

S&DS 265 / 565

# Unsupervised Learning: PCA

Tuesday, October 12

Yale

# Checkpoint

- Assn 3 due Thursday; Assn 4 out
- Quiz 2 available from 12:01pm today for 24 hours
- Midterm one week from today, in class
- Practice midterm posted on Canvas
- Review sessions scheduled
- Questions?

# Unsupervised Learning

Supervised learning is about being able to predict a  $Y$  using predictors  $X_1, X_2, \dots, X_p$

Unsupervised learning deals with data that do not have labels  $Y$

We are not trying to predict anything. So what else might we hope to do?

# Unsupervised Learning

Supervised learning is about being able to predict a  $Y$  using predictors  $X_1, X_2, \dots, X_p$

Unsupervised learning deals with data that do not have labels  $Y$

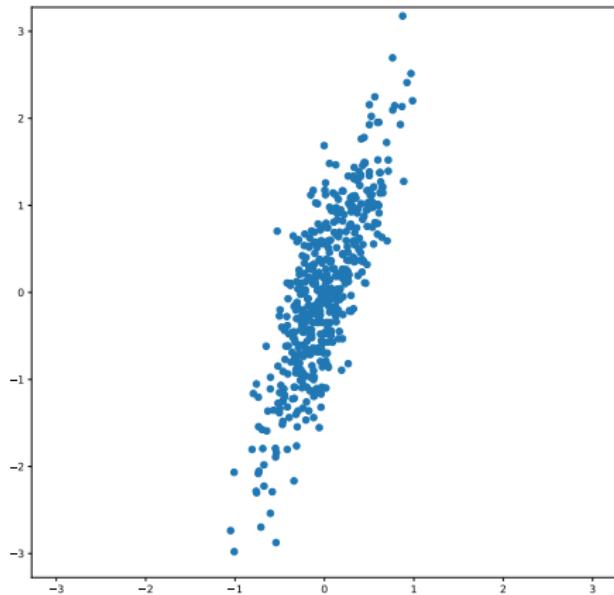
We are not trying to predict anything. So what else might we hope to do?

Consider:

- Are there interesting ways to visualize/summarize the data?
- Are there natural subgroups in the data?
- What are the important types of variation in the data?

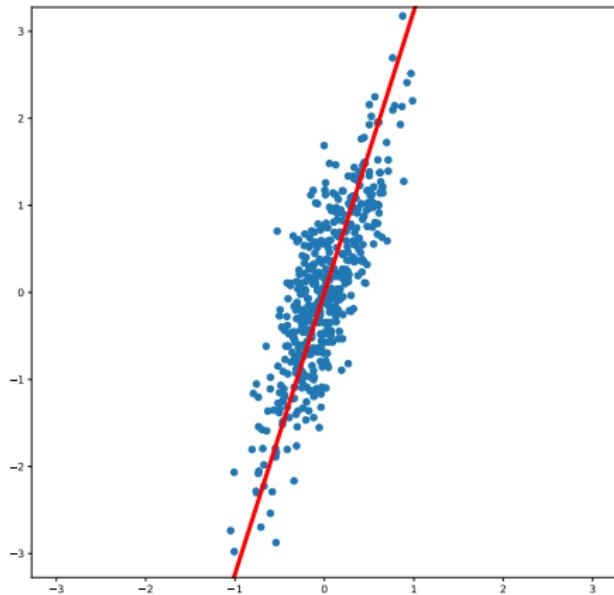
# Principal Component Analysis (PCA)

PCA finds the directions of greatest variability in the data.



# Principal Component Analysis (PCA)

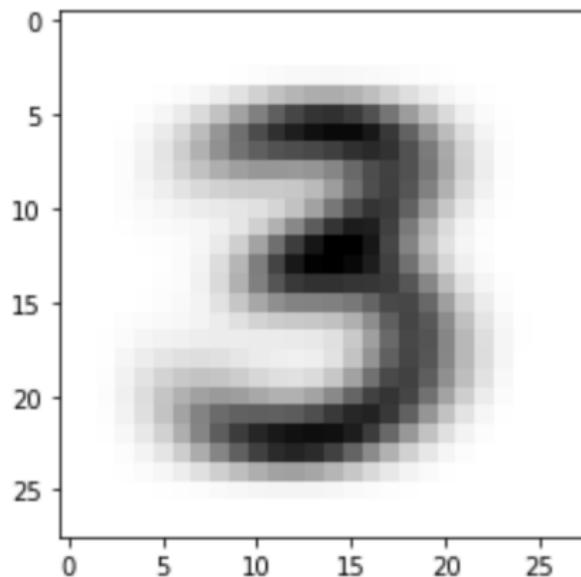
PCA finds the directions of greatest variability in the data.



