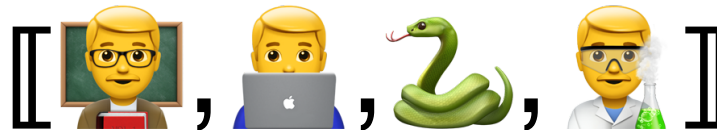


Lecture Notes for **Machine Learning in Python**



Professor Eric Larson
Visualization and Dimensionality Reduction

Class Logistics and Agenda

- Dimensionality Reduction
 - PCA
 - Randomized PCA
 - Images Representation with PCA

Class Overview, by topic

Table Data
Visualization

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

Dimension
Reduction and
Image Processing

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

Linear and
Logistic
Regression

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

Neural Networks
and Back Prop.

Numpy
Detailed mathematics for NN operations

Wide and Deep
Networks

Convolutional
Networks

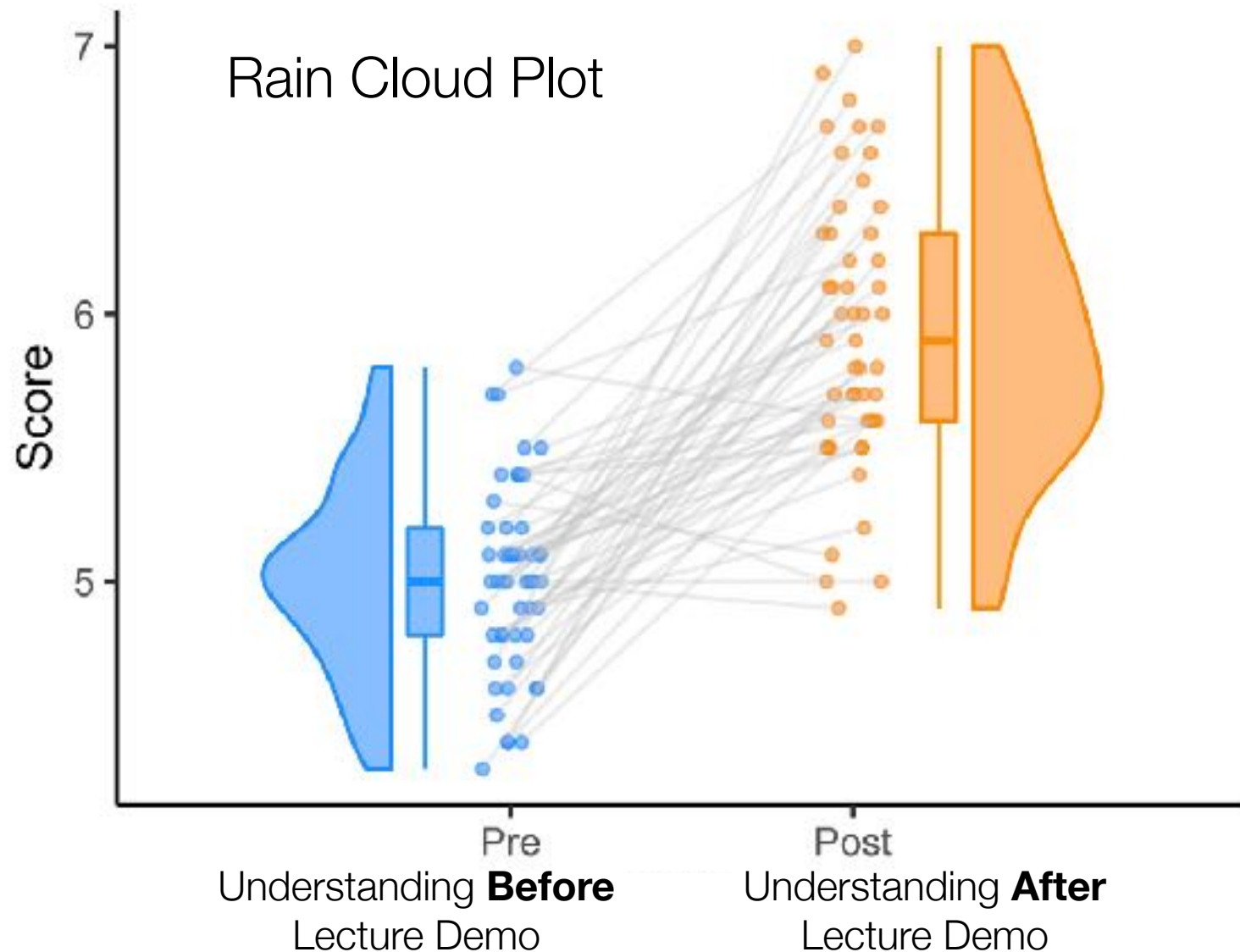
Recurrent
Networks

Keras, Tensorflow
Intuition, Detailed implement.

