

Updates

- ◆ HW 5 published
- ◆ Midterm regrades: later this week

Unsupervised Learning

- ◆ **Spectral methods**
 - Eigenvector/singular vector decomposition (SVD)
 - PCA, CCA
- ◆ **Reconstruction methods**
 - PCA, ICA, auto-encoders
- ◆ **Clustering and Probabilistic methods**
 - K-means
 - Gaussian mixtures
 - Latent Dirichlet Allocation (LDA)

PCA

Lyle Ungar

Learning objectives

PCA as change of basis

PCA minimizes reconstruction error

PCA maximizes variance

PCA relation to eigenvalues/vectors

PCR: PCA for feature creation

Based in part on slides by Jia Li
(PSU) and Barry Slaff (Upenn)

PCA

- ◆ Express a vector x in terms of coefficients on an (orthogonal) basis vector (eigenvectors v_k)

$$x_i = \sum_k z_{ik} v_k$$

- We can describe how well we approximate x in terms of the eigenvalues

- ◆ PCA is used for dimensionality reduction

- visualization
- semi-supervised learning
- eigenfaces, eigenwords, eigengrasps

PCA

- ◆ Express a vector x in terms of coefficients on an (orthogonal) basis vector (eigenvectors v_k)

$$x_i = \sum_k z_{ik} v_k$$

- ◆ Find z_{ik} by projection

$$x_i \cdot v_j = \sum_k z_{ik} v_k \cdot v_j$$

$$x_i \cdot v_j = z_{ij}$$

PCA

- ◆ PCA can be viewed as

- minimizing distortion $\|X - ZV^T\|_F$
 - Or the square of the above: $\sum_i \|x_i - \sum_k z_{ik}v_k\|_2^2$
 - Note that either definition gives the same result
- A rotation to a new coordinate system to maximize the variance in the new coordinates

- ◆ Generally done by mean centering first

- Sometimes standardize

Nomenclature

$$\mathbf{X} = \mathbf{Z}\mathbf{V}^T$$

- ◆ \mathbf{Z} ($n \times k$)
 - principal component **scores**
- ◆ \mathbf{V} ($m \times k$)
 - **Loadings**
 - Principal component **coefficients**
 - Principal components

In PCA world, \mathbf{X} is $n \times m$

PCA minimizes Distortion

- ◆ First subtract off the average $\bar{\mathbf{x}}$ from all the \mathbf{x}_i
 - From here, we'll assume this has been done
- ◆ Approximate \mathbf{x} in terms of an orthonormal basis \mathbf{v}
 - $\hat{\mathbf{x}}_i = \sum_k z_{ik} \mathbf{v}_k$ or $\mathbf{X} = \mathbf{ZV}^T$
- ◆ Distortion (this is the square of the earlier definition)

$$\sum_{i=1}^n \|\mathbf{x}^i - \hat{\mathbf{x}}^i\|_2^2 = \sum_{i=1}^n \sum_{j=1}^m (x_j^i - \hat{x}_j^i)^2.$$

PCA minimizes distortion

$$\begin{aligned}\text{Distortion}_k &: \sum_{i=1}^n \sum_{j=k+1}^m \mathbf{u}_j^\top (\mathbf{x}^i - \bar{\mathbf{x}})(\mathbf{x}^i - \bar{\mathbf{x}})^\top \mathbf{u}_j \\ &= \sum_{j=k+1}^m \mathbf{u}_j^\top \left(\sum_{i=1}^n (\mathbf{x}^i - \bar{\mathbf{x}})(\mathbf{x}^i - \bar{\mathbf{x}})^\top \right) \mathbf{u}_j \\ &= n \sum_{j=k+1}^m \mathbf{u}_j^\top \Sigma \mathbf{u}_j = n \sum_{j=k+1}^m \lambda_j\end{aligned}$$

See the course wiki!

PCA maximizes variance

$$\begin{aligned}\text{Variance}_k &: \sum_{i=1}^n \sum_{j=1}^k (\mathbf{u}_j^\top \mathbf{x}^i - \mathbf{u}_j^\top \bar{\mathbf{x}})^2 \\ &= \sum_{j=1}^k \mathbf{u}_j^\top \left(\sum_{i=1}^n (\mathbf{x}^i - \bar{\mathbf{x}})(\mathbf{x}^i - \bar{\mathbf{x}})^\top \right) \mathbf{u}_j \\ &= n \sum_{j=1}^k \mathbf{u}_j^\top \Sigma \mathbf{u}_j.\end{aligned}$$

See the course wiki!

