# IAML: Dimensionality Reduction

Victor Lavrenko and Charles Sutton

School of Informatics

Semester 1

## Overview
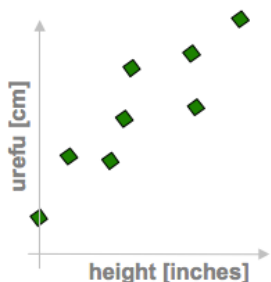
- Curse of dimensionality
- Different ways to reduce dimensionality
- Principal Components Analysis (PCA)
- Examples: Eigen Faces, Topics in Text
- PCA for classification
- Witten & Frank section 7.3
  - only the PCA section required

## True vs. observed dimensionality

- Get a population, predict some property
  - instances represented as {urefu, height} pairs
  - what is the dimensionality of this data?



urefu [cm]

height [inches]

"height" = "urefu" in Swahili

- Data points over time from different geographic areas over time:
  - $X_1$: # of traffic accidents
  - $X_2$: # of burst water pipes
  - $X_3$: snow-plow expenditures
  - $X_4$: # of forest fires
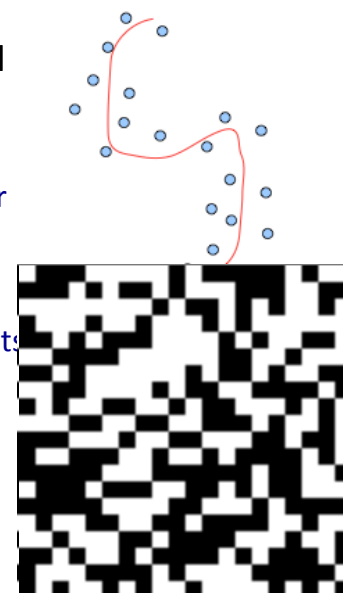  - $X_5$: # patients with heat stroke

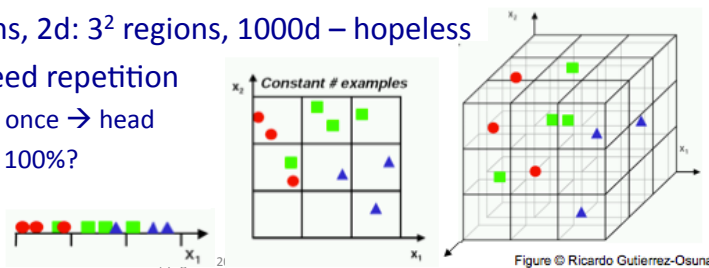Temperature below freezing?

## Curse of dimensionality

- Datasets typically high dimensional
  - vision: $10^4$ pixels, text: $10^6$ words
    - the way we observe / record them
  - true dimensionality often much lower
    - a manifold (sheet) in a high-d space
- Example: handwritten digits
  - 20 x 20 bitmap: $\{0,1\}^{400}$ possible events
    - will never see most of these events
    - actual digits: tiny fraction of events
  - true dimensionality:
    - possible variations of the pen-stroke

# Curse of dimensionality (2)

- Machine learning methods are statistical by nature
  - count observations in various regions of some space
  - use counts to construct the predictor f(x)
  - e.g. decision trees: $p_+/p_-$ in {o=rain,w=strong,T>28°}
  - text: #documents in {"hp" and "3d" and not "$" and …)
- As dimensionality grows: fewer observations per region
  - 1d: 3 regions, 2d: $3^2$ regions, 1000d – hopeless
  - statistics need repetition
    - flip a coin once → head
    - P(head) = 100%?