Ethics in
Language Models

ConceptNet
Case studies

Last time: visualization



Dimensionality Reduction: PCA



Kyle 🚀 🐬 🪐 🦖 @KyleMorgens... · 1d ...

eigenvalues are just the TLDR for a matrix

💬 38

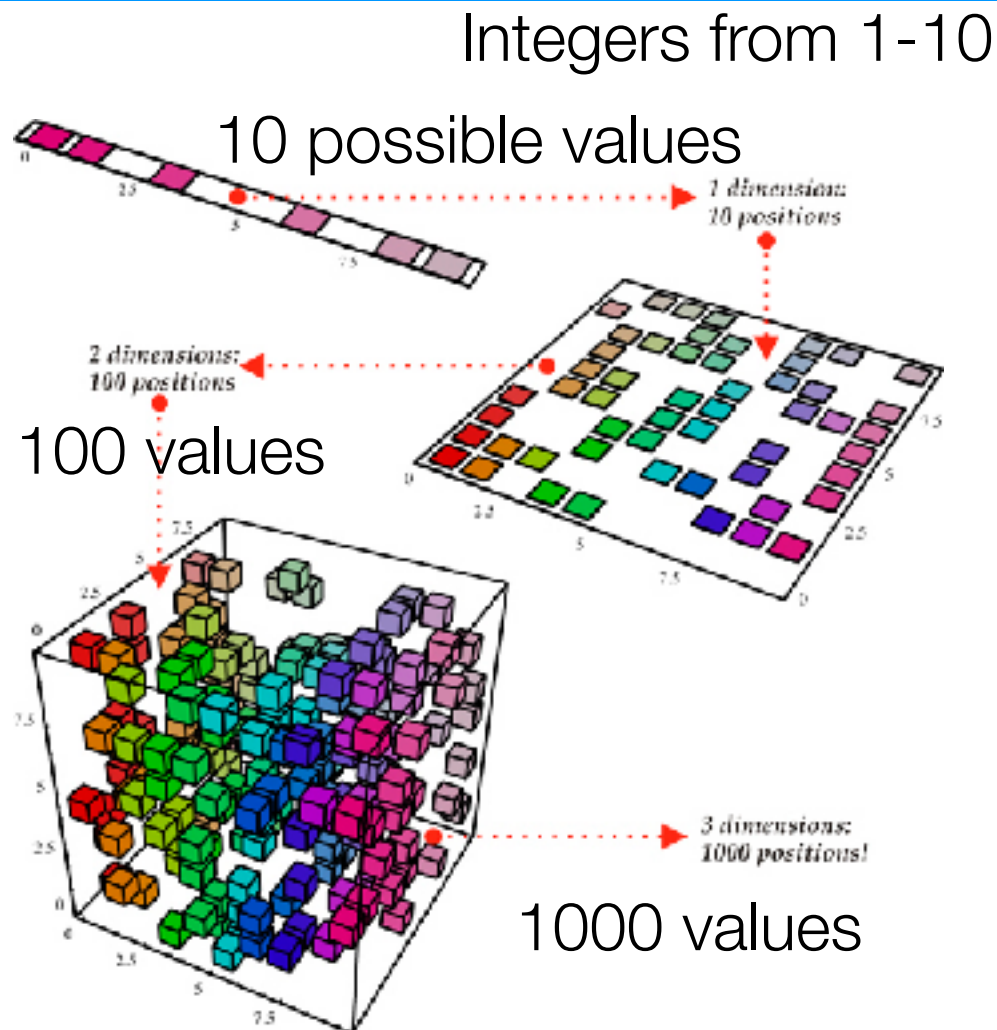
↻ 602

❤️ 6,046



Curse of Dimensionality

- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful



Dimensionality Reduction

- Purpose:
 - Avoid curse of dimensionality
 - Select subsets of independent features
 - Reduce amount of time and memory required by data mining algorithms
 - Allow data to be more easily visualized
 - May help to eliminate irrelevant features or reduce noise
- Techniques
 - Principle Component Analysis
 - Non-linear PCA
 - Stochastic Neighbor Embedding



I invented PCA...
and *Social Darwinism*



Aside: Eigen Vectors are your friend!

Three Blue One Brown:

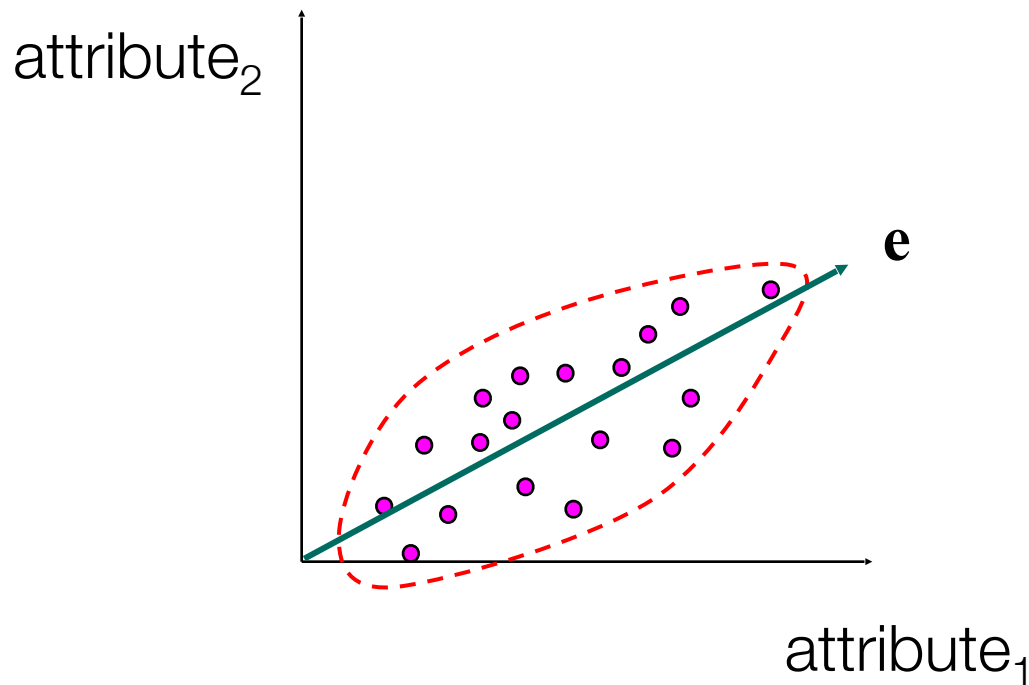
<https://www.youtube.com/watch?v=PFDu9oVAE-g>