PCA - Summary

$$\hat{\mathbf{x}}^i = \mathbf{x}^i = \bar{\mathbf{x}} + \sum_{j=1}^m z_j^i \mathbf{u}_j$$

$$\text{Variance}_k + \text{Distortion}_k = n \sum_{j=1}^m \lambda_j$$

See the course wiki!

Principal Component Analysis

$$\mathbf{X} \rightarrow \mathbf{X}_c = \mathbf{UDV}^T = \mathbf{ZV}^T$$

\mathbf{X}_c is $(n \times p)$, \mathbf{Z} is $(n \times p)$, \mathbf{V} is $(p \times p)$.

\mathbf{Z} is the transformation of \mathbf{X} into “PC space”

Column vector \mathbf{z}_i is the i'th *PC score vector*.

Column vector \mathbf{v}_i is the i'th *PC direction or loading*.

Since \mathbf{V} is orthogonal, $\mathbf{X}_c \mathbf{V} = \mathbf{ZV}^T \mathbf{V} = \mathbf{Z}$, and therefore:

$$\mathbf{z}_i = \mathbf{X}_c \mathbf{v}_i = \mathbf{u}_i D_{ii}$$

Hence \mathbf{z}_i is the projection of the row vectors of \mathbf{X}_c on the (unit) direction \mathbf{v}_i , scaled by D_{ii} .

Principal Component Analysis

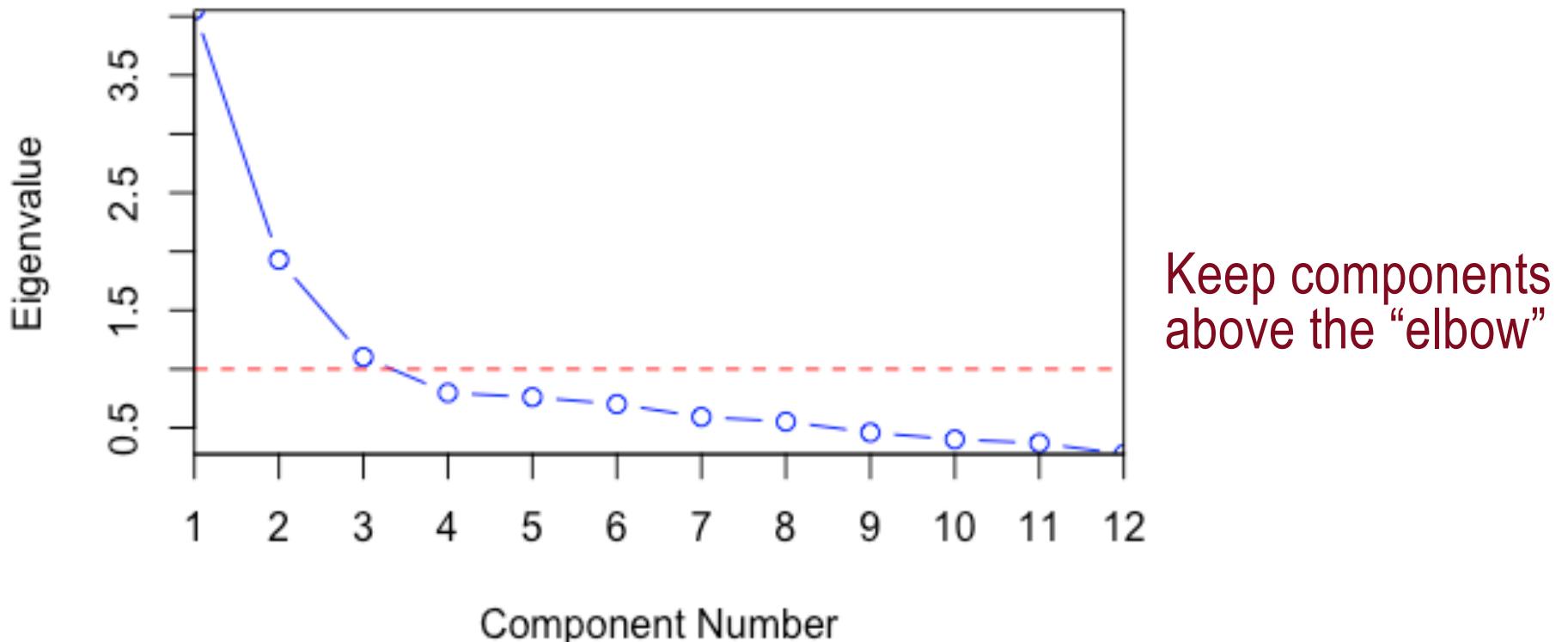
$$\mathbf{X} \rightarrow \mathbf{X}_c = \mathbf{UDV}^T = \mathbf{ZV}^T$$

$$\mathbf{X}_c^T \mathbf{X}_c = \sum_{i=1}^p (D_{ii})^2 \mathbf{v}_i \mathbf{v}_i^T$$

