

# Data Mining

## S4240 Section 001

Giovanni Motta

Columbia University, Department of Statistics

September 23, 2015

# Outline

Today:

1. High dimensional data
2. Principal components analysis (PCA) overview
3. (Review: eigenvalues and eigenvectors)
4. PCA computation

Next time: more PCA, doing PCA with R

# High-Dimensional Data

**Problem:** want to describe student performance and compare students

Data ( $n = 175$  students,  $p = 8$  scores):

- ▶ scores on six homeworks
- ▶ midterm score
- ▶ final score

How can I summarize data for a student in a way that will be useful for comparisons?

# Higher Dimensional Data

If data *truly* high dimensional:

- ▶ scores for each of 7 items have low correlation  $\leftrightarrow$  need all items to summarize student
- ▶ example: good on HWs 1,2,3, middling on midterm, poor on HW 4, 5, and 6, great on final

But: we would not expect to see a student like the one above.  
Why?

Any suggestions for summarizing student performance in a way that is useful for comparing students? Why would this be reasonable?

# Dimensionality Reduction

Data with  $n$  observations and  $p$  dimensions:

$$\mathbf{X}_{n \times p} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ x_{21} & \dots & x_{2p} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p]$$

Here  $\mathbf{x}^\top$  means the transpose of  $\mathbf{x}$ .

Today (and often, but not always), **bold** means matrix or vector, and *plain* means scalar.

Note that  $\mathbf{X} \in \mathbb{R}^{n \times p}$

# Dimensionality Reduction : $\kappa < p$

Two approaches:

- ▶ **Feature Selection:** choose a subset of features for prediction

$$\begin{array}{ccc} [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p] & \rightarrow & [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_\kappa] \\ p \text{ variables} & & \kappa \text{ variables} \end{array}$$

But you need to know your prediction task before you do feature selectionn.

- ▶ **Feature Extraction:** create new features by combining existing ones

$$\begin{aligned} [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p] &\rightarrow [f_1(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p), \dots, f_\kappa(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p)] \\ &= [\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_\kappa] \end{aligned}$$

Feature extraction is often a data preprocessing step since you do not need to know which prediction methods you will use.

# Dimensionality Reduction

Today, we will look at *linear feature extraction*

$$\begin{bmatrix} \mathbf{X}_1^\top \\ \mathbf{X}_2^\top \\ \vdots \\ \mathbf{X}_j^\top \\ \vdots \\ \mathbf{X}_p^\top \end{bmatrix}_{p \times n} \rightarrow \begin{bmatrix} \mathbf{Y}_1^\top \\ \mathbf{Y}_2^\top \\ \vdots \\ \mathbf{Y}_\kappa^\top \end{bmatrix}_{\kappa \times n} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1j} & \dots & w_{1p} \\ w_{21} & w_{22} & \dots & w_{2j} & \dots & w_{2p} \\ \vdots & \vdots & & \vdots & & \vdots \\ w_{\kappa 1} & w_{\kappa 2} & \dots & w_{\kappa j} & \dots & w_{\kappa p} \end{bmatrix}_{\kappa \times p} \begin{bmatrix} \mathbf{X}_1^\top \\ \mathbf{X}_2^\top \\ \vdots \\ \mathbf{X}_j^\top \\ \vdots \\ \mathbf{X}_p^\top \end{bmatrix}_{p \times n}$$

$$\mathbf{Y}_{\kappa \times n}^\top = \mathbf{W}_{\kappa \times p}^\top \mathbf{X}_{p \times n}^\top \quad \text{or} \quad \mathbf{Y}_{n \times \kappa} = \mathbf{X}_{n \times p} \mathbf{W}_{p \times \kappa}$$

Terminology:

- ▶  $\mathbf{Y}_1, \dots, \mathbf{Y}_\kappa$  are the *scores*,
- ▶  $\mathbf{w}_1, \dots, \mathbf{w}_\kappa$  are the *loadings*.

# Linear Feature Extraction

Student score averages are a form of linear feature extraction:

$$y_{i1} = [0.05, 0.05, 0.05, 0.05, 0.05, 0.05, 0.3, 0.4] \begin{bmatrix} HW1 \\ HW2 \\ HW3 \\ HW4 \\ HW5 \\ HW6 \\ Midterm \\ Final \end{bmatrix}$$

The student score average is the *score* and the assignment weights are the *loadings* (here  $p = 8$  and  $\kappa = 1$ ).



