# EE16B
# Designing Information Devices and Systems II

Lecture 10B
PCA

–Last Time:

– Uniqueness and Geometry of SVD

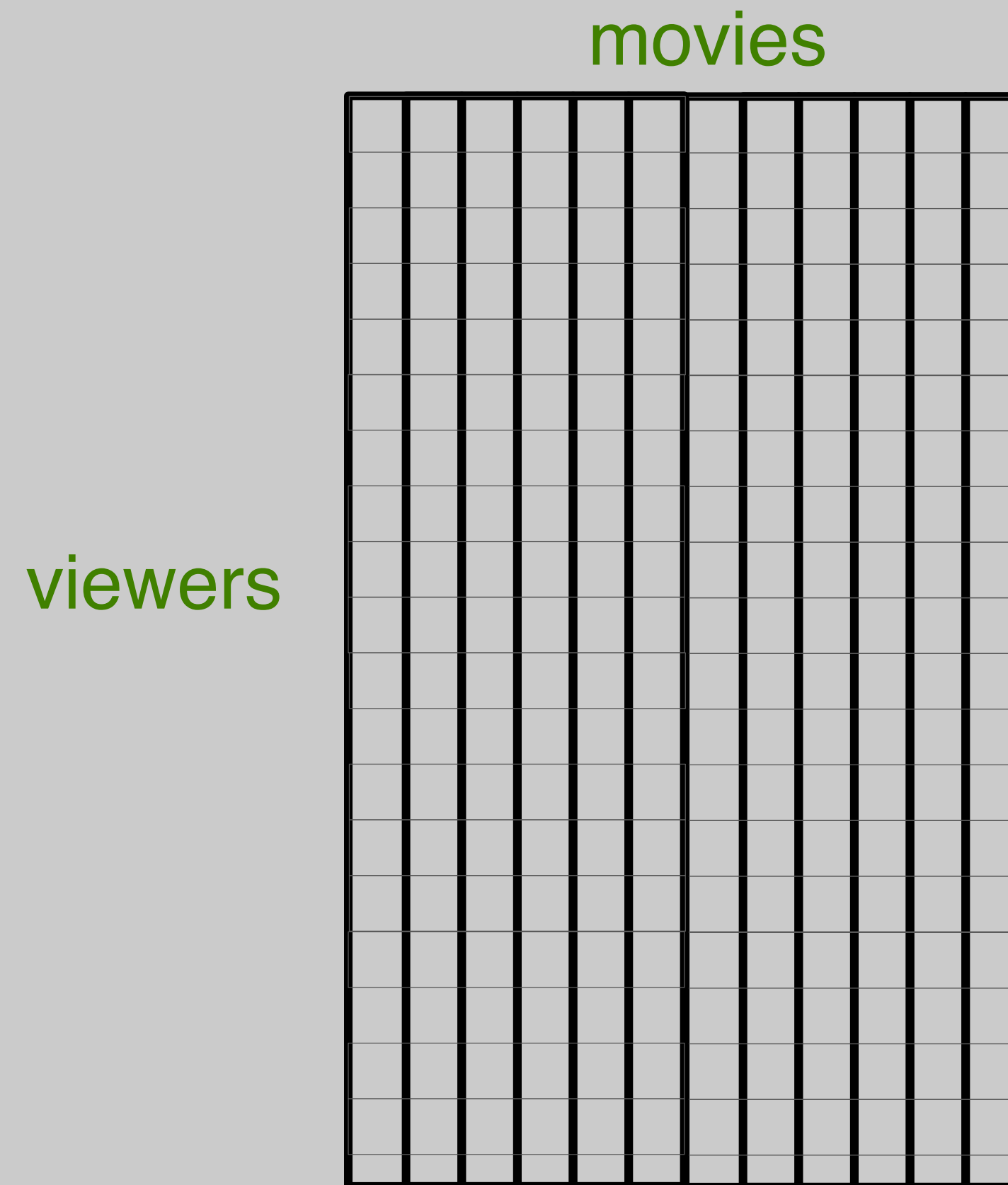– Finished proofs

– Started PCA

–Today:

– Continue PCA

–Examples of PCA

– K-means (maybe)

# Principal Component Analysis
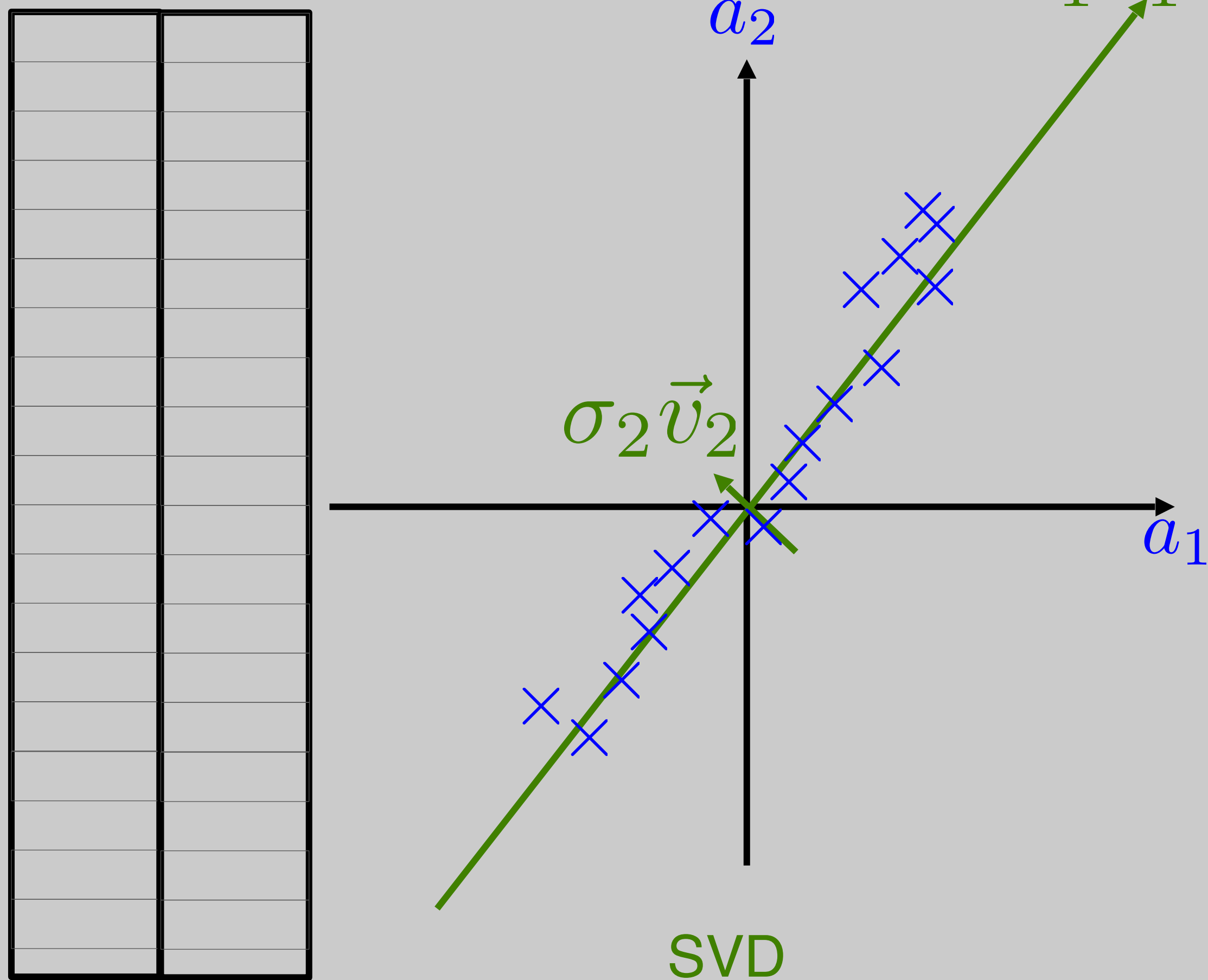
Application of the SVD to datasets to learn features

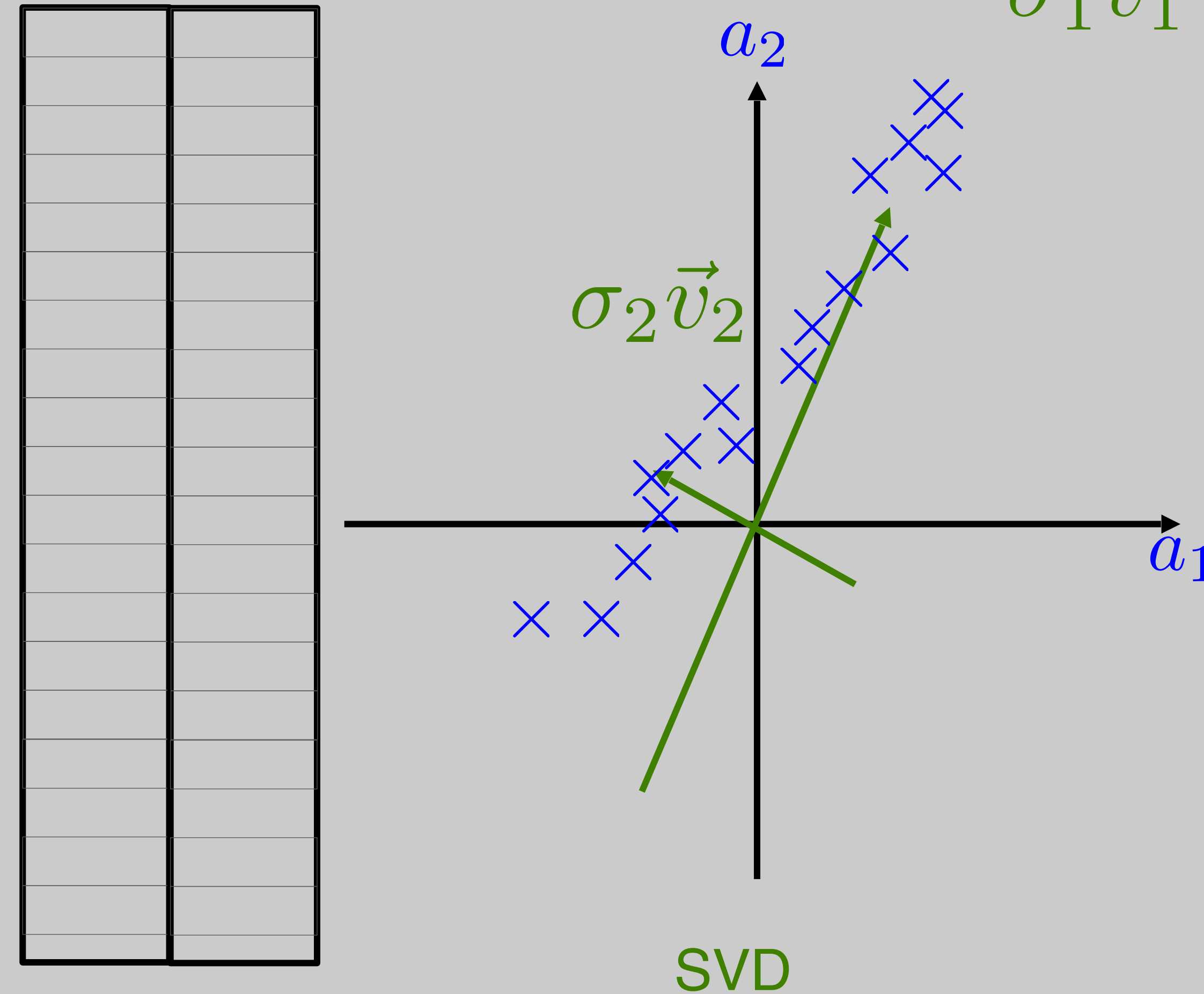PCA is a tool in statistics and machine learning, which can be computed using SVD

movies

viewers

# Example -- PCA

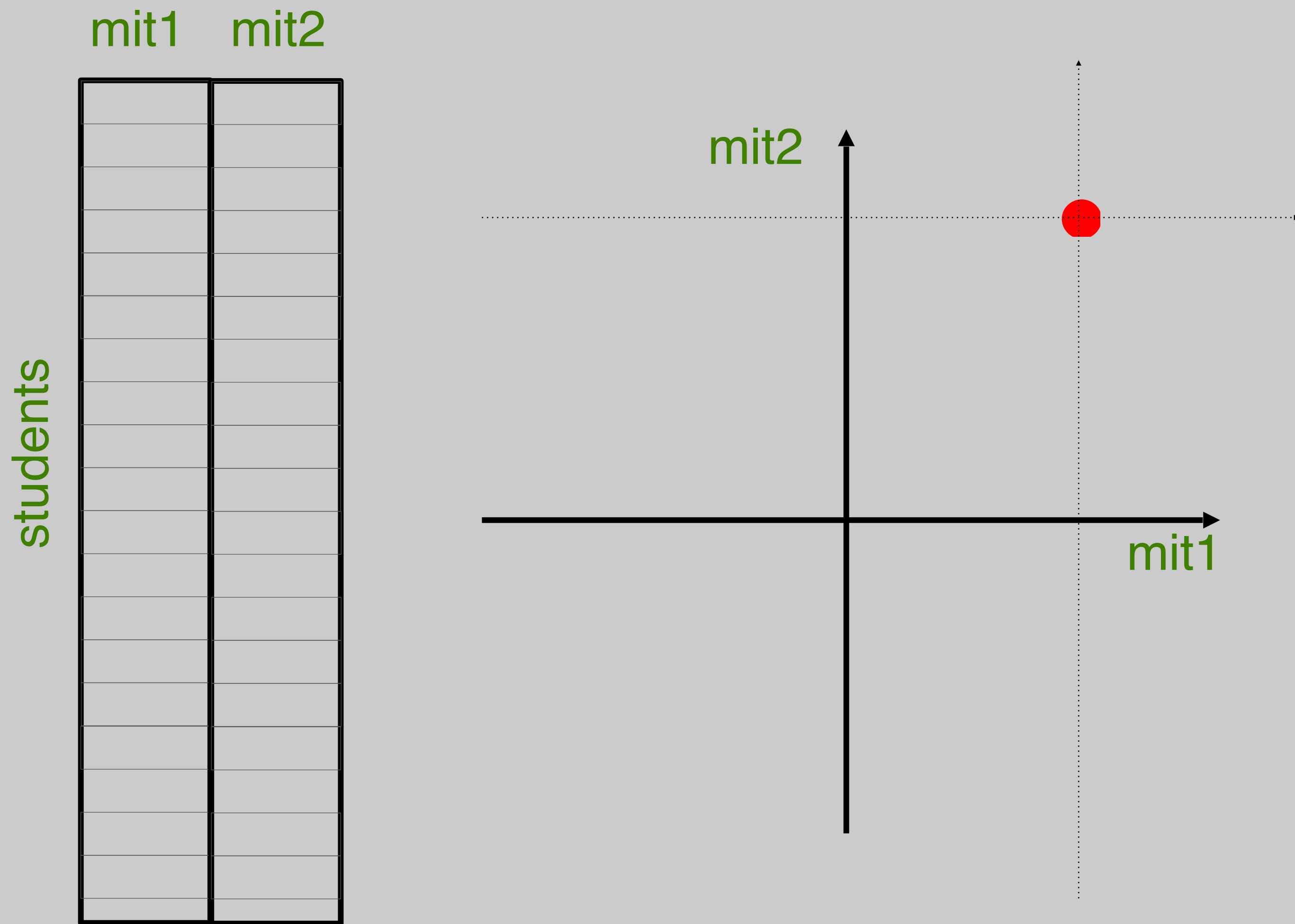Consider data s.t.