Figure © Ricardo Gutierrez-Osuna

# Dealing with high dimensionality

- Use domain knowledge
  - feature engineering: SIFT, MFCC
- Make assumption about dimensions
  - independence: count along each dimension separately
  - smoothness: propagate class counts to neighboring regions
  - symmetry: e.g. invariance to order of dimensions: x1 ⇔ x2
- Reduce the dimensionality of the data
  - create a new set of dimensions (variables)
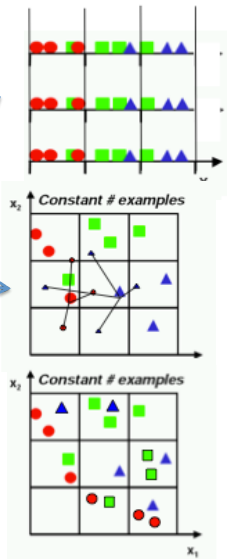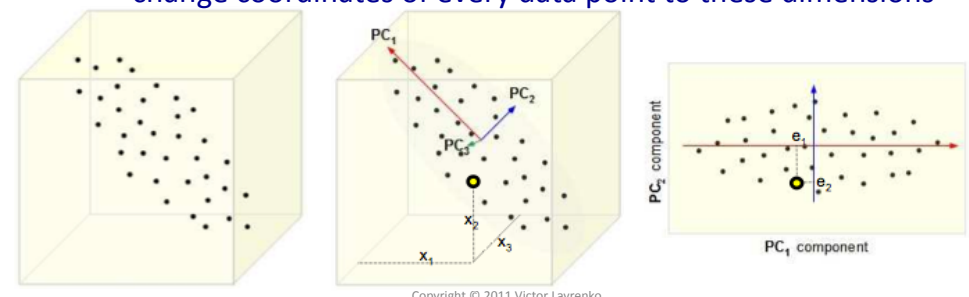
# Dimensionality reduction

- Goal: represent instances with fewer variables
  - try to preserve as much structure in the data as possible
  - discriminative: only structure that affects class separability
- Feature selection
  - pick a subset of the original dimensions $X_1 X_2 X_3 … X_{d-1} X_d$
  - discriminative: pick good class "predictors" (e.g. gain)
- Feature extraction
  - construct a new set of dimensions $E_1 E_2 … E_m$
    $$E_i = f(X_1…X_d)$$
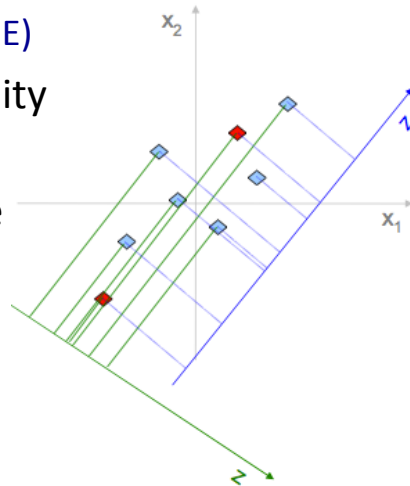  - (linear) combinations of original $X_1 X_2 X_3 … X_d$

# Principal Components Analysis

- Defines a set of principal components
  - 1st: direction of the greatest variability in the data
  - 2nd: perpendicular to 1st, greatest variability of what's left
  - … and so on until d (original dimensionality)
- First *m* components become *m* new dimensions
  - change coordinates of every data point to these dimensions
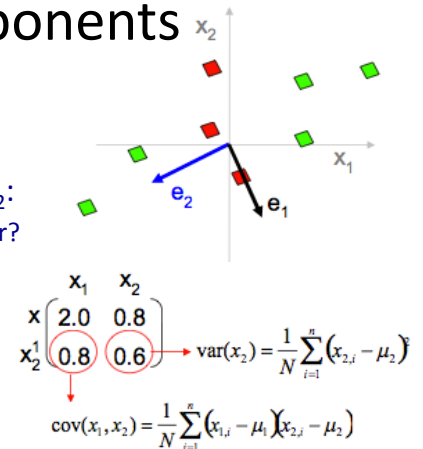
# Why greatest variability?

- Example: reduce 2-dimensional data to 1-d
  - $\{x_1, x_2\} \rightarrow e$ (along new axis E)
- Pick E to maximize variability
- Reduces cases when two points are close in e-space but very far in (x,y)-space
- Minimizes distances between original points and their projections
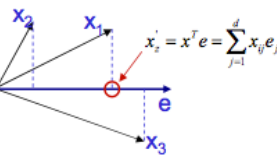
# Principal components

- Compute covariance matrix $\Sigma$
  - covariance of dimensions $x_1$ and $x_2$:
    - do $x_1$ and $x_2$ tend to increase together?
    - or does $x_2$ decrease as $x_1$ increases?
  - covariance: measure of variability

$$x \begin{matrix} x_1 & x_2 \\ 2.0 & 0.8 \\ 0.8 & 0.6 \end{matrix}$$

$$var(x_2) = \frac{1}{N} \sum_{i=1}^{n} (x_{2,i} - \mu_2)^2$$

$$cov(x_1, x_2) = \frac{1}{N} \sum_{i=1}^{n} (x_{1,i} - \mu_1)(x_{2,i} - \mu_2)$$

