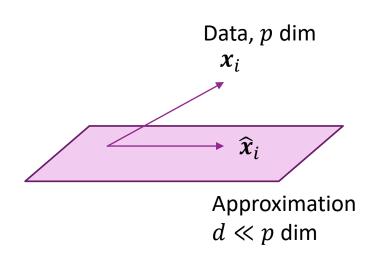


## Low-Dimensional Representations

- $\square$ Given data  $x_i$ , i = 1, ..., N. Each  $x_i \in \mathbb{R}^p$
- $\square$  Problem: Find basis vectors  $v_i$ , j = 1, ..., d such that:

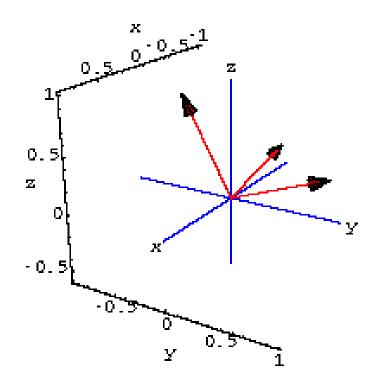
$$\mathbf{x}_i \approx \widehat{\mathbf{x}}_i = \overline{\mathbf{x}} + \sum_{j=1}^d \alpha_{ij} \mathbf{v}_j$$

- Sample mean + linear combination of basis vectors
- $lpha \ lpha_i = (lpha_{i1}, ..., lpha_{id})$  is an approximate coordinates of  $m{x}_i$  in basis  $(m{v}_1, ..., m{v}_d)$
- □ Dimensionality reduction:
  - $\circ$  If  $d \ll p$  we have represented  $oldsymbol{v}_i$  with a smaller number of coefficients.



#### Orthonormal Sets and Bases

- lacktriangle Definition: A set of vectors  $v_1, ..., v_d$  are an orthonormal set if:
  - $||v_j|| = 1$  for all j (unit length)
  - $\boldsymbol{v}_j^T \boldsymbol{v}_k = 0$  if  $j \neq k$  (perpendicular to one another)
- lacksquare Matrix form: If  $V = [v_1 \dots v_d]$ , then  $V^T V = I_d$
- $\square$  If d=p then  $oldsymbol{v_1},...,oldsymbol{v_p}$  is called an orthonormal basis
  - $\circ$   $\emph{V}$  is an orthogonal matrix
- ☐ Key property: the PCs form an orthonormal basis



#### Coefficients in an Orthonormal Basis

- $\square$  Suppose  $v_1, ..., v_p$  is an orthonormal basis
- $\Box$  Given a vector z, can write

$$oldsymbol{z} = \sum_{j=1}^p lpha_j oldsymbol{v}_j$$
 ,  $lpha_j = oldsymbol{v}_j^T oldsymbol{z}$ 

- Simple expression for computing coefficients in an orthonormal basis
- Matrix form:

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

# Approximating the Data Matrix

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

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

 $\square$  Approximation with d coefficients:

$$\widehat{\boldsymbol{x}}_i = \overline{\boldsymbol{x}} + \sum_{j=1}^d \alpha_{ij} \boldsymbol{v}_j$$

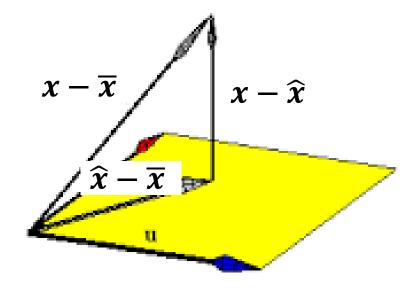


# Geometry of Approximations

- □ Approximation can be interpreted geometrically
- $\square$  Let V be set of all linear combinations

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

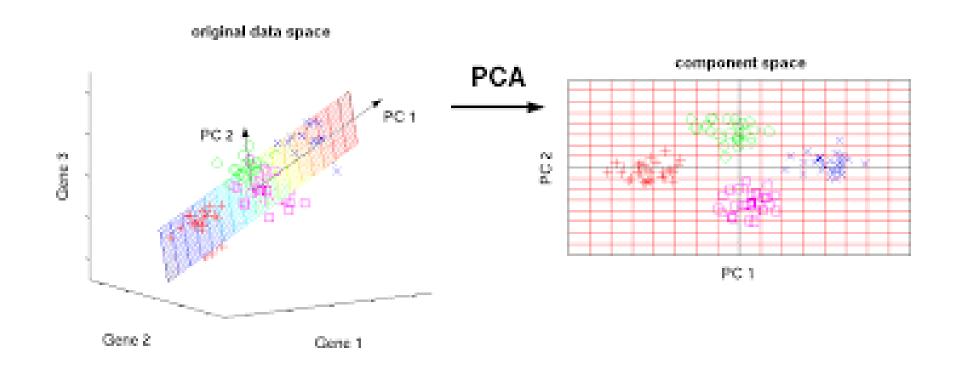
- ∘ *V* is a vector space
- $\circ$  Called the span of  $v_1, ..., v_d$
- $\square \widehat{x} \overline{x}$  is the closest vector in V to  $x \overline{x}$ 
  - Note the subtraction of the mean