$\vec{a}_1 \quad \vec{a}_2 \approx 3\vec{a}_1$

$\sigma_1 \vec{v}_1$

$a_2$

$\sigma_2 \vec{v}_2$

$a_1$

SVD

$\vec{a}_1 \quad \vec{a}_2 \approx 3\vec{a}_1 + 1$

$\sigma_1 \vec{v}_1$

$a_2$

$\sigma_2 \vec{v}_2$

$a_1$

SVD

# Example -- PCA

Consider miterm data

mit1    mit2

students



mit2

mit1

# Example -- PCA

## Consider miterm data



mit1   mit2

students

mit2

mit1

# Example -- PCA

## Consider miterm data

mit1   mit2

students

consistency

up-down mobility

mit2

mit1

# PCA Procedure

Remove averages from column of A

From AᵀA, find $\sigma_i, \quad \vec{v_i}$

$\vec{v_i}$ are principal components!

# AᵀA as sample covariance matrix

$$A = \vec{a} \qquad a_\mu = \frac{1}{N} \sum_{i=0}^{N-1} a_i \qquad \tilde{A} = \vec{a} - a_\mu \vec{1}$$

$$\tilde{A}^T \tilde{A} = (\vec{a} - a_\mu \vec{1})^T (\vec{a} - a_\mu \vec{1})$$

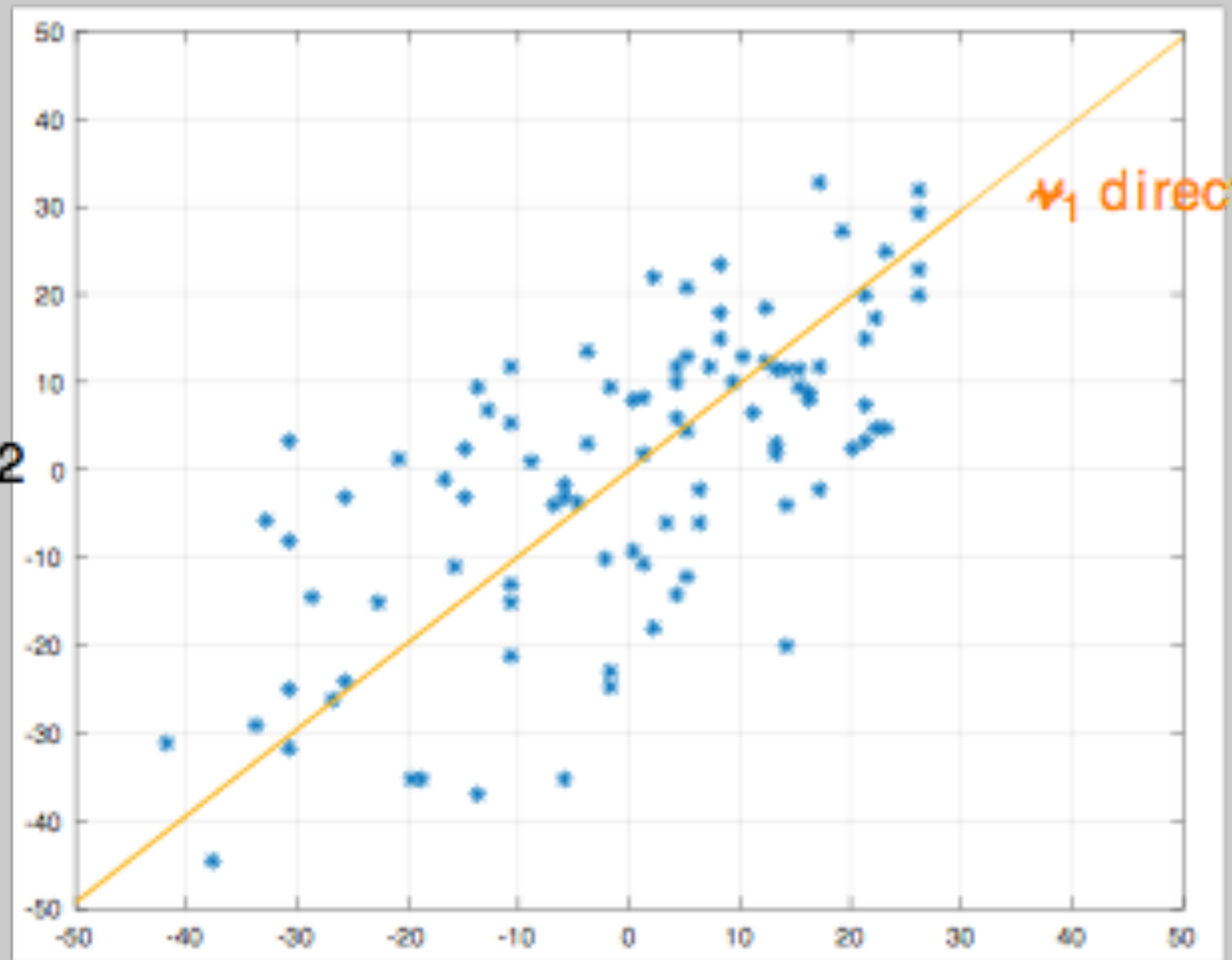$$= \vec{a}^T \vec{a} - 2N a_\mu^2 + N a_\mu^2 \quad = \vec{a}^T \vec{a} - N a_\mu^2$$

$$\frac{1}{N} \tilde{A}^T \tilde{A} = \frac{1}{N} \vec{a}^T \vec{a} - a_\mu^2 = \frac{1}{N} \sum_{i=0}^{N-1} a_i^2 - a_\mu^2 \quad = a_\sigma^2$$

Sample
Variance!

# Example midterm



$$\frac{1}{93} A^T A = \begin{bmatrix} 297.69 & 202.53 \\ 202.53 & 292.07 \end{bmatrix}$$
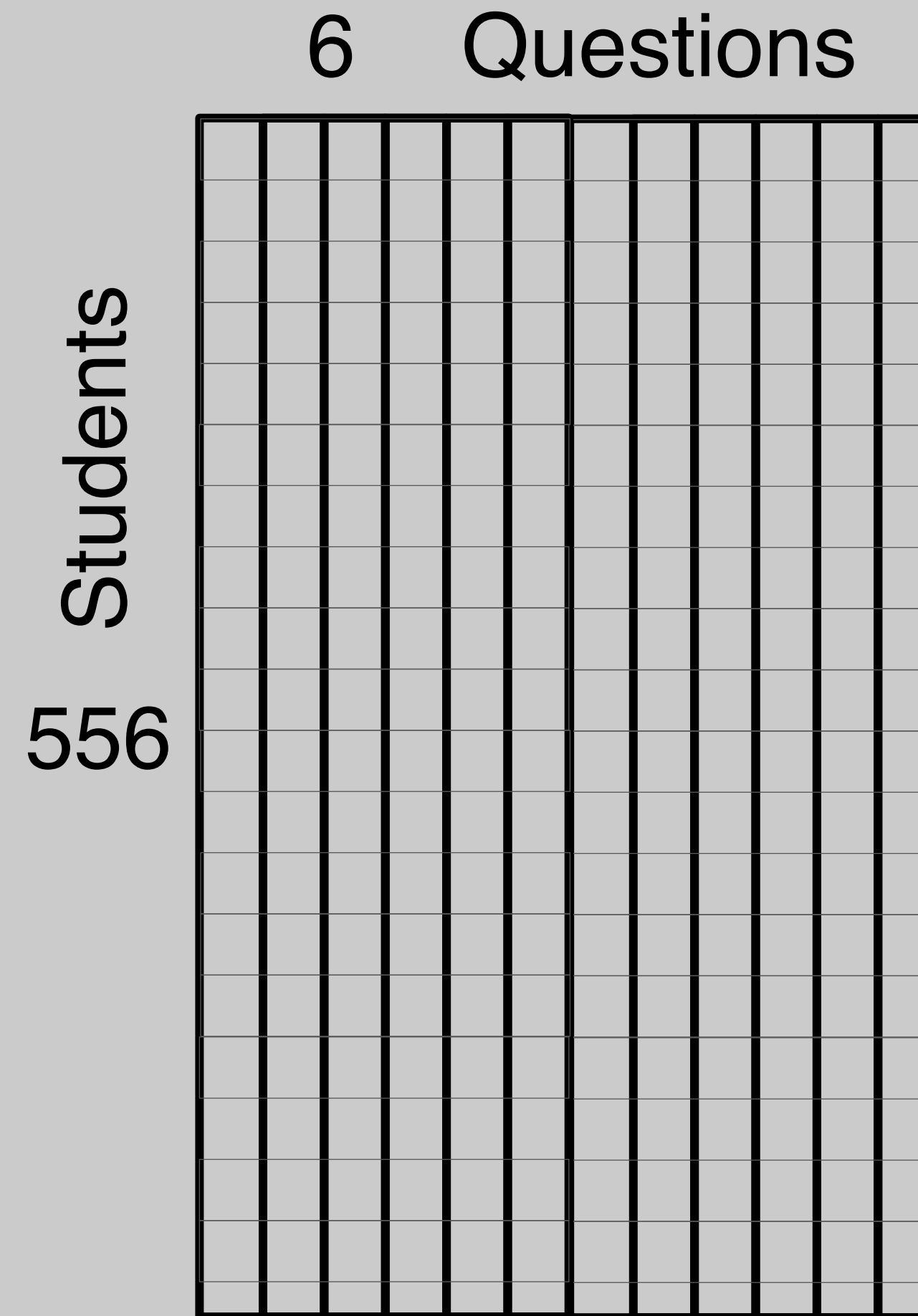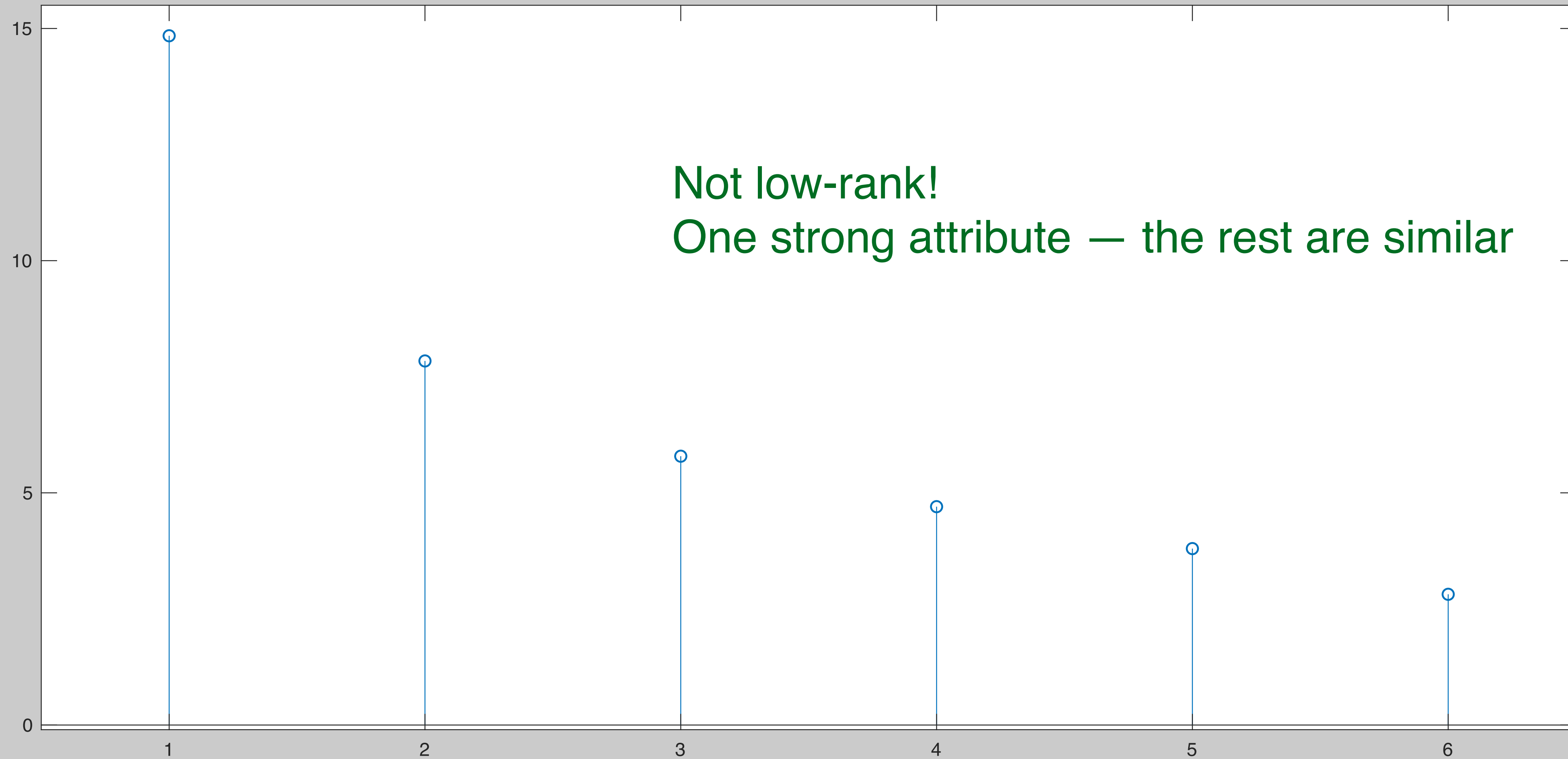
# Mid Semester Survey Results

