

EECS16B

Designing Information  
Devices and Systems II

Lecture 14B

Principal Component Analysis

# Data Analysis with SVD

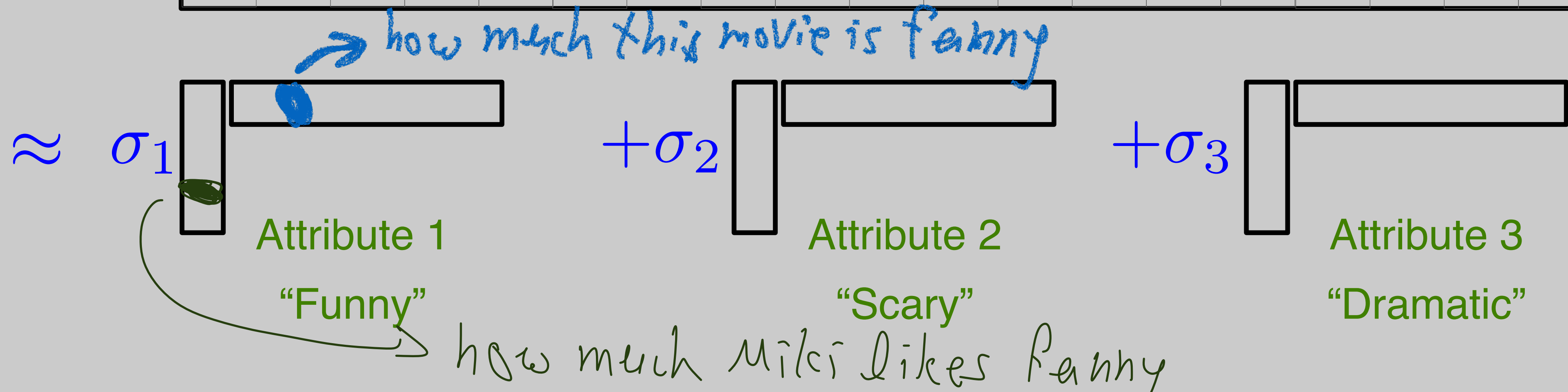
$$A \approx \sigma_1 \vec{u}_1 \vec{v}_1^T + \sigma_2 \vec{u}_2 \vec{v}_2^T + \dots + \sigma_{\hat{r}} \vec{u}_{\hat{r}} \vec{v}_{\hat{r}}^T$$

n movies ratings

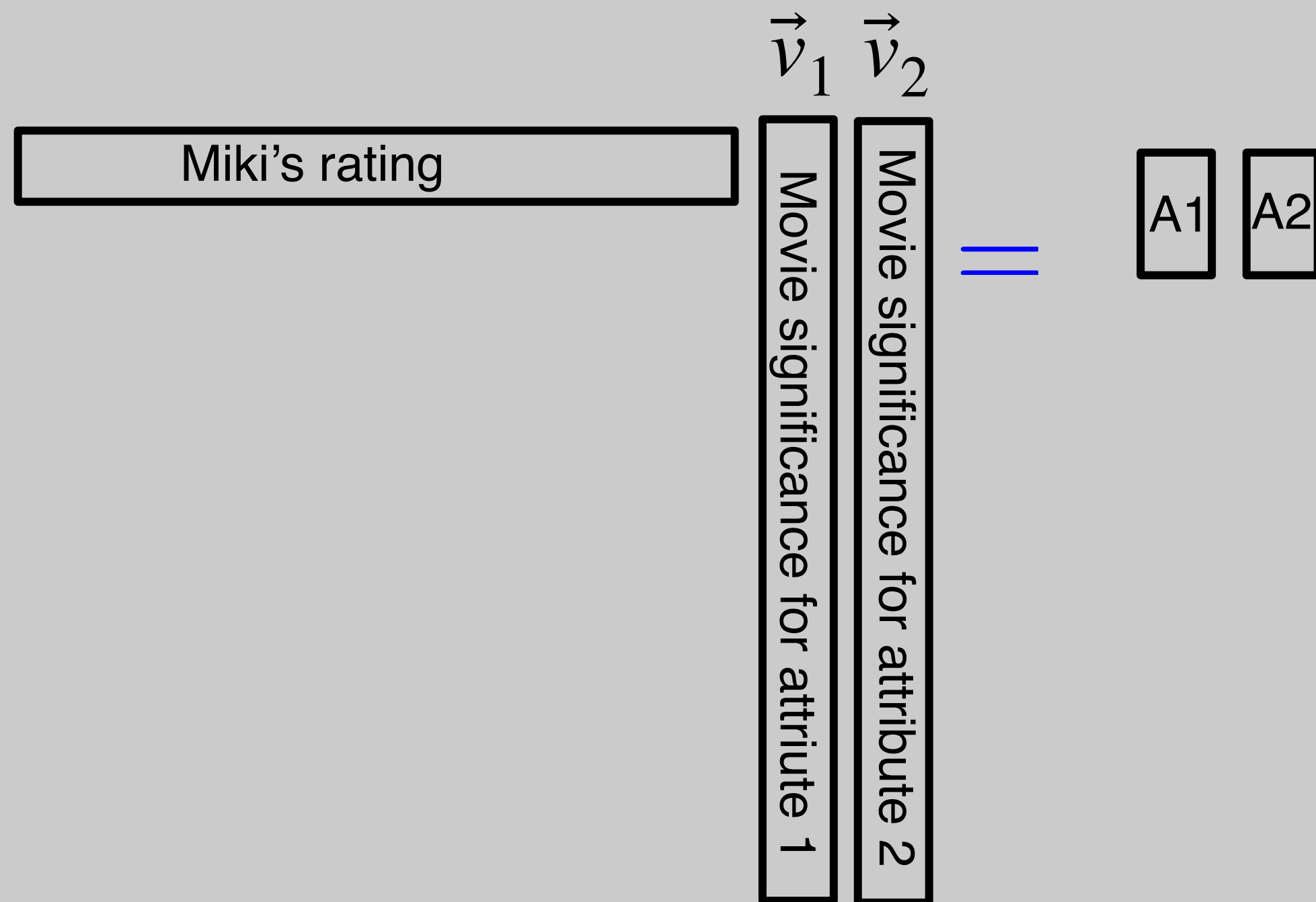
effective

m viewers

1	2	5	5	1	3	3	3	2	5	5	4	4	3	2	2	5	1	1
5	3	2	1	1	3	3	3	1	1	2	4	5	5	4	4	1	3	2
5	1	2	1	1	2	3	3	1	1	2	1	5	5	3	5	1	1	2
5	3	2	1	1	3	2	3	1	1	2	4	1	1	4	4	5	1	5
1	2	3	2	1	3	2	3	2	1	2	1	1	1	4	4	5	1	5
1	1	1	1	5	3	3	3	1	5	2	4	4	4	2	5	5	1	1