Space spanned by  $v_1, \dots, v_d$ 

# Visualizing the Representation

☐ Finds a low-dimensional representation



## **Example Calculation**

#### □ Problem:

$$\circ$$
 Let  $v_1 = \frac{1}{\sqrt{2}}[1,1,0], v_2 = \frac{1}{\sqrt{6}}[1,-1,2]$ 

 $\circ$  Show  $v_1$  and  $v_2$  are orthogonal

#### □ Solution:

$$v_1^T v_1 = \frac{1}{2} (1^2 + 1^2 + 0^2)$$

$$v_2^T v_2 = \frac{1}{6} (1^2 + (-1)^2 + 2^2) = 1$$

$$v_1^T v_2 = \frac{1}{\sqrt{2(3)}} (1(1) + 1(-1) + 0(2)) = 0$$

## **Example Calculation Continued**

#### □ Problem:

- $\circ$  Let  $v_1 = \frac{1}{\sqrt{2}}[1,1,0], v_2 = \frac{1}{\sqrt{6}}[1,-1,2]$  be two PCs
- Let  $\overline{x} = [0,1,2]$  be the mean of the data
- Find the approximation of x = [2,4,4] with the two PCs

#### ■ Solution:

- Subtract mean:  $x \overline{x} = [2,3,2]$
- Coeff on PC1:  $\alpha_1 = v_1^T(x \overline{x}) = \frac{1}{\sqrt{2}}[2 + 3 + 0] = \frac{5}{\sqrt{2}}$
- Coeff on PC2:  $\alpha_2 = v_2^T (x \overline{x}) = \frac{1}{\sqrt{6}} [2 3 + 4] = \frac{3}{\sqrt{6}}$
- Approximation:  $\hat{x} = \overline{x} + \sum_{j=1}^{d} \alpha_j v_j = [0,1,2] + \frac{5}{2}[1,1,0] + \frac{3}{6}[1,-1,2] \approx [3,3,3]$



## Average Approximation Error

- $\square$  Let  $\widehat{x}_i$  = approximation with d PCs
- $\square$  Error in sample i:

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

 $\Box$ Theorem: Average error with a d PC approximation is:

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

 $\circ$  Sum of the smallest p-d eigenvalues

# Proportion of Variance (PoV)

☐ Total variance of data set:

$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \overline{x}||^2 = \sum_{j=1}^{p} \lambda_j$$

☐ Average approximation error:

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

 $\Box$  The proportion of variance explained by d PCs is:

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

 $\circ$  Measure of approximation error in using d PCs

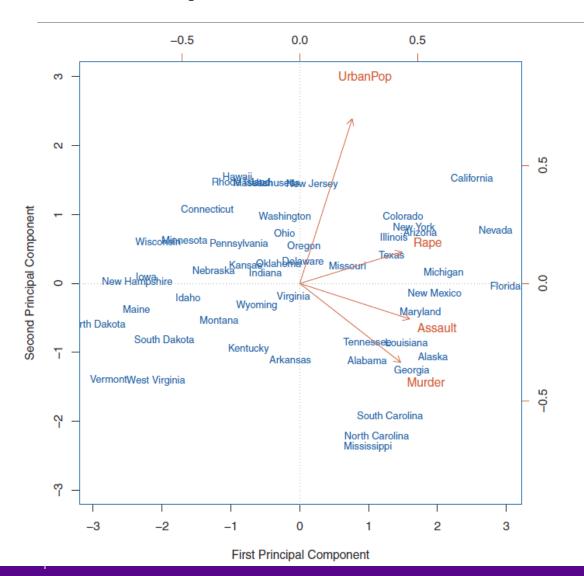
#### Example

PC index	$\lambda_i$	POV(i)
1	10	10/14.3≈ 0.70
2	4	14/14.3≈ 0.98
3	0.2	14.2/14.3≈ 0.99
4	0.1	14.3/14.3= 1