1) HW difficulty

2) HW Length

3) Lab hour/week

4) Current rating

5) Previous Rating

6) Comfortable attending OH

6    Questions

Students

556
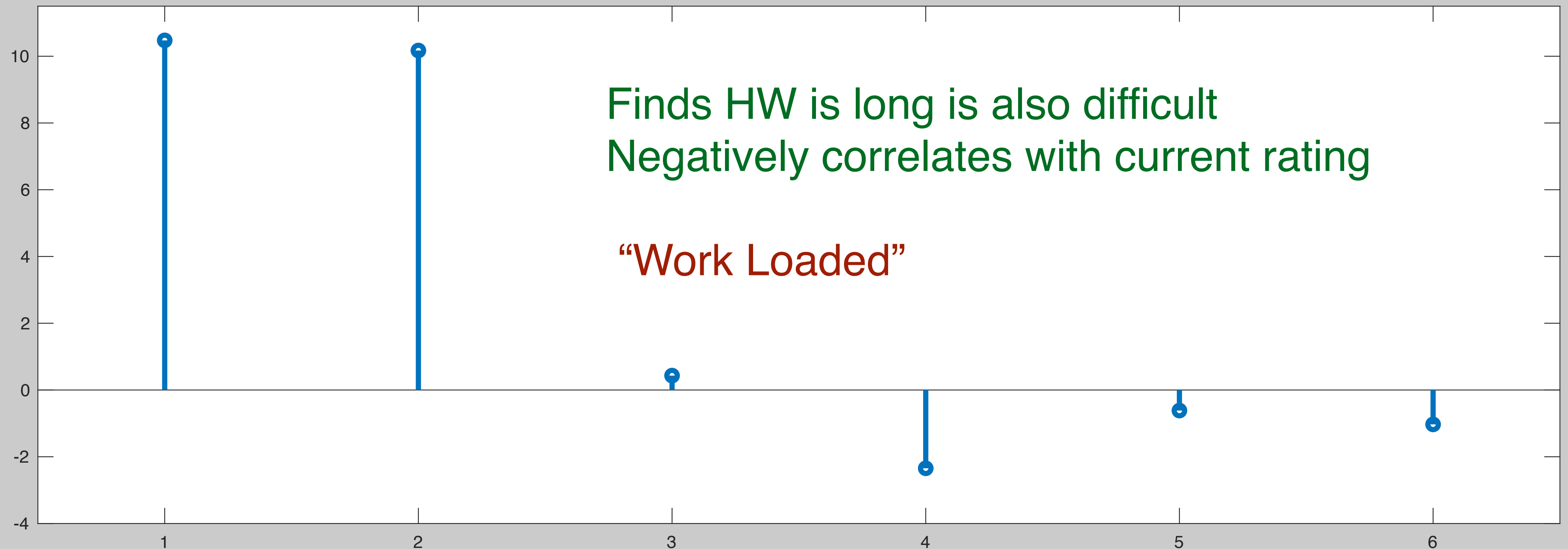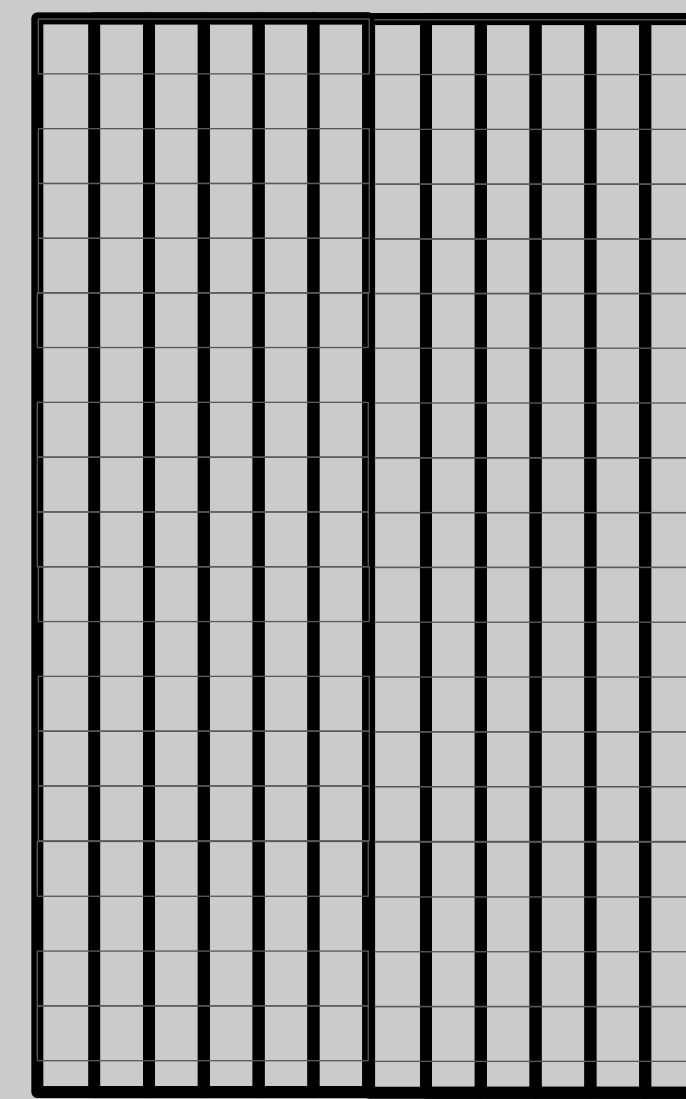
# Data Science

## Singular values



Not low-rank!
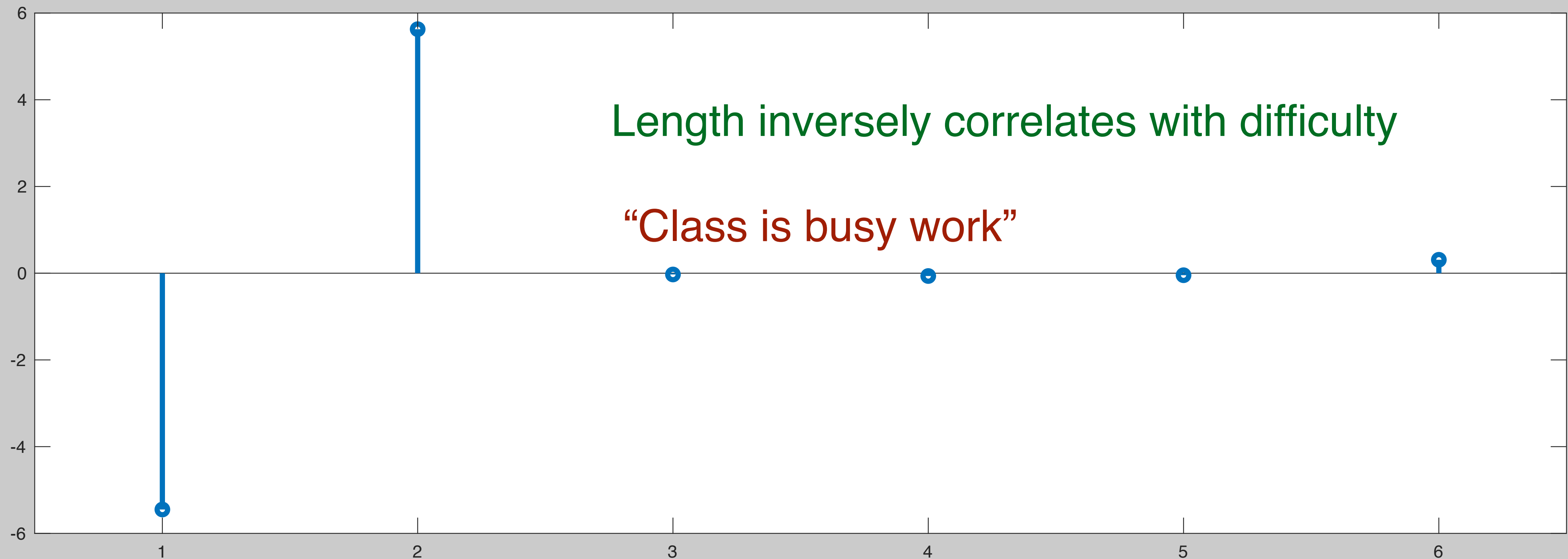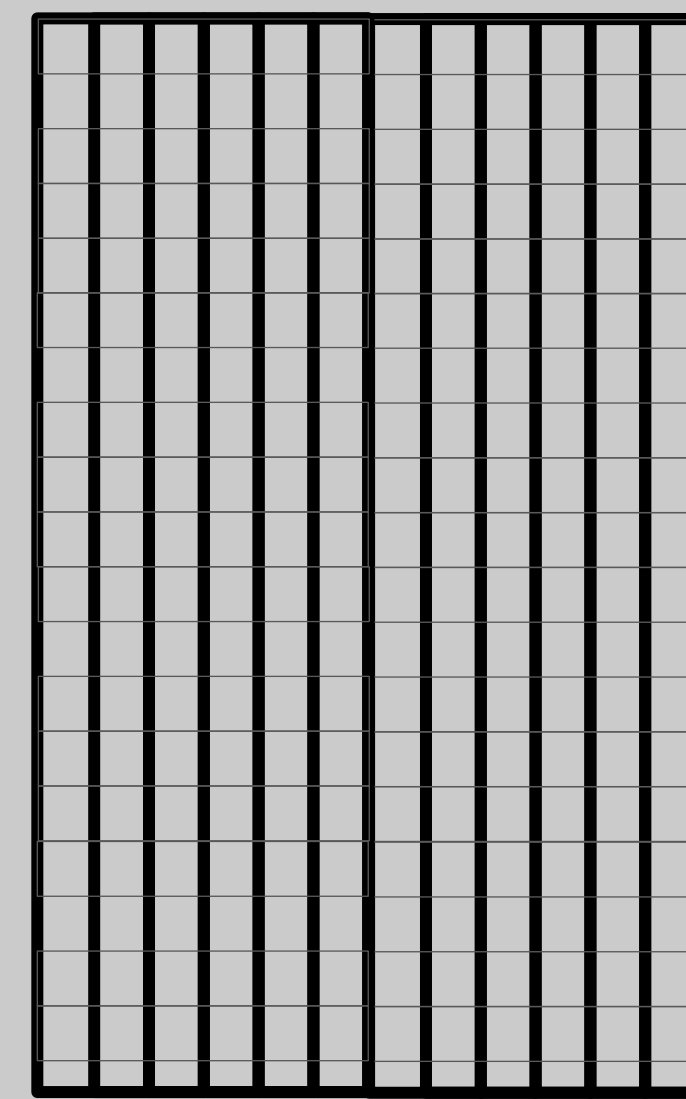One strong attribute — the rest are similar

# Data Science

$$A^T \vec{u}_1$$

1)  HW difficulty
2)  HW Length
3)  Lab hour/week
4)  Current rating
5)  Previous Rating
6)  Comfortable attending OH