“% Variance explained by the i'th principal component:”

$$= 100 \cdot \frac{(D_{ii})^2}{\sum_{j=1}^p (D_{jj})^2} = 100 \lambda_i / \sum_i \lambda_i$$

Scree plot



[https://en.wikipedia.org/
wiki/Scree plot](https://en.wikipedia.org/wiki/Scree_plot)

PCA

True or false:

If X is any matrix, and X has singular value decomposition $X = UDV^T$

then the principal component scores for X are the columns of

$$Z = UD$$

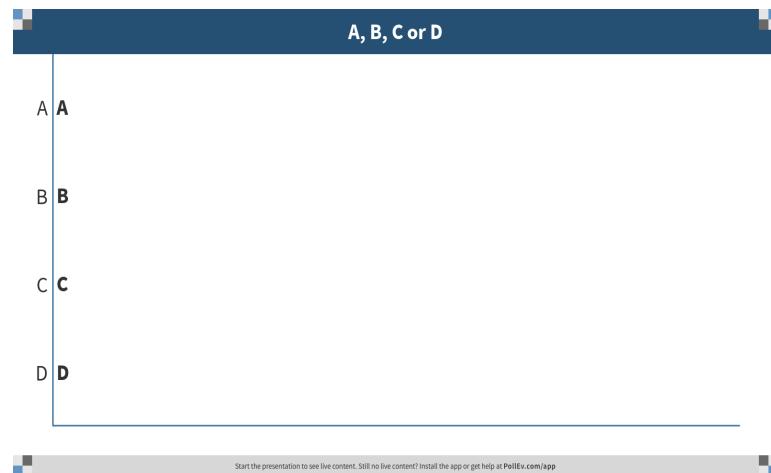


- a) True
- b) False

PCA

If X is mean-centered, then PCA finds...?

- (a) Eigenvectors of $X^T X$
- (b) Right singular vectors of X
- (c) Projection directions of maximum covariance of X
- (d) All of the above



PCA: Reconstruction Problem

PCA can be viewed as an L₂ optimization, minimizing distortion, the reconstruction error.

$$Z^*, V^* = \underset{\substack{Z \in \mathbb{R}^{n \times k}, V \in \mathbb{R}^{p \times k}, \\ v_i^T v_j = \delta_{ij}}}{\operatorname{argmin}} |X_c - ZV^T|_F$$

Here we have constrained Z , V by dimension:

X_c is still $(n \times p)$.

Z is $(n \times k)$, with $k \leq p$.

V is $(p \times k)$.

If $k=p$ then the reconstruction is perfect. $k < p$, not.

PCA via SVD

◆ $X = ZV^T = UDV^T$

- X n x p U n x k D k x k V^T k x p

◆ $Z = UD$ - component scores or "factor scores"

- the transformed variable values corresponding to a particular data point

◆ V^T - loadings

- the weight by which each standardized original variable should be multiplied to get the component score

PCA via SVD

- ◆ $x_i = \sum_k z_{ik} v_k$
- ◆ What is z_{ik} ?
 - $x_i = \sum_k u_{ik} d_{kk} v_k$

Sparse PCA

- ◆ $\operatorname{argmin}_{Z,V} \|X - ZV^T\|_F$
 - $v_i'v_j = \delta_{ij}$ (orthonormality)
 - Eigenvectors give the optimal solution
- ◆ you can add an L_1 penalty
- ◆ $\operatorname{argmin}_{Z,V} \|X - ZV^T\|_F + \lambda_1 \|Z\|_1 + \lambda_2 \|V\|_1$
 - No longer eigenvectors
 - Convex in Z given V and in V given Z
 - Solve by alternating gradient descent

What you should know

- ◆ PCA as minimum reconstruction error ('distortion')
- ◆ PCA as finding direction of maximum covariance
- ◆ Sensitivity of PCA to standardizing
- ◆ Nomenclature: scores, coefficients/loadings
- ◆ Coming next: autoencoders, eigenfaces, eigenwords