#### Latent Representations

- □ Each record is of the form:  $x_i \approx \overline{x} + \sum_{j=1}^d \alpha_{ij} v_j$
- $\square$  Variance in  $x_i$  explained by small number of "latent components"
  - $\circ$  Coefficients  $\alpha_{ij}$  are the latent representations of  $x_i$
- ■Example:
  - $x_i$  = 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

#### In-Class Exercise

#### Exercise 2: Computing the Approximation Error

We now verify the approximation errors. For each k:

- . Compute Xhat the PC approximation of X using k PC coefficients
- Compute the approximation error, err[k] the average error 1/ns sum\_{ij} (X[i,j]-Xhat[i,j])\*\*2)
- Compute the expected approximation error, err\_pred[k] the expected approximation error based on the eigenvalues lam.

Remember you will need to sort the eigenvalues in descending order. You can use the command

```
I = np.argsort(lam)[::-1]
```

#### Outline

- □ Dimensionality reduction
- ☐ Principal components and directions of variance
- □ Approximation with PCs
- Computing PCs via the SVD
  - ☐ Face example in python
  - ☐ Training models from PCs
  - □Low rank approximations and recommender systems



# Singular Value Decomposition

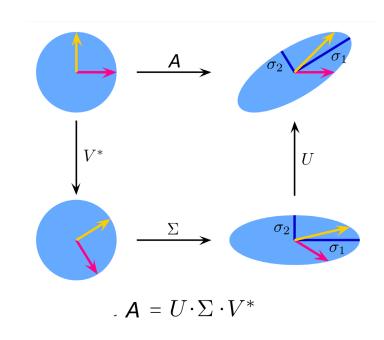
□SVD: Powerful method in linear algebra

#### $\square$ Given a matrix A:

- $\circ$  Decomposes the matrix into a product:  $A = USV^T$
- Provides orthonormal bases of the input and output spaces
- Multiplication of *A* is equivalent to scaling in that basis

#### For PCA:

- Identifies low rank subspaces for data
- Computes coefficients in that subspace



# Singular Value Decomposition Defined

- $\square$  Given matrix  $A \in \mathbb{F}^{n \times d}$ 
  - For PCA, this will be a scaled version of the data matrix
- $\square$ SVD is  $A = U\Sigma V^T$ , where
  - $\cdot$   $U \in \mathbb{F}^{n \times r}$ , columns are orthonormal
  - $\circ \ \emph{\textbf{V}} \in \mathbb{F}^{d imes r}$ , columns are orthonormal
  - $\circ$   $\Sigma = \operatorname{diag}(s_1, ..., s_r)$ , sorted  $s_1 \ge s_2 \ge ... \ge s_r \ge 0$ .
  - Called the singular values
- □ All matrices have an SVD
  - Matrices do not have to be square.
- Number of singular values  $r \le \min(n, d)$



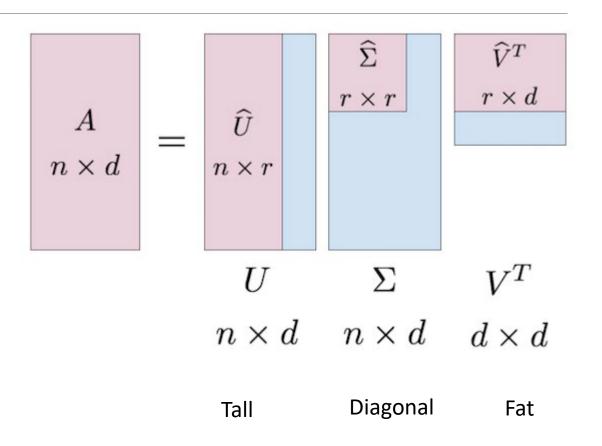
# Economy vs. Full SVD

- □Suppose  $A \in \mathbb{R}^{n \times d}$  with rank  $r \leq \min\{n, d\}$
- ☐ Two types of SVDs
- $\square$  Economy SVD:  $A = USV^*$ 
  - $U \in \mathbb{F}^{n \times r}$ , columns are orthonormal
  - $V \in \mathbb{F}^{d \times r}$ , columns are orthonormal
  - $\Sigma \in \mathbb{F}^{r \times r}$  diagonal  $\Sigma = diag(s_1, ..., s_r)$ ,
- $\square$  Full SVD:  $A = USV^*$ 
  - $\circ~U \in \mathbb{F}^{n \times n}$ , columns are an orthonormal basis of  $\mathbb{R}^n$
  - $V \in \mathbb{F}^{d \times d}$ , columns are an orthonormal basis of  $\mathbb{R}^d$
  - $\Sigma \in \mathbb{F}^{n \times d}$  with diagonal upper left  $\Sigma = \begin{bmatrix} \widehat{\Sigma} & 0 \\ 0 & 0 \end{bmatrix}$