Finds HW is long is also difficult
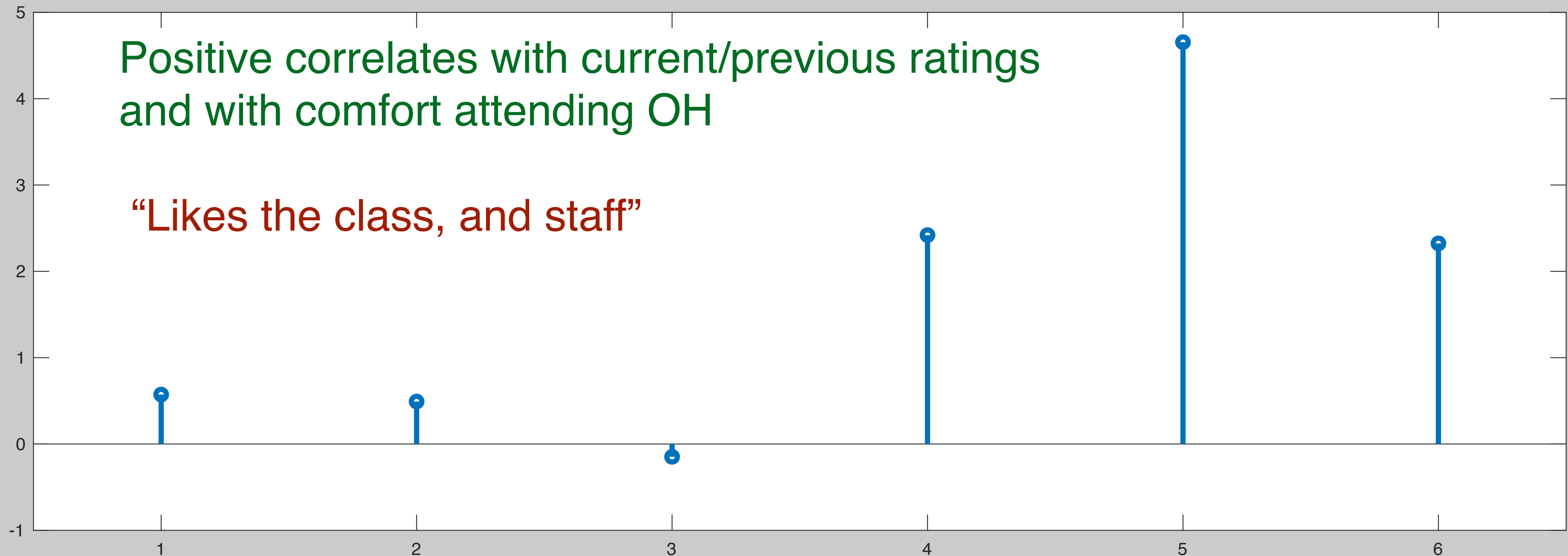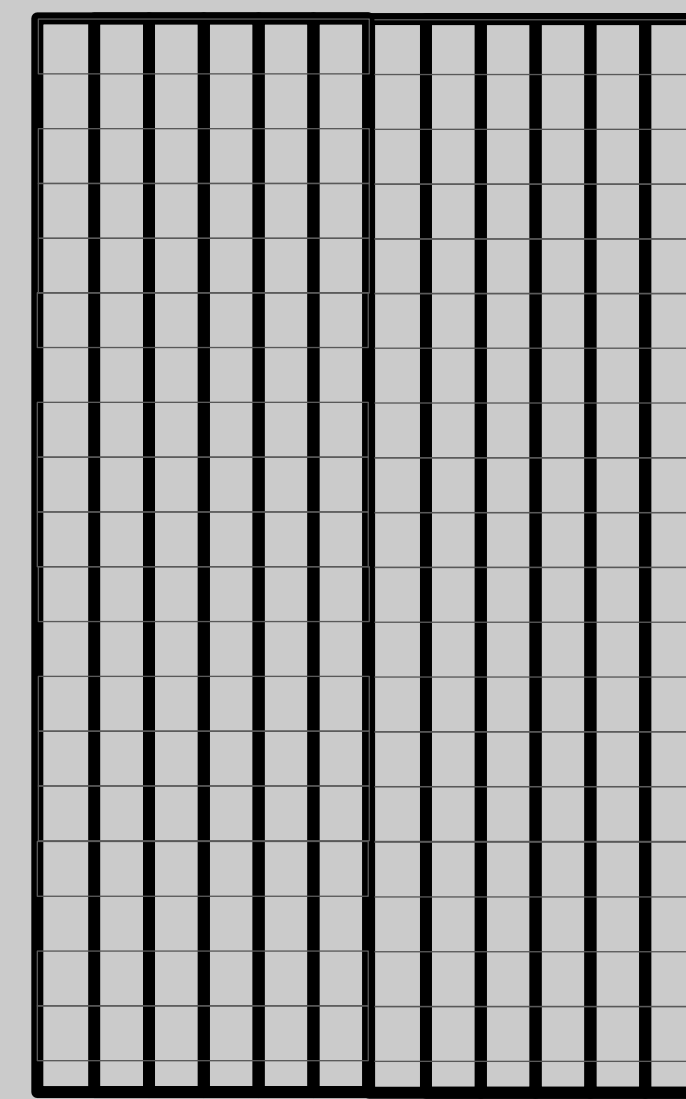Negatively correlates with current rating

"Work Loaded"

# Data Science

$$A^T \vec{u}_2$$
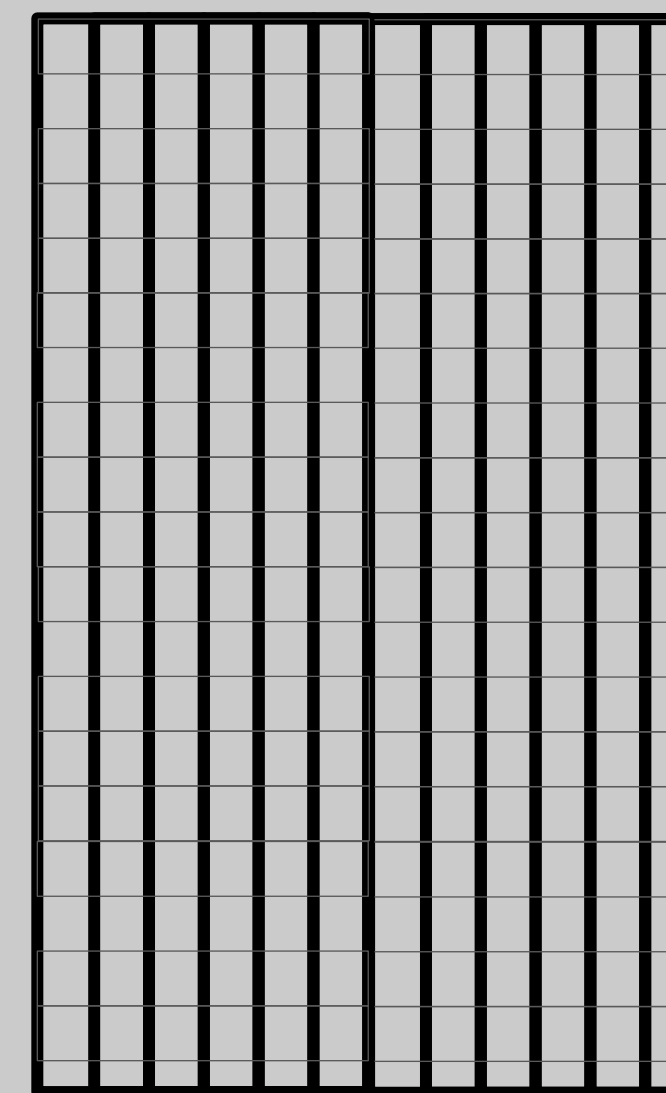
1) HW difficulty
2) HW Length
3) Lab hour/week
4) Current rating
5) Previous Rating
6) Comfortable attending OH

Length inversely correlates with difficulty

"Class is busy work"

# Data Science

$$A^T \vec{u}_3$$

Positive correlates with current/previous ratings
and with comfort attending OH
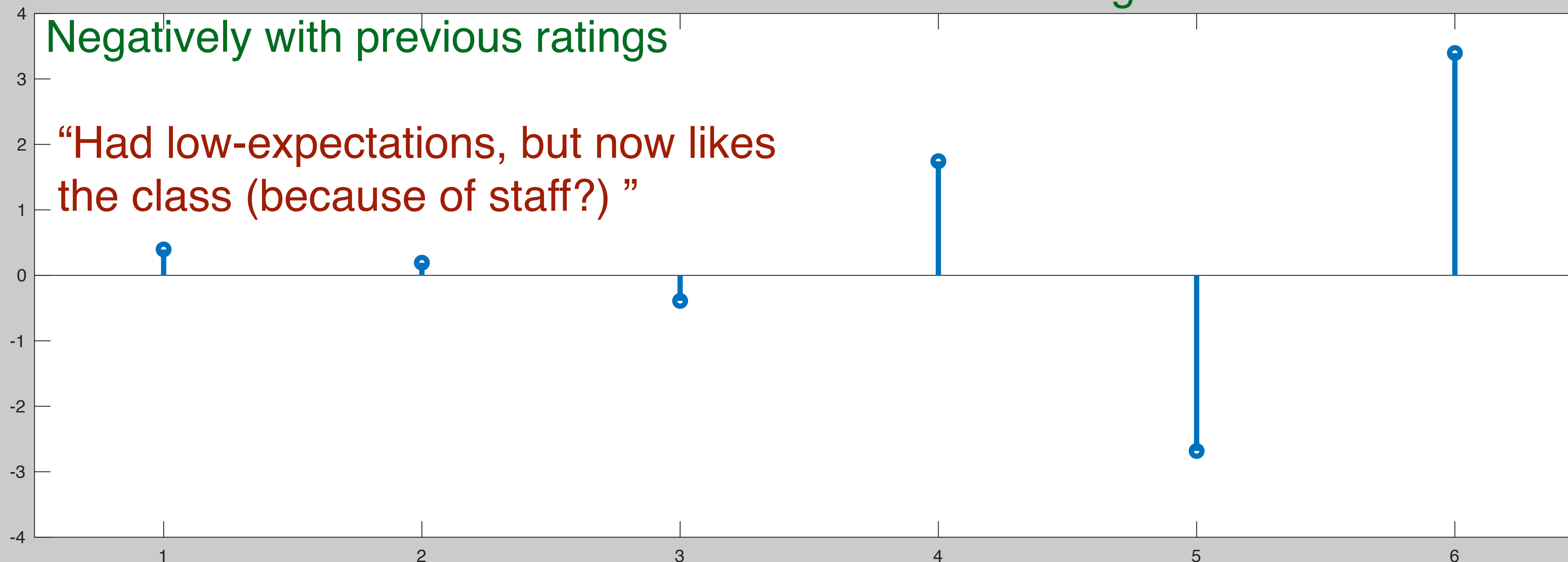
"Likes the class, and staff"

# Data Science

$$A^T \vec{u}_4$$

1) HW difficulty
2) HW Length
3) Lab hour/week
4) Current rating
5) Previous Rating
6) Comfortable attending OH

Positive correlates with current & with comfort attending OH
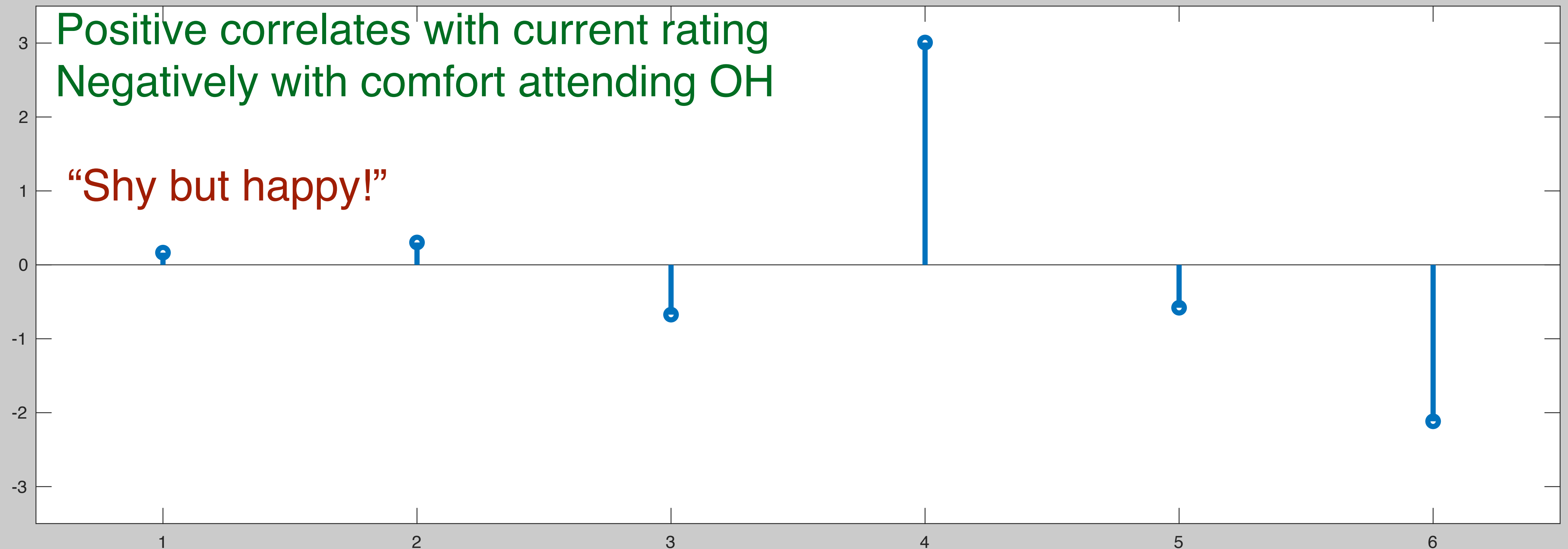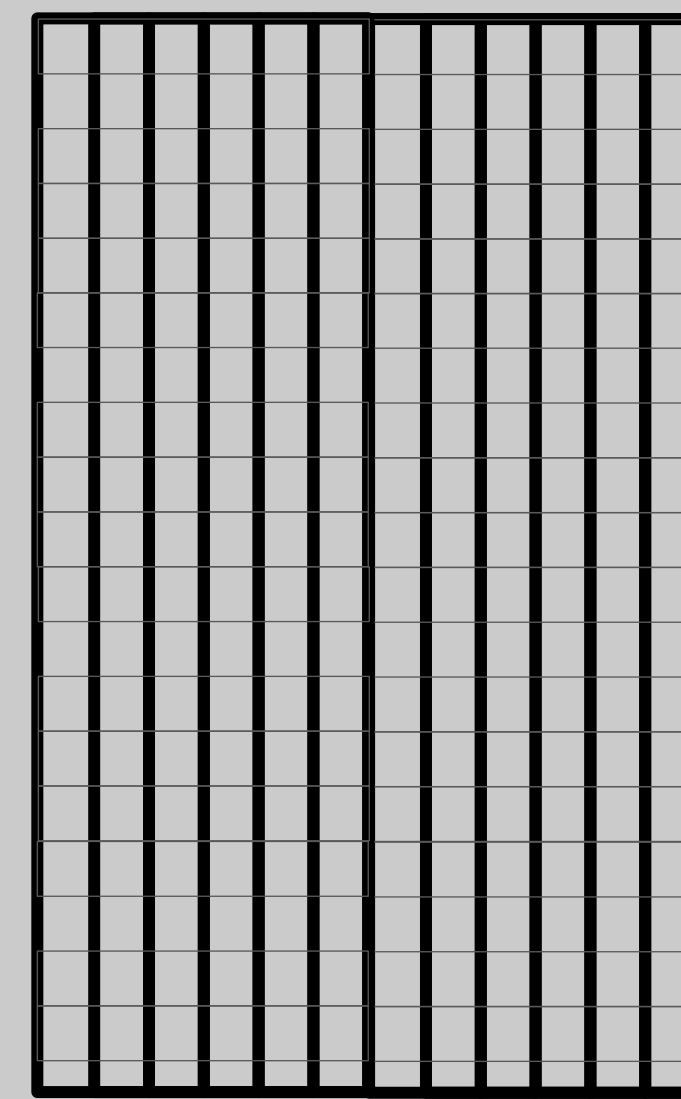Negatively with previous ratings