- Find the basis of $\Sigma$
  - find vectors $e_i$ which aren't turned by $\Sigma$
    - $\Sigma e_i = \lambda_i e_i$: eigenvalue / eigenvector

$$\lambda_1 \begin{bmatrix} 0.26 \end{bmatrix} \quad x \begin{matrix} e_1 & e_2 \\ 0.4 & -0.9 \end{matrix}$$
$$\lambda_2 \begin{bmatrix} 2.42 \end{bmatrix} \quad y \begin{matrix} -0.9 & -0.4 \end{matrix}$$

  - $1^{st}$ PC: "longest" $e_i$ (has largest $\lambda_i$), $2^{nd}$ PC: next longest, ...

# Direction of greatest variability

- Select dimension e which maximizes the variance
- Points x "projected" onto vector e:

$$x_z' = x^T e = \sum_{j=1}^{d} x_{ij} e_j$$

- Variance of projections:

$$\frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{d} x_{ij} e_j - \mu \right)^2 = \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{d} x_{ij} e_j \right)^2$$

- Maximize variance
  - want unit length: $||e||=1$
  - add Lagrange multiplier

$$L = \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{d} x_{ij} e_j \right)^2 - \lambda \left( \left( \sum_{k=1}^{d} e_j^2 \right) - 1 \right)$$

$$\frac{\partial L}{\partial e_a} = \frac{2}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{d} x_{ij} e_j \right) x_{ia} - 2\lambda e_a = 0$$
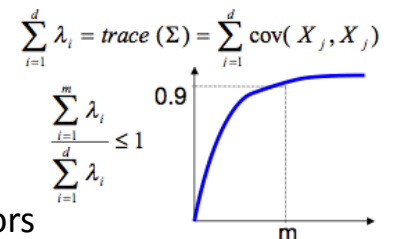
$$\sum_{j=1}^{d} cov(a,j) e_j = \lambda e_a \quad \text{for a = 1...d} \quad \Leftarrow \quad 0 = 2 \sum_{j=1}^{d} e_j \left( \frac{1}{n} \sum_{i=1}^{n} x_{ia} x_{ij} \right) - 2\lambda e_a$$

$$\underbrace{\phantom{\frac{1}{n} \sum_{i=1}^{n} x_{ia} x_{ij}}}_{covariance}$$

$$\Sigma e = \lambda e \quad \rightarrow \text{ e must be an eigenvector}$$
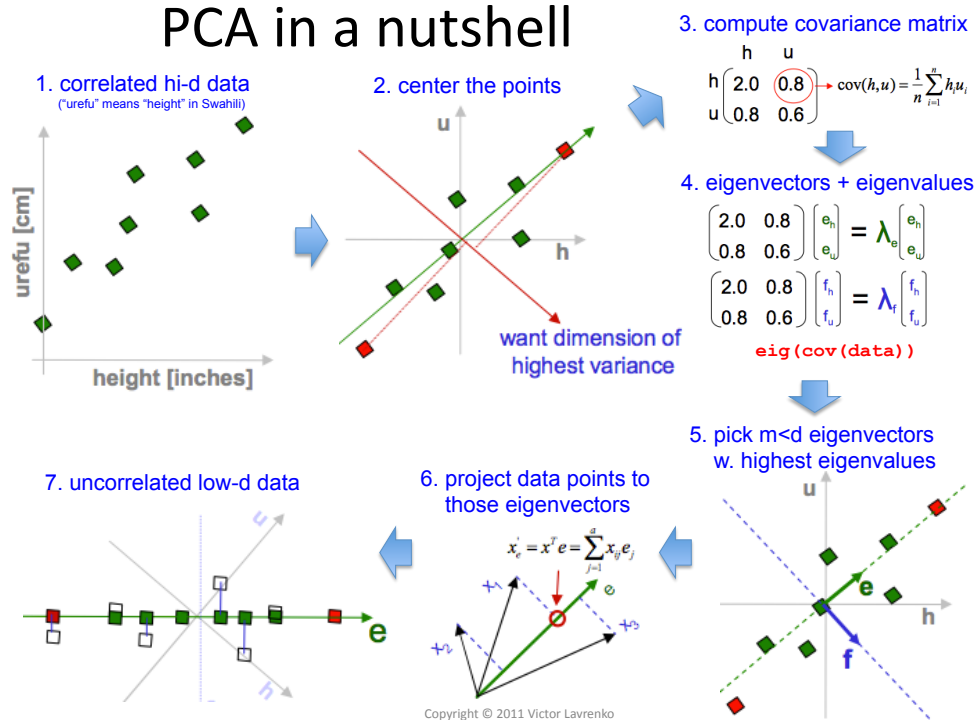
# Projecting to new dimensions