#### **SVD** Visualized

- ☐ Pink matrices represent "economy" SVD
- ☐Blue represent "full SVD"



## Example

 $\Box$ Then can check that  $A = U\Sigma V^*$ 

$$\mathbf{U} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

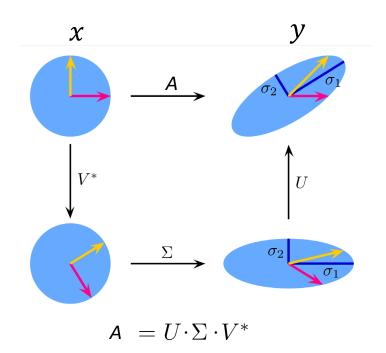
$$\Sigma = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\mathbf{U} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \qquad \mathbf{\Sigma} = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & \mathbf{0} & 0 \end{bmatrix} \qquad \mathbf{V}^* = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$

- Also verify that  $UU^* = I_5$  and  $VV^* = I_5$
- This can be found by (cleverly) permute the rows of A
- But, in general, use a computer to compute SVD

### Geometric Interpretation

- $\square$  Let  $A = U\Sigma V^*$  and y = Ax
- ☐ Consider a transformed space
  - $\mathbf{w} = \mathbf{V}^* \mathbf{x} = [w_1, ..., w_N]$  coefficients in input basis  $V = [v_1, ..., v_N]$
  - $\mathbf{z} = \mathbf{U}^* \mathbf{y} = [z_1, ..., z_M]$ : coefficients in output basis  $U = [u_1, ..., u_M]$
- □ Then:  $\mathbf{z} = \mathbf{\Sigma} \mathbf{w}$  so  $z_i = \sigma_i w_i$
- lacktriangle Each input direction  $oldsymbol{v}_i$  is mapped to  $\sigma_i oldsymbol{u}_i$
- ☐Consequence:
  - SVD finds orthonormal bases *U*, *V* such that matrix *A* is a linear scaling in each basis vector



### Example Problem

- □Suppose that  $A = U\Sigma V^* \in \mathbb{R}^{3\times 4}$  with  $\Sigma = diag(3,0.2,0,0)$
- $\Box$  If  $x = 2v_1 + 3v_2 + 4v_3 + 5v_4$  find y = Ax in terms of basis  $u_1, u_2, u_3$
- Solution:
  - $\mathbf{v} \cdot \mathbf{A} \mathbf{v}_i = \sigma_i \mathbf{u}_i$  for all i
  - Therefore,

$$y = Ax = 2Av_1 + 3Av_2 + 4Av_3 + 5Av_4$$
  
= 2(3) $u_1 + 3(0.2)u_2 + 4(0)u_3$   
=  $6u_1 + 0.6 u_2$ 



### Computing the SVD in Python

☐ Random matrix

```
# Create some random matrix
A = np.random.normal(0,1,(100,10))
```

☐ Full SVD

```
# Full SVD
U,s,Vtr = np.linalg.svd(A)
```

```
A.shape = (100, 10)
U.shape = (100, 100)
s.shape = (10,)
Vtr.shape = (10, 10)
```

☐ Economy SVD

```
# Economy SVD
U,s,Vtr = np.linalg.svd(A, full_matrices=False)
U.shape = (100, 10)
s.shape = (10,)
Vtr.shape = (10, 10)
```

☐ Reconstruction:

```
# Recovers back A
Ahat = (U*s[None,:]).dot(Vtr)
```

## Computing the PCA via SVD

- Let  $A = \frac{1}{\sqrt{N}}\widetilde{X}$  = scaled data matrix with sample mean removed.
- $\Box$  Take SVD:  $A = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$
- ☐ Properties:
  - Sample covariance matrix is  $Q = \frac{1}{N} \widetilde{X}^T \widetilde{X} = A^T A = V \Sigma \mathbf{U}^T U \Sigma V^T = V \Sigma^2 V^T$
  - $\circ$  Eigenvalues of  $\mathbf{Q}=$  squared singular values of  $\mathbf{A}$
  - $\circ$  PCs are  $v_i$ , columns of V
  - $\circ$  Coefficients are  $oldsymbol{Z} = \widetilde{oldsymbol{X}} oldsymbol{V} = \sqrt{N} oldsymbol{A} oldsymbol{V} = \sqrt{N} oldsymbol{U} oldsymbol{\Sigma}$
- $\square$  Hence, SVD provides PCs, eigenvalues coefficients Z in the PCA representation.