# Linear Feature Extraction

How can we find a good set of loadings and scores for a general data set?

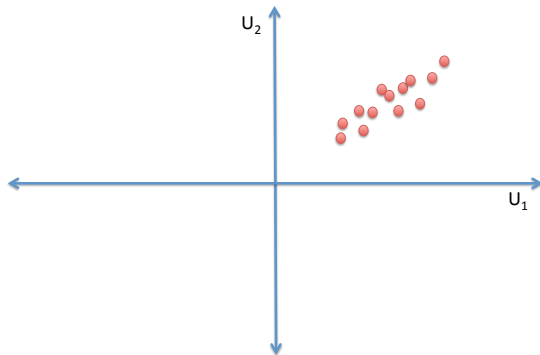
Principal components analysis (PCA): a covariance matrix singular value decomposition method for unsupervised data

Principal components: a set of linearly uncorrelated variables—these will be the loadings

Multiply loadings with original data to get scores!

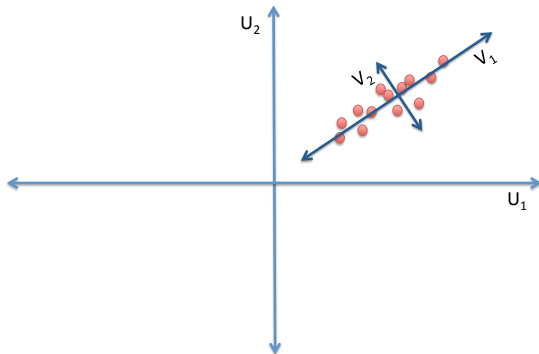
# Principal Components Analysis

Basic idea: consider a dataset



# Principal Components Analysis

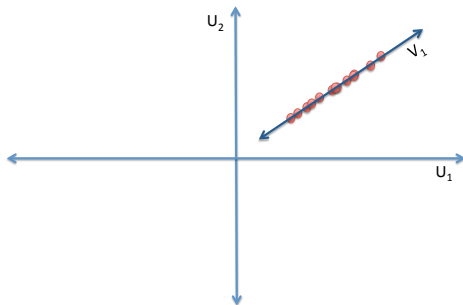
Basic idea: find a rotation that best describes data



A coordinate (linear) transformation from  $U$  to  $V$ !

# Principal Components Analysis

Basic idea: given linear transformation, we can throw out the less descriptive dimensions and still have a decent representation



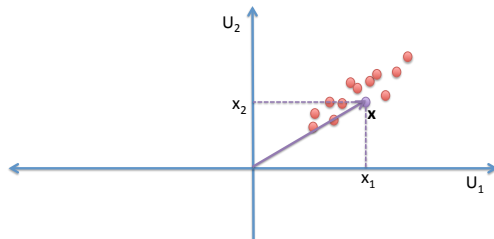
Consider the tree heights and weights. What is an interesting dimension? What is less descriptive?

# Principal Components Analysis

Mathematically, how do we do this?

Simple case:  $\mathbf{x} \in \mathbb{R}^2$  (i.e.,  $p = 2$ ), and we want a good projection  $y \in \mathbb{R}$  (i.e.,  $\kappa = 1$ )

- ▶  $\mathbf{x} = [x_1, x_2]^\top$ , where  $x_1, x_2$  are scalars
- ▶ the axes vectors  $\mathbf{u}_1$  and  $\mathbf{u}_2$  are orthonormal
- ▶  $\mathbf{x} = x_1 \mathbf{u}_1 + x_2 \mathbf{u}_2 = \mathbf{I}_2 \mathbf{x}$
- ▶ here  $\mathbf{x}$  is one among the  $p$ -dimensional vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n$



# Principal Components Analysis

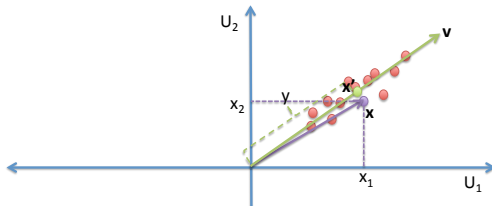
Mathematically, how do we do this?