Eigen-things aren't
actually so bad



Dimensionality Reduction: PCA

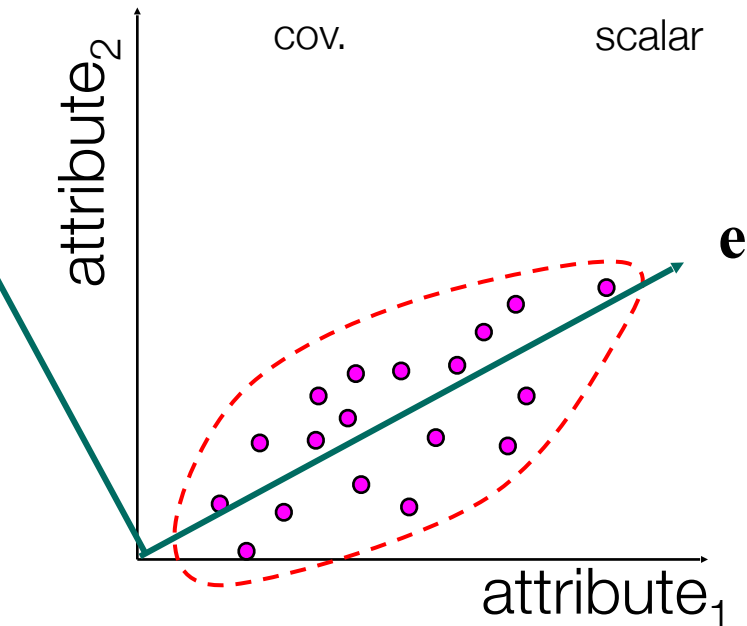
- Goal is to find a projection that captures the largest amount of variation in data



Dimensionality Reduction: PCA

- Find the **eigenvectors** of the **covariance** matrix
- keep the “k” **largest** eigenvectors

$$\underset{\text{cov.}}{\mathbf{C}} \cdot \underset{\text{scalar}}{\mathbf{e}} = \lambda \mathbf{e}$$



| $E1$ | $E2$ |
|-----------------|----------------|
| 0.749 | 0.662 |
| 0.662 | -0.749 |
| $\lambda=268.3$ | $\lambda=1.57$ |

covariance

| | |
|-------|-------|
| 151.5 | 132.4 |
| 132.4 | 118.3 |

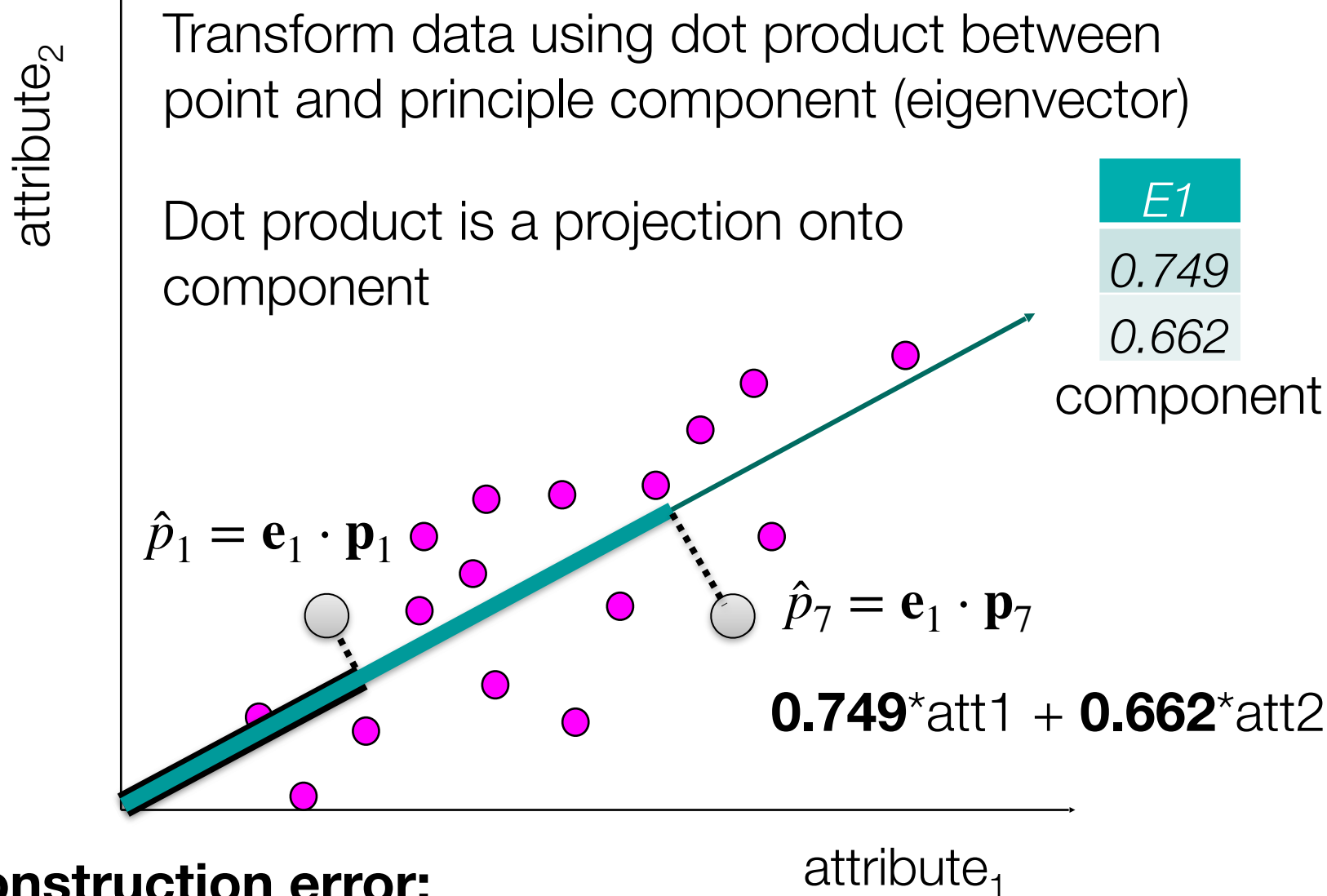
| | A1 | A2 |
|---|------|------|
| 1 | 14 | 12.6 |
| 2 | 26 | 26.6 |
| 3 | 36.3 | 33.3 |
| 4 | 2.5 | 3.6 |
| 5 | 15 | 17.4 |
| 6 | 8 | 11 |