#### **In-Class Exercise**

#### Exercise 3: Computing the PCA via the SVD ¶

- Compute the matrix A = 1/np.sqrt(ns)\*Xm
- Compute the economy SVD of A with the np.linalg.svd() command. Use full\_matrices=False option. Print the dimensions of the components of the SVD



#### Outline

- ☐ Dimensionality reduction
- ☐ Principal components and directions of variance
- □ Approximation with PCs
- ☐ Computing PCs via the SVD
- Face recognition using PCA in python
- ☐ Training models from PCs
- □ Low rank approximations and recommender systems



### Example: Face Recognition



Labeled Faces in the Wild Home



- ☐ Face recognition challenges:
  - Face images can be high-dimensional
  - We will use 50 x 37 = 1850 pixels
- □ Applying PCA:
  - Should be few degrees of freedom
  - Can transform to lower dimensional representations
- ☐ 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 20 years ago



## Loading the Data

- ■Built-in routines to load data is sciket-learn
- ☐ Can take several minutes the first time (Be patier

```
Image size = 50 \times 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 14:22:13,640 Loading face #00001 / 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 PCA

```
☐ Manually compute the PCs with SVD
 npix = h*w

    Remove the mean

 Xmean = np.mean(X,0)

    Use broadcasting

 Xs = X - Xmean[None,:]
U,S,Vtr = np.linalg.svd(Xs, full matrices=False)
                                                           Compute the SVD
                                                       ☐ Use sklearn builtin PCA function
from sklearn.decomposition import PCA

    Construct a PCA object

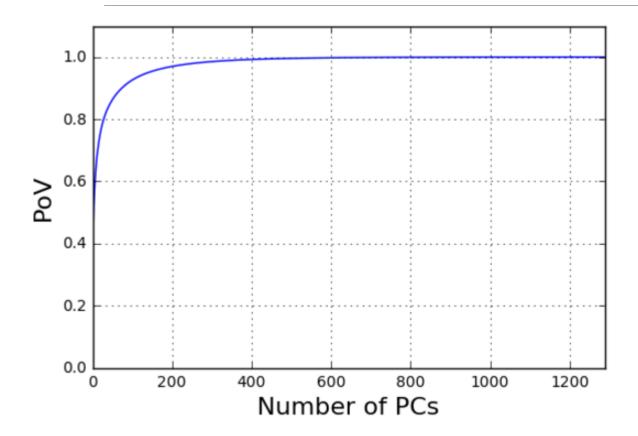
# Construct the PCA object
pca = PCA(n components=ncomp,
         svd solver='randomized', whiten=True)
                                                         Call fit: Computes mean and PC components