Simple case:  $\mathbf{x} \in \mathbb{R}^2$ , and we want a good projection  $y \in \mathbb{R}$

- ▶ we want to find a new direction vector  $\mathbf{v}$
- ▶ let  $\mathbf{w}_1$  be the linear transformation of  $\mathbf{u}_1$ ,  $\mathbf{u}_2$  into  $\mathbf{v}$ :

$$\mathbf{v} = w_{11}\mathbf{u}_1 + w_{12}\mathbf{u}_2 \quad (= \mathbf{I}_2 \mathbf{w}_1 = \mathbf{w}_1)$$

- ▶ Remember:  $y_{ik} = \mathbf{x}_i^\top \mathbf{w}_k$ , with  $1 \leq i \leq n$  and  $1 \leq k \leq \kappa$  Here  $\kappa = 1$ ;  $y := y_{i1} = \mathbf{x}^\top \mathbf{v} = \mathbf{x}^\top \mathbf{w}_1$
- ▶ project  $\mathbf{x}$  onto  $\mathbf{v}$  to make  $\mathbf{x}' = y \mathbf{w}_1$  ( $y$  is the length of  $\mathbf{x}'$ )
- ▶ two coordinates  $x_1, x_2$  become  $y$  (what have we lost?)



# Principal Components Analysis

To find best  $\mathbf{w}_1$ :

- ▶ “closeness” is based on squared error between original points and new points
- ▶ usual notion of distance is Euclidean:

$$d(\mathbf{x}_i, \mathbf{x}'_i) = \sqrt{\sum_{j=1}^p (x_{ij} - x'_{ij})^2}$$

- ▶ we will measure distance as Euclidean distance *squared*  $\rightarrow$  big deviations heavily penalized, smaller deviations less so
- ▶ objective is to minimize squared errors

$$\hat{\mathbf{w}}_1 = \arg \min_{\mathbf{w}_1} \left\{ \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{x}'_i\|_2^2 = \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - x'_{ij})^2 \right\}$$

# Principal Components Analysis

To find best  $\mathbf{w}_1$ :

- ▶ also have constraint:

$$\|\mathbf{w}_1\|_2 = \sqrt{\sum_{j=1}^p w_{1j}^2} = 1$$

this means one step in old coordinates  $\mathbf{u}$  is equal to one step in new coordinates  $\mathbf{v}$  (this is called the L2 norm... we will use this and the L1 norm  $\|\mathbf{x}\|_1 = \sum_{j=1}^p |x_j|$  a lot in this class)



# Principal Components Analysis

To find the best single feature  $\mathbf{w}_1$ :

1. center the data (subtract mean)... in R?
2. find best single feature,  $\mathbf{w}_1$  by

$$\begin{aligned}\hat{\mathbf{w}}_1 &= \arg \min_{\mathbf{w}_1: \|\mathbf{w}_1\|_2=1} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - x'_{ij})^2 \\&= \arg \min_{\mathbf{w}_1: \|\mathbf{w}_1\|_2=1} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - y_{i1} w_{1j})^2 \\&= \arg \min_{\mathbf{w}_1: \|\mathbf{w}_1\|_2=1} \sum_{i=1}^n \sum_{j=1}^p \left( x_{ij} - w_{1j} \mathbf{x}_i^\top \mathbf{w}_1 \right)^2 \\&= \dots \quad (\text{using some matrix algebra} - \text{ see next slide}) \\&= \arg \max_{\mathbf{w}_1: \|\mathbf{w}_1\|_2=1} \left\{ \mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1 = \frac{\mathbf{w}_1^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{w}_1}{\mathbf{w}_1^\top \mathbf{w}_1} \right\}\end{aligned}$$