| | A1 | A2 |
|---|--------|--------|
| 1 | -2.96 | -4.82 |
| 2 | 9.03 | 9.18 |
| 3 | 19.33 | 15.88 |
| 4 | -14.46 | -13.82 |
| 5 | -1.96 | -0.02 |
| 6 | -8.96 | -6.42 |

normalize: zero mean
optional: unit std

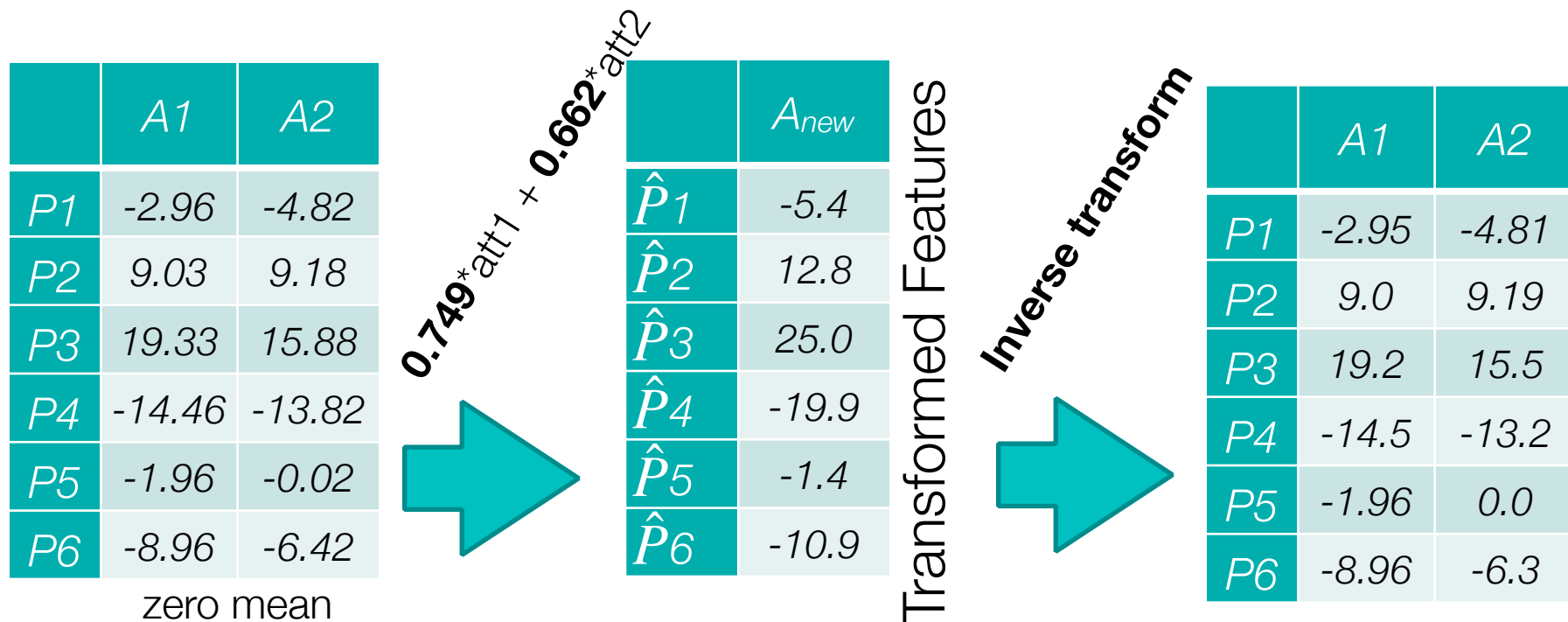
Dimensionality Reduction: PCA



Reconstruction error:

difference between projection and original point in 2D space

Dimensionality Reduction: PCA



This projection is called a **Transform**
known as the **Karhunen-Loève Transform (KLT)**