    Stores values internally in the pca class

# Fit the PCA components on the entire dataset
```

pca.fit(X)

# 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
                     # True=Use sklearn reconstruction, else use SVD
use pca = True
# Loop over figures
iplt = 0
for ind in inds:
   for d in ds:
       plt.subplot(nplt,nd+1,iplt+1)
       if use pca:
           # Zero out coefficients after d.
                                                                                □ Reconstruction using sklearn method
           # Note, we need to copy to not overwrite the coefficients
           Zd = np.copy(Z[ind,:])

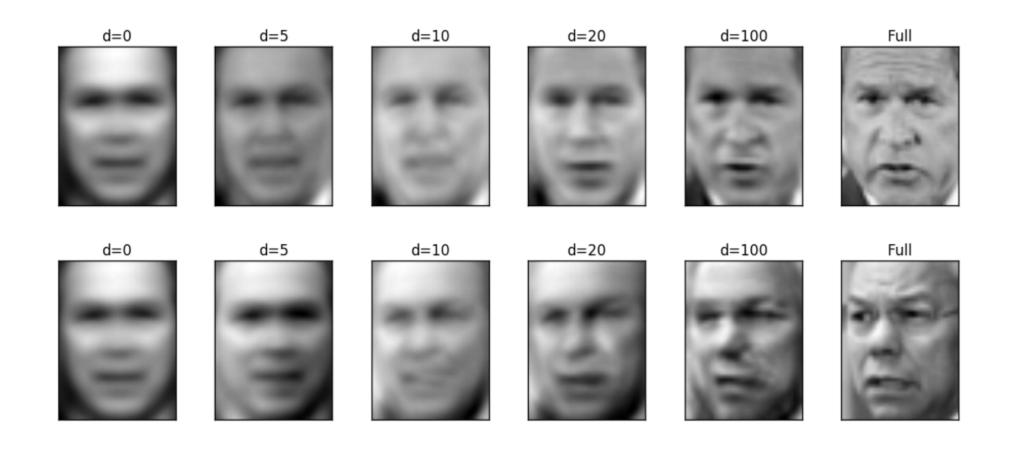
    Uses the inverse transform method

           Zd[d:] = 0
           Xhati = pca.inverse transform(Zd)
       else:
           # Reconstruct with SVD
          Xhati = (U[ind,:d]*S[None,:d]).dot(Vtr[:d,:]) + Xmean
                                                                                Reconstruction using SVD
       plt face(Xhati)
       plt.title('d={0:d}'.format(d))
       iplt += 1

    Note use of broadcasting

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

# Plotting the Approximations





# Plotting the PCs

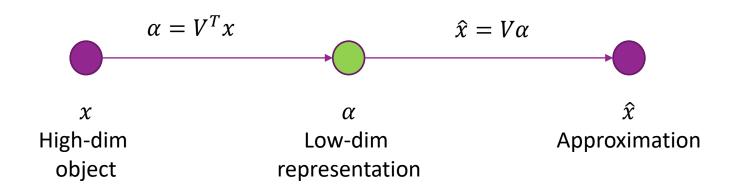
☐ The PCs can be plotted as well





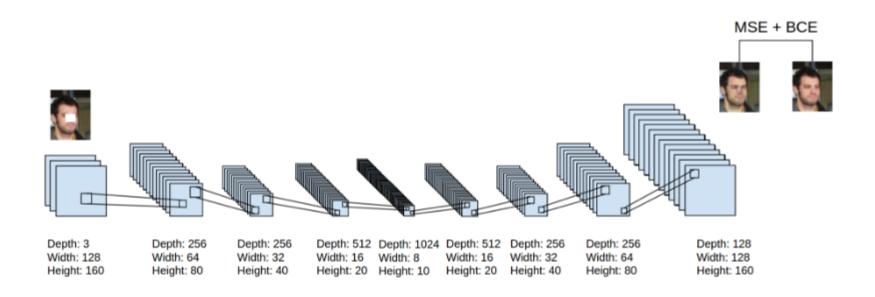
#### State-of-the-Art: Auto-Encoders

- □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/
  - <a href="https://swarbrickjones.wordpress.com/2016/01/13/enhancing-images-using-deep-convolutional-generative-adversarial-networks-dcgans/">https://swarbrickjones.wordpress.com/2016/01/13/enhancing-images-using-deep-convolutional-generative-adversarial-networks-dcgans/</a> (Code in Theano not tensorflow)

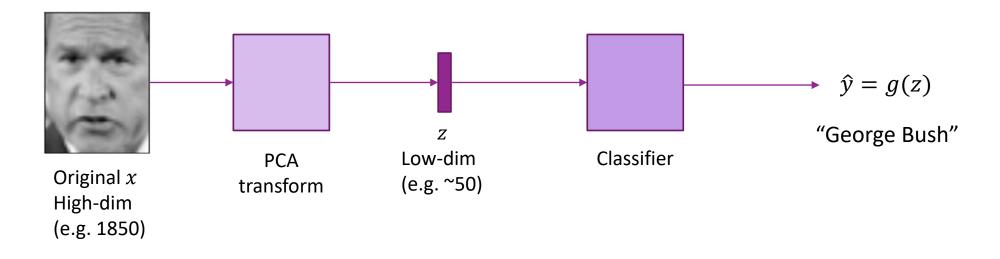


#### Outline

- □ Dimensionality reduction
- ☐ Principal components and directions of variance
- □ Approximation with PCs
- ☐ Computing PCs via the SVD
- ☐ Face example in python