# Principal Components Analysis

$$\sum_{i=1}^n \sum_{j=1}^p (x_{ij} - x'_{ij})^2 = \sum_{i=1}^n \sum_{j=1}^p \left( x_{ij} - w_{1j} \mathbf{x}_i^\top \mathbf{w}_1 \right)^2$$

$$\left\| \mathbf{X} - \mathbf{X} \mathbf{w}_1 \mathbf{w}_1^\top \right\|_2^2 = \left\| \mathbf{X} (\mathbf{I}_p - \mathbf{w}_1 \mathbf{w}_1^\top) \right\|_2^2$$

$$\text{tr}[(\mathbf{I}_p - \mathbf{w}_1 \mathbf{w}_1^\top)^\top \mathbf{X}^\top \mathbf{X} (\mathbf{I}_p - \mathbf{w}_1 \mathbf{w}_1^\top)] = \text{tr}[(\mathbf{X}^\top \mathbf{X} - \mathbf{w}_1 \mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X}) (\mathbf{I}_p - \mathbf{w}_1 \mathbf{w}_1^\top)]$$

$$= \text{tr}[\mathbf{X}^\top \mathbf{X}] - \text{tr}[\mathbf{X}^\top \mathbf{X} \mathbf{w}_1 \mathbf{w}_1^\top] - \text{tr}[\mathbf{w}_1 \mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X}] + \text{tr}[\mathbf{w}_1 \mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1 \mathbf{w}_1^\top]$$

(using the cyclical property of the trace and  $\mathbf{w}_1^\top \mathbf{w}_1 = 1$ )

$$= \text{tr}[\mathbf{X}^\top \mathbf{X}] - 2\text{tr}[\mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1] + \text{tr}[\mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1]$$

$$= \text{tr}[\mathbf{X}^\top \mathbf{X}] - \text{tr}[\mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1]$$

$$= \text{tr}[\mathbf{X}^\top \mathbf{X}] - \mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1$$

# Principal Components Analysis

An important fact:  $\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$  is the maximum likelihood estimate of the covariance matrix for  $[x_1, \dots, x_p]^\top$

The scalar

$$\frac{\mathbf{w}_1^\top \left( \frac{1}{n} \mathbf{X}^\top \mathbf{X} \right) \mathbf{w}_1}{\mathbf{w}_1^\top \mathbf{w}_1} \quad \left( = \frac{1}{n} \mathbf{Y}_1^\top \mathbf{Y}_1 = \frac{1}{n} \sum_{i=1}^n y_{i1}^2 \right)$$

is called a *Rayleigh quotient*

Another fact: since  $\frac{1}{n} \mathbf{X}^\top \mathbf{X}$  is symmetric,

- ▶ the value of  $\mathbf{w}_1$  that maximizes the Rayleigh quotient is the eigenvector associated with the largest eigenvalue of  $\frac{1}{n} \mathbf{X}^\top \mathbf{X}$
- ▶ the corresponding maximum of the Rayleigh quotient is the largest eigenvalue of  $\frac{1}{n} \mathbf{X}^\top \mathbf{X}$

# Principal Components Analysis

- ▶ center the data (column-wise) and redefine  $\mathbf{X} = \mathbf{H}\mathbf{X}$ , where  $\mathbf{H} = \mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top$
- ▶ Let  $\mathbf{Y}_1 = \mathbf{X}\mathbf{w}_1$ , and rewrite the problem as

$$\max_{\mathbf{w}_1^\top \mathbf{w}_1 = 1} \|\mathbf{Y}_1\| = \max_{\mathbf{w}_1^\top \mathbf{w}_1 = 1} \mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1$$

- ▶ Lagrange function  
 $L_1(\mathbf{w}_1, \mathbf{X}, \lambda) = \mathbf{w}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_1 - \lambda(\mathbf{w}_1^\top \mathbf{w}_1 - 1)$
- ▶ Our problem is equivalent to

$$\max_{\mathbf{w}_1} L_1(\mathbf{w}_1, \mathbf{X}, \lambda)$$

- ▶ solution

$$\frac{\partial L_1}{\partial \mathbf{w}_1} = 2\mathbf{X}^\top \mathbf{X} \mathbf{w}_1 - 2\lambda \mathbf{w}_1 = 0 \iff \mathbf{X}^\top \mathbf{X} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1$$

How to find the second new direction?