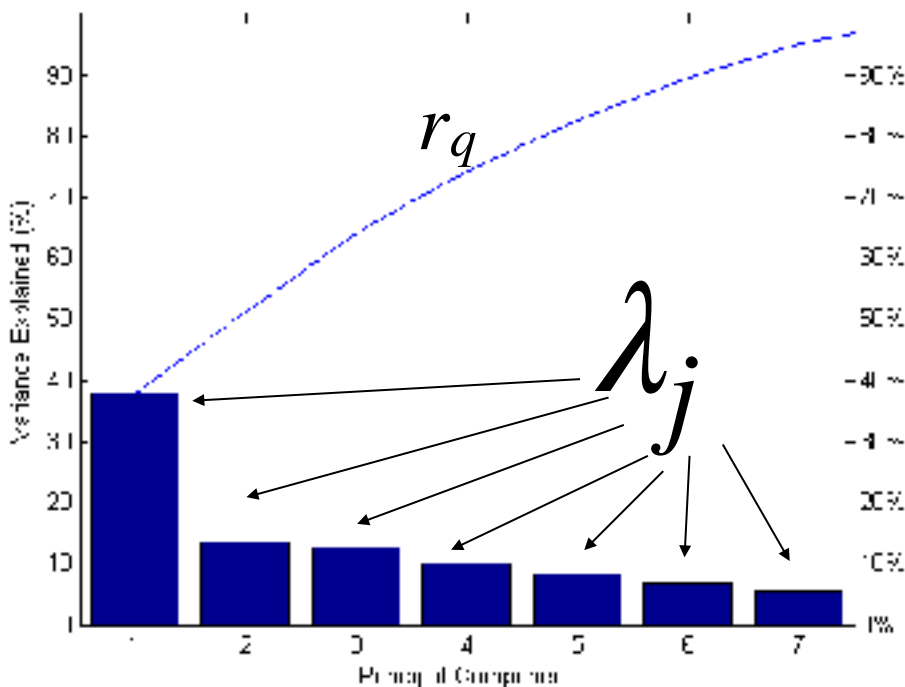
Explained Variance

- Each principle component **explains** a certain **amount of variation** in the data.
- This explained variation is **encoded** in the **eigenvalues** of each **eigenvector**

sum of q largest eigenvalues

$$r_q = \frac{\sum_{j=1}^q \lambda_j}{\sum_{\forall i} \lambda_i}$$

sum of all eigenvalues



Dimensionality Reduction: PCA

Genetic profiles distilled to 2 components

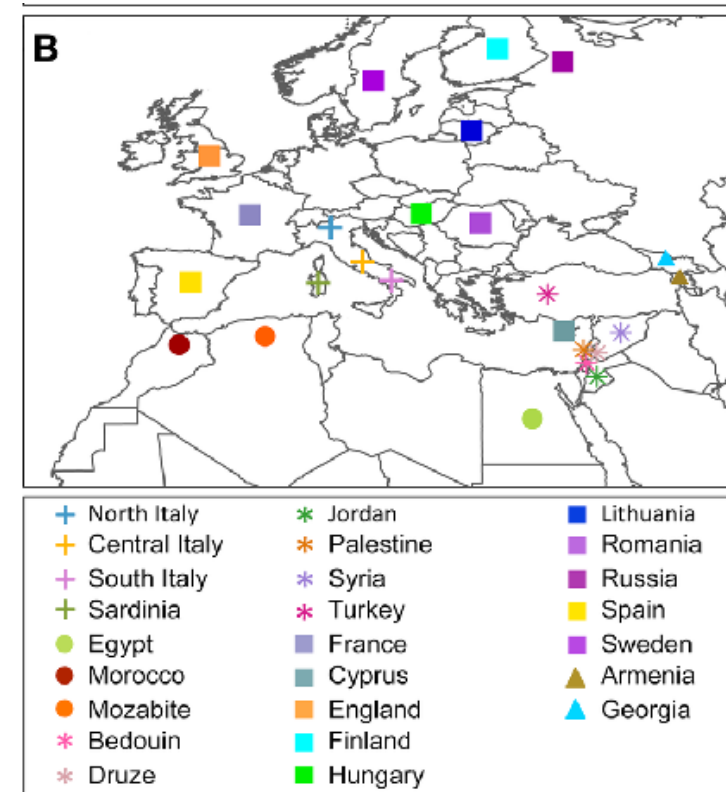
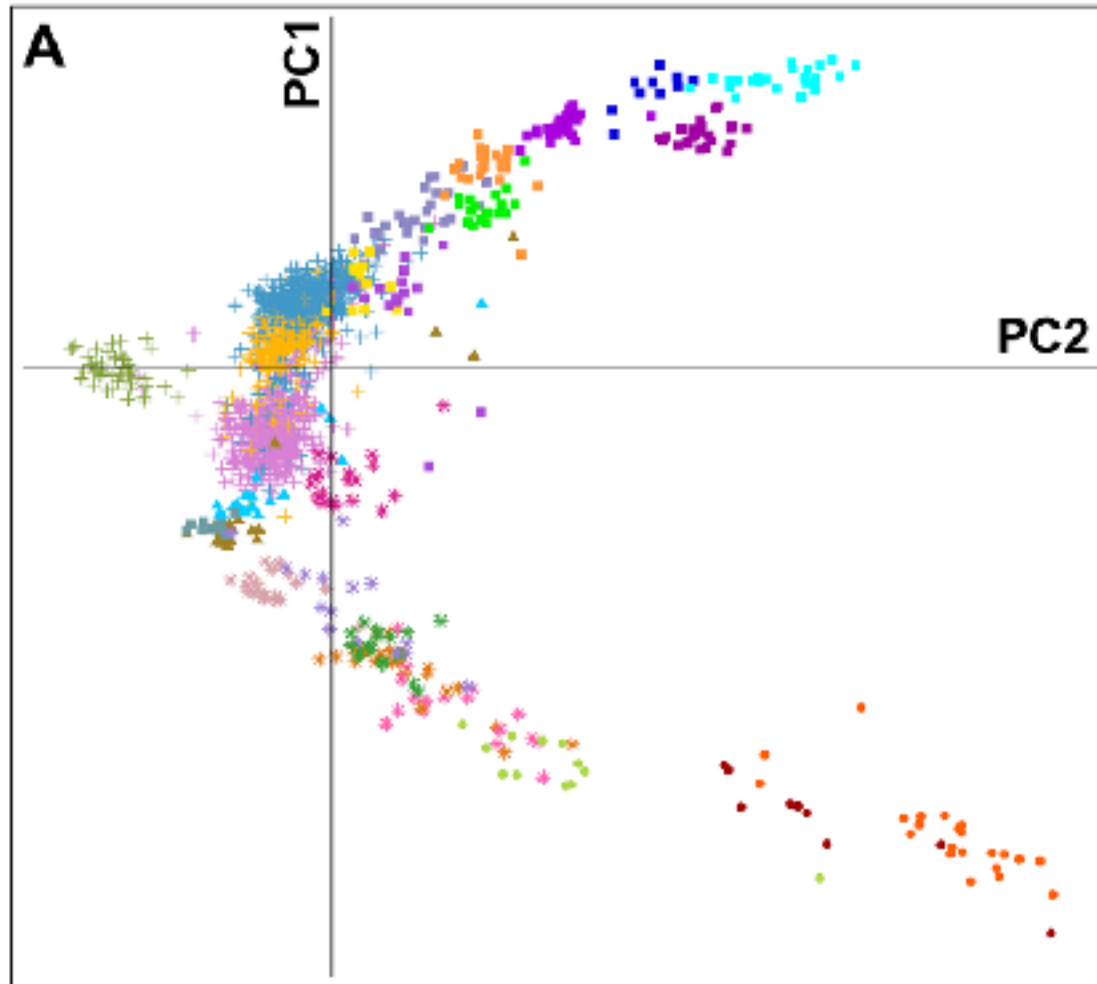