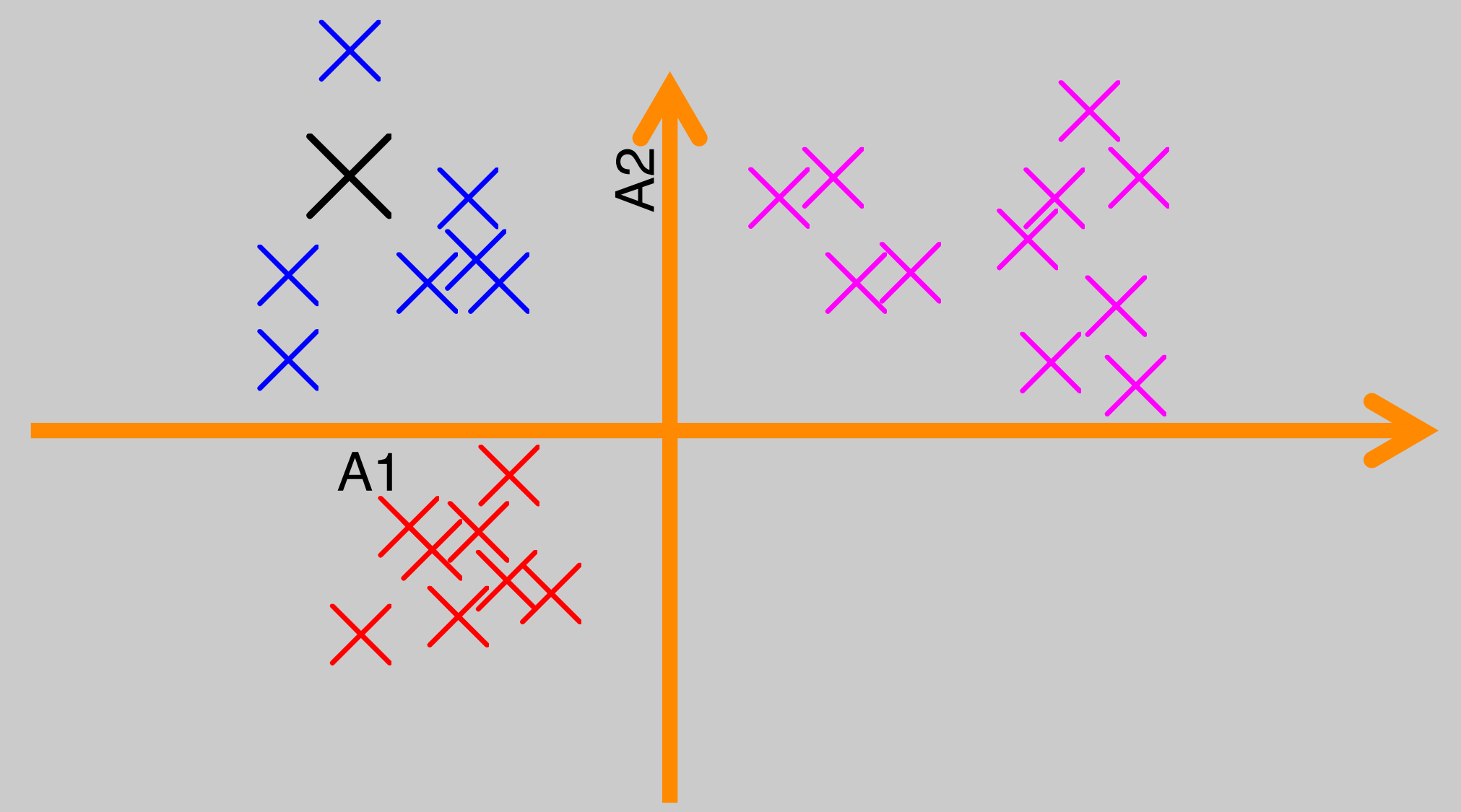
# Classification with SVD



m viewers

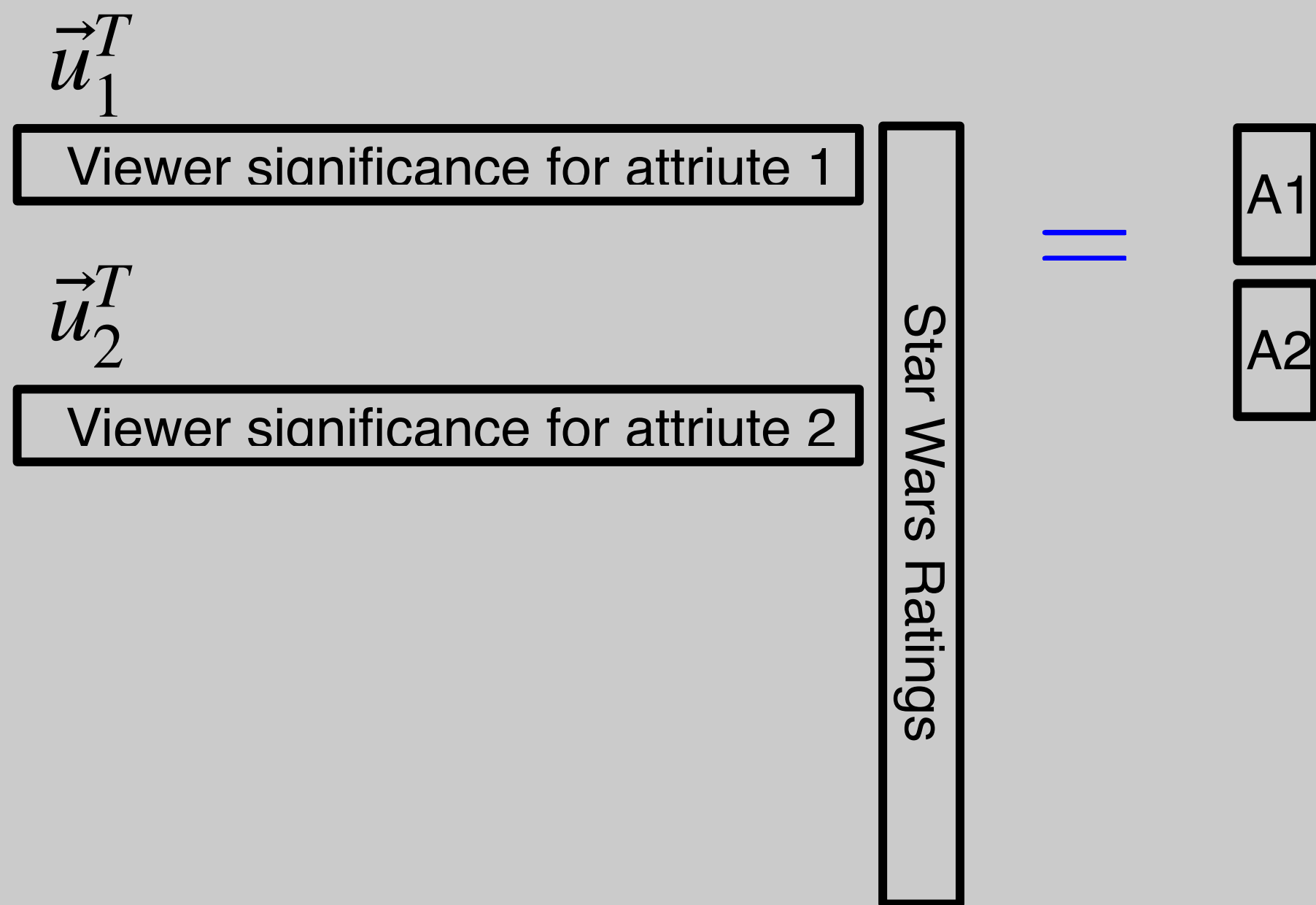
n movies

1	2	5	5	1	3	3	3	2	5	5	4	4	3	2	2	5	1	1
5	3	2	1	1	3	3	3	1	1	2	4	5	5	4	4	1	3	2
5	1	2	1	1	2	3	3	1	1	2	1	5	5	3	5	1	1	2
5	3	2	1	1	3	2	3	1	1	2	4	1	1	4	4	5	1	5
1	2	3	2	1	3	2	3	2	1	2	1	1	1	4	4	5	1	5
1	1	1	1	5	3	3	3	1	5	2	4	4	4	2	5	5	1	1