# Principal Components Analysis

- ▶ Let  $\mathbf{Y}_2 = \mathbf{X}\mathbf{w}_2$ , and rewrite the problem as

$$\max_{\mathbf{w}_2^\top \mathbf{w}_2=1, \mathbf{w}_2^\top \mathbf{w}_1=0} \|\mathbf{Y}_2\| = \max_{\mathbf{w}_2^\top \mathbf{w}_2=1, \mathbf{w}_2^\top \mathbf{w}_1=0} \mathbf{w}_2^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_2$$

- ▶ Lagrange function

$$L_2(\mathbf{w}_2, \mathbf{w}_1, \mathbf{X}, \lambda) = \mathbf{w}_2^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_2 - \lambda(\mathbf{w}_2^\top \mathbf{w}_2 - 1) - \mu(\mathbf{w}_2^\top \mathbf{w}_1)$$

- ▶ Our problem is equivalent to

$$\max_{\mathbf{w}_2} L_2(\mathbf{w}_2, \mathbf{w}_1, \mathbf{X}, \lambda)$$

- ▶ solution

$$\frac{\partial L_2}{\partial \mathbf{w}_2} = 2\mathbf{X}^\top \mathbf{X} \mathbf{w}_2 - 2\lambda \mathbf{w}_2 - \mu \mathbf{w}_1 = 0$$

Pre-multiply by  $\mathbf{w}_2^\top$

$$2\mathbf{w}_2^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}_2 - 2\lambda \mathbf{w}_2^\top \mathbf{w}_2 - \mu \mathbf{w}_2^\top \mathbf{w}_1 = 0 \iff \mathbf{X}^\top \mathbf{X} \mathbf{w}_2 = \lambda_2 \mathbf{w}_2$$

How to find the third new direction? Solution:  $\mathbf{X}^\top \mathbf{X} \mathbf{w}_3 = \lambda_3 \mathbf{w}_3$   
... and so on ... up to the  $p$ -th eigenvector (the one corresponding to the smallest eigenvalues)

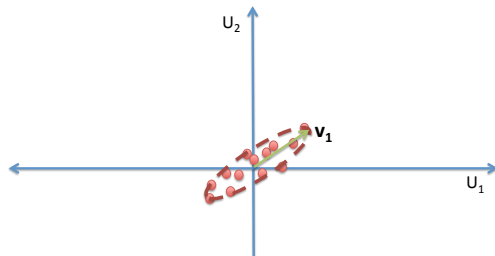
# Principal Components Analysis

Geometric interpretation.

- ▶ objective is to minimize squared errors

$$\hat{\mathbf{w}}_1 = \arg \min_{\mathbf{w}_1 : \|\mathbf{w}_1\|_2=1} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - x'_{ij})^2$$

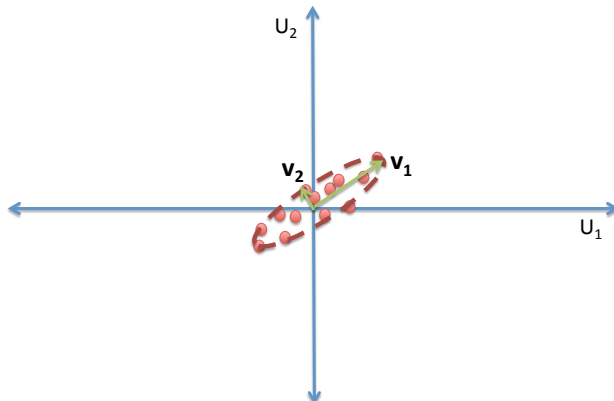
- ▶ center data
- ▶ fit a multivariate Gaussian distribution fit to the data
- ▶ Gaussian has mean  $\mathbf{0}$ , covariance  $\Sigma$
- ▶  $\mathbf{w}_1$  is the direction of the covariance ellipse with maximum variance



# Principal Components Analysis

This generalizes to multiple dimensions

- ▶ suppose we have  $p$  original dimensions and we would like  $\kappa$  new ones
- ▶ the optimal vectors have the same direction as the  $\kappa$  ellipse directions with maximum variance



# Principal Components Analysis

How do we find the  $\kappa$  most descriptive directions?