"Had low-expectations, but now likes the class (because of staff?)"
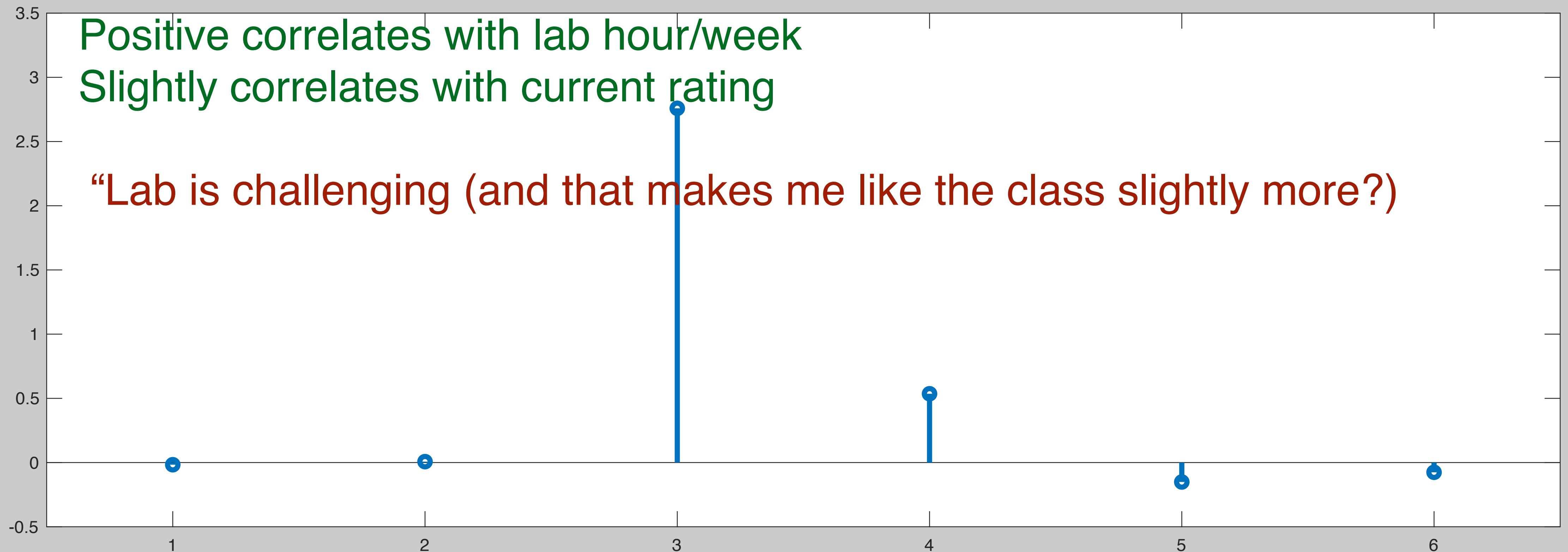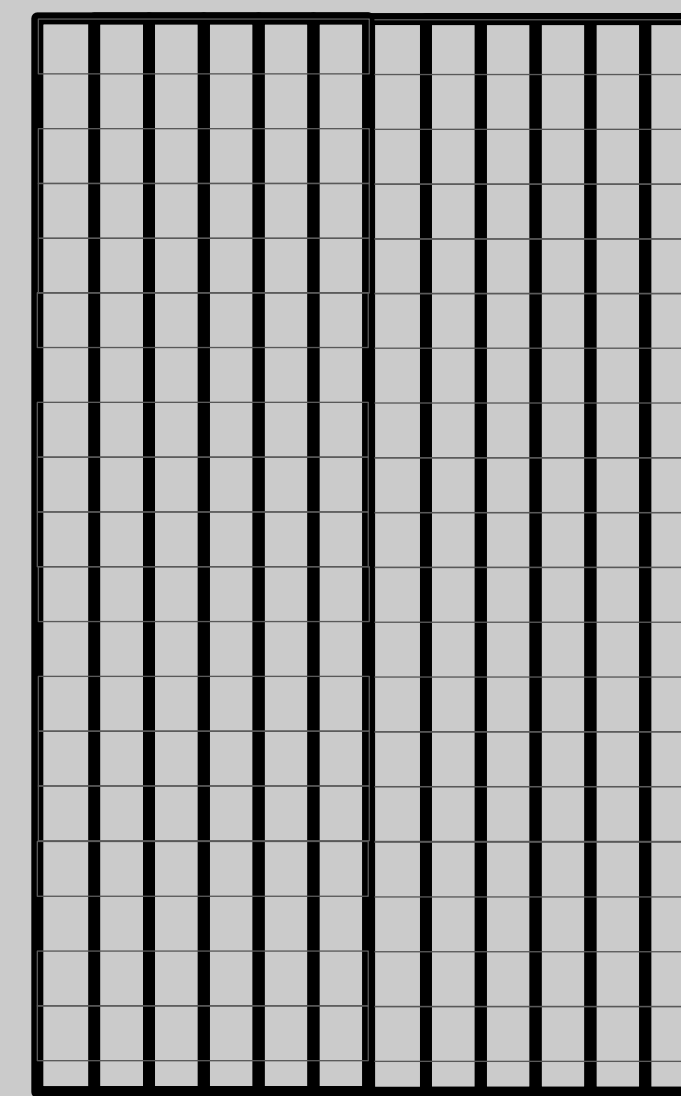
# Data Science

$$A^T \vec{u}_5$$

1) HW difficulty
2) HW Length
3) Lab hour/week
4) Current rating
5) Previous Rating
6) Comfortable attending OH

Positive correlates with current rating
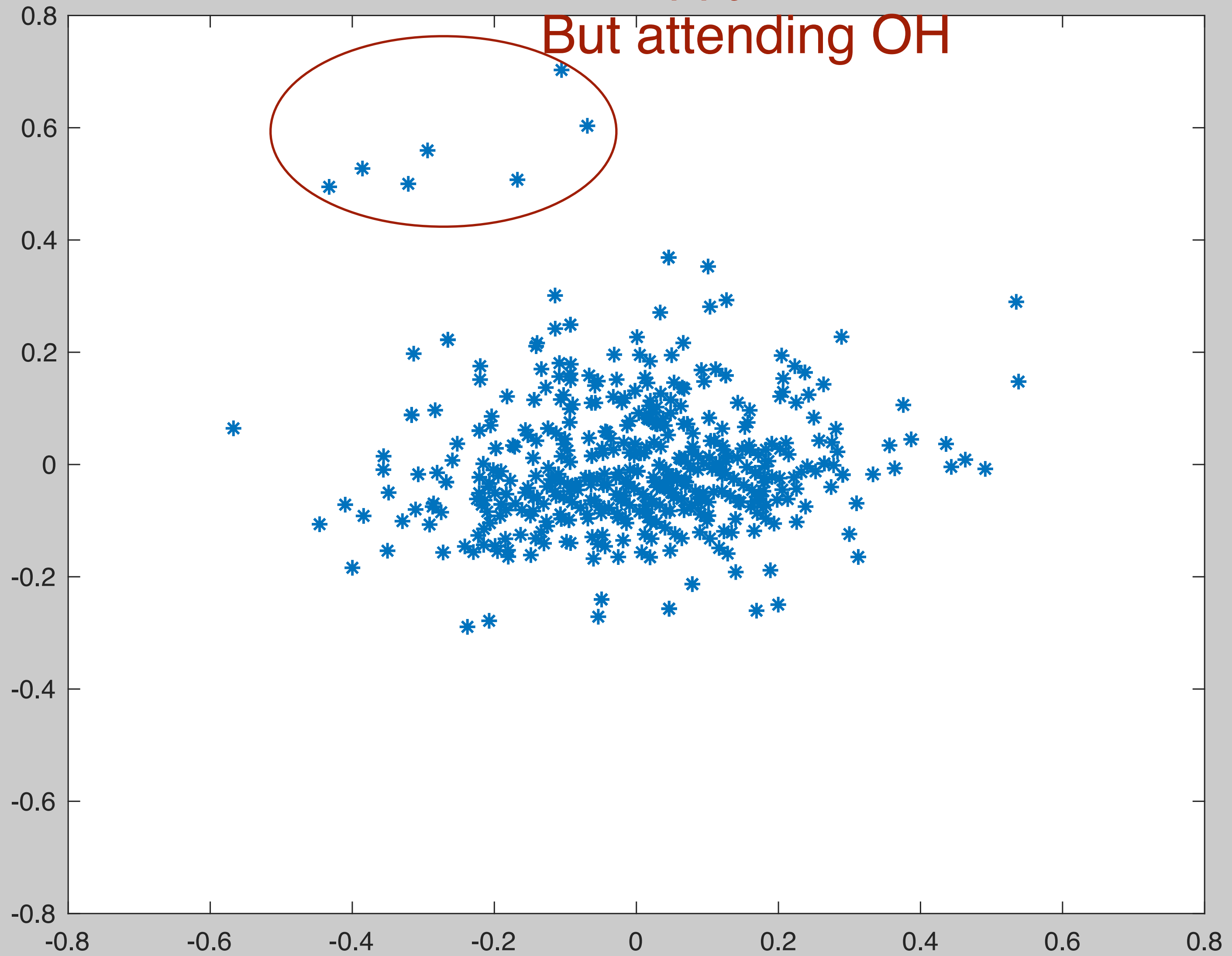Negatively with comfort attending OH

"Shy but happy!"

# Data Science

$$A^T \vec{u}_6$$

1) HW difficulty
2) HW Length
3) Lab hour/week
4) Current rating
5) Previous Rating
6) Comfortable attending OH

Positive correlates with lab hour/week
Slightly correlates with current rating

"Lab is challenging (and that makes me like the class slightly more?)

# Data Science

# Data Science

$A\vec{v}_4$    Had low-expectations, but now likes the class



Students

# Data Science

Histogram of $A\vec{v}_4$



Had low-expectations, but now like the class

# PCA in Genetics Reveals Geography

Study:

