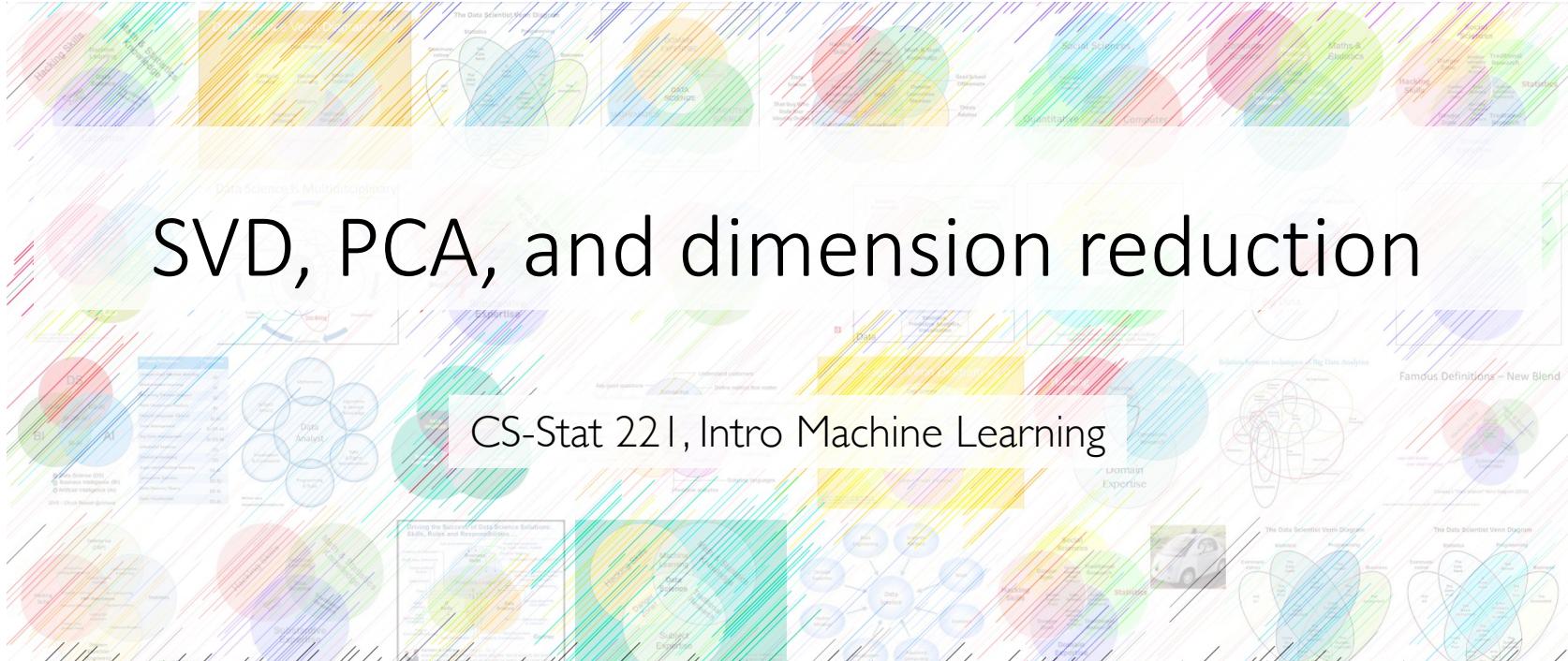


# SVD, PCA, and dimension reduction

CS-Stat 221, Intro Machine Learning



THE UNIVERSITY OF  
CHICAGO

## Consider genotype data

## Individuals (few 1000s)

# Consider genotype data

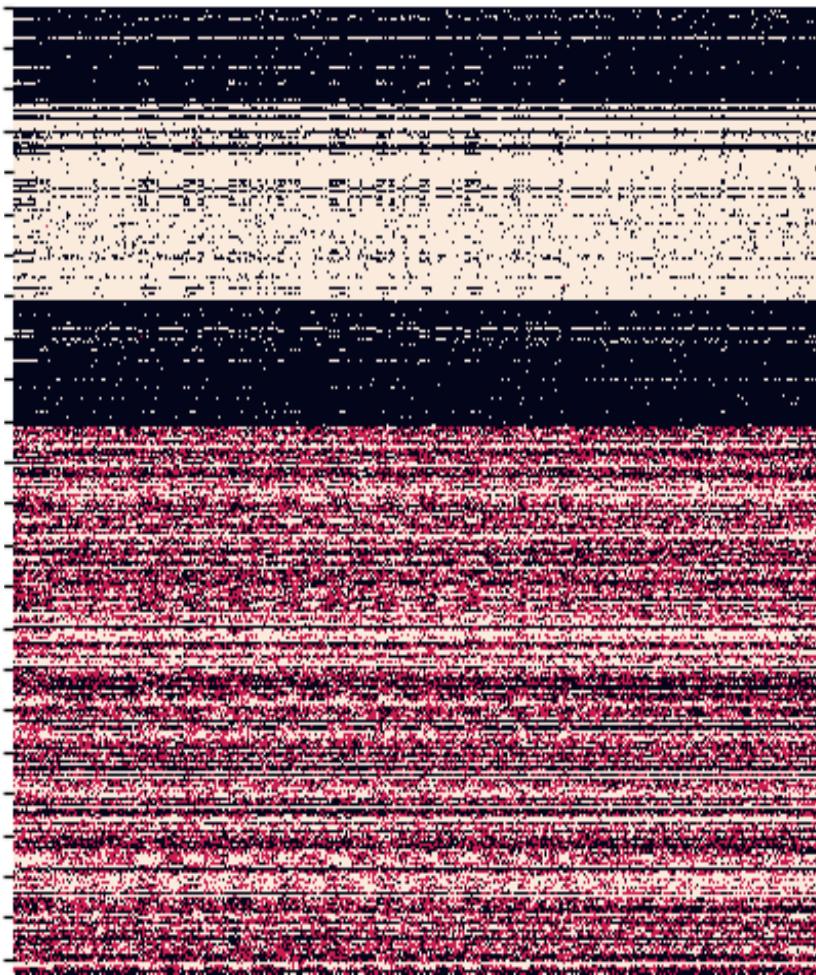
## Individuals (few 1000s)

Suppose your data looks like genotype data

## Individuals (few 1000s)

Genetic loci

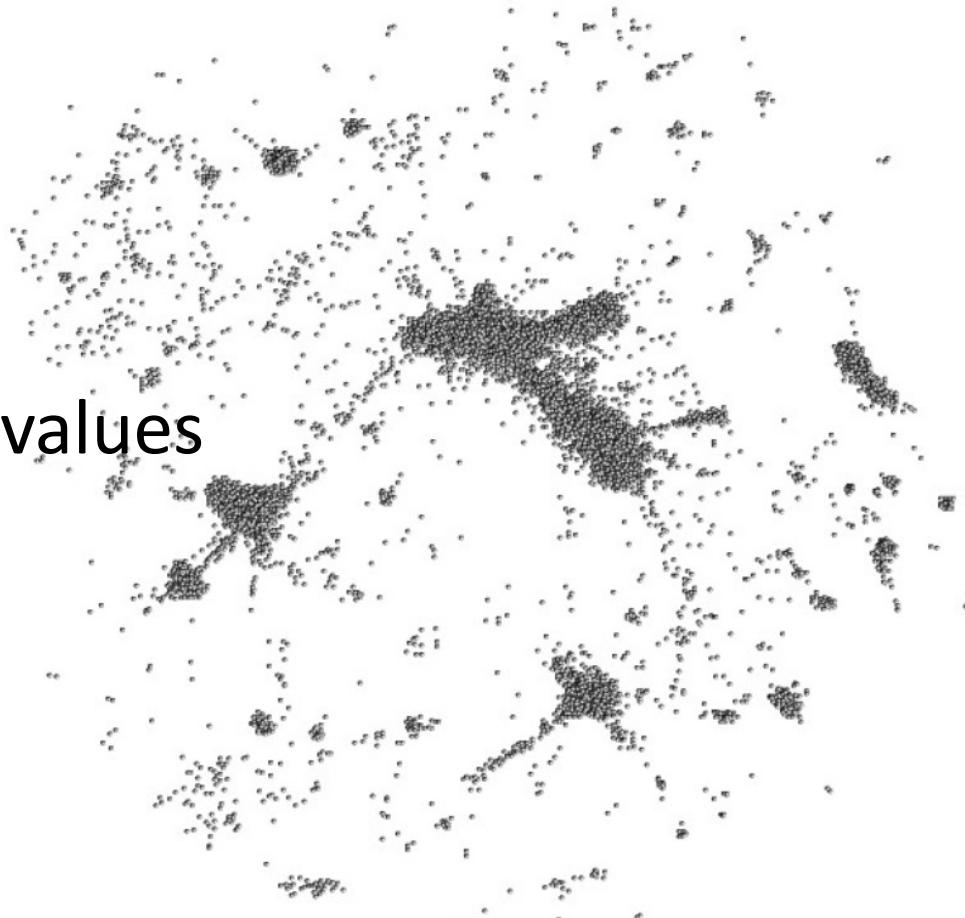
Individuals



- Some rows are similar to each other, some columns are similar to each other
- The data are in some way “predictable”
- Can we build a low-complexity object that explains most of the data matrix.