- ▶ center the  $p$  variables so that they all have mean 0
- ▶ estimate the covariance matrix

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^{\top}$$

- ▶ find the *eigenvalues* and *eigenvectors* of the covariance matrix
- ▶ the  $\kappa$  most descriptive directions are the eigenvectors associated with the  $\kappa$  largest eigenvalues



# The Truth: Matrices, Eigenvalues, and Eigenvectors

Matrices can be used to describe transformations:

- ▶ under transformations, some of the original directions may be preserved

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$$

- ▶ the direction  $\mathbf{x}$  is an *eigenvector* of  $\mathbf{A}$
- ▶ the scalar  $\lambda$  is an *eigenvalue*  $\mathbf{A}$
- ▶ the eigenvalue describes the magnitude of the change in  $\mathbf{x}$  under  $\mathbf{A}$

Let's play with an animated gif [eigshow in MATLAB]. What is the matrix? What are the eigenvectors? What are the eigenvalues?

# The Truth: Matrices, Eigenvalues, and Eigenvectors

Matrices can be used to describe transformations:

- ▶  $\mathbf{x} \rightarrow \mathbf{Ax}$  is a transformation of  $\mathbf{x}$

- ▶ Examples (look at these graphically):

- ▶  $\mathbf{x} =$

- $(1, 1), (1, 0), (1, -1), (0, 1), (-1, 1), (-1, 0), (-1, -1), (0, -1)$

- ▶  $\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$

- ▶  $\mathbf{A} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

# The Truth: Matrices, Eigenvalues, and Eigenvectors

To find the eigenvectors of  $\mathbf{A}$ , find the roots of the polynomial

$$\det(\mathbf{A} - \lambda I) = 0$$

Let's do this with our examples:

►  $\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$

►  $\mathbf{A} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

# The Truth: Matrices, Eigenvalues, and Eigenvectors

We can use the eigenvalues to find the eigenvectors

- ▶ start with eigenvalue  $\lambda_i$
- ▶ solve the set of linear equations

$$\mathbf{A}\mathbf{x} = \lambda_i\mathbf{x}$$

- ▶  $\mathbf{x}$  is the eigenvector associated with  $\lambda_i$

Let's do this with our examples:

- ▶  $\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$

- ▶  $\mathbf{A} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

# The Truth: Matrices, Eigenvalues, and Eigenvectors

We have 2 matrices:

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix},$$

$$\mathbf{B} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}.$$

Let's find the eigenvalues:

$$\begin{aligned} 0 &= \det(\mathbf{A} - \lambda I) \\ &= \det\left(\begin{bmatrix} 2 - \lambda & 1 \\ 1 & 2 - \lambda \end{bmatrix}\right) \\ &= (2 - \lambda)^2 - 1 \\ &= (\lambda - 3)(\lambda - 1) \end{aligned}$$

$$\begin{aligned} 0 &= \det(\mathbf{B} - \lambda I) \\ &= \det\left(\begin{bmatrix} 2 - \lambda & 0 \\ 0 & 2 - \lambda \end{bmatrix}\right) \\ &= (2 - \lambda)^2 \end{aligned}$$

So the eigenvalues are 3 and 1 for  $\mathbf{A}$ ; 2 and 2 for  $\mathbf{B}$ .

# The Truth: Matrices, Eigenvalues, and Eigenvectors

Now use the eigenvalues to find the eigenvectors:  $\mathbf{Ax} = \lambda\mathbf{x}$ .

Solve for  $\mathbf{A}$ :

$$\mathbf{Ax}_1 = 3\mathbf{x}_1$$

$$\begin{bmatrix} 2x_{11} + x_{12} \\ x_{11} + 2x_{12} \end{bmatrix} = \begin{bmatrix} 3x_{11} \\ 3x_{12} \end{bmatrix}$$

$$\mathbf{x}_1 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}$$

$$\mathbf{Ax}_2 = 1\mathbf{x}_2$$

$$\begin{bmatrix} 2x_{21} + x_{22} \\ x_{21} + 2x_{22} \end{bmatrix} = \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix}$$

$$\mathbf{x}_2 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} \end{bmatrix}$$

Solve for  $\mathbf{B}$ :

$$\mathbf{Bx}_1 = 2\mathbf{x}_1$$

$$\begin{bmatrix} 2x_{11} \\ 2x_{12} \end{bmatrix} = \begin{bmatrix} 2x_{11} \\ 2x_{12} \end{bmatrix}$$

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

# Principal Components Analysis

Back to PCA...