Miki belongs to class: like A2 don't like A1

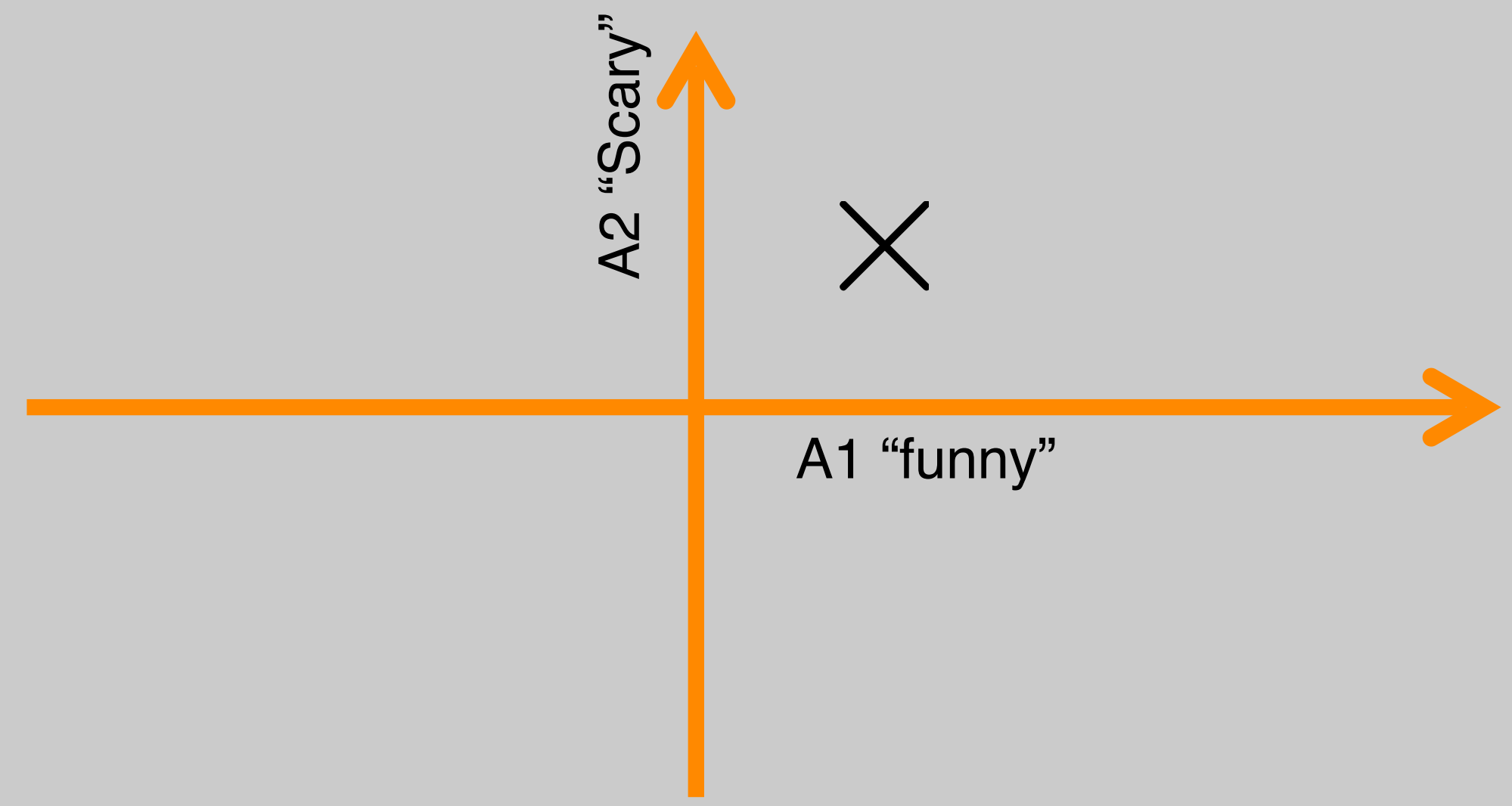
# Classification with SVD



m viewers

n movies

1	2	5	5	1	3	3	3	2	5	5	4	4	3	2	2	5	1	1
5	3	2	1	1	3	3	3	1	1	2	4	5	5	4	4	1	3	2
5	1	2	1	1	2	3	3	1	1	2	1	5	5	3	5	1	1	2
5	3	2	1	1	3	2	3	1	1	2	4	1	1	4	4	5	1	5
1	2	3	2	1	3	2	3	2	1	2	1	1	1	4	4	5	1	5
1	1	1	1	5	3	3	3	1	5	2	4	4	4	2	5	5	1	1