# Handwritten Digits (3s)

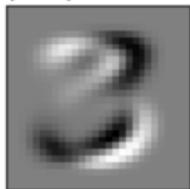
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3  
3 3 3 3 3 3 3

# Handwritten Digits (3s) – Average

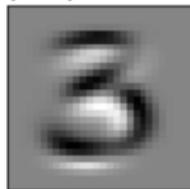


# Handwritten Digits (3s) – Principal vectors

principal vector 1



principal vector 2



principal vector 3



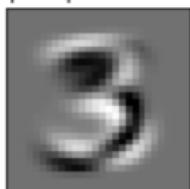
principal vector 4



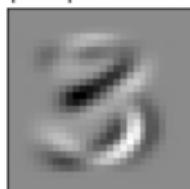
principal vector 5



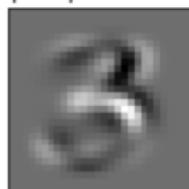
principal vector 6



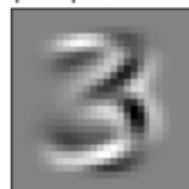
principal vector 7



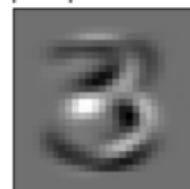
principal vector 8



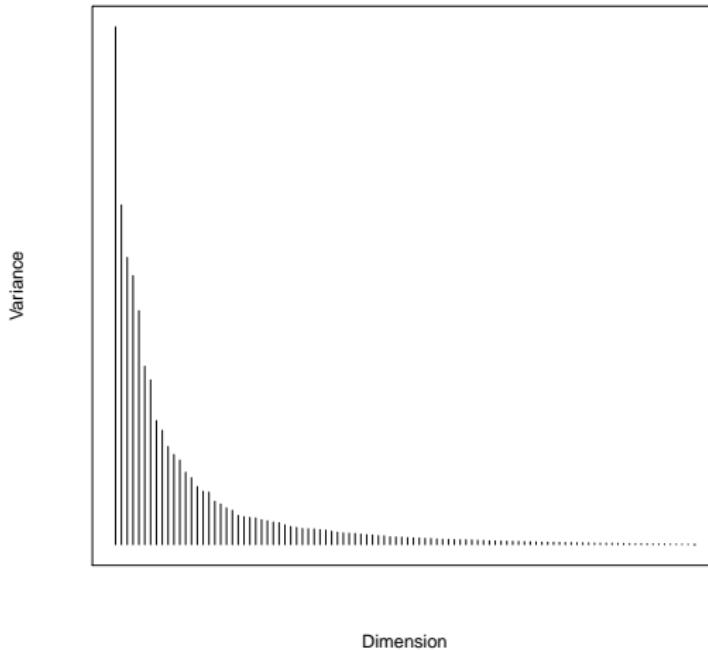
principal vector 9



principal vector 10



# Handwritten Digits (3s) – PCA variance

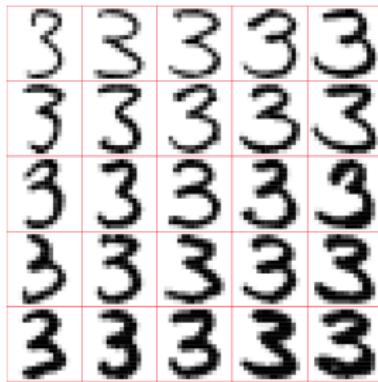
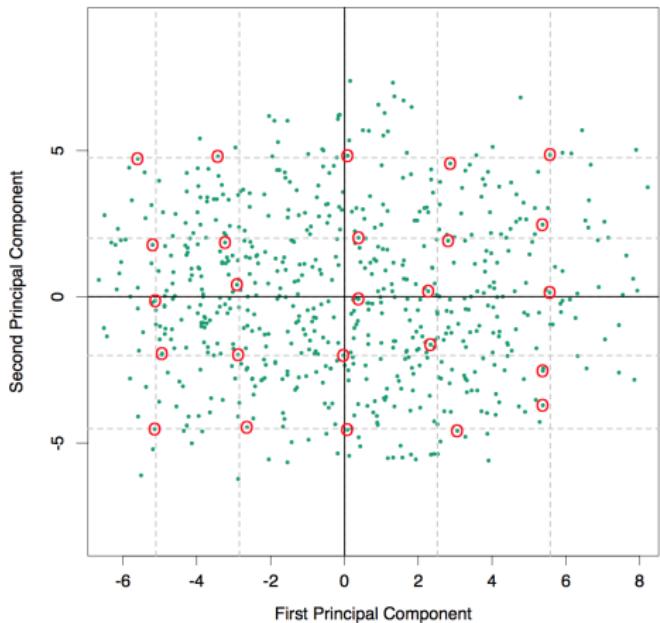