Characterized genetic variatios in 3,000 Europeans from 36 Countries
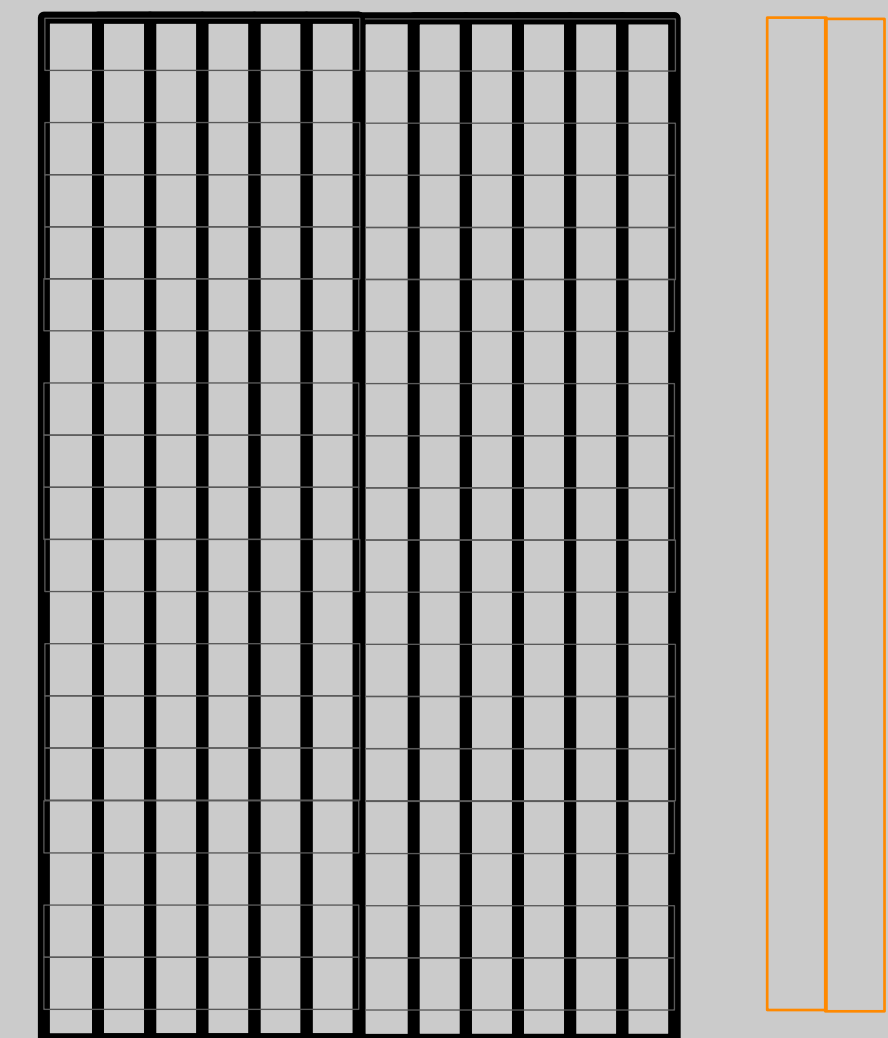
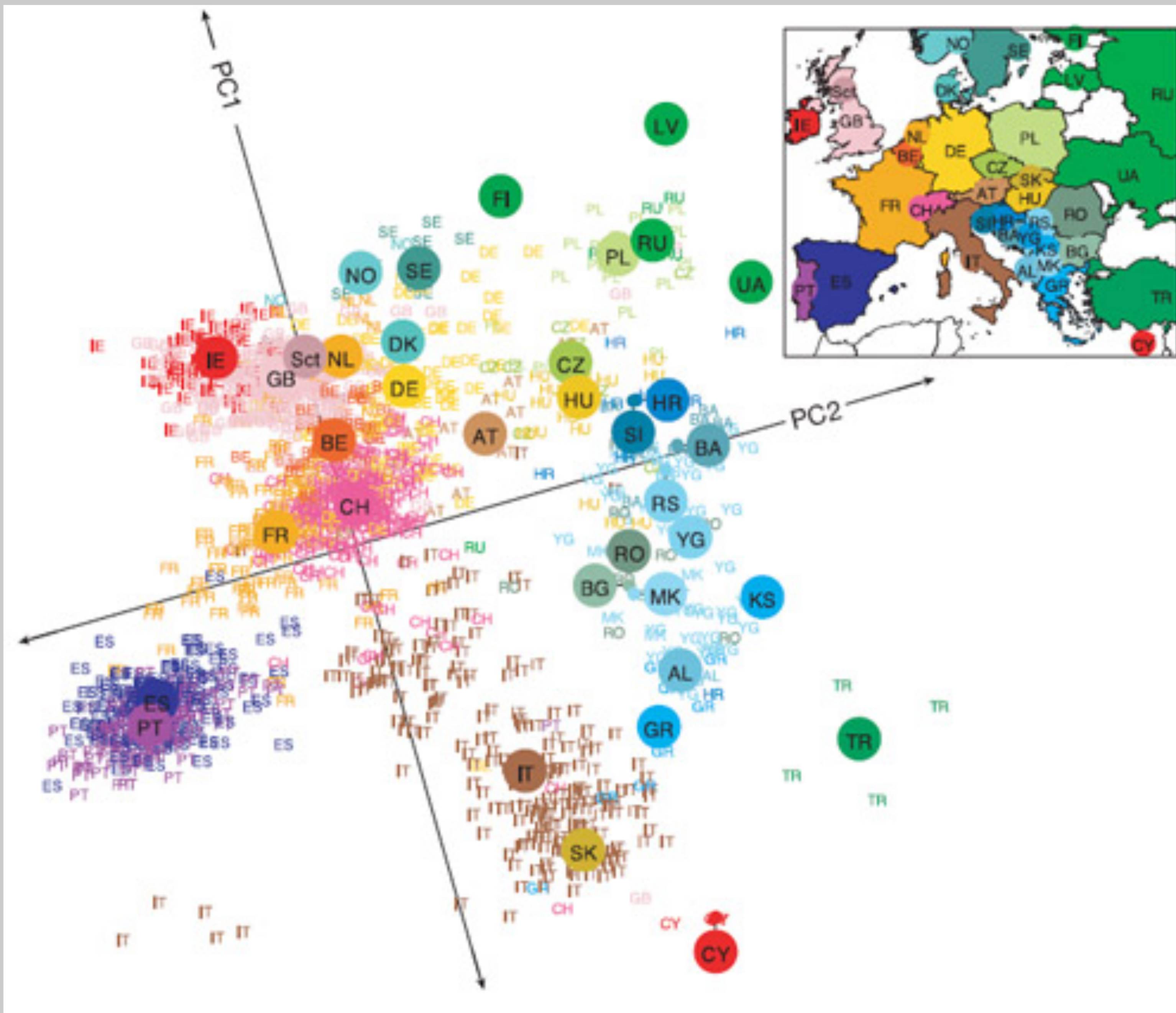Built a matrix of 200K SNPs (single nucleotide polymorphisms)

Computed largest 2 principle components

Projected subjects on 2 dimentional data

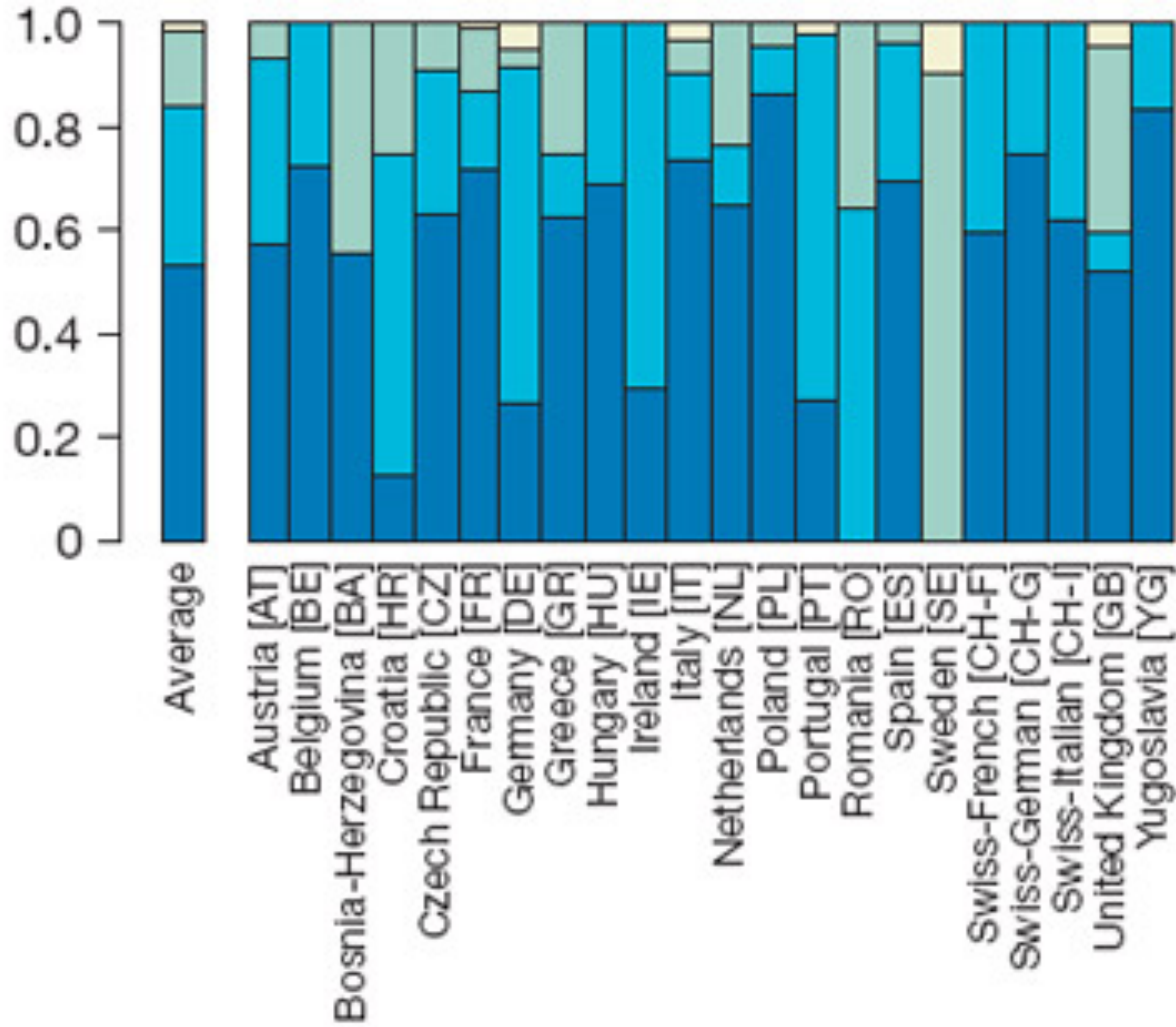Overlayed the result on the map of Europe

$$A\vec{v}_1 \quad A\vec{v}_2$$

# Interesting conclusions

"The results have implications for a lot of biomedical research. Many scientists are scanning entire genomes on a hunt for SNPs that affect a person's risk of diseases like cancer or their reaction to drugs. Novembre says that researchers who are running these "whole-genome studies" need to bear in mind where their sample has come from. Even if a study looks at a small and seemingly related parts of Europe, it would have to adjust for any geographical influences in the genetic variations it uncovers."

http://phenomena.nationalgeographic.com/2008/09/01/european-genes-mirror-european-geography/

# 23 and me

# Physical features

# Labeled VS non labeled Classification

Word1

Word2

Word3

Word4

Word5

Word6

Word7

Word8

# Labeled VS non labeled Classification

Word1

Word2

Word3

Word4

Word5

Word6

Word7

Word8

# Labeled VS non labeled Classification

"Banana"

"Banana"

"Banana"

"Mango"

"Mango"

"Mango"

"Chop"

"Chop"

"Chop"

PC2

PC1

# k-means

Given: $\vec{x}_1, \vec{x}_2, \cdots, \vec{x}_m \in \mathrm{R}^n$
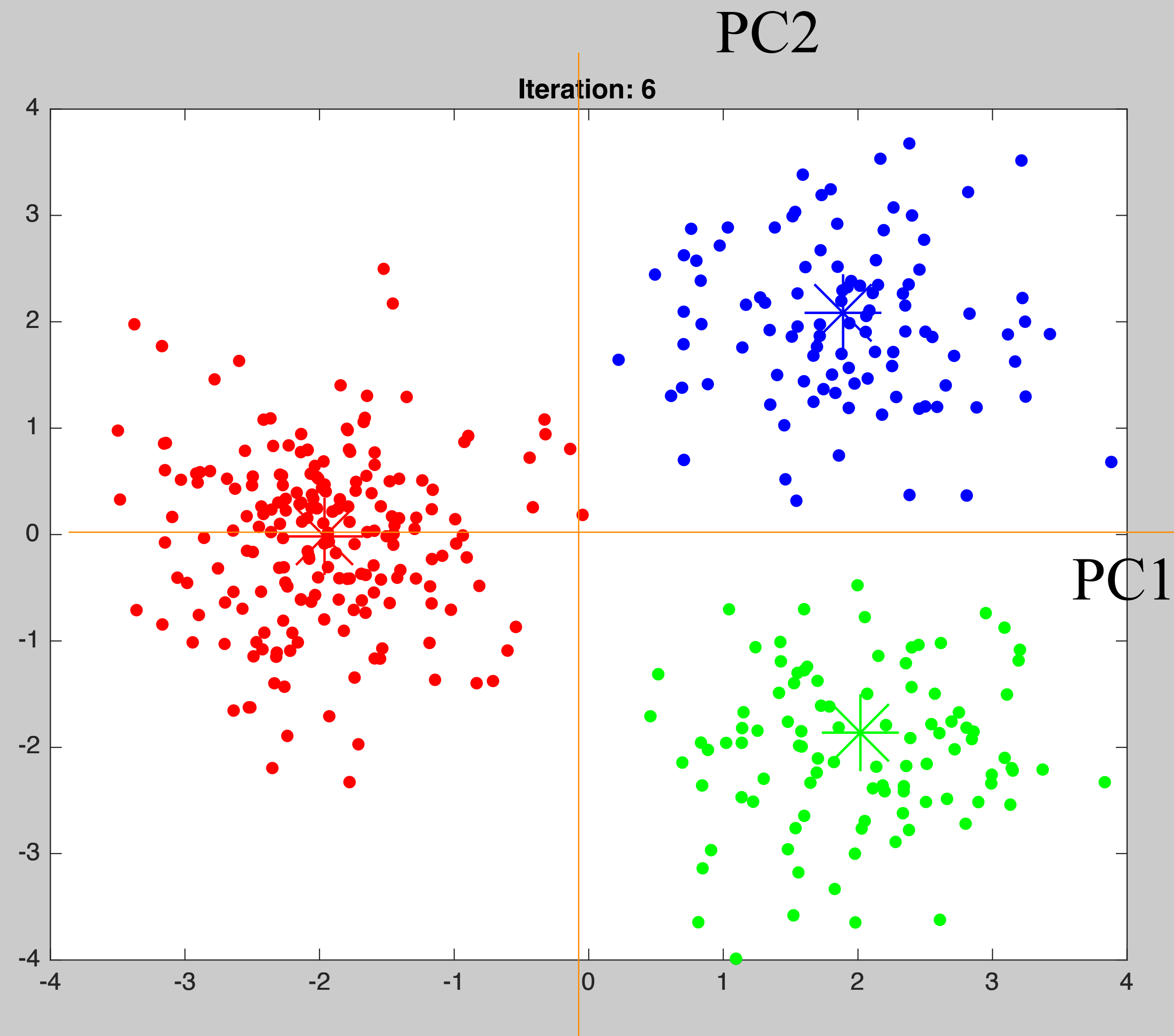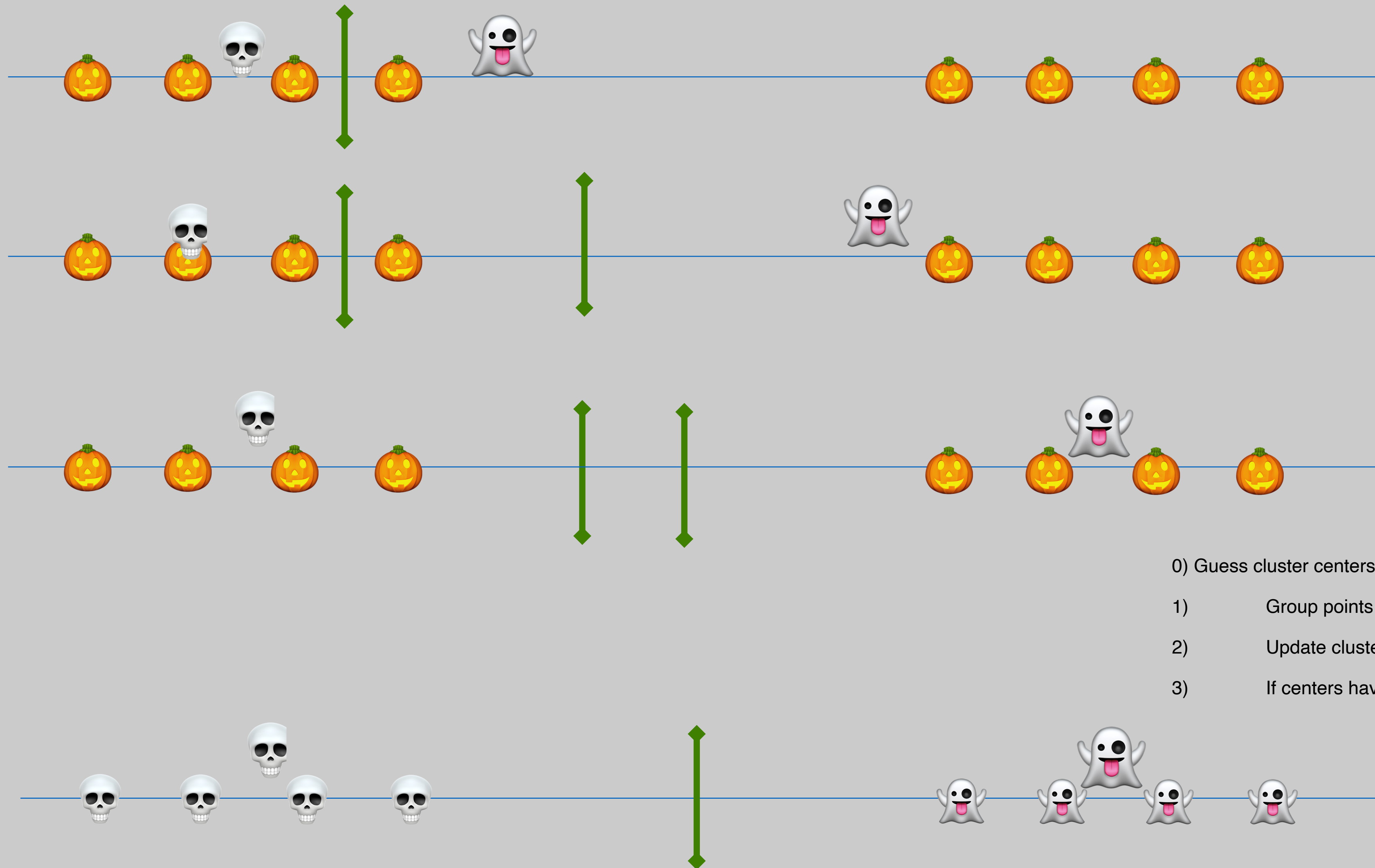
Partition them into k << m groups