Star Wars is somewhat funny and somewhat scary

# Prediction with SVD

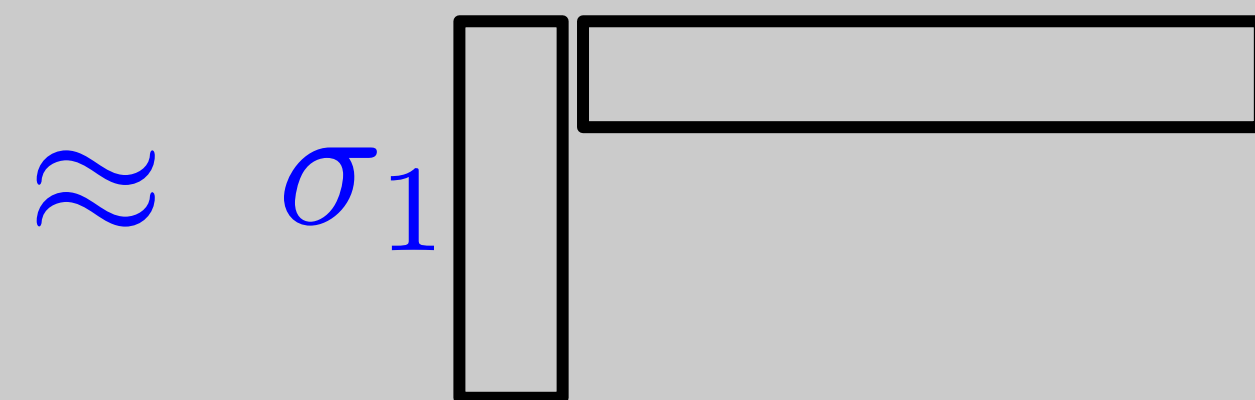
Can try to predict preferences of a new customer with few ratings

See homework!

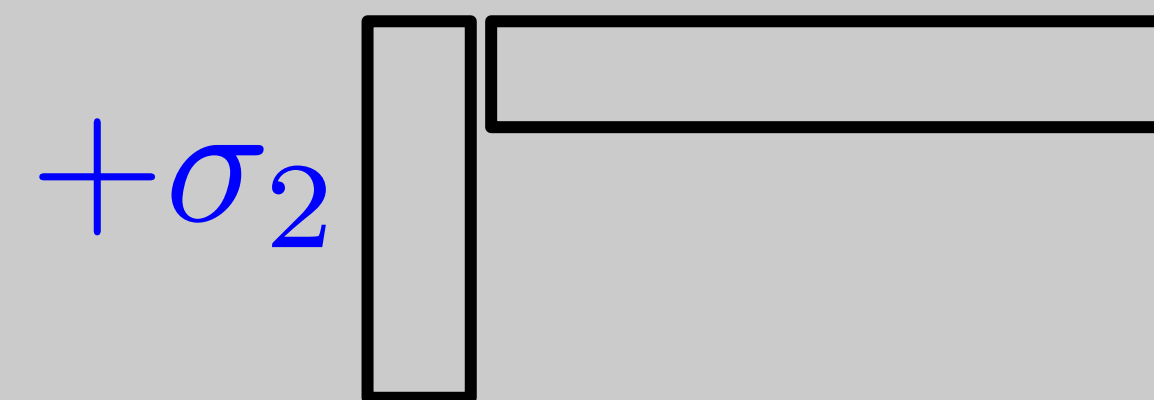
n movies

1	2	5	5	1	3	3	3	2	5	5	4	4	3	2	2	5	1	1	
5	3	2	1	1	3	3	3	1	1	2	4	5	5	4	4	1	3	2	
5	1	2	1	1	2	3	3	1	1	2	1	5	5	3	5	1	1	2	
5	3	2	1	1	3	2	3	1	1	2	4	1	1	4	4	5	1	5	
1	2	3	2	1	3	2	3	2	1	2	1	1	1	4	4	5	1	5	
1	1	1	1	5	3	3	3	1	5	2	4	4	4	2	5	5	1	1	
1	?	?	2	?	?	?	?	3	5	1	?	?	?	?	5	2	?	3	

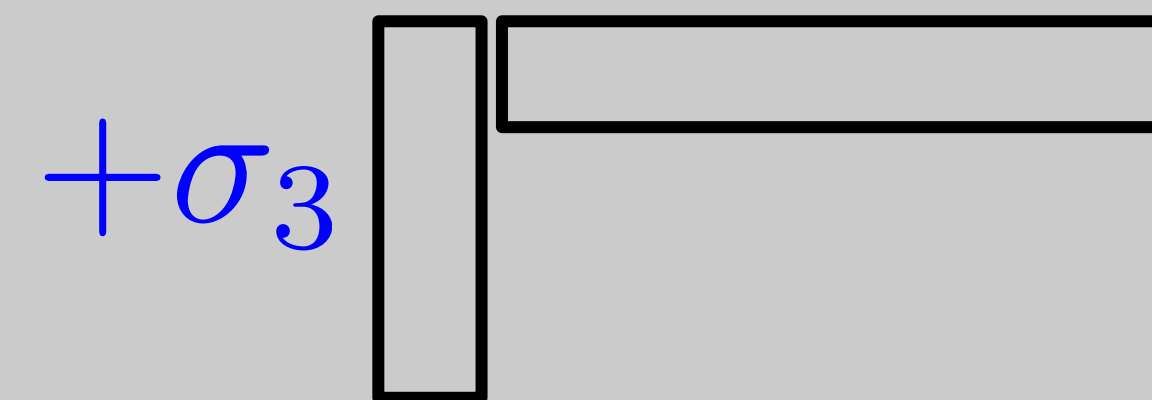
m views



Attribute 1



Attribute 2



Attribute 3

# Low-rank Completion

What if my database is full of “holes”?

Should be still low-rank!

n movies

1		5	5	1	3		3	2	5	5	4		3	2	2	5		1	
5	3		1	1	3	3		1	1	2		5		4	4	1	3	2	
	1	2	1	1		3	3		1	2	1	5	5	3	5		1	2	
5	3		1		3	2		1		2	4	1		4		5		5	
1		3	2	1	3	2	3	2	1		1		1	4		5	1	5	
1	1			5	3	3		5	2	4			2	5	5	1	1		

m views

Q) Can we complete missing data?