- Training models from PCs
  - □ Low rank approximations and recommender systems



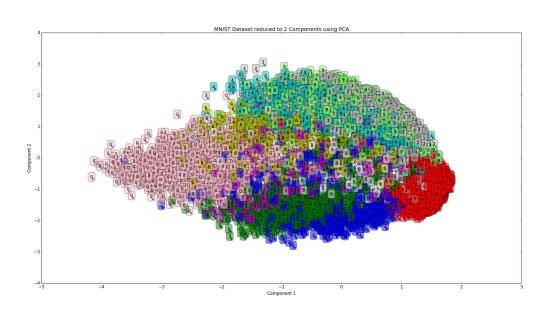
## Classification Using PCs



- $\blacksquare$  Many problems: Dimensionality of data x is too large
  - Classifier in original space will have too many parameters
- ☐ Key idea:
  - Learn a dimension reducing transform via PCA: z = f(x)
  - Train classifier on low-dim transform  $\hat{y} = g(z)$



## Why This Would Work?



- □PCA works if: classes are separable in transformed domain
- ☐ Example to right:
  - MNIST digits plotted in two PCs
  - Can mostly separate the classses



## Training and Testing

- $\square$ Split data in training and test:  $X_{tr}, y_{tr}, X_{ts}, y_{ts}$
- $\square$  Fit PCA transform on Z = g(X) on training data  $X_{tr}$ 
  - Do not include test data in PCA fit!
  - Many students make this mistake.
- ☐ Transform training and test:

$$Z_{tr} = g(X_{tr}), Z_{ts} = g(X_{ts})$$

- $\square$  Fit classifier  $\hat{y} = f(z)$  on transformed training data  $(Z_{tr}, y_{tr})$
- $\square$  Predict classifier on transformed test data:  $\hat{y}_{ts} = f(Z_{ts})$
- $\square \text{Score error rate / MSE on test data: } \epsilon = \frac{1}{N} \# \{ \hat{y}_{ts}^i \neq y_{ts}^i \}$



#### **Cross-Validation**

- ☐ To find number of PCs and other parameters use cross-validation
- $\square$ Split data in training and test:  $X_{tr}, y_{tr}, X_{ts}, y_{ts}$
- ☐ For each set of parameters:
  - Fit PCA transform on Z = g(X, numPCs) on training data  $X_{tr}$
  - Transform training and test:  $Z_{tr} = g(X_{tr})$ ,  $Z_{ts} = g(X_{ts})$
  - Fit classifier  $\hat{y} = f(z)$  on transformed training data  $(Z_{tr}, y_{tr})$
  - Predict classifier on transformed test data:  $\hat{y}_{ts} = f(Z_{ts})$
  - Score (e.g. error rate / MSE) on test data:  $\epsilon = \frac{1}{N} \# \{ \hat{y}_{ts}^i \neq y_{ts}^i \}$
- Select the parameters with lowest score

### Example: SVM classification with PCAs

