

S&DS 265 / 565  
**Introductory Machine Learning**

# **PCA and Review**

October 13

# Plan for today

- Reminders
- Quick recap of PCA
- No new material
- Demo notebook
- Brief review for midterm

# Quiz 3

## Quiz Summary

Section Filter ▾

 Student Analysis

 Item Analysis

 Average Score

**93%**

 High Score

**100%**

 Low Score

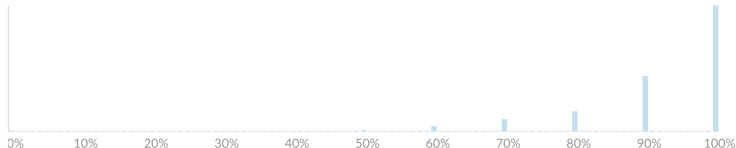
**50%**

 Standard Deviation

**1.04**

 Average Time

**13:25**

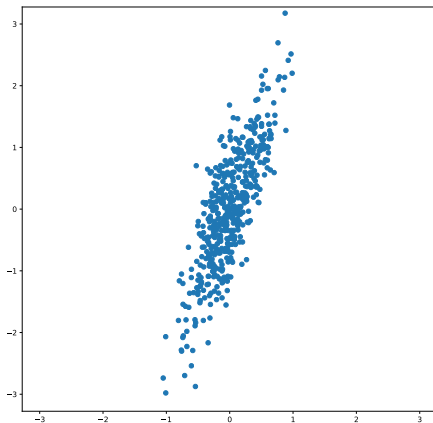


# Reminders

- Assn 2 due today at midnight; Assn 3 out
- Midterm next Tuesday, October 18, in class
- “Closed book, notes, computer...”
- $8\frac{1}{2} \times 11$  sheet of notes, handwritten double-sided
- Practice midterms posted on Canvas (with solutions)
- Will go over practice exams in review sessions
- Questions?

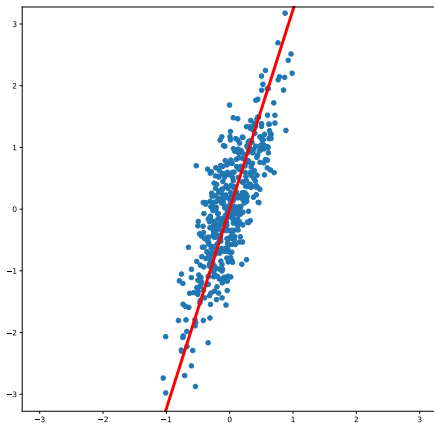
# Principal Component Analysis (PCA)

PCA finds the directions of greatest variability in the data.

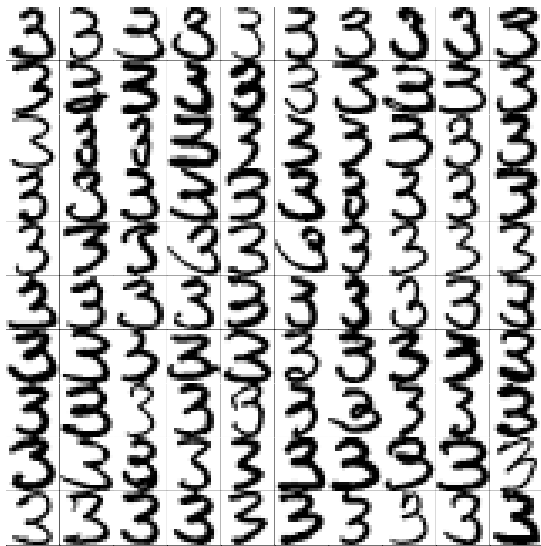


# Principal Component Analysis (PCA)

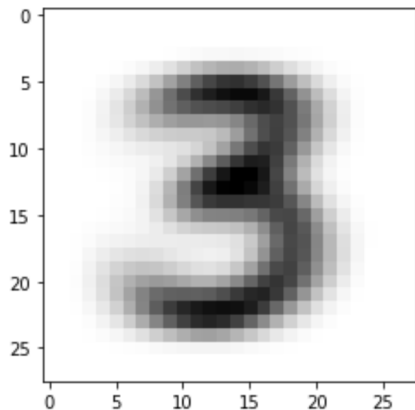
PCA finds the directions of greatest variability in the data.