- Got a set of principal components $e_1 \ldots e_d$
  - orthogonal, unit length
  - corresponding eigenvalues $\lambda_1 \ldots \lambda_d$:

$$\sum_{i=1}^{d} \lambda_i = trace\ (\Sigma) = \sum_{i=1}^{d} cov(X_j, X_j)$$

  - fraction of variation explained by first $m$ principal components

$$\frac{\sum_{i=1}^{m} \lambda_i}{\sum_{i=1}^{d} \lambda_i} \leq 1$$

  - typical threshold values: 0.9 or 0.95
- $e_1 \ldots e_m$ are new dimension vectors
- Change coordinates: $x_{1..d} \rightarrow x'_{1..m}$
  - subtract mean from old dimensions
  - dot product each dimension with $e_1 \ldots e_m$
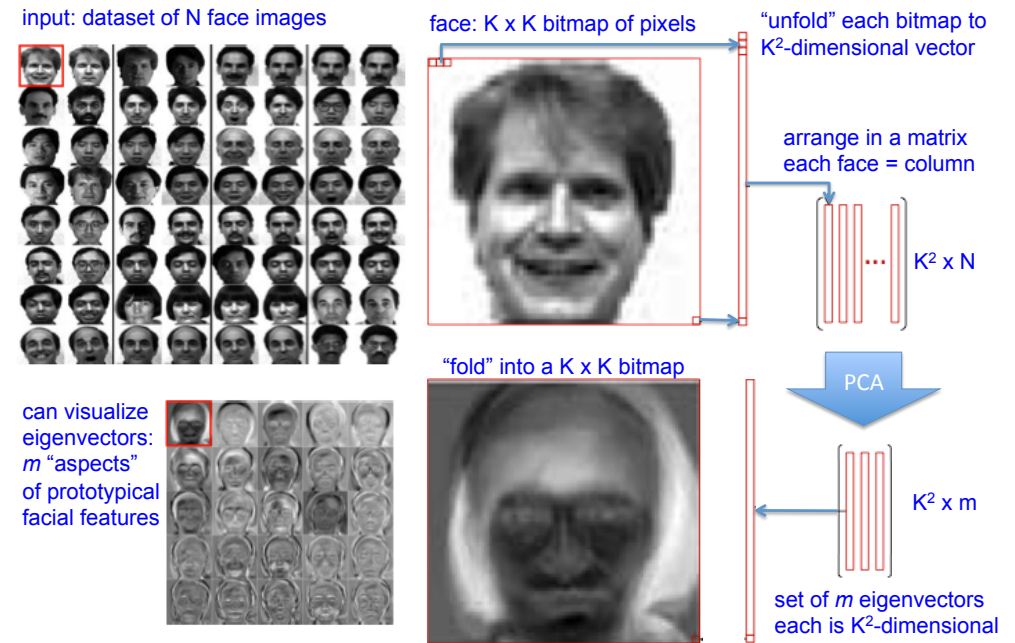
$$\begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_m \\ \cdots \\ x_d \end{bmatrix} \Rightarrow \begin{bmatrix} x_1' = \sum_{j=1}^{d} e_{i,j} x_j \\ x_2' = \sum_{j=1}^{d} e_{2,j} x_j \\ \cdots \\ x_m' = \sum_{j=1}^{d} e_{m,j} x_j \end{bmatrix}$$