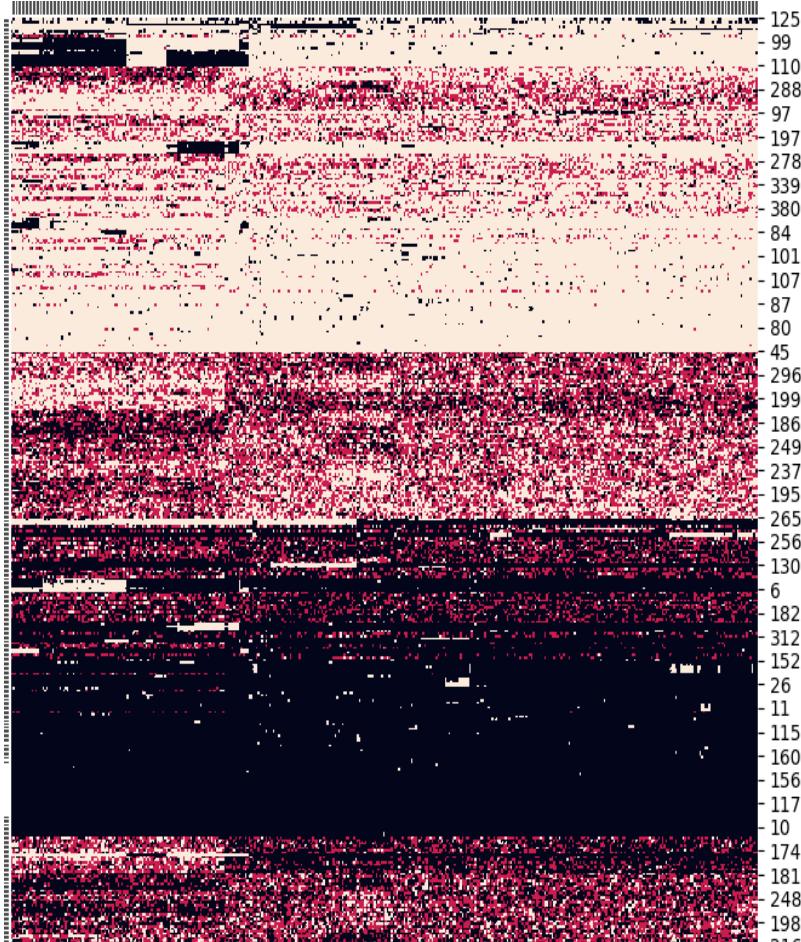
# Roadmap

- Dimension reduction
- Absolute essentials of eigenvectors and eigenvalues
- SVD and PCA
- Examples



Individuals

Genetic loci



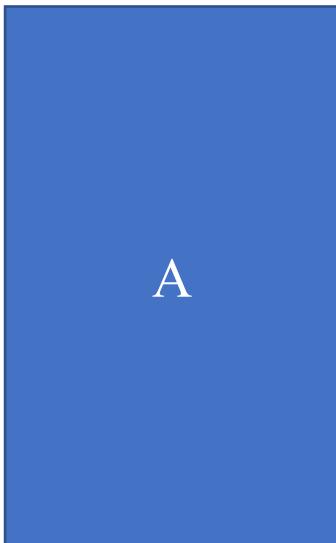
# Now that looks like structure

- Reorder rows and columns to put similar rows next to each other?
- Libraries to make exactly this sort of eye-candy are mature. This one is `sanborn.clustermap`
- Organizing data like this is deeply satisfying—it makes us think we have discovered truth.
- What if I told you linear algebra could do this?

# Why Reduce dimensions?

n rows  
"samples"

m columns  
"features"



- Reducing the number of dimensions can make ML easier

- May reduce noise.

- Can exploit unlabeled data to make nicer inputs for ML

- We need dimension reduction and categorization to understand high-dimensional data

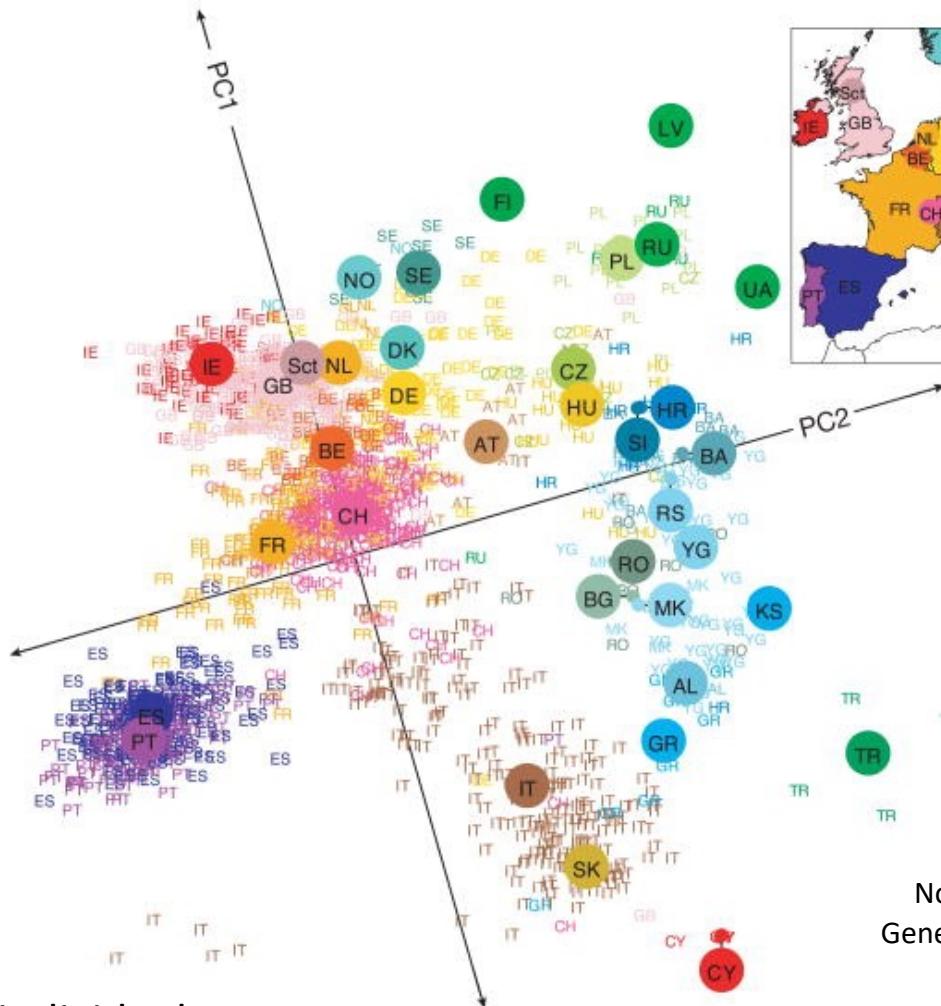
- Useful for making visualizations of high-dimensional data.

r  
"reduced features"

n samples



Example:  
Genes mirror  
geography within  
Europe



Novembre et al, 2008  
Genes mirror geography  
within Europe

$m \times n = 200k$  genetic loci  $\times 1400$  individuals

10.1038/nature07331

# How many dimensions again?

$n \times m$  can be large (compared to our computers and our displays):

Twitter (500M tweets / day, May 2020)

Youtube (2.5 billion views / day)

Natural Language processing ( $10^4$ - $10^6$  words)

Netflix prize was 500K users x 20K movies

Biochemical or genetic assays ( $10^4$ - $10^7$  features of interest x sample size)

Image processing ( $10^7$  pixels easy)

# Definition of eigenvectors and eigenvalues

There exists a set of vectors  $\{e_i\}$

That has the following nice  
property

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

square matrix  
 $n \times n$

ith eigenvalue

ith eigenvector  
 $n \times 1$

square  
matrix

$$\begin{matrix} & n \\ \mathbf{A} & \end{matrix} \quad = \quad \begin{matrix} & \lambda_i \\ e & \end{matrix}$$

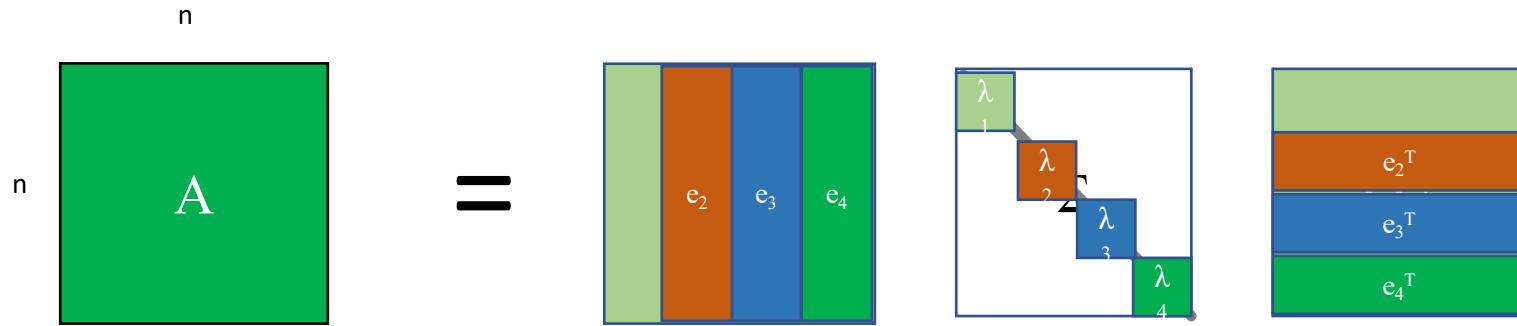
$n \times n$        $n \times 1$        $n \times 1$