A) Sometimes! Very recent mathematical and practical results show you can.  
Keywords: Compressed Sensing, Low-rank completion, robust PCA

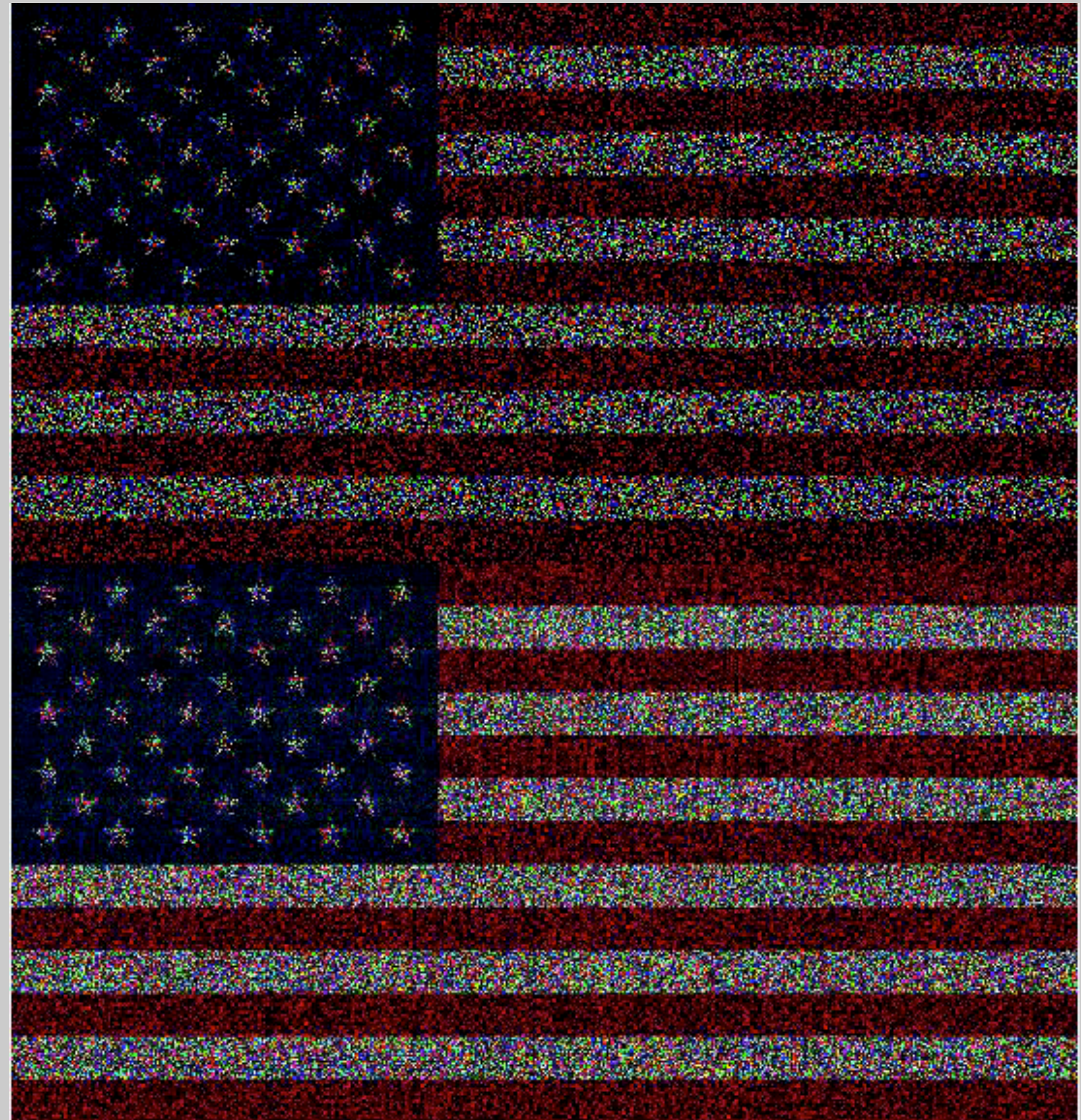
E. Candes and B. Recht, *Foundation of Computational Mathematics*, 2009;9:717