# Handwritten Digits (3s)

$$\begin{aligned}\hat{f}(\lambda) &= \bar{x} + \lambda_1 v_1 + \lambda_2 v_2 \\ &= \boxed{3} + \lambda_1 \cdot \boxed{3} + \lambda_2 \cdot \boxed{3}.\end{aligned}$$

# Handwritten Digits (3s) – Top 2 components

$$\begin{aligned}\hat{f}(\lambda) &= \bar{x} + \lambda_1 v_1 + \lambda_2 v_2 \\ &= \boxed{\text{3}} + \lambda_1 \cdot \boxed{\text{3}} + \lambda_2 \cdot \boxed{\text{3}}.\end{aligned}$$



# Handwritten Digits (3s) – PCA reconstruction



# Faces

test face 0



test face 1



test face 2



test face 3



test face 4



test face 5



test face 6



test face 7



test face 8



test face 9



test face 10



test face 11



# Eigenfaces (principal vectors)

eigenface 0



eigenface 1



eigenface 2



eigenface 3



eigenface 4



eigenface 5



eigenface 6



eigenface 7



eigenface 8



eigenface 9



eigenface 10



eigenface 11



## Genes mirror geography within Europe

John Novembre,<sup>1,2</sup> [Toby Johnson](#),<sup>4,5,6</sup> [Katarzyna Bryc](#),<sup>7</sup> [Zoltán Kutalik](#),<sup>4,6</sup> [Adam R. Boyko](#),<sup>7</sup> [Adam Auton](#),<sup>7</sup> Amit Indap,<sup>7</sup> [Karen S. King](#),<sup>8</sup> [Sven Bergmann](#),<sup>4,6</sup> [Matthew R. Nelson](#),<sup>8</sup> [Matthew Stephens](#),<sup>2,3</sup> and [Carlos D. Bustamante](#)<sup>7</sup>

[Author information ▶](#) [Copyright and License information ▶](#)

The publisher's final edited version of this article is available at [Nature](#)

This article has been corrected. See the correction in volume 456 on page 274.

See commentary "[Editorial comment should accompany hot papers online](#)." in *Nature*, volume 455 on page 861.

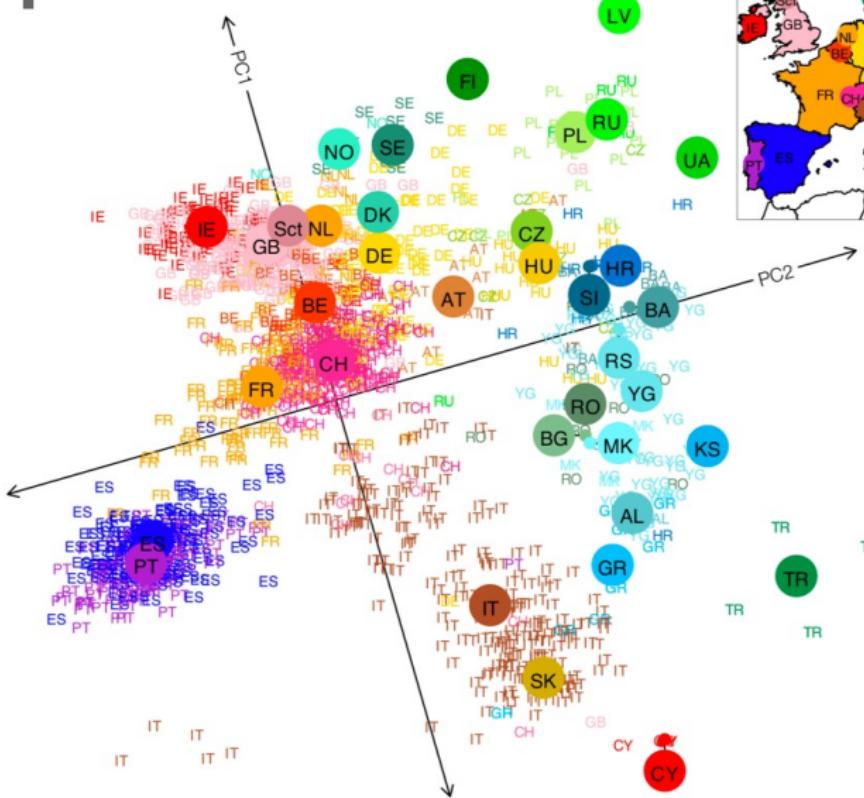
See other articles in PMC that [cite](#) the published article.