# Eigenvalue decomposition

These vectors can be constructed  
orthogonal (zero inner products)

and normal ( $x \cdot x = 1$ )

$$A = \sum_i e_i^T \lambda_i e_i$$



square  
matrix

eigenvectors  
(orthogonal)

eigenvalues  
(diagonal)

eigenvectors  
(orthogonal)

$A$

$Q^{-1}$

$\Sigma$

$Q$

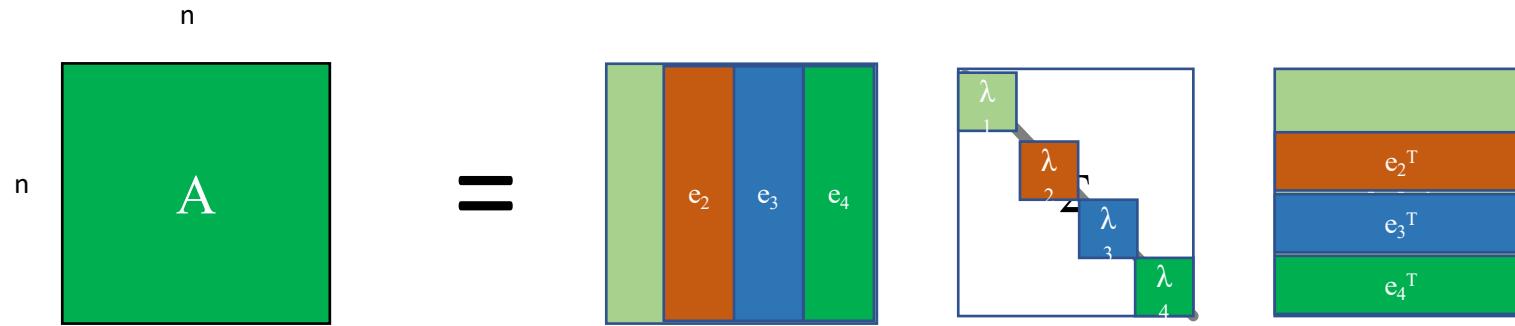
# Eigenvalue decomposition

```
import numpy as np
from numpy import linalg as LA

# A library function will take a square matrix and give
# back a vector of eigenvalues and a matrix of eigenvectors.

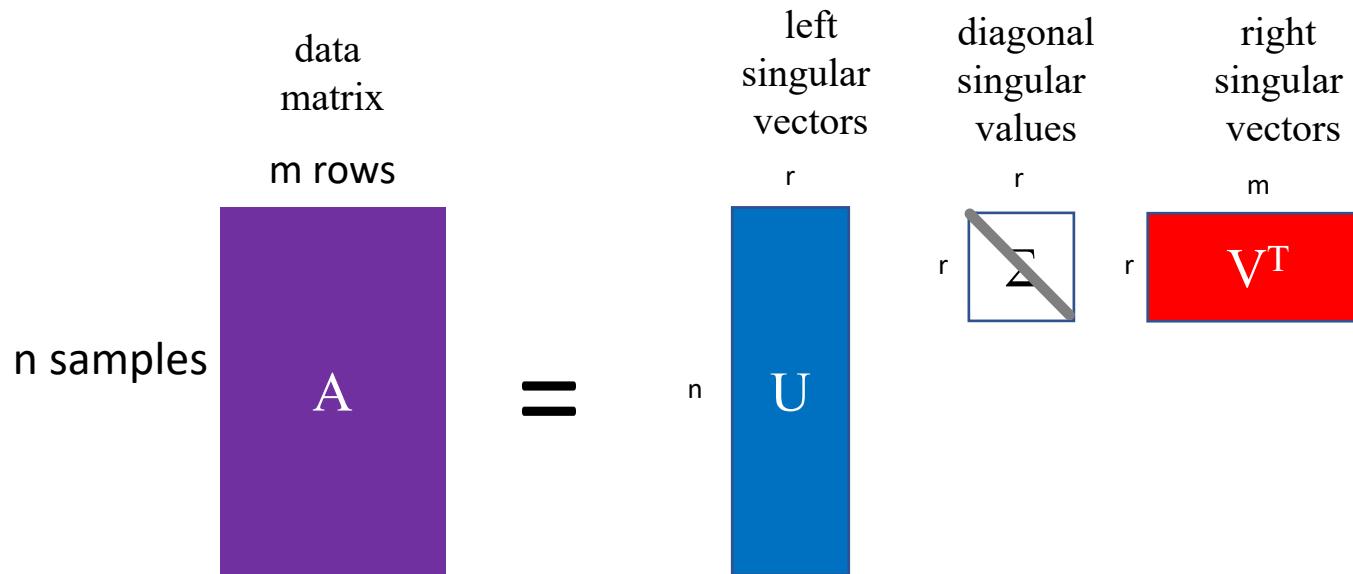
A=np.array( [ [16, 2, 3, 13 ],[5, 11, 10, 8], [9, 7, 6, 12],
[4, 14, 15, 1] ] )
evals, evectors = LA.eig(A)
print(evals)
print(evectors)
B=np.dot(np.dot(evectors , evals * np.eye(4)) ,
LA.inv(evectors) )
```

# Singular Value Decomposition (SVD)



- Eigenvalue decomposition is very nice, but only applies to square data matrices.
- Singular Value Decomposition (SVD) is a generalization for rectangular matrices

# Singular Value Decomposition (SVD)



- Eigenvalue decomposition is very nice, but only applies to square data matrices.
- Singular Value Decomposition (SVD) is a generalization for rectangular matrices

# Constructing SVD - rows

- Recipe for building SVD:

$$\begin{matrix} m \\ n \end{matrix} \quad A \quad \begin{matrix} m \\ n \end{matrix} \quad A^T \quad = \quad \begin{matrix} n \\ n \end{matrix} \quad AA^T$$

The diagram illustrates the matrix multiplication  $A \cdot A^T$ . On the left, there is a dark purple rectangular box labeled 'A' with dimensions  $m \times n$  above it. To its right is a dark purple rectangular box labeled  $A^T$  with dimensions  $n \times m$  above it. Between them is a black double equals sign. To the right of the multiplication is a light blue rectangular box labeled  $AA^T$  with dimensions  $n \times n$  above it.

This is the matrix of inner products between rows;  
summed over all the samples

# Constructing SVD - rows

$$\begin{matrix} n \\ \text{AA}^T \end{matrix} = \begin{matrix} n \\ U \end{matrix} \begin{matrix} \sigma \\ \vdots \end{matrix} \begin{matrix} U^T \end{matrix}$$

- Now take the eigenvalue decomposition of  $\text{AA}^T$

# Constructing SVD - rows

$$\begin{matrix} & n \\ \text{AA}^T & = \end{matrix} \begin{matrix} U & 0 \\ \cancel{\sigma} & 0 \end{matrix} \begin{matrix} U^T \\ 0 \end{matrix}$$

- Now take the eigenvalue decomposition of  $\text{AA}^T$
- I have a bunch of zero eigenvalues, so I have a bunch of eigenvectors that can't contribute to the sum

# Constructing SVD - rows

$$\begin{matrix} n \\ \text{AA}^T \end{matrix} = \begin{matrix} r \\ U \end{matrix} \quad \begin{matrix} r \\ \sigma \end{matrix} \quad \begin{matrix} n \\ U^T \end{matrix}$$

The diagram illustrates the Singular Value Decomposition (SVD) of a matrix  $\mathbf{A}$ . On the left, a blue rectangle labeled  $\mathbf{AA}^T$  represents an  $n \times n$  matrix. An equals sign follows. To its right is another blue rectangle labeled  $\mathbf{U}$ , representing an  $n \times r$  matrix. To the right of the equals sign is a smaller blue square labeled  $\sigma$ , representing an  $r \times r$  diagonal matrix with a single non-zero element  $\sigma$ . A diagonal line through the  $\sigma$  matrix indicates it is a scalar multiple of the identity matrix. To the right of the  $\sigma$  matrix is a final blue rectangle labeled  $\mathbf{U}^T$ , representing an  $r \times n$  matrix.