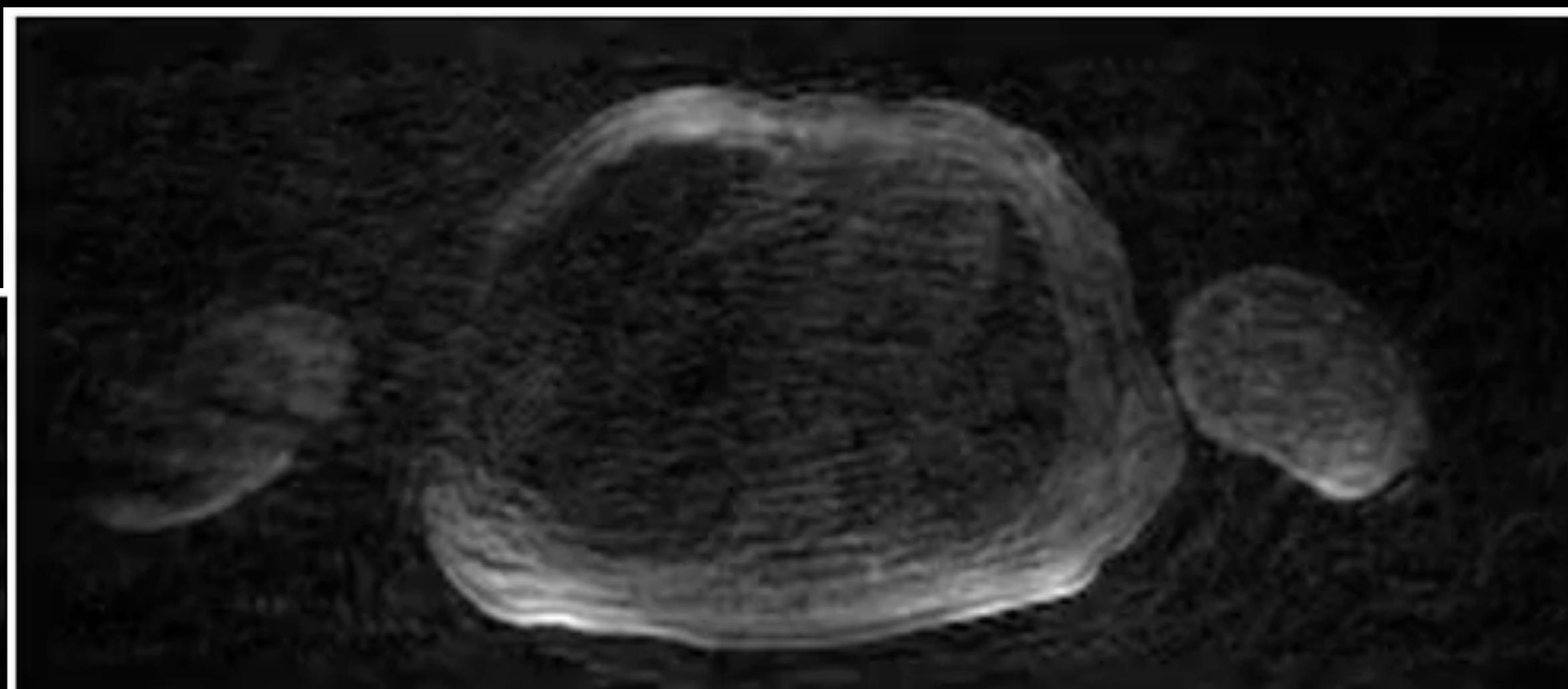
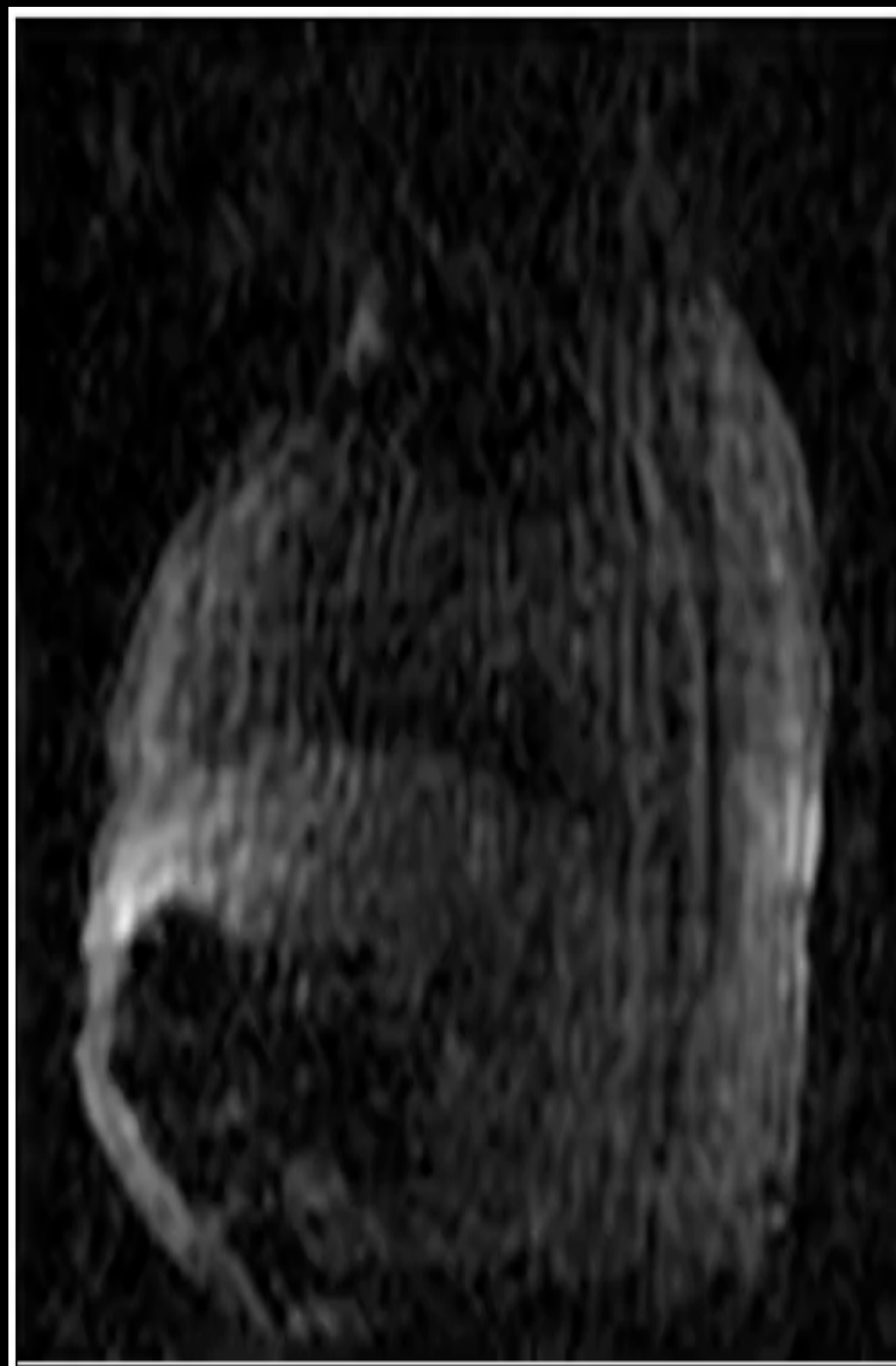
# Low-Rank Recovery from 20% pixels

- Algorithm for low-rank completion:
  - `flag_hat = flag`
  - Compute  $[U, S, V] = \text{svd}(\text{flag\_hat})$
  - $$\text{flag}_{\text{hat}} = \sum_{i=0}^6 \sigma_i \vec{u}_i \vec{v}_i^T$$
  - update missing pixels in `flag` from `flag_hat`
  - repeat (250 times here)