# PCA in a nutshell

1. correlated hi-d data
("urefu" means "height" in Swahili)



2. center the points



want dimension of highest variance

3. compute covariance matrix

$$h \quad u$$
$$h \begin{bmatrix} 2.0 & 0.8 \\ u & 0.8 & 0.6 \end{bmatrix} \rightarrow cov(h,u) = \frac{1}{n}\sum_{i=1}^{n} h_i u_i$$

4. eigenvectors + eigenvalues

$$\begin{bmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{bmatrix}\begin{bmatrix} e_h \\ e_u \end{bmatrix} = \lambda_e \begin{bmatrix} e_h \\ e_u \end{bmatrix}$$

$$\begin{bmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{bmatrix}\begin{bmatrix} f_h \\ f_u \end{bmatrix} = \lambda_f \begin{bmatrix} f_h \\ f_u \end{bmatrix}$$

**eig(cov(data))**

5. pick m<d eigenvectors w. highest eigenvalues



6. project data points to those eigenvectors

$$x'_e = x^T e = \sum_{j=1}^{d} x_{ij} e_j$$

7. uncorrelated low-d data



Copyright © 2011 Victor Lavrenko

# PCA example: Eigen Faces

input: dataset of N face images



face: K x K bitmap of pixels

"unfold" each bitmap to K²-dimensional vector

arrange in a matrix each face = column

K² x N

PCA

K² x m

set of m eigenvectors each is K²-dimensional

"fold" into a K x K bitmap

can visualize eigenvectors: m "aspects" of prototypical facial features

# Eigen Faces: Projection



= 0.9 * ... - 0.2 * ... + 0.4 * ... + ...

- Project new face to space of eigen-faces
- Represent vector as a linear combination of principal components
- How many do we need?

Copyright © 2011 Victor Lavrenko

# (Eigen) Face Recognition

- Face similarity
  - in the reduced space
  - insensitive to lighting expression, orientation
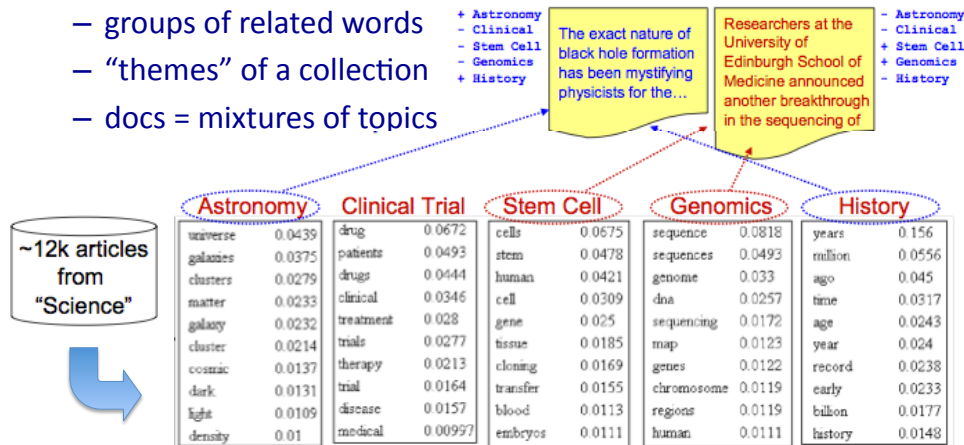- Projecting new "faces"
  - everything is a face



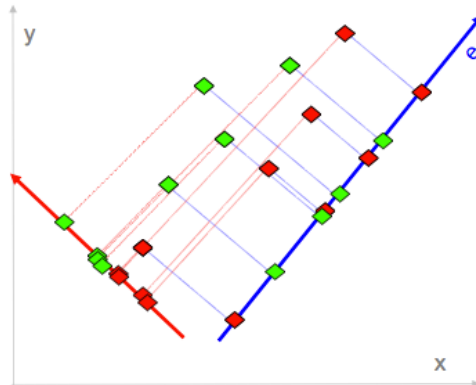new face (not in training)

projected to eigenfaces

Copyright © 2011 Victor Lavrenko

# PCA example: Topics in Text

- Can run variants of PCA on news, scientific papers
- Eigenvectors can be interpreted as "topics"
  - groups of related words
  - "themes" of a collection
  - docs = mixtures of topics