- Now take the eigenvalue decomposition of  $\mathbf{AA}^T$
- $\mathbf{U}$  ( $n \times r$ ) ,  $\sigma$  ( $r$ ),  $\mathbf{U}^T$  ( $r \times n$ )

# Constructing SVD - columns

$$\begin{matrix} & n \\ m & \end{matrix} \quad A^T \quad \begin{matrix} m \\ n \end{matrix} \quad A \quad = \quad \begin{matrix} m \\ m \end{matrix} \quad A^T A$$

- This matrix has dimensions  $m\_features \times m\_features$

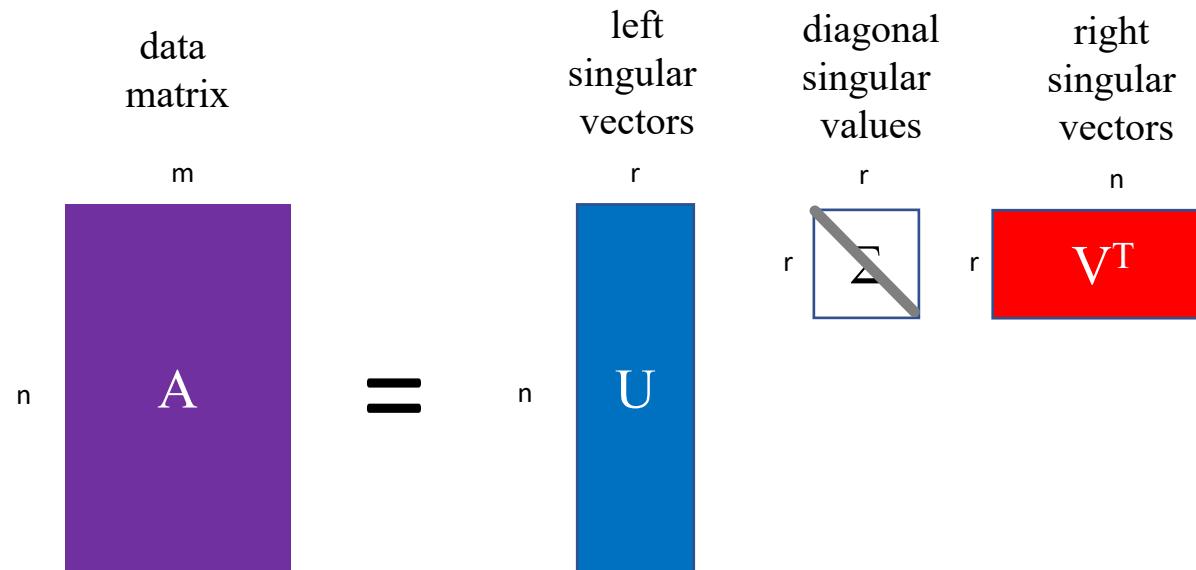
# Constructing SVD - columns

Eigenvalue decomposition again

$$\begin{matrix} & m \\ \begin{matrix} & m \\ A^T A \end{matrix} & \end{matrix} = \begin{matrix} & m \\ \begin{matrix} & r \\ V \end{matrix} & \end{matrix} \begin{matrix} & r \\ \begin{matrix} & r \\ \sigma \end{matrix} & \end{matrix} \begin{matrix} & r \\ \begin{matrix} & m \\ V^T \end{matrix} & \end{matrix}$$

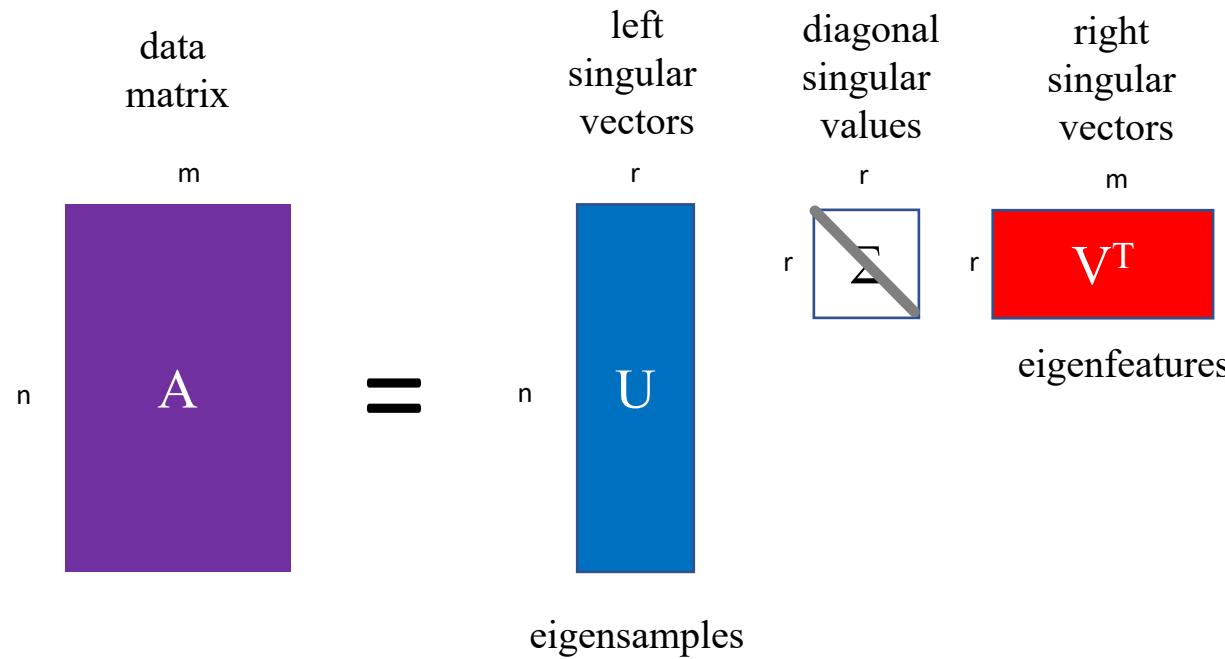
- Because of the relationship between  $A^T A$  and  $AA^T$  these two matrices have the same eigenvalues. (But one is  $m \times m$  and one is  $n \times n$  !!)
- Some of the eigenvalues have to be zero in the larger matrix

# Singular Value Decomposition (SVD)



- I have two sets of basis vectors and one set of  $r$  singular values.
- Two powers of  $A$  -> elements of  $\Sigma$  are square roots of eigenvalues

# Singular Value Decomposition (SVD)



- I have two sets of basis vectors and one set of  $r$  singular values.
- Two powers of  $A$  -> elements of  $\Sigma$  are square roots of eigenvalues

# SVD

```
import numpy as np
from numpy import linalg as LA
#
https://numpy.org/doc/stable/reference/generated/numpy.linalg.svd.html
```

```
U, SIGMA VT = LA.svd(A)
```

We could calculate SVD  
ourselves in two lines –  
but we shouldn't

# Approximate, truncated SVD

```
import numpy as np
from numpy import linalg as LA
from sklearn.utils.extmath import
randomized_svd
```

```
# https://scikit-learn.org/stable/modules/generated/sklearn.utils.extmath.randomized\_svd.html
```

```
U, SIGMA, VT = randomized_svd(A,
n_components=50)
```

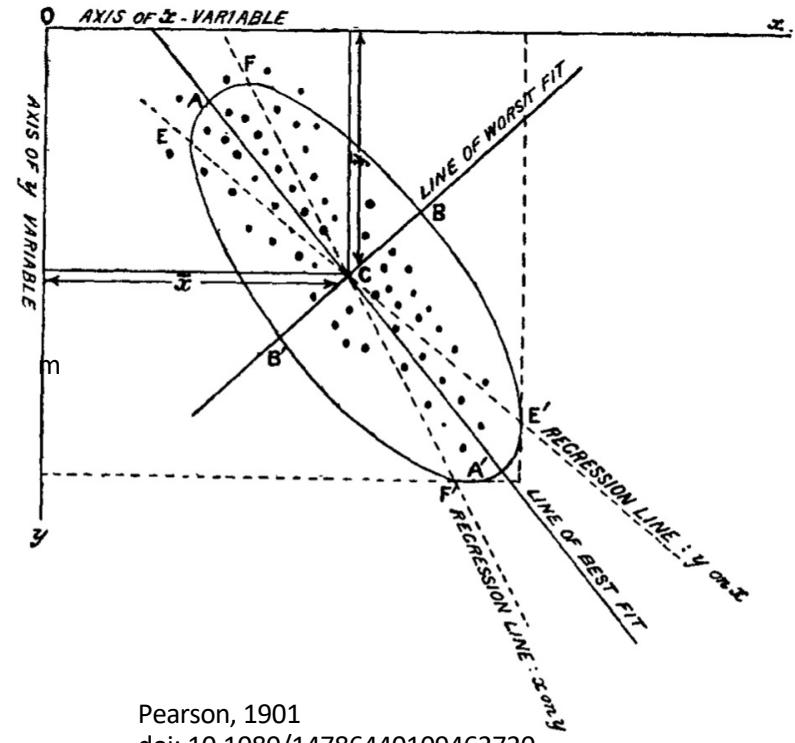
The library functions use convergence or optimization to find approximate matrix factorizations.

Does not actually build or store  $A^T A$

Approximate and truncated decompositions are cheaper to compute.

# Principal Component analysis

- SVD is a just linear algebra + geometry interpretation of the columnwise & rowwise inner products.
- PCA is the doctrine of using vectors derived from SVD to interpret data.
- The data points are rotated onto the singular vectors, renamed PCA coordinates
- $XV^T = T$  (new coordinates)  
 $(n \times m) * (m \times r) = (n \times r)$
- Can make pretty pictures



Pearson, 1901

doi: 10.1080/14786440109462720

**Table 1**  
*Principal Components Analysis of BIS-11 Items (Oblique Rotation)*

| BIS-11 items                                  | First-order factors |      |      |      |      |      |
|---|---------------------|------|------|------|------|------|
|   | 1                   | 2    | 3    | 4    | 5    | 6    |
| 11. I "squirm" at plays or lectures.          | .84                 | .17  | -.08 | -.03 | .03  | .02  |
| 32. I am restless at the theater or lectures. | .84                 | .19  | -.12 | -.06 | -.00 | -.03 |
| 5. I don't "pay attention."                   | .57                 | .04  | .16  | -.02 | .27  | .02  |
| 17. I act "on impulse."                       | .15                 | .74  | .08  | -.02 | -.20 | .06  |
| 20. I act on the spur of the moment.          | .12                 | .72  | .19  | -.10 | -.19 | .01  |
| 23. I buy things on impulse.                  | -.08                | .59  | -.04 | .28  | .10  | .11  |
| 12. I am a careful thinker. <sup>a</sup>      | .17                 | -.13 | .64  | .17  | -.18 | .05  |
| 1. I plan tasks carefully. <sup>a</sup>       | -.05                | .16  | .64  | -.04 | .11  | -.10 |
| 8. I am self-controlled. <sup>a</sup>         | .10                 | .00  | .63  | -.24 | .08  | -.17 |