### Abstract

Go to:

Understanding the genetic structure of human populations is of fundamental interest to medical, forensic and anthropological sciences. Advances in high-throughput genotyping technology have markedly improved our understanding of global patterns of human genetic variation and suggest the potential to use large samples to uncover variation among closely spaced populations<sup>1–5</sup>. Here we characterize genetic variation in a sample of 3,000 European individuals genotyped at over half a million variable DNA sites in the human genome. Despite low average levels of genetic differentiation among Europeans, we find a close correspondence between genetic and geographic distances; indeed, a geographical map of Europe arises naturally as an efficient two-dimensional summary of genetic variation in Europeans. The results emphasize that when mapping the genetic basis of a disease phenotype, spurious associations can arise if genetic structure is not properly accounted for. In addition, the results are relevant to the prospects of genetic ancestry testing<sup>6</sup>; an individual's DNA can be used to infer their geographic origin with surprising accuracy—often to within a few hundred kilometres.

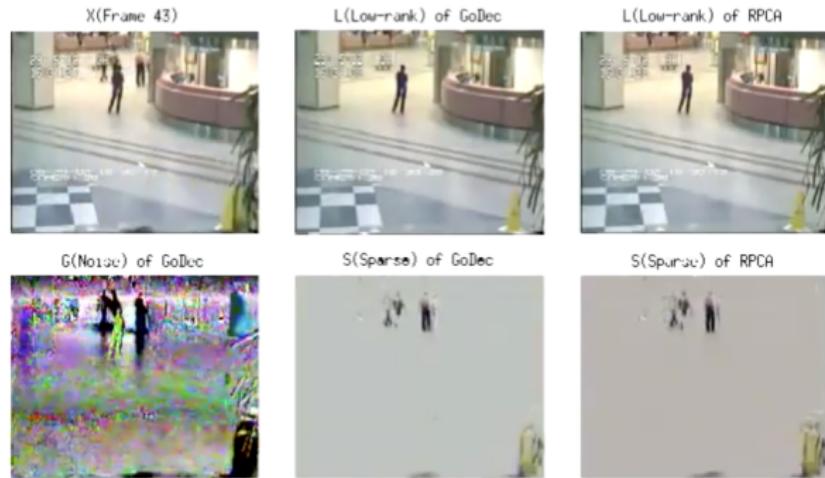
1





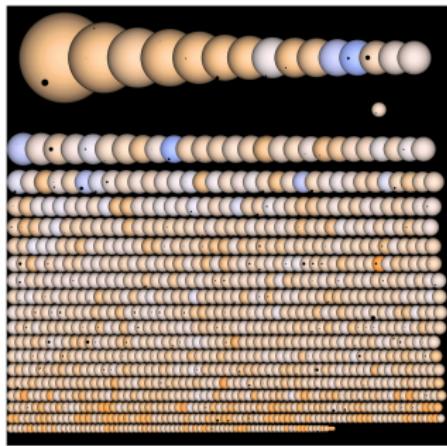
# Robust PCA

Robust PCA (low rank plus sparse) can be used for background subtraction in video.