1. center the data
2. compute  $\mathbf{X}^\top \mathbf{X}$
3. compute the  $\kappa$  eigenvectors corresponding to the largest  $\kappa$  eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  to get loadings
4. for the eigenvectors with the  $\kappa$  largest eigenvalues, make factor scores by setting

$$\mathbf{y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{i\kappa} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1p} \\ w_{21} & w_{22} & \dots & w_{2p} \\ \vdots & & & \vdots \\ w_{\kappa 1} & w_{\kappa 2} & \dots & w_{\kappa p} \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ \vdots \\ \vdots \\ x_{ip} \end{bmatrix}$$

# Principal Components Analysis

- do computations in the (smaller,  $\kappa$ -dimensional) space, the space of factor scores  $\mathbf{y}$
- transform results back to original space (if needed)

$$\mathbf{y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{i\kappa} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1p} \\ w_{21} & w_{22} & \dots & w_{2p} \\ \vdots & & & \vdots \\ w_{\kappa 1} & w_{\kappa 2} & \dots & w_{\kappa p} \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ \vdots \\ \vdots \\ x_{ip} \end{bmatrix}$$



# Principal Components Analysis

Let's do an example:

$$\mathbf{X} = \begin{bmatrix} -4 & -4 \\ -1 & 1 \\ 1 & -1 \\ 4 & 4 \end{bmatrix}$$

$$[\mathbf{w}_1, \mathbf{w}_2] =$$

$$[\mathbf{y}_1, \mathbf{y}_2] =$$

Let's find the eigenvalues and eigenvectors of the empirical covariance matrix of  $\mathbf{X}$  for PCA.

## PCA: steps to find the optimal rotation

1. Is the data centered? Yes, the sum of all of the columns of  $\mathbf{X}$  is 0.)
2. Find the empirical covariance matrix  $\frac{1}{n}\mathbf{X}^\top\mathbf{X}$ :

$$\frac{1}{n}\mathbf{X}^\top\mathbf{X} = \frac{1}{4} \begin{bmatrix} 34 & 30 \\ 30 & 34 \end{bmatrix}$$

3. Now find the eigenvalues:

$$\begin{aligned} 0 &= \det \left( \begin{bmatrix} \frac{34}{4} - \lambda & \frac{30}{4} \\ \frac{30}{4} & \frac{34}{4} - \lambda \end{bmatrix} \right) \\ &= (\lambda - 16)(\lambda - 1) \end{aligned}$$

Therefore  $\lambda = 16$  is the eigenvalue associated with the most expressive linear rotation,  $\lambda = 1$  is the second the eigenvalue associated with the second most expressive.

## PCA: steps to find the optimal rotation

4. Find the eigenvectors associated with those eigenvalues:

$$\begin{array}{l|l} \frac{1}{4}\mathbf{X}^\top\mathbf{X}\mathbf{w}_1 = 16\mathbf{w}_1 & \frac{1}{4}\mathbf{X}^\top\mathbf{X}\mathbf{w}_2 = 1\mathbf{w}_2 \\ \left[ \begin{array}{c} \frac{34}{4}w_{11} + \frac{30}{4}w_{12} \\ \frac{30}{4}w_{11} + \frac{34}{4}w_{12} \end{array} \right] = \left[ \begin{array}{c} 16w_{11} \\ 16w_{12} \end{array} \right] & \left[ \begin{array}{c} \frac{34}{4}w_{21} + \frac{30}{4}w_{22} \\ \frac{30}{4}w_{21} + \frac{34}{4}w_{22} \end{array} \right] = \left[ \begin{array}{c} w_{21} \\ w_{22} \end{array} \right] \\ \Leftrightarrow \mathbf{w}_1 = \left[ \begin{array}{c} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{array} \right] & \Leftrightarrow \mathbf{w}_2 = \left[ \begin{array}{c} -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{array} \right] \end{array}$$

Note that the eigenvectors have magnitude 1:

$$\sum_{j=1}^p w_{jk}^2 = 1, \quad 1 \leq k \leq p.$$

This means that one step in the new directions is equivalent to one step in the old directions.

## PCA: steps to find the optimal rotation

5. The eigenvectors make up your loading matrix,  $\mathbf{W}$ :

$$\mathbf{W} = \begin{bmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}$$

Since eigenvectors form an orthonormal basis and have magnitude 1,  $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$ .

6. Find the scores for each data point by  $\mathbf{Y} = \mathbf{XW}$ .