Patton et al. 1995 Factor structure of the Barratt Impulsiveness Scale  
doi: 10.1002/1097-4679(199511)51:6<768::AID-JCLP2270510607>3.0.CO;2-1

**Table 1**  
*Principal Components Analysis of BIS-11 Items (Oblique Rotation)*

| BIS-11 items                                  | First-order factors |      |      |      |      |      |
|---|---------------------|------|------|------|------|------|
|   | 1                   | 2    | 3    | 4    | 5    | 6    |
| 11. I “squirm” at plays or lectures.          | .84                 | .17  | -.08 | -.03 | .03  | .02  |
| 32. I am restless at the theater or lectures. | .84                 | .19  | -.12 | -.06 | -.00 | -.03 |
| 5. I don’t “pay attention.”                   | .57                 | .04  | .16  | -.02 | .27  | .02  |
| 17. I act “on impulse.”                       | .15                 | .74  | .08  | -.02 | -.20 | .06  |
| 20. I act on the spur of the moment.          | .12                 | .72  | .19  | -.10 | -.19 | .01  |
| 23. I buy things on impulse.                  | -.08                | .59  | -.04 | .28  | .10  | .11  |
| 12. I am a careful thinker. <sup>a</sup>      | .17                 | -.13 | .64  | .17  | -.18 | .05  |
| 1. I plan tasks carefully. <sup>a</sup>       | -.05                | .16  | .64  | -.04 | .11  | -.10 |
| 8. I am self-controlled. <sup>a</sup>         | .10                 | .00  | .63  | -.24 | .08  | -.17 |

Patton et al. 1995 Factor structure of the Barratt Impulsiveness Scale

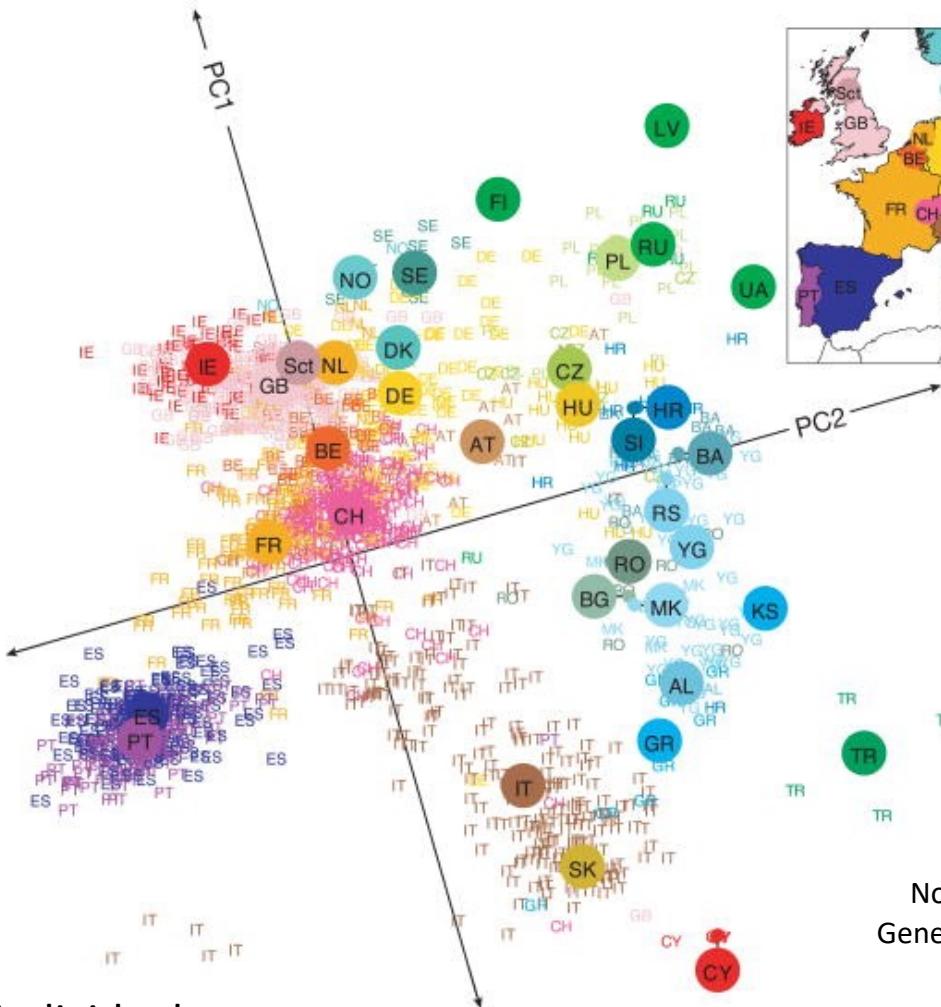
doi: 10.1002/1097-4679(199511)51:6<768::AID-JCLP2270510607>3.0.CO;2-1

**Table 1**  
*Principal Components Analysis of BIS-11 Items (Oblique Rotation)*

| BIS-11 items                                  | First-order factors |      |      |      |      |      |
|---|---------------------|------|------|------|------|------|
|   | 1                   | 2    | 3    | 4    | 5    | 6    |
| 11. I “squirm” at plays or lectures.          | .84                 | .17  | -.08 | -.03 | .03  | .02  |
| 32. I am restless at the theater or lectures. | .84                 | .19  | -.12 | -.06 | -.00 | -.03 |
| 5. I don’t “pay attention.”                   | .57                 | .04  | .16  | -.02 | .27  | .02  |
| 17. I act “on impulse.”                       | .15                 | .74  | .08  | -.02 | -.20 | .06  |
| 20. I act on the spur of the moment.          | .12                 | .72  | .19  | -.10 | -.19 | .01  |
| 23. I buy things on impulse.                  | -.08                | .59  | -.04 | .28  | .10  | .11  |
| 12. I am a careful thinker. <sup>a</sup>      | .17                 | -.13 | .64  | .17  | -.18 | .05  |
| 1. I plan tasks carefully. <sup>a</sup>       | -.05                | .16  | .64  | -.04 | .11  | -.10 |
| 8. I am self-controlled. <sup>a</sup>         | .10                 | .00  | .63  | -.24 | .08  | -.17 |

Data correlations here used to inform reduction of dimension from 35-question survey to 6 scores.

# Genes mirror geography within Europe



Novembre et al, 2008  
Genes mirror geography  
within Europe

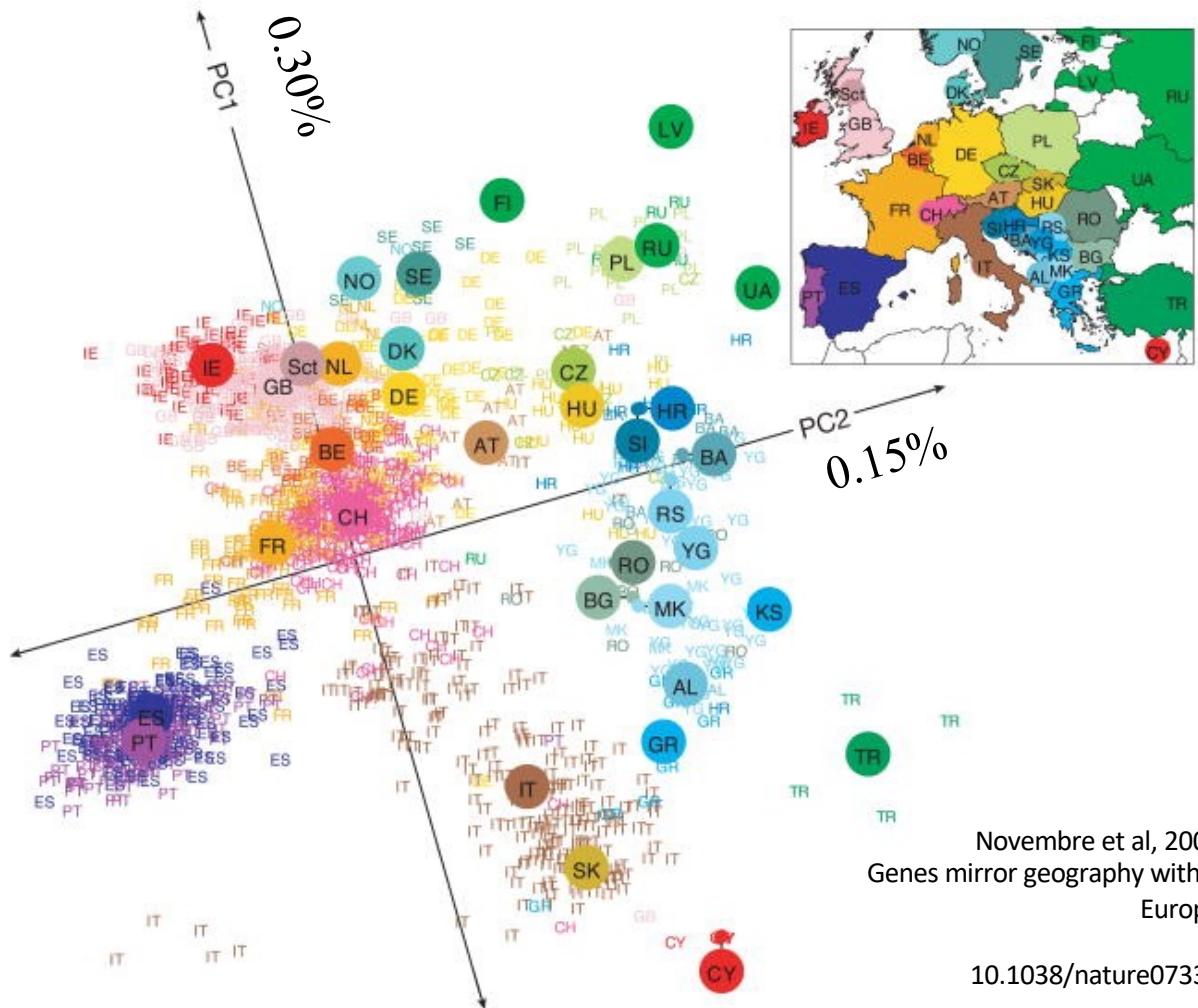
$m \times n = 200k$  genetic loci  $\times 1400$  individuals

10.1038/nature07331

# Genes mirror geography within Europe

Principal components and eigenvalues come back from largest to smallest magnitude.