image source: Wikipedia 14

04.Dimension Reduction and Images.ipynb

PCA
biplots



Other Tutorials:

http://scikit-learn.org/stable/auto_examples/decomposition/plot_pca_vs_lda.html#example-decomposition-plot-pca-vs-lda-py

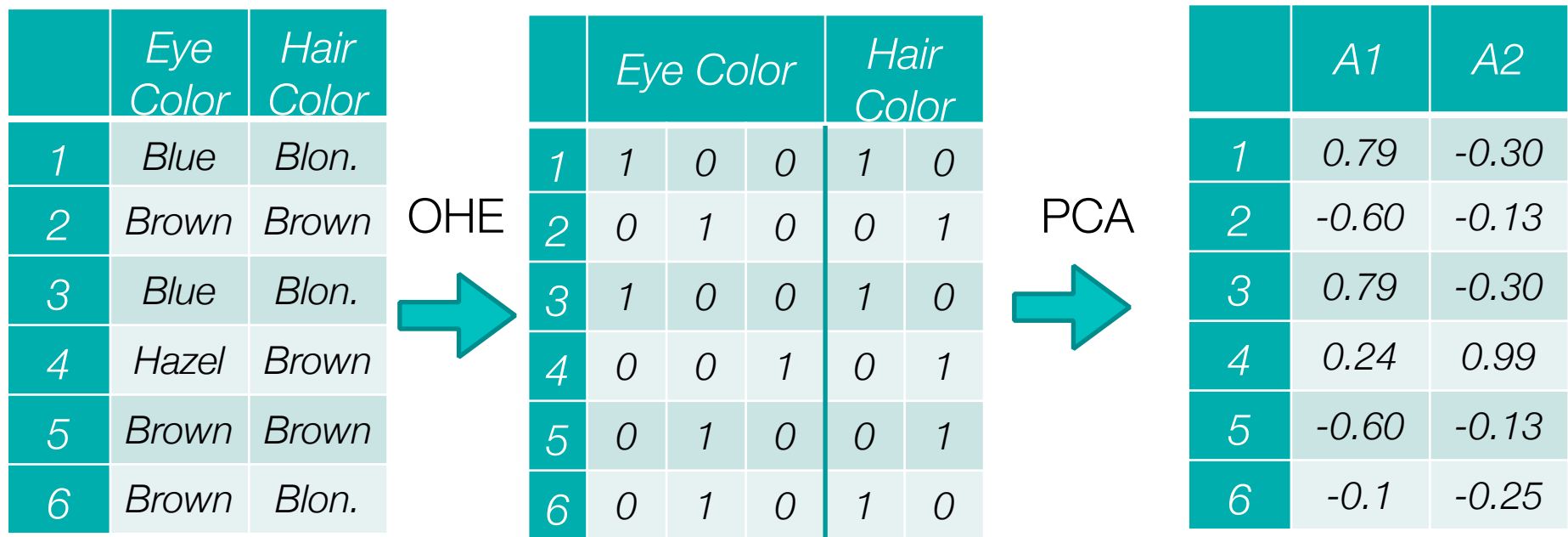
<http://nbviewer.ipython.org/github/ogrisel/notebooks/blob/master/Labeled%20Faces%20in%20the%20Wild%20recognition.ipynb>

Self Test ML2b.1

Principal Components Analysis works well for categorical data by design.