0) Guess cluster centers to initialize

1) Group points around nearest center

2) Update cluster centers by averaging within group

3) If centers have changed, repeat 1-3

# k-means 1D example

0) Guess cluster centers to initialize

1) Group points around nearest center

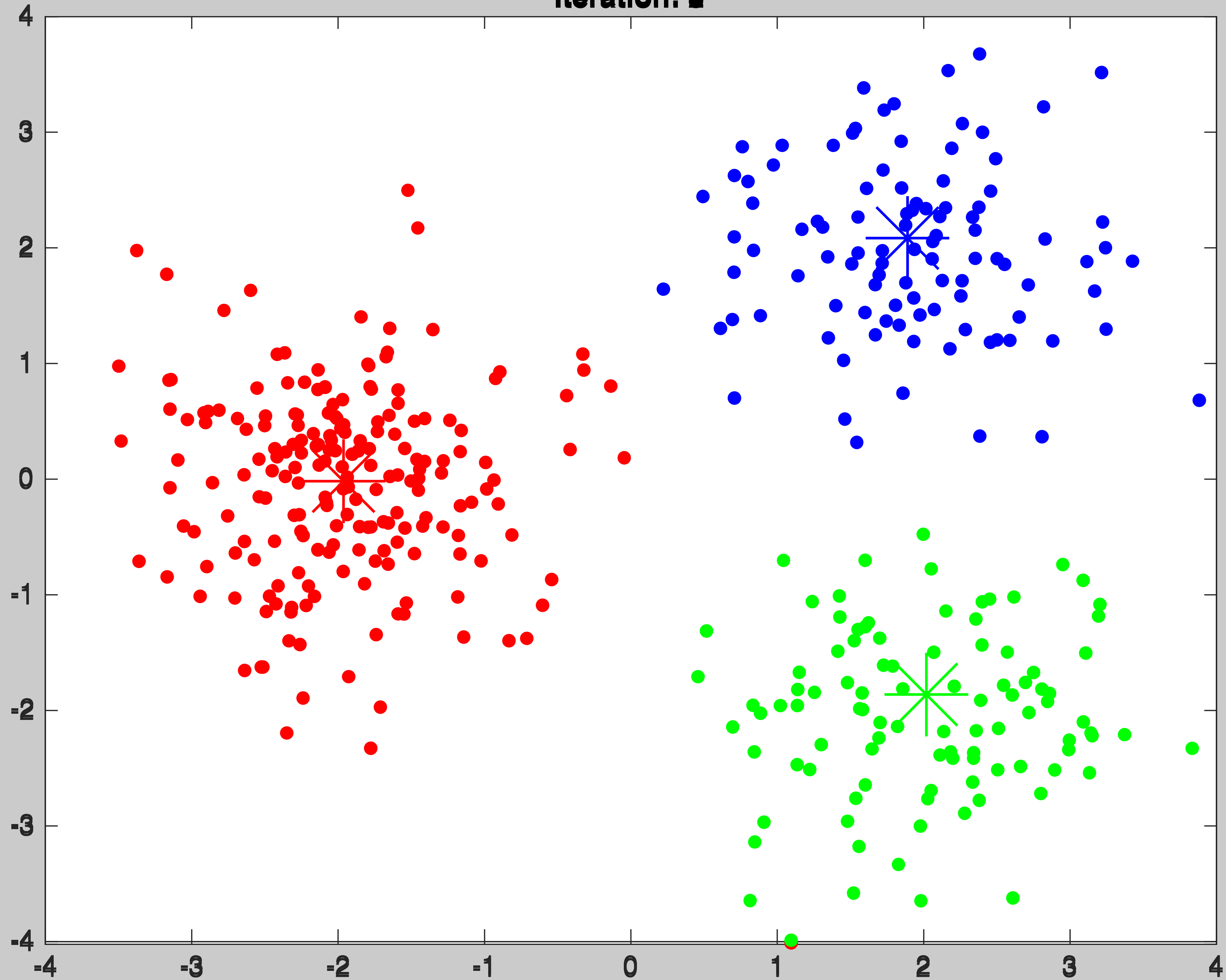2) Update cluster centers by averaging within group

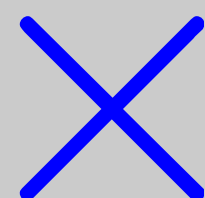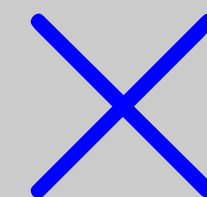3) If centers have changed, repeat 1-3

# General k-means Algorithm

0) Initialize k cluster centers $\vec{m}_1, \vec{m}_2, \cdots, \vec{m}_k$

1) Assign points to cluster: point $\vec{x}$ goes to cluster $i$
   if,
   $$\|\vec{x} - \vec{m}_i\| < \|\vec{x} - \vec{m}_j\| \quad \forall j \neq i$$

2) Let $S_i$ be the set of samples in cluster $i$
   recompute cluster centers:
   $$\vec{m}_i = \frac{1}{|S_i|} \sum_{\vec{x} \in S_i} \vec{x}$$

3) If any $m_i$ has changed, repeat 1-3

Iteration: 0

EE16B M. Lustig,  EECS UC Berkeley

# Objective Function

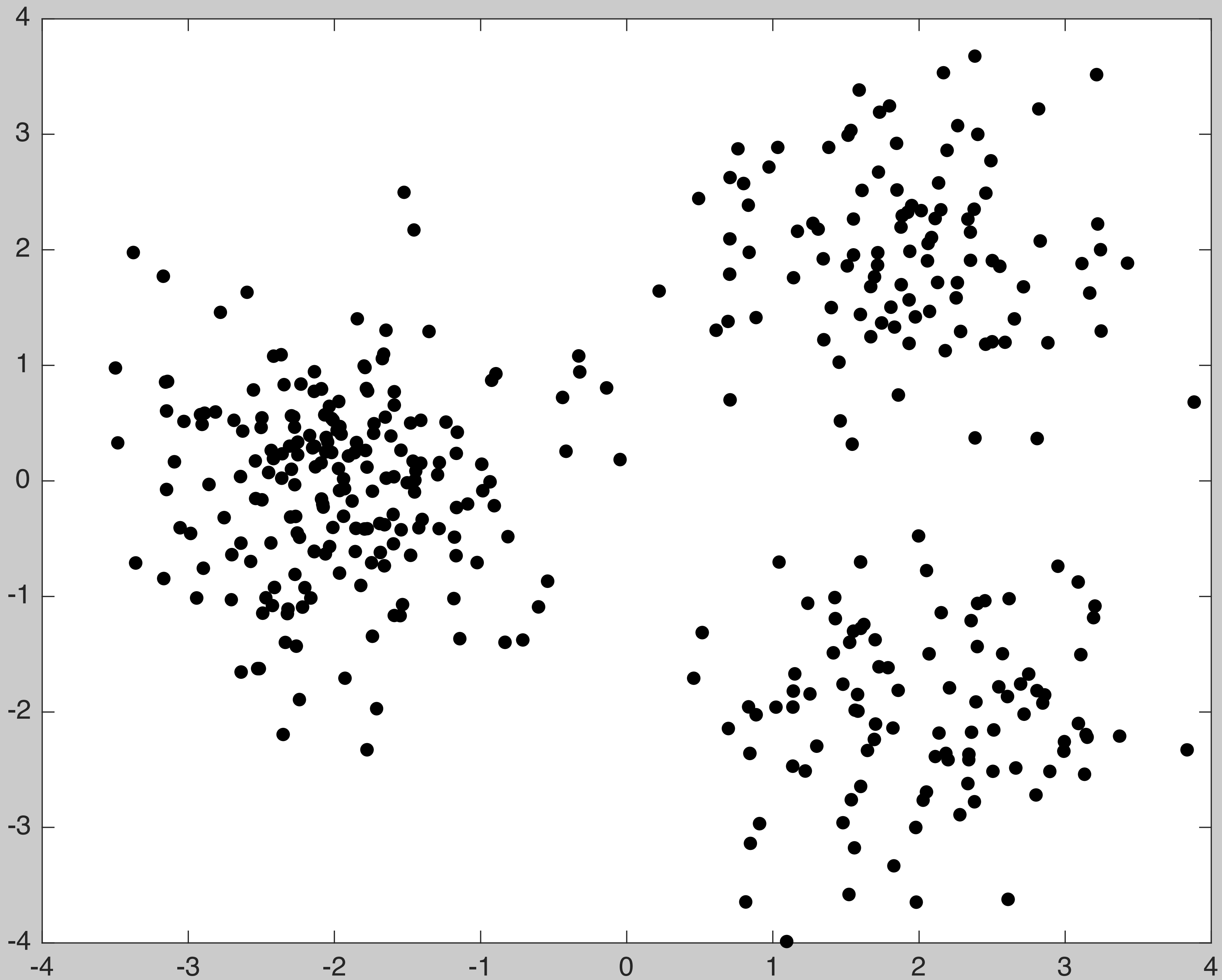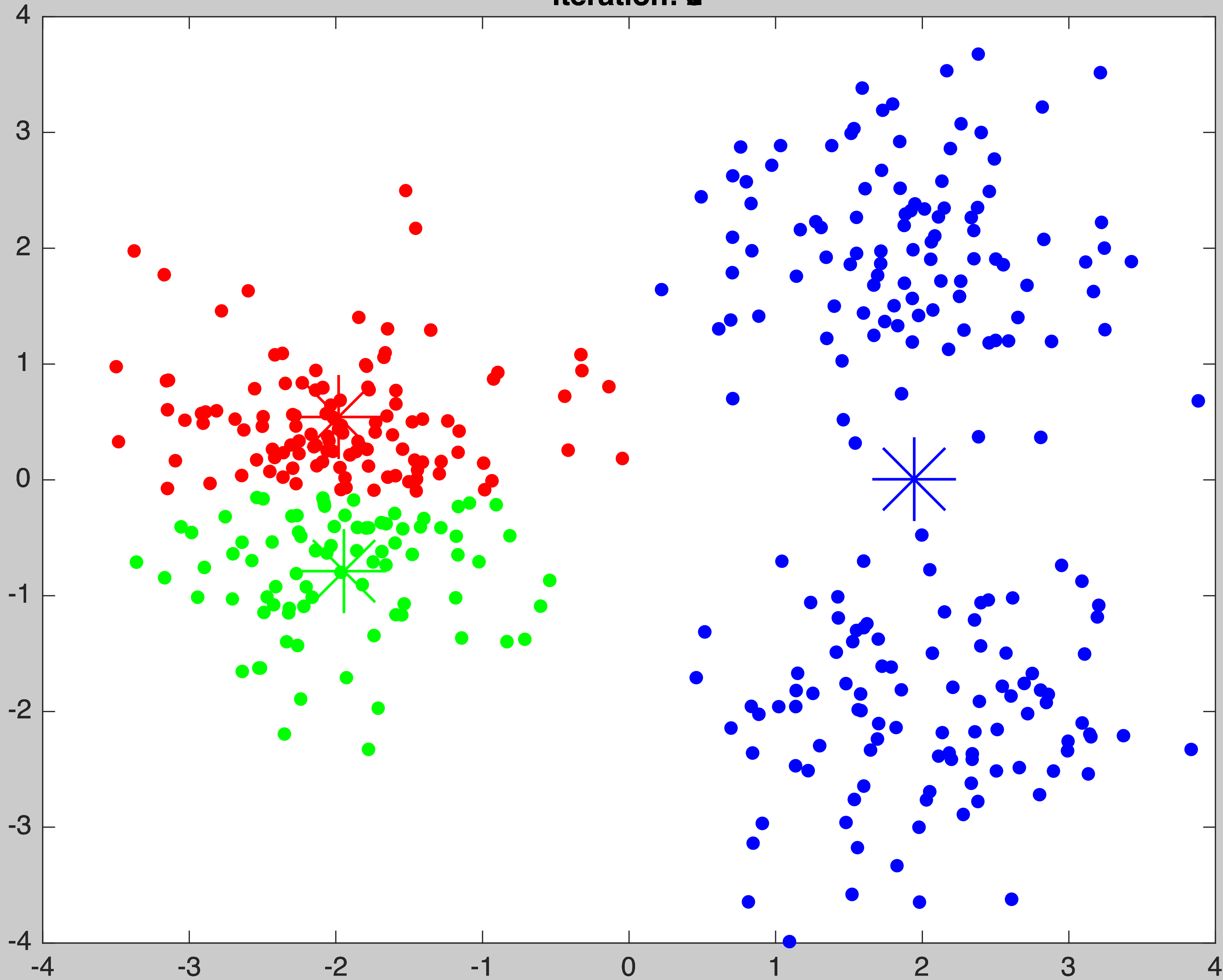Find the clustering of $\vec{x}_1, \cdots, \vec{x}_m$ into sets $S_1, \cdots, S_k$
which minimizes:

$$D = \sum_{i=1}^{k} \sum_{\vec{x} \in S_i} \|\vec{x} - \mu_i\| \qquad \mu_i = \frac{1}{|S_i|} \sum_{x \in S_i} \vec{x}$$

While the algorithm decreases the objective, the objective is non-convex and can be stuck on local mimima.