It is conventional to report principal components along with the fraction of variance explained by each



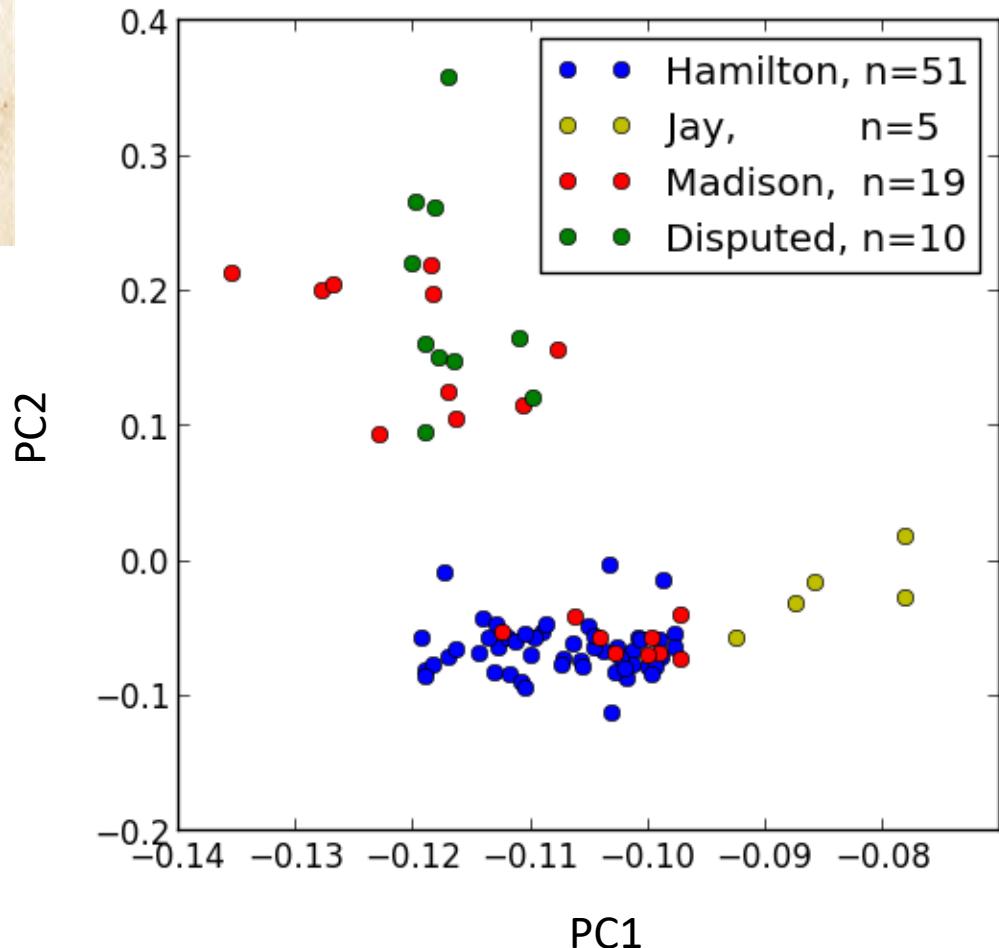
FEDERALIST:  
A COLLECTION OF  
ESSAYS,

- 85 documents, initially published anonymously 1787-1788
  - Authors' reminiscences disagree about authorship of some of the essays
  - 85 documents parsed into vocabulary of ~10000 “words”

# What if your data looks like this?

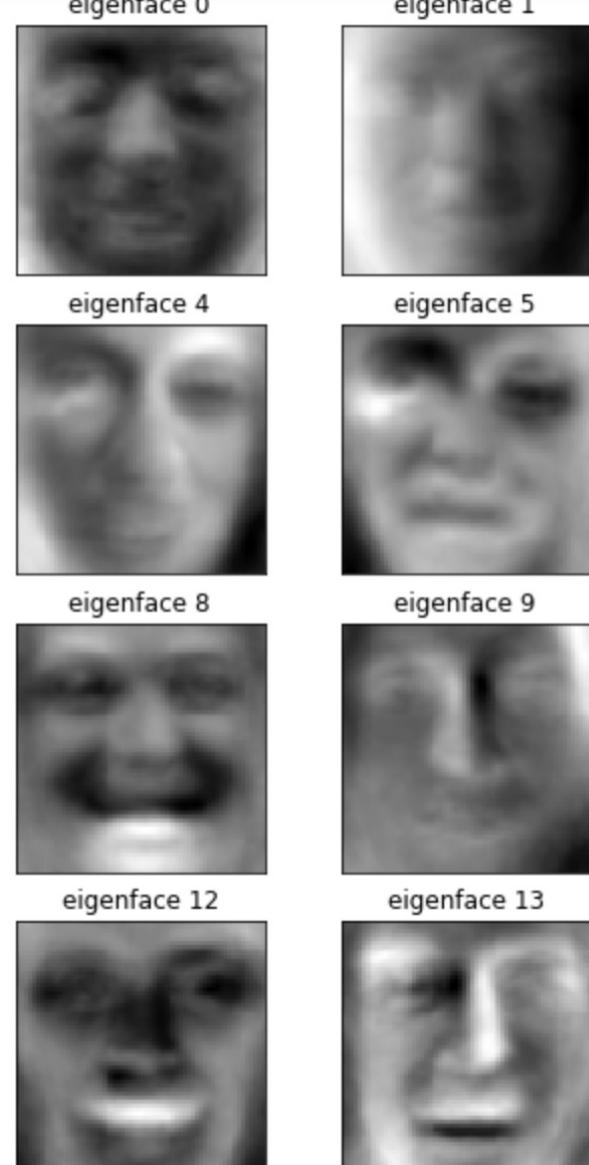
THE  
FEDERALIST:  
A COLLECTION OF  
ESSAYS,

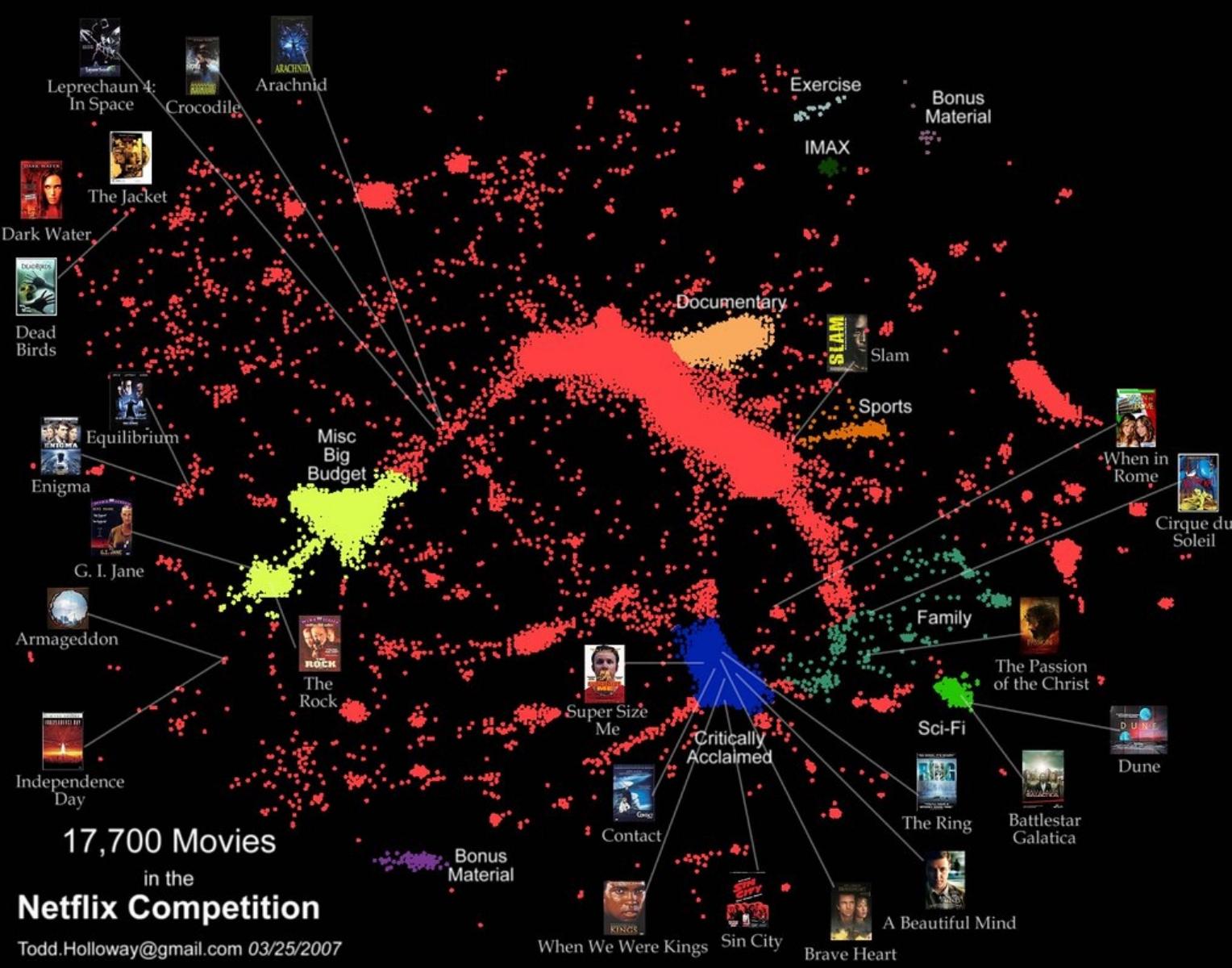
- Sometimes  $\sum(x - \bar{x})(y - \bar{y})$  is not an appropriate distance score for your data.
- You can run PCA on an  $n \times n$  distance matrix evaluated using a domain-specific columnwise comparison.
- Resulting low-dimensional coordinates visualize vocabulary/style relationships