# PCA: practical issues

- Covariance extremely sensitive to large values
  - multiply some dimension by 1000
    - dominates covariance
    - becomes a principal component
  - normalize each dimension to zero mean and unit variance:
    $x' = (x - \text{mean}) / \text{st.dev}$
- PCA assumes underlying subspace is linear
  - 1d: straight line
    2d: flat sheet
  - transform to handle non-linear spaces (manifolds)



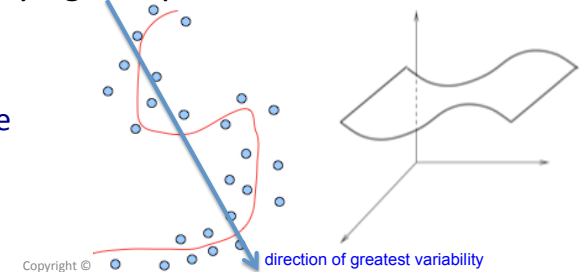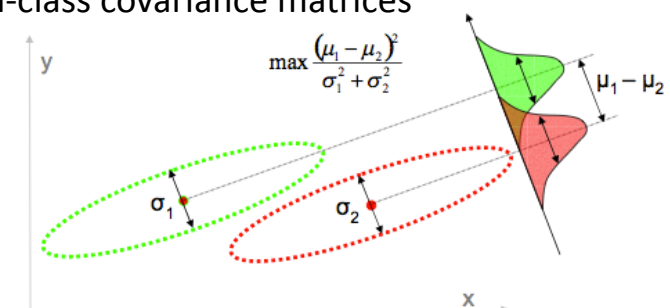direction of greatest variability

# PCA and classification

- PCA is unsupervised
  - maximizes overall variance of the data along a small set of directions
  - does not know anything about class labels
  - can pick direction that makes it hard to separate classes
- Discriminative approach
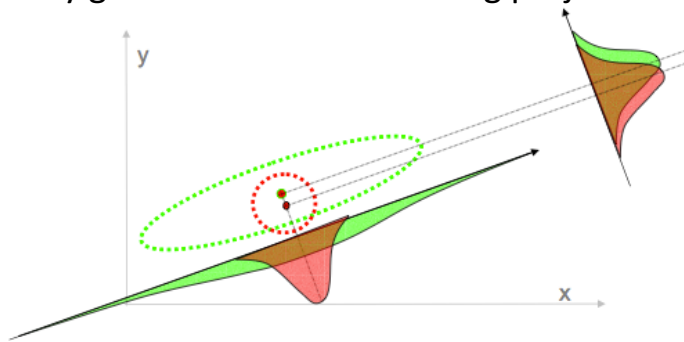  - look for a dimension that makes it easy to separate classes

# Linear Discriminant Analysis

- LDA: pick a new dimension that gives:
  - maximum separation between means of projected classes
  - minimum variance within each projected class
- Solution: eigenvectors based on between-class and within-class covariance matrices



$$\max \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$$

$\mu_1 - \mu_2$

$\sigma_1$   $\sigma_2$

# PCA vs. LDA

- LDA not always good for classification
  - assumes classes are unimodal Gaussians
  - fails when discriminatory information is not in the mean, but in the variance of the data
- PCA may give a more discriminating projection



# Dimensionality reduction

- Pros
  - reflects our intuitions about the data
  - allows estimating probabilities in high-dimensional data
    - no need to assume independence etc.
  - dramatic reduction in size of data
    - faster processing (as long as reduction is fast), smaller storage
- Cons
  - too expensive for many applications (Twitter, web)
  - disastrous for tasks with fine-grained classes
  - understand assumptions behind the methods (linearity etc.)
    - there may be better ways to deal with sparseness

# Summary

- True dimensionality << observed dimensionality
- High dimensionality ➜ sparse, unstable estimates
- Dealing with high dimensionality:
  - use domain knowledge
  - make an assumption: independence / smoothness / symmetry
  - dimensionality reduction: feature selection / feature extraction
- Principal Components Analysis (PCA)
  - picks dimensions that maximize variability
    - eigenvectors of the covariance matrix
  - examples: Eigen Faces, Topics in Text
  - variant for classification: Linear Discriminant Analysis