- A. True
- B. False
- C. It doesn't but people do it anyway

Mutual Correspondence Analysis



Dimensionality Reduction: Randomized PCA

- **Problem:** PCA on all that data can take a while to compute
 - What if the number of instances is gigantic?
 - What if the number of dimensions is gigantic?
- What if we partially construct the covariance matrix with a lower rank matrix?
 - By **transforming** our table data, \mathbf{A} , with another orthogonal matrix, \mathbf{Q} , we can **approximate** the **covariance matrix**, but with **lower rank**
 - Gives a matrix with typically good enough precision of actual eigenvectors, like using SVD. $\mathbf{Q}\mathbf{Q}^T\mathbf{A}$ is surrogate

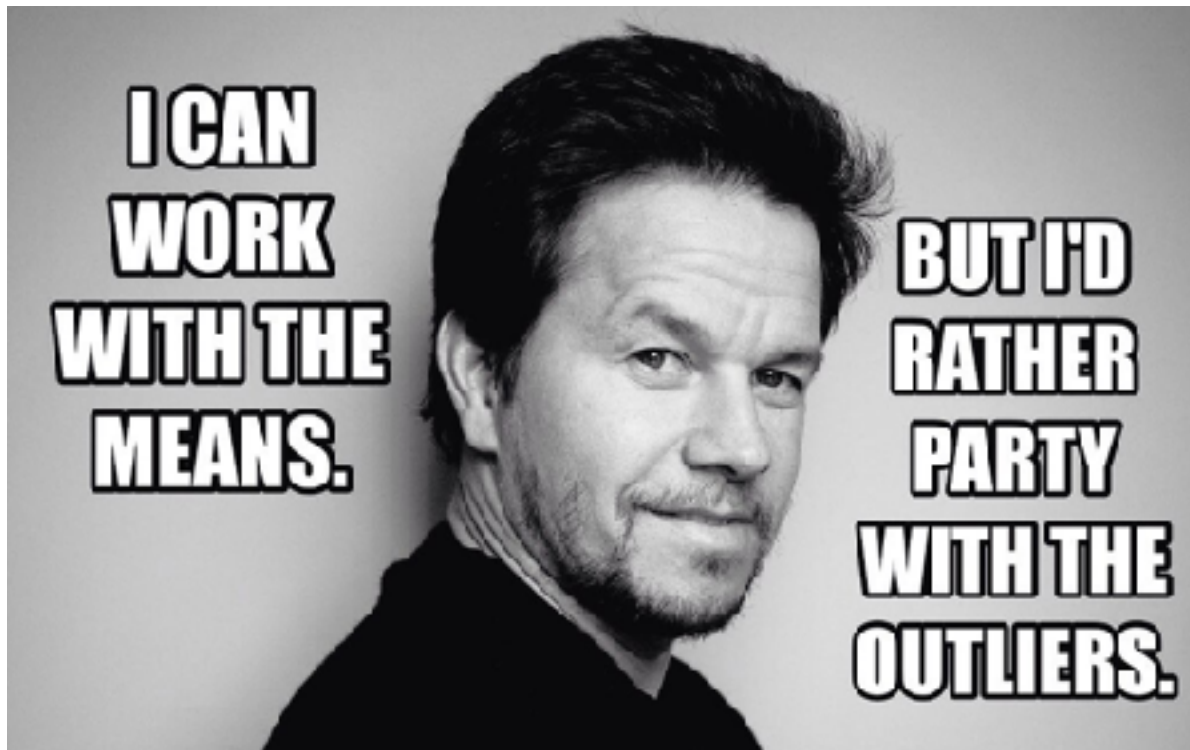
Example
Objective

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \left[1 + 11\sqrt{k+p} \cdot \sqrt{\min\{m,n\}}\right] \sigma_{k+1}$$

Halko, et al., (2009) Finding structure with randomness: Stochastic algorithms for constructing approximate matrix decompositions. <https://arxiv.org/pdf/0909.4061.pdf>

Image Processing and Representation

Our first @ResearchMark meme



Images as data

- an image can be represented in many ways
- most common format is a matrix of pixels
 - each “pixel” is BGR(A)
- used for capture and display

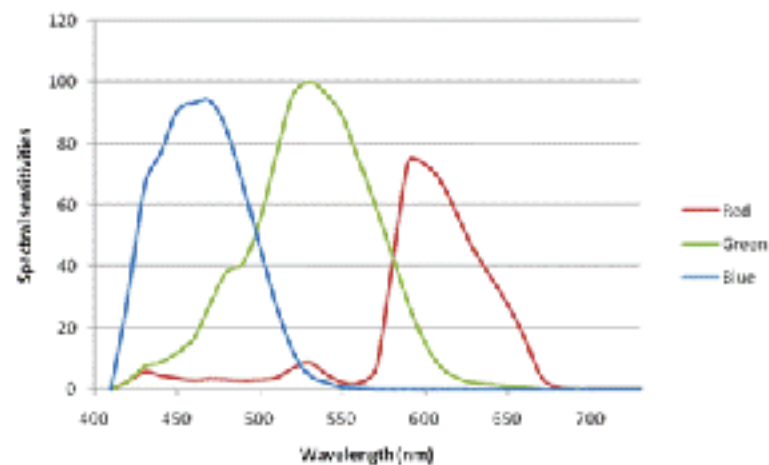
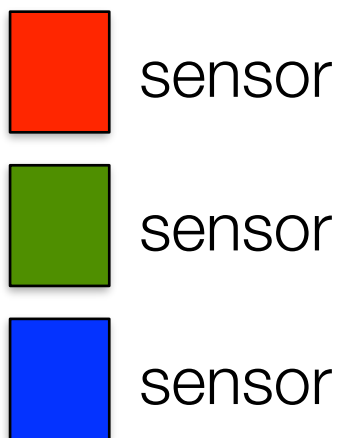
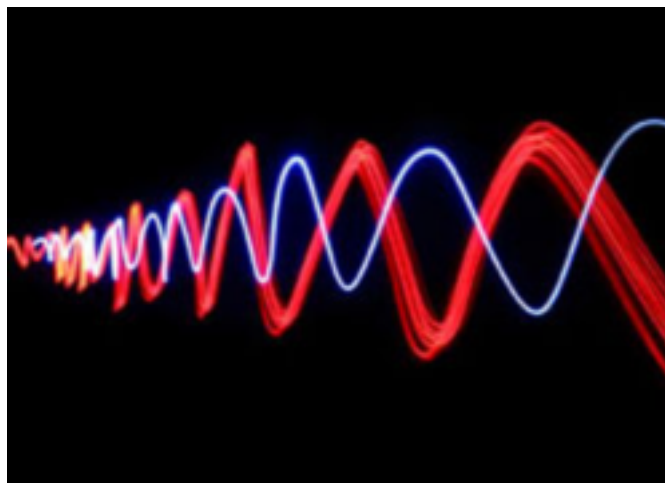
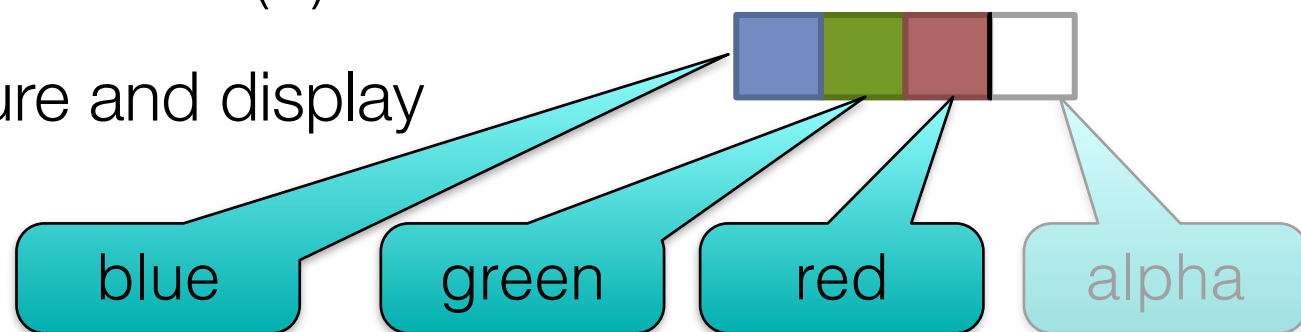