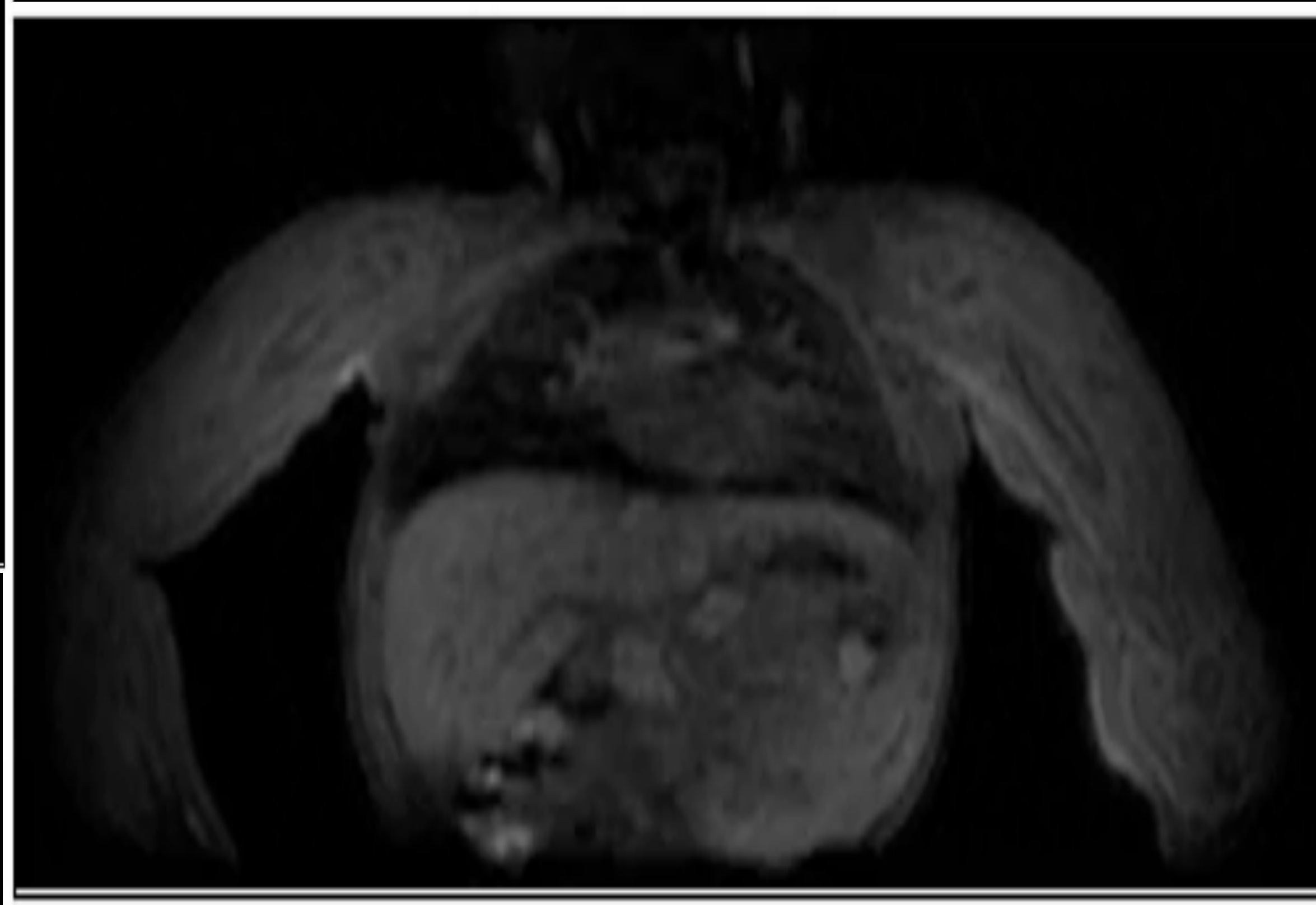
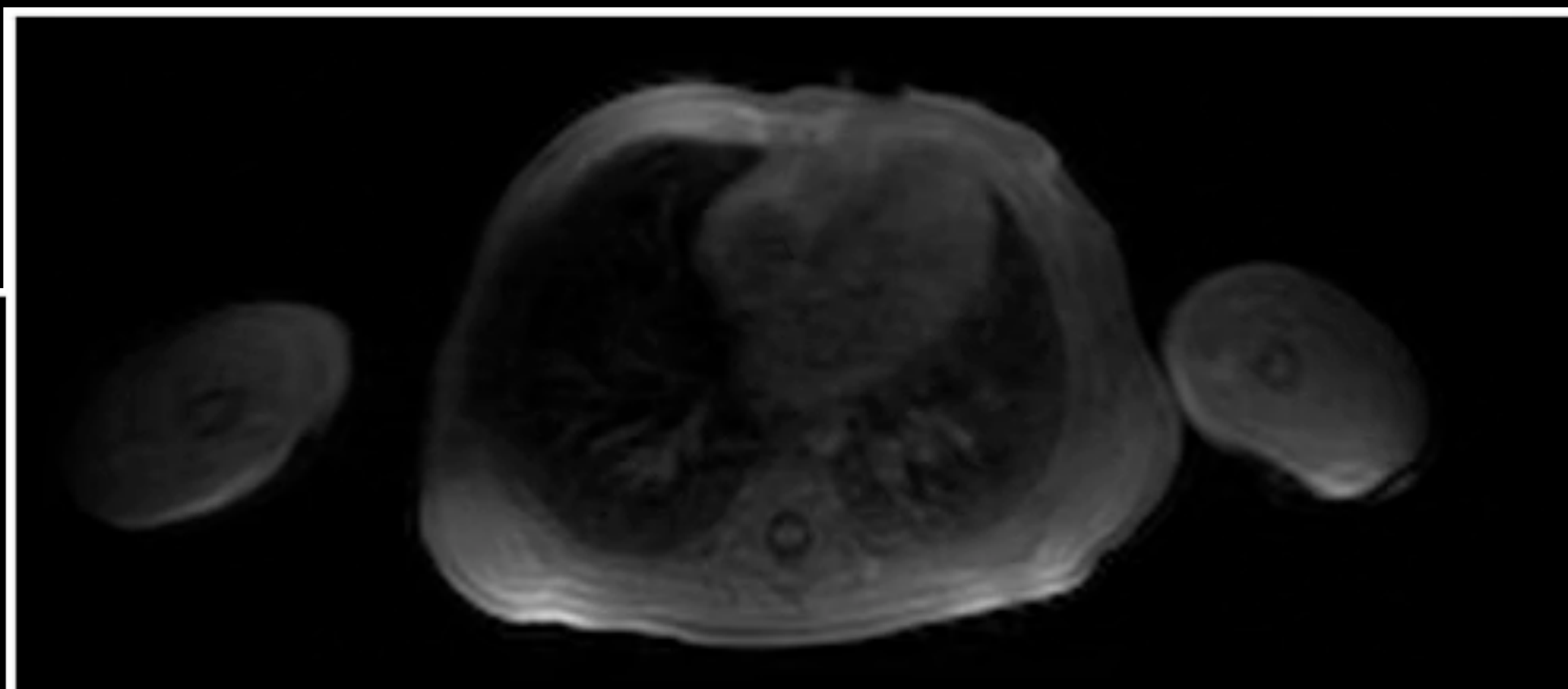
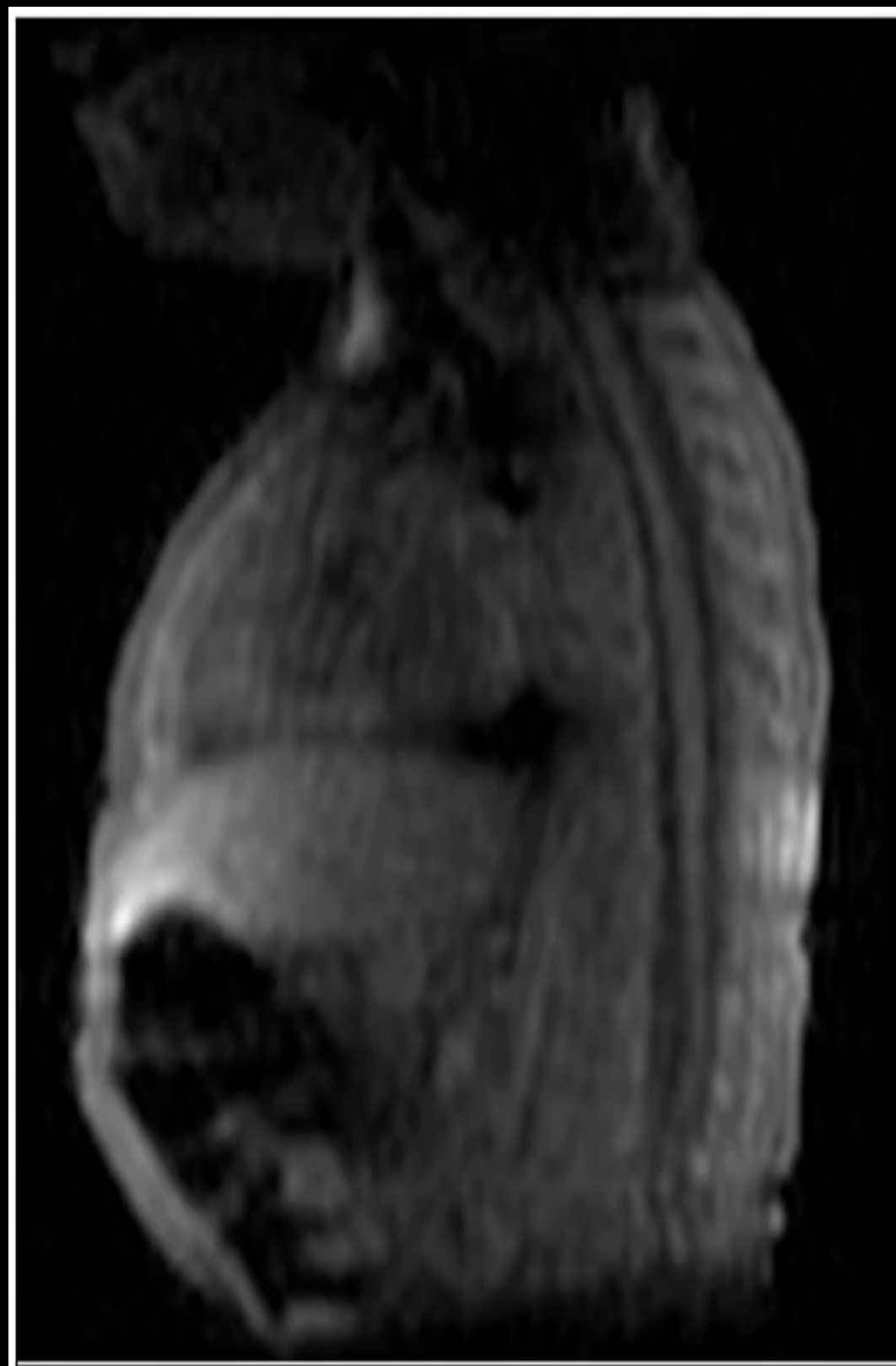
# Dynamic MRI Missing Data