## Handwritten Digits (3s)



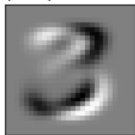
## Handwritten Digits (3s) – Average



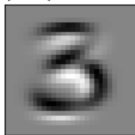


# Handwritten Digits (3s) – Principal vectors

principal vector 1



principal vector 2



principal vector 3



principal vector 4



principal vector 5



principal vector 6



principal vector 7



principal vector 8



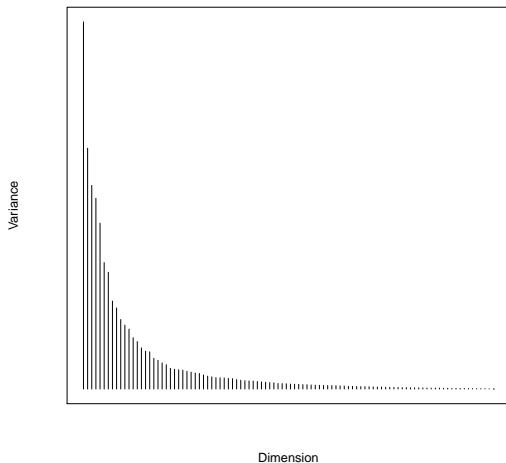
principal vector 9



principal vector 10



# Handwritten Digits (3s) – PCA variance

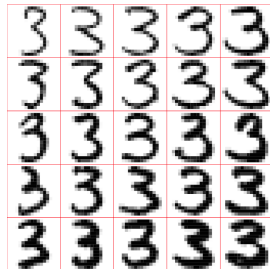
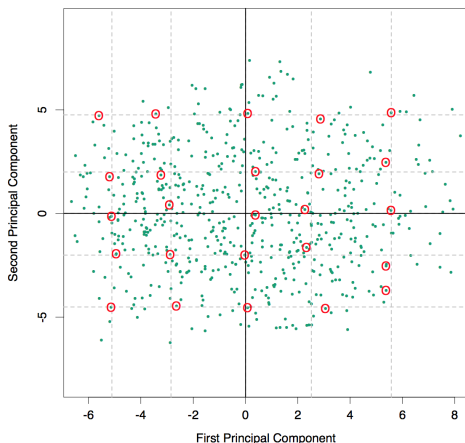


# Handwritten Digits (3s)

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

# Handwritten Digits (3s) – Top 2 components

$$\begin{aligned}\hat{f}(\lambda) &= \bar{x} + \lambda_1 v_1 + \lambda_2 v_2 \\ &= \text{[Image of mean digit 3]} + \lambda_1 \cdot \text{[Image of } v_1 \text{]} + \lambda_2 \cdot \text{[Image of } v_2 \text{]}.\end{aligned}$$



# 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  $S$ ,

$$v_1, \dots, v_k \in \mathbb{R}^d, \quad S v_j = \lambda_j v_j$$

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

$$x_i \mapsto \bar{x} + (v_1^T x_i) v_1 + \dots + (v_k^T x_i) v_k$$

# PCA: Algorithm

- ① We can compute everything directly
- ② Except for the eigenvectors
- ③ Let's illustrate this in the demo notebook

# Let's go to the notebook

The number  $x^T v_k$  is the amount of  $x$  that lies in the direction of the principal vector  $v_k$ . This is easily translated into Python:

```
In [26]: v = principal_vectors.reshape(num_components, height*width)

xhat = avgimg
for k in np.arange(num_components):
    xhat = xhat + np.dot(x, v[k]) * v[k]
    plot_face_reconstruction(x, xhat, 'George W. Bush', 'Reconstruction using %d vectors' % (k+1))
```

George W. Bush



Reconstruction using 100 vectors





# Using PCA for classification or regression

- A combination of supervised learning and unsupervised learning
- Given data  $\{x\}$  extract principal vectors and components
- Map each data point  $x_i$  to its principal components

$$z_i \equiv (x_i^T v_1, \dots, x_i^T v_K)$$

- For labeled data  $\{(x_i, y_i)\}$ , now train a supervised learning algorithm using the transformed data  $\{(z_i, y_i)\}$ .

# Example notebook

## Flower Power: PCA and classification (30 points)



In this problem you will carry out principal components analysis and classification on the iris data. The task will be to reduce the dimension from four to two using PCA, and then to train logistic regression models on the projected data.

# 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

# Review

Week	Dates	Topics	Demos & Tutorials	Lecture Slides	Readings and Notes	Assignments & Exams
1	Sept 1	Course overview		Sept 1: <a href="#">Course overview</a>		
2	Sept 6, 8	Python and background concepts	<a href="#">Python elements</a> <a href="#">Covid trends</a>	Sept 6: <a href="#">Python elements</a> Sept 8: <a href="#">Pandas and linear regression</a>	Data8 Chapters 3, 4, 5	Thu: <a href="#">Quiz 1</a>
3	Sept 13, 15	Linear regression and classification	<a href="#">Covid trends (revisited)</a> <a href="#">Classification examples</a>	Sept 13: <a href="#">Regression concepts</a> Sept 15: <a href="#">Classification</a>	ISL Sections 3.1, 3.2, 3.5 Notes on <a href="#">regression</a> ISL Sections 4.3, 4.4 <a href="#">Notes on classification</a>	Thu: <a href="#">Covid Assn1 out</a>
4	Sept 20, 22	Stochastic gradient descent	<a href="#">SGD examples</a>	Sept 20: <a href="#">Classification (continued)</a> Sept 22: <a href="#">Stochastic gradient descent</a>	ISL Section 6.2.2 ISL Section 10.7.2	Thu: <a href="#">Quiz 2</a>
5	Sept 27, 29	Bias and variance, cross-validation	<a href="#">Bias-variance tradeoff</a> <a href="#">Covid trends (revisited)</a> <a href="#">California housing</a>	Sept 27: <a href="#">Bias and variance</a> Sept 29: <a href="#">Cross-validation</a>	ISL Section 2.2 ISL Section 5.1	Thu: Assn 1 in <a href="#">Covid Assn2 out</a>
6	Oct 4, 6	Tree-based methods	<a href="#">Trees and forests</a> <a href="#">Visualizing trees</a> <a href="#">Bagging operations</a>	Oct 4: <a href="#">Trees</a> Oct 6: <a href="#">Forests</a>	ISL Sections 8.1, 8.2	Thu: <a href="#">Quiz 3</a>
7	Oct 11, 13	PCA and dimension reduction	<a href="#">PCA examples</a> <a href="#">PCA revisited</a> <a href="#">Used for regression</a>	Oct 11: <a href="#">PCA</a> Oct 13: <a href="#">PCA and review</a>	ISL Section 12.2	Thu: Assn 2 in <a href="#">Covid Assn3 out</a>

**Questions?**