```
npc_test = [25,50,75,100,200]
gam test = [1e-3,4e-3,1e-2,1e-1]
                                                                        ☐ Parameters to search
C = 100
n0 = len(npc test)

    Number of PCs and gamma

n1 = len(gam test)
acc = np.zeros((n0,n1))
acc max = 0
for i0, npc in enumerate(npc test):
                                                                        ☐ Fit on the training data.
   # Fit PCA on the training data
   pca = PCA(n components=npc, svd solver='randomized', whiten=True)
                                                                          This is in the loop!
   pca.fit(Xtr)
   # Transform the training and test
                                                                        ☐ Transform the data
   Ztr = pca.transform(Xtr)
   Zts = pca.transform(Xts)
   for i1, gam in enumerate(gam test):
       # Fiting on the transformed training data
                                                                        ·□Fit classifier on transformed training data
       svc = SVC(C=c, kernel='rbf', gamma = gam)
       svc.fit(Ztr, ytr)
       # Predict on the test data
                                                                        ☐ Test on the transformed test data
       yhat = svc.predict(Zts)
       # Compute the accuracy
                                                                        □ Score on test data
       acc[i0,i1] = np.mean(yhat == yts)
       print('npc=%d gam=%12.4e acc=%12.4e' % (npc,gam,acc[i0,i1]))
```

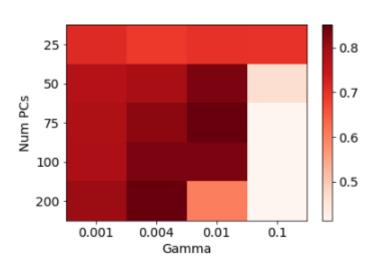
#### Example: Parameter Search

#### ■Search over:

- Number of PCs  $\in \{25,50,75,100,200\}$
- $\gamma \in \{0.001, 0.004, 0.01, 0.1\}$

- ☐ Plotted is the test accuracy
- □ Best test accuracy  $\approx 85\%$

#### **Test Accuracy**



Optimal num PCs = 75 Optimal gamma = 0.010000

## Examples

□Correct images

George W Bush George W Bush





George W Bush George W Bush





Original Reduced

☐ Error images





Gerhard Schroeder George W Bush





Original

Reduced

#### Outline

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

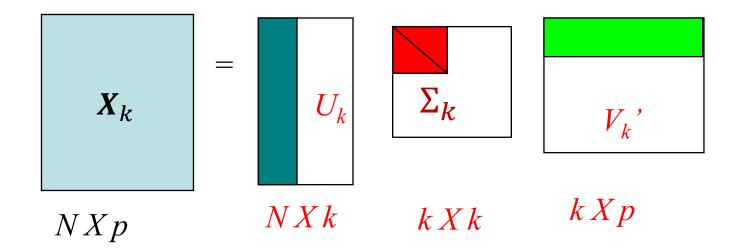
Low rank approximations and recommender systems



## Low-Rank Approximations

- □SVD can be used for a low-rank approximation
- $\square$ SVD can be written:  $X = \mathbf{U} \Sigma \mathbf{V}^T = \sum_{j=1}^r \alpha_j \mathbf{u}_j \mathbf{v}_j^T$
- $\square$  Consider k —term approximation:  $X_k = \sum_{j=1}^k \alpha_j u_j v_j^T$
- ☐Properties:
  - $\circ$   $X_k$  is rank k
  - $\circ X_k = U_k \Sigma_k V_k^T$
  - Error is  $||X X_k||_F^2 = \sum_i \sum_j (X_{ij} X_{k,ij})^2 = \sum_{j=k+1}^r \alpha_j^2$
  - $\circ$  If  $s_{k+1}, ..., s_r$  is small then matrix is well approximated by rank k matrix

## Low-Rank Approximation Visualized



 $\square$  Can show: Reconstructed matrix  $X_k$  is optimal rank k approximation

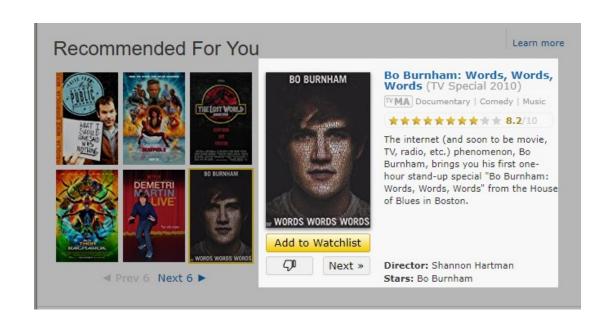


#### Recommender Systems

- ☐ How do you recommend a movie to a user?
- MovieLens dataset:
  - Get past ratings from users
  - Make recommendations for future

t[3]:

movield		title	genres	
0	1	Toy Story (1995)	Adventure   Animation   Children   Comedy   Fantasy	
1	2	Jumanji (1995)	Adventure Children Fantasy	
2	3	Grumpier Old Men (1995)	Comedy Romance	
3	4	Waiting to Exhale (1995)	Comedy Drama Romance	
4	5	Father of the Bride Part II (1995)	Comedy	

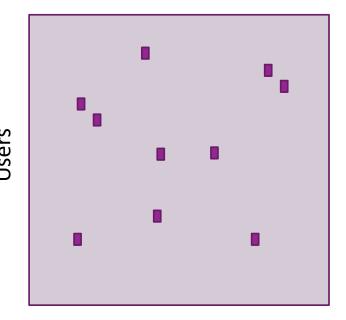


#### Ratings Matrix

- □ Data can be represented as ratings matrix
  - Users x movies
- ☐ Problem: Most users have only rated a small fraction
- Need to estimate unseen entries
  - Very sparse
- ☐ How can we do complete this matrix

Name	Dates	Users	Movies	Ratings	Density
ML Latest	'95 <b>–</b> '16	247,753	34,208	22,884,377	0.003%
ML Latest Small	'96 <b>–</b> '16	668	10,329	105,339	0.015%

#### Movies





### Latent Factor Model for Ratings

- □ Idea: Ratings for movies dependent on small number of latent factors
  - E.g. Action, famous actors, genre, ...
- ☐ Mathematically model as:

$$R_{ij} \approx \widehat{R}_{ij} = b_i^u + b_j^m + \sum_{k=1}^K A_{ik} B_{jk}$$

- $R_{ij}$  =Rating of movie j by user i
- $b_i^u = \text{Bias of user } i$
- $b_i^m = \text{Bias of movie } j$
- K = number of latent factors. Typically small  $K \ll N_{user}$ ,  $N_{movies}$
- $A_{ik}$  = Preference of user i to factor k
- $B_{jk} = \text{Component of factor } k \text{ in movie } j$



#### More to be added

- ☐ These slides are still under construction.
- ☐ More will be added on low rank approximations and embedding layers.