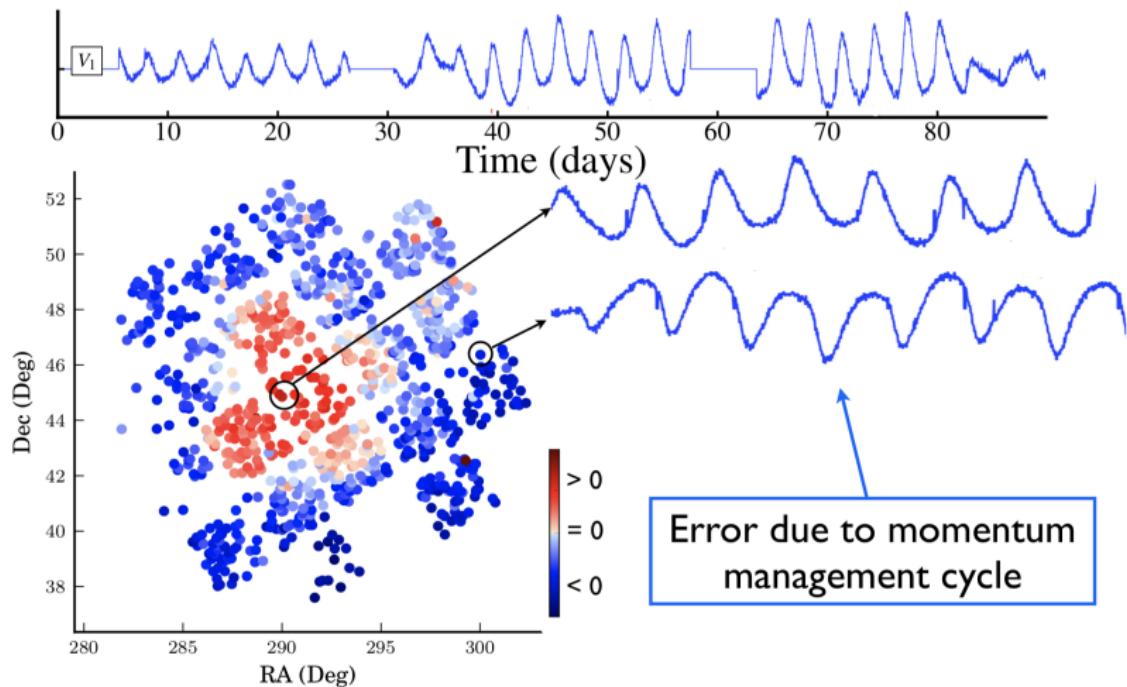
<https://www.youtube.com/watch?v=BTrbow8u4Cw>

# Kepler telescope

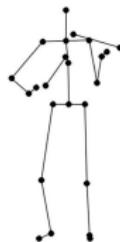


# Kepler telescope

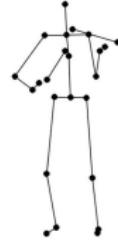
## ***Identification of correlated errors using robust PCA***



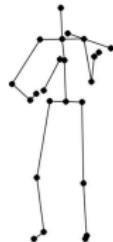
## RESULTS: PC VISUALISATION, TECHNICAL VS. NON-TECHNICAL MOTION



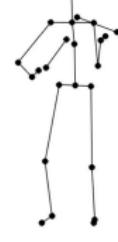
Original



PC1 (45% of variance accounted)  
**Side-to-side swaying (non-technical)**



PC2 (25% of variance accounted)  
**Right arm string crossing (technical)**



PC3 (15% of variance accounted)  
**Left-right rotation (non-technical)**

# PCA: Algorithm

- ① Center the data:  $x_i \mapsto x_i - \bar{x}$
- ② Compute the  $d \times d$  sample covariance  $S = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$
- ③ Find the first  $k$  eigenvectors of  $S$
- ④ Project the data onto those  $k$  vectors

# PCA: Algorithm

- ① Center the data:  $x_i \mapsto x_i - \frac{1}{n} \sum_{j=1}^n x_j = x_i - \bar{x}$
- ② Compute the  $d \times d$  sample covariance  $S = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ . Note that

$$\frac{1}{n} \sum_i (x_{ij} - \bar{x})^2$$

is the sample variance of  $j$ th coordinate of data.

- ③ Find the first  $k$  eigenvectors of  $S$ ,

$$\phi_1, \dots, \phi_k \in \mathbb{R}^d, \quad S\phi_j = \lambda_j \phi_j$$

- ④ Project the data onto those  $k$  vectors:

$$x_i \mapsto \bar{x} + (\phi_1^T x_i) \phi_1 + \dots + (\phi_k^T x_i) \phi_k$$

# Let's go to the notebook

```
pca = PCA(num_components).fit(cimages)
principal_vectors = pca.components_
principal_vectors = principal_vectors.reshape((num_components, height, width))
pcs = pca.fit_transform(cimages)
capprox = pca.inverse_transform(pcs)
labels = ['principal vector %d' % (i+1) for i in np.arange(num_components)]
plot_images(principal_vectors, labels, height, width, int(num_components/5.), 5)
ratio = pca.explained_variance_ratio_.sum()
print('Variance explained by first %d principal vectors: %.2f%%' % (num_components, ratio*100))
```

Variance explained by first 25 principal vectors: 72.46%

principal vector 1    principal vector 2    principal vector 3    principal vector 4    principal vector 5



principal vector 6    principal vector 7    principal vector 8    principal vector 9    principal vector 10



# PCA: Summary

- PCA is an unsupervised method
- Finds directions of greatest variation in the data
- The directions are called the *principal vectors*; the weightings on the vectors are called the *principal components*
- The first few vectors may be interpretable
- Orthogonality makes interpretation difficult for the higher components
- Can be used for visualization or dimensionality reduction