General problem is N-P Complete

# Management of intersections with multi-modal high-resolution data[☆,☆☆]

Ajith Muralidharan[1], Samuel Coogan[2], Christopher Flores, Pravin Varaiya[*]

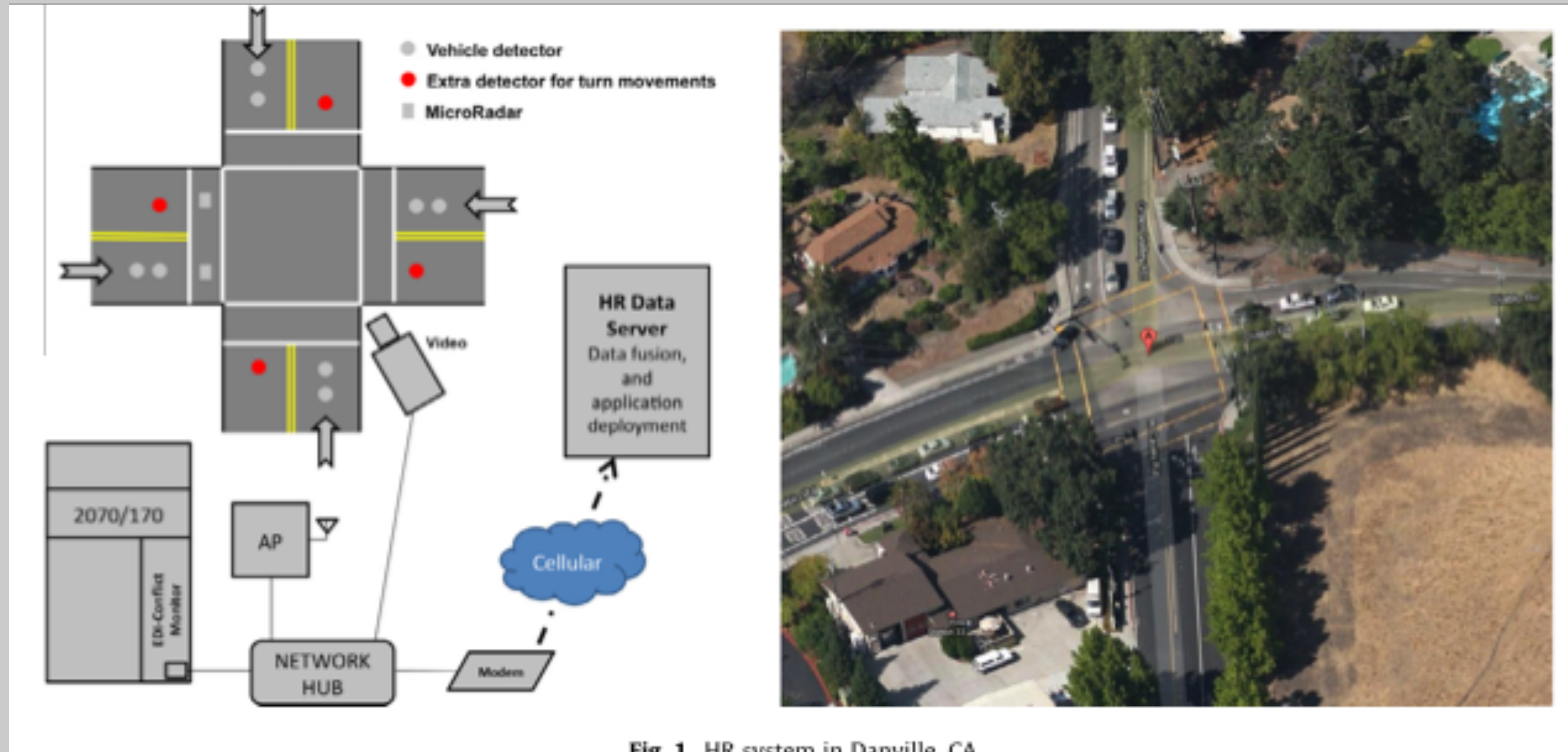*Sensys Networks, Inc, Berkeley, CA 94710, United States*

**Fig. 1.** HR system in Danville, CA.

# Traffic Patterns

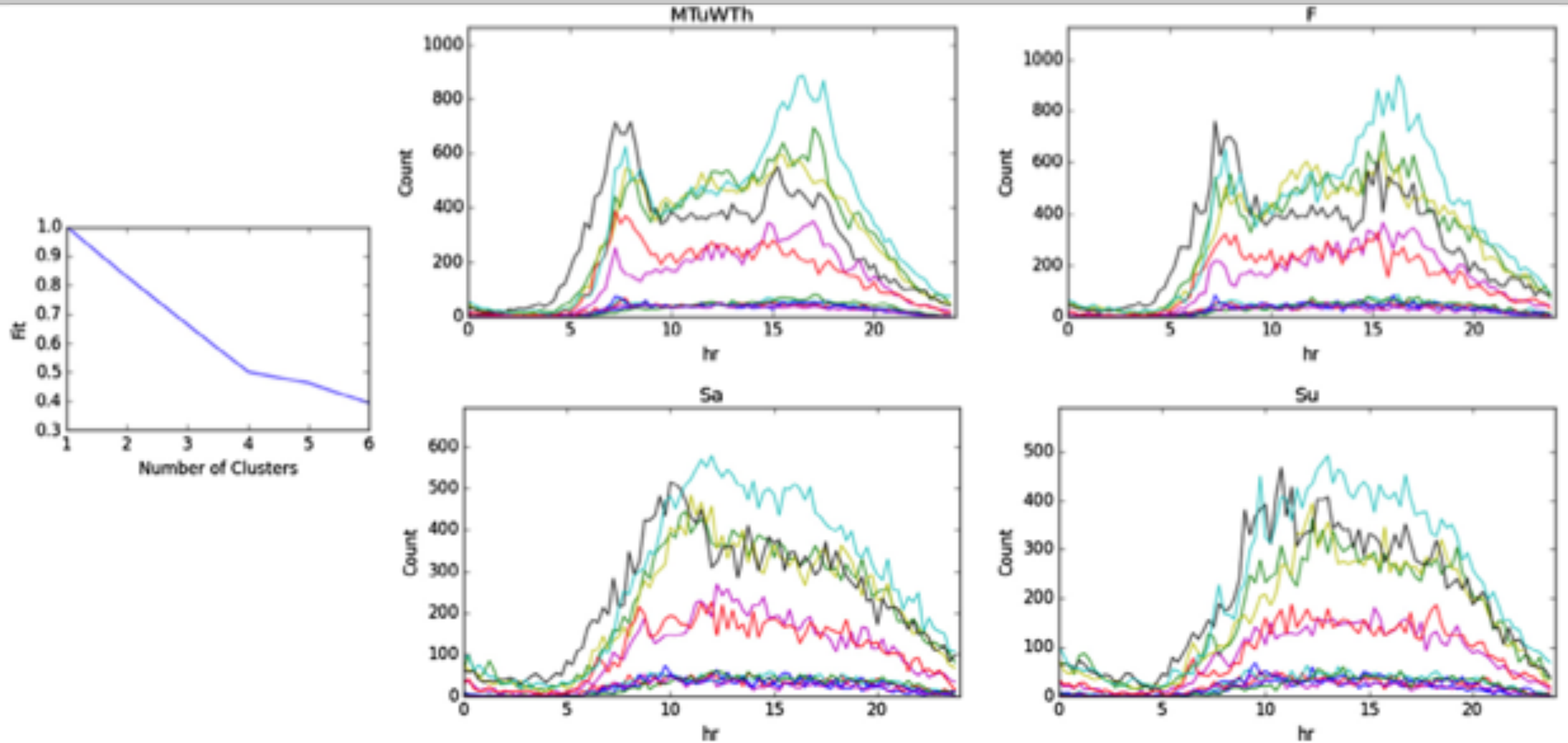| 1-2 | 2-3 | 3-4 | 4-5 | 5-6 | 6-7 | 7-8 | 8-9 | 9-10 | 10-11 | ... | 12-1 |
|-----|-----|-----|-----|-----|-----|-----|-----|------|-------|-----|------|
| | | | | | | | | | | | |

days

What would k-means cluster to?

**Fig. 5.** Clustering of daily data for Dec 2014 to May 2015 in an intersection in Beaufort, SC.
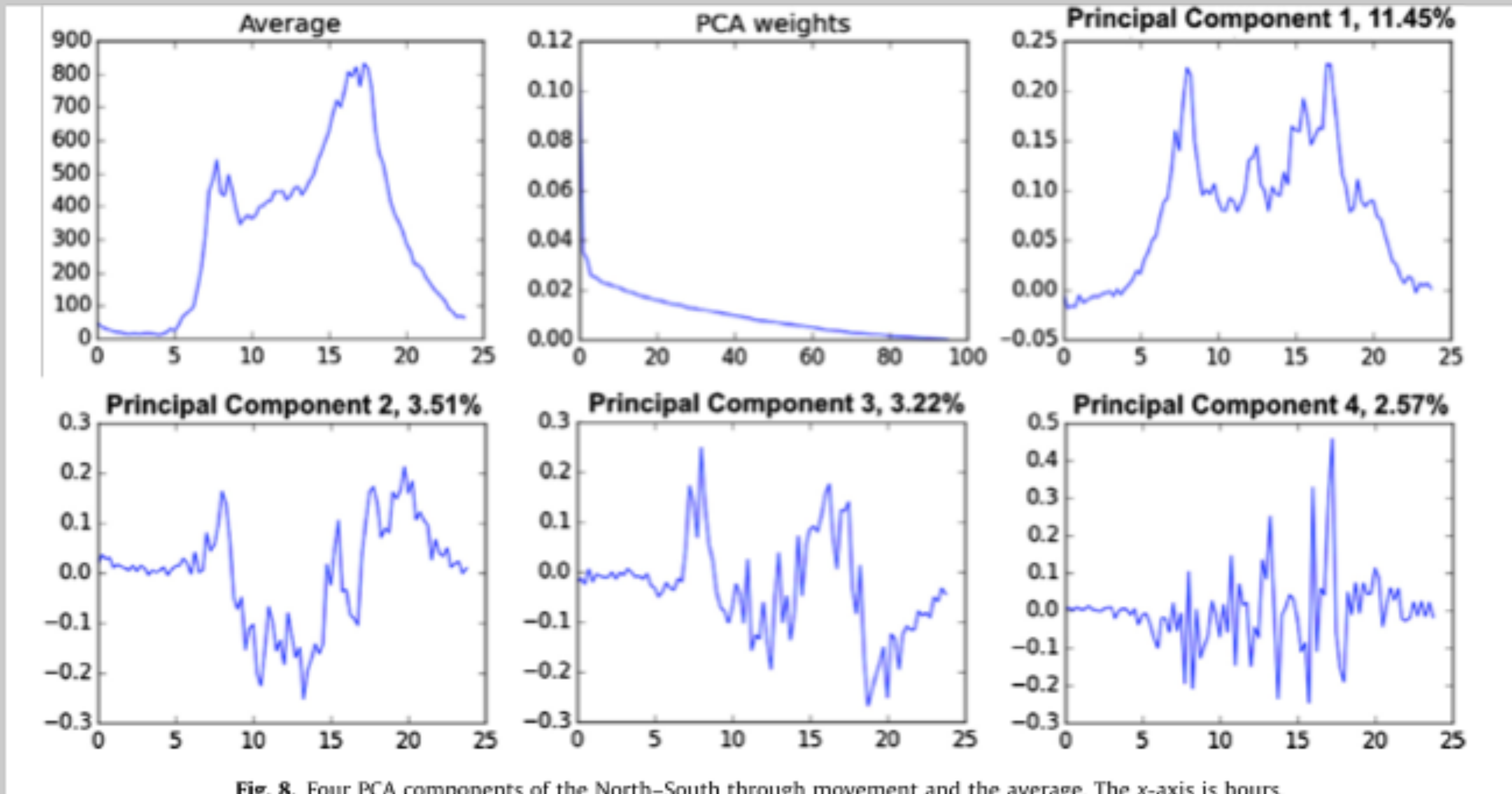
**Fig. 8.** Four PCA components of the North–South through movement and the average. The x-axis is hours.

EE16B M. Lustig, EECS UC Berkeley