# Caveats

- SVD/PCA don't know anything but the cloud of points and its correlations.
- Each new sample builds a different set of basis vectors
- Hard to interpret: positive and negative values in the eigenvectors
- Does not handle non-numerical values—have to change the objective function for that (Netflix)





# Dimension reduction approaches

- **Hierarchical clustering (built into clustermap):** takes pairwise distances as input; outputs a tree.
- **PCA/SVD (Principal Components Analysis, Singular Value Decomposition)**  
Linear algebra hammer. Finds vectors with maximum projection of covariance. This is a linear transformation of the data.
- **MDS Multidimensional Scaling:** Finds an embedding (low-dimensional projection) that optimally preserves pairwise distances.
- **tSNE (t-distributed stochastic neighbor embedding )** : more elaborate transformations into an embedding space optimized to put same-label points near each other.
- The APIs to find the **three embeddings** look the same, and the lore says for small problems the three approaches give similar results.

# These library functions work by optimization

- Even though PCA is just math, the just-math calculation approach is often really expensive. ( $n^3$  complexity, for instance, in finding determinant.)
- Particularly for large problems, close-enough algorithms can be made to work faster and on larger datasets than we could otherwise handle.
- These close-enough approximation algorithms often use optimization; don't be surprised when your matrix inversion says it is "searching for a good fit."



# Linear algebra objects

$\lambda$

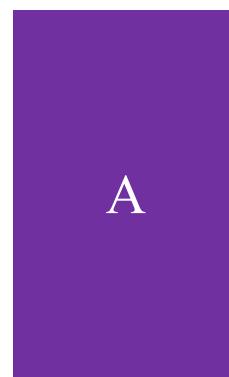
number  
“scalar”