Image Representation

- need a compact representation

- **grayscale**

$$0.3 \cdot R + 0.59 \cdot G + 0.11 \cdot B,$$

“luminance”

gray

| | | | | | |
|---|---|---|---|---|---|
| 1 | 4 | 2 | 5 | 6 | 9 |
| 1 | 4 | 2 | 5 | 5 | 9 |
| 1 | 4 | 2 | 8 | 8 | 7 |
| 3 | 4 | 3 | 9 | 9 | 8 |
| 1 | 0 | 2 | 7 | 7 | 9 |
| 1 | 4 | 3 | 9 | 8 | 6 |
| 2 | 4 | 2 | 8 | 7 | 9 |

Numpy Matrix
`image[rows, cols]`

on

R

G

B

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|--|--|
| | | 1 | 4 | 2 | 5 | 6 | 9 | | |
| | 1 | 4 | 2 | 5 | 6 | 9 | 9 | | |
| 1 | 4 | 2 | 5 | 6 | 9 | 9 | 7 | | |
| 1 | 4 | 2 | 5 | 5 | 9 | 7 | 8 | | |
| 1 | 4 | 2 | 8 | 8 | 7 | 8 | 9 | | |
| 3 | 4 | 3 | 9 | 9 | 8 | 9 | 6 | | |
| 1 | 0 | 2 | 7 | 7 | 9 | 6 | 9 | | |
| 1 | 4 | 3 | 9 | 8 | 6 | 9 | | | |
| 2 | 4 | 2 | 8 | 7 | 9 | | | | |

Numpy Matrix
`image[rows, cols, channels]`

Image Representation, Features

Problem: need to represent image as table data

- need a compact representation

| | | | | | |
|---|---|---|---|---|---|
| 1 | 4 | 2 | 5 | 6 | 9 |
| 1 | 4 | 2 | 5 | 5 | 9 |
| 1 | 4 | 2 | 8 | 8 | 7 |
| 3 | 4 | 3 | 9 | 9 | 8 |
| 1 | 0 | 2 | 7 | 7 | 9 |
| 1 | 4 | 3 | 9 | 8 | 6 |
| 2 | 4 | 2 | 8 | 7 | 9 |

Image Representation, Features

Problem: need to represent image as table data

- need a compact representation

Solution: row concatenation (also, vectorizing)

| | | | | | | | | | | | | | | | | | | | |
|-------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Row 1 | 1 | 4 | 2 | 5 | 6 | 9 | 1 | 4 | 2 | 5 | 5 | 9 | 1 | 4 | 2 | 8 | 8 | 7 | 3 |
| Row 2 | 1 | 4 | 2 | 8 | 8 | 7 | 3 | 4 | 3 | 9 | 9 | 8 | 1 | 4 | 2 | 5 | 5 | 9 | 1 |
| ... | | | | | | | | | | | | | | | | | | | |
| Row N | 9 | 4 | 6 | 8 | 8 | 7 | 4 | 1 | 3 | 9 | 2 | 1 | 1 | 5 | 2 | 1 | 5 | 9 | 1 |

Self test: 3a-1

- When vectorizing images into table data, each “feature column” corresponds to:
 - a. the value (color) of pixel
 - b. the spatial location of a pixel in the image
 - c. the size of the image
 - d. the spatial location and color channel of a pixel in an image

Images Representation
in PCA and
Randomized PCA



04.Dimension Reduction and Images.ipynb

For Next Lecture

- Next Lecture:
 - Finish Dimension Reduction Demo
 - Crash-course Image Feature Extraction