$m \times 1$

vector

$n$



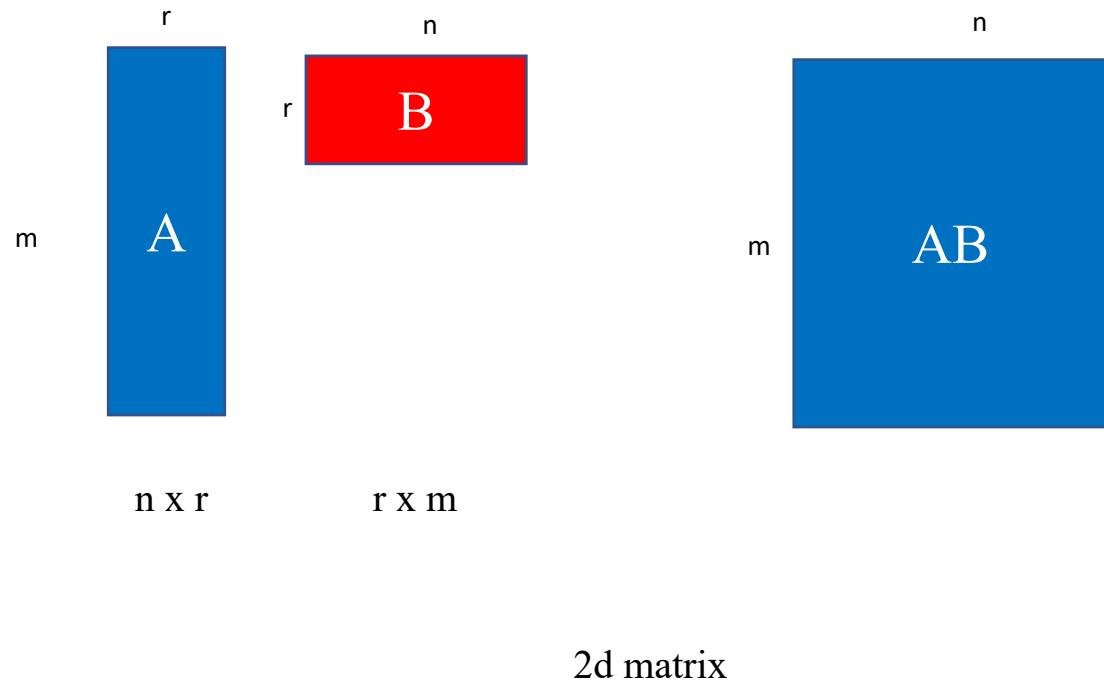
$m$

$A$

$m \times n$

2d matrix

# Matrix multiplication

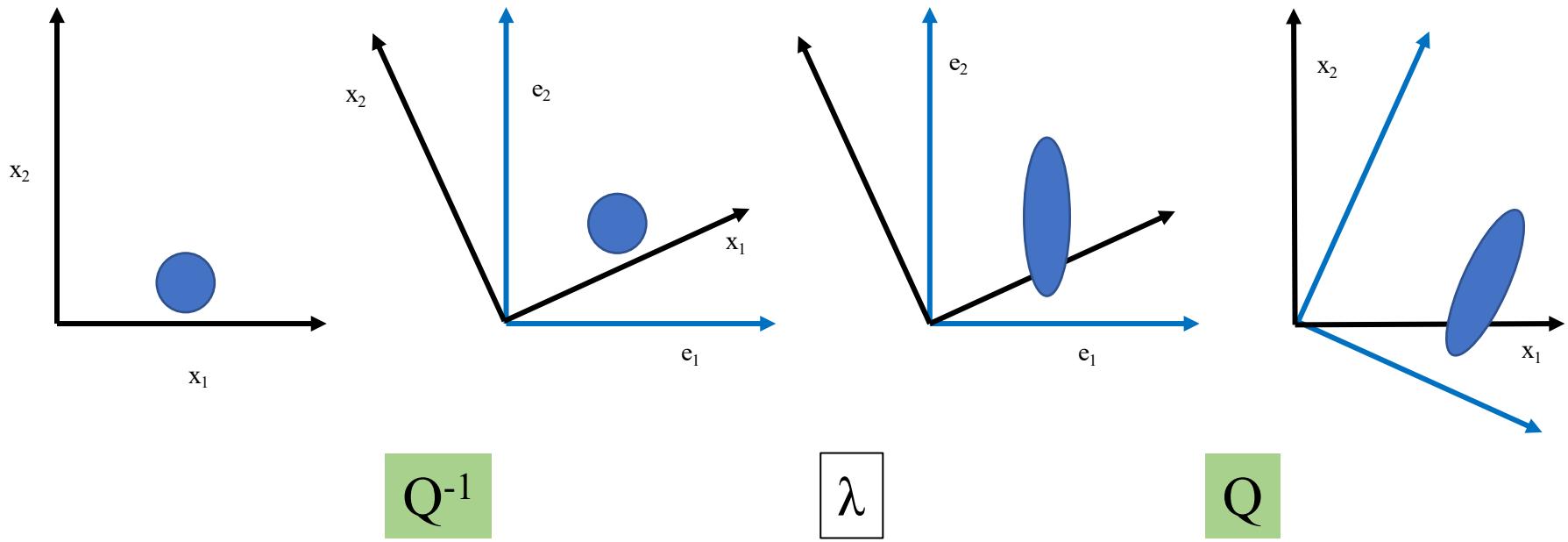


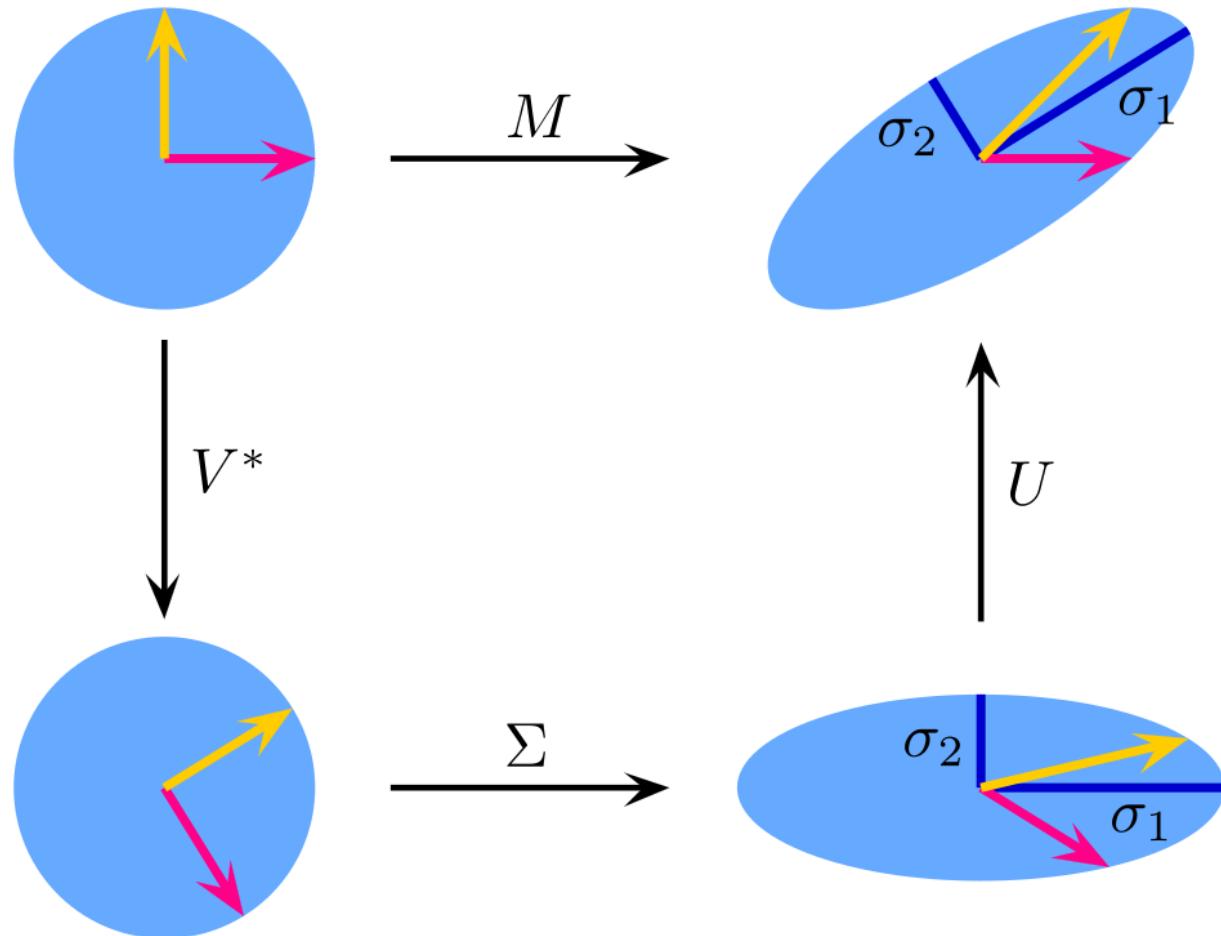
# Principal Component analysis

$$\begin{matrix} & m \\ r & \boxed{U^T} \end{matrix} \quad \begin{matrix} & n \\ m & \boxed{A} \end{matrix} = \begin{matrix} & n \\ r & \boxed{T} \end{matrix}$$

This is the dimension reduction we might have asked for

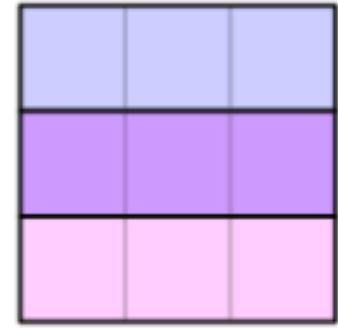
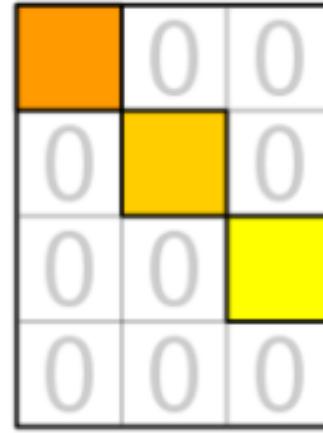
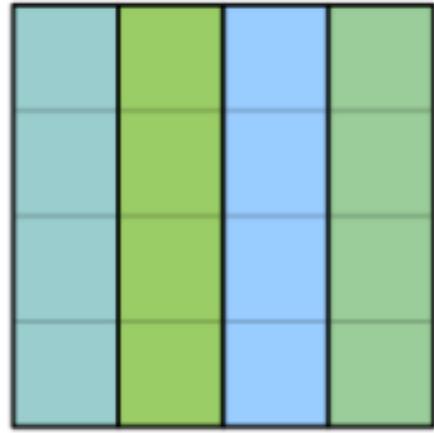
# Transformation of $x_1, x_2$ using $A$





$$M = U \cdot \Sigma \cdot V^*$$

<https://towardsdatascience.com/simple-svd-algorithms-13291ad2eef2>



$$\mathbf{M}_{m \times n} = \mathbf{U}_{m \times m} \mathbf{\Sigma}_{m \times n} \mathbf{V}^*_{n \times n}$$

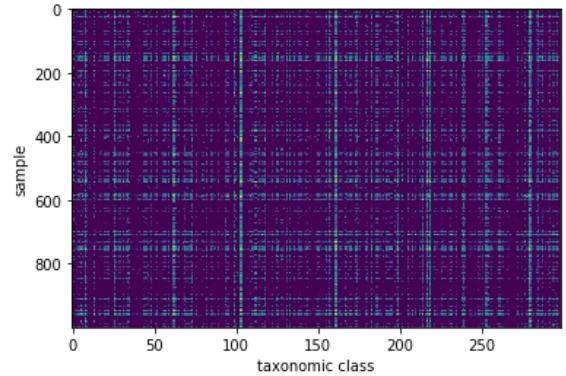
# Dimension reduction

For a general matrix representing n observations of m traits:

$$A \text{ ( } m \times n \text{ )}$$

It is unlikely that all of the elements in this matrix are independent; there is likely some “structure”

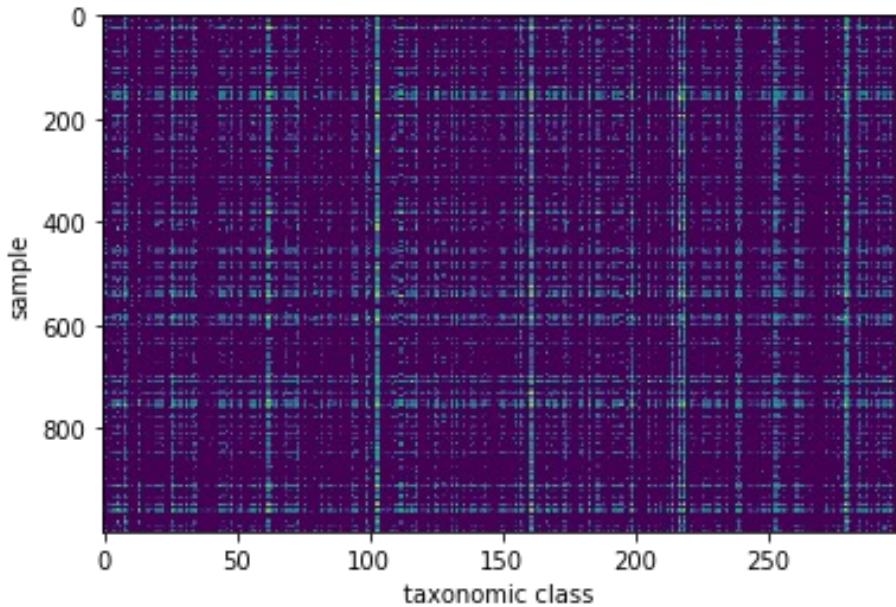
If I could replace A with an approximation of A that has fewer than  $n \times m$  numbers, A will be cheaper to store, faster to copy, can run on computers with less memory.



# Eigenvalue decomposition

$$A = Q \Sigma Q^{-1}$$

$A$  ( $n \times m$ )  
 (rows x columns)  
 ( $n$  “samples” by  $m$  “features”)



|                 | Hydrozoa | Sphingobacteriia | Marattiopsida | Pedinophyceae | Chrysiogenetes<br>(class) | Fusobacteria<br>(class) | Amphibia |
|-----------------|----------|------------------|---------------|---------------|---------------------------|-------------------------|----------|
| (Index,)        |          |                  |               |               |                           |                         |          |
| mgm4703051.3.mg | 2        | 354              | 0             | 0             | 0                         | 0                       | 0        |
| mgm4755207.3.mg | 3430     | 0                | 25            | 14            | 33111                     | 101765                  | 5094     |
| mgm4473196.3.mg | 150      | 0                | 0             | 10            | 1265                      | 487614                  | 11598    |
| mgm4679127.3.mg | 0        | 245              | 0             | 0             | 0                         | 0                       | 0        |
| mgm4529798.3.mg | 2        | 4588             | 0             | 0             | 25                        | 0                       | 56       |

# Now this looks like “structure”

Organizing data like this is deeply satisfying—it makes us think we have discovered truth.

Mature libraries to make exactly this sort of eye-candy graphs. This one is `sanborn.clustermap`

What if I told you linear algebra could do this?