**$\sim 1.5 \times 1.5 \times 3 \text{ mm}^3$**



# Dynamic MRI Low-Rank



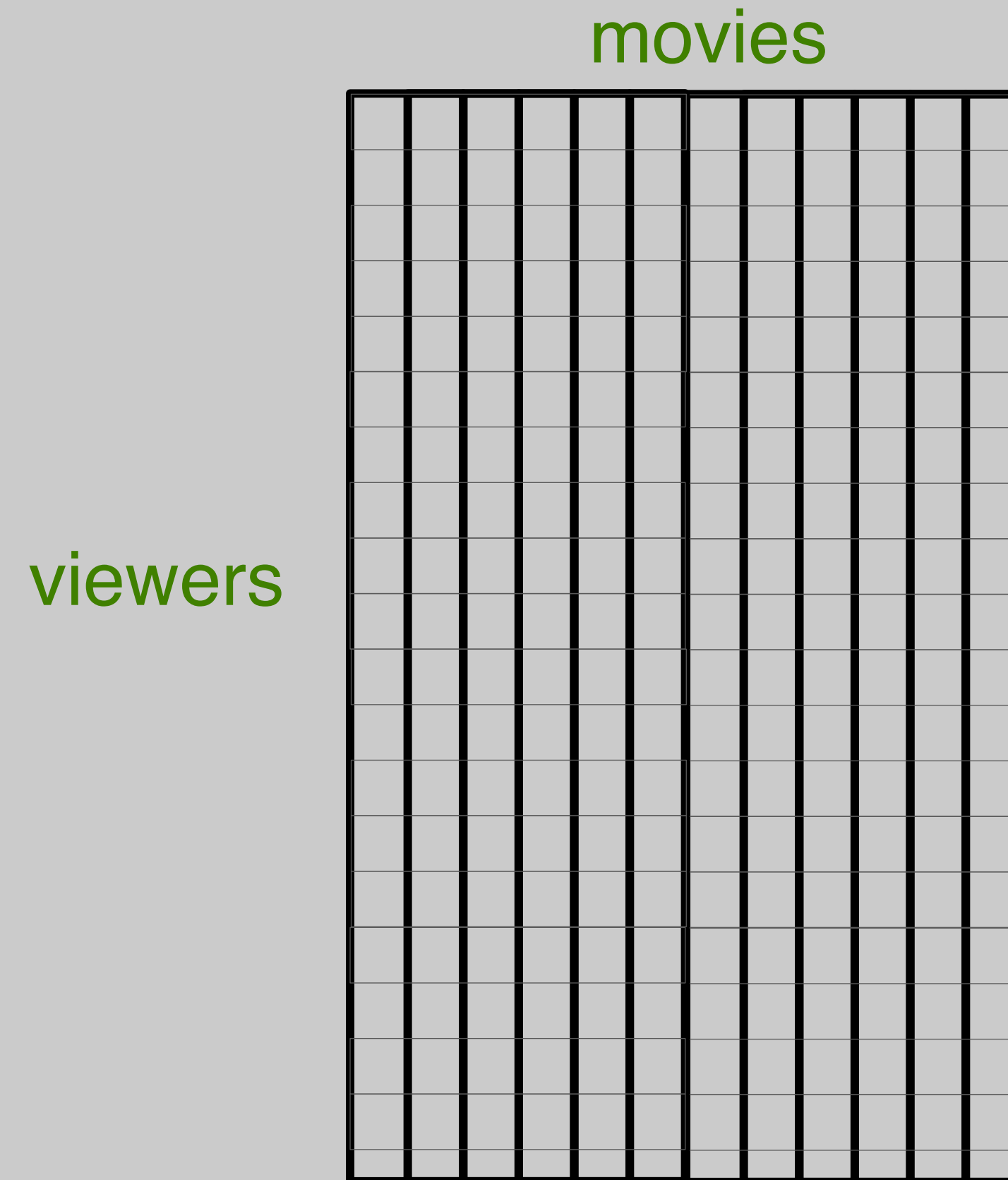
**$\sim 1.5 \times 1.5 \times 3 \text{ mm}^3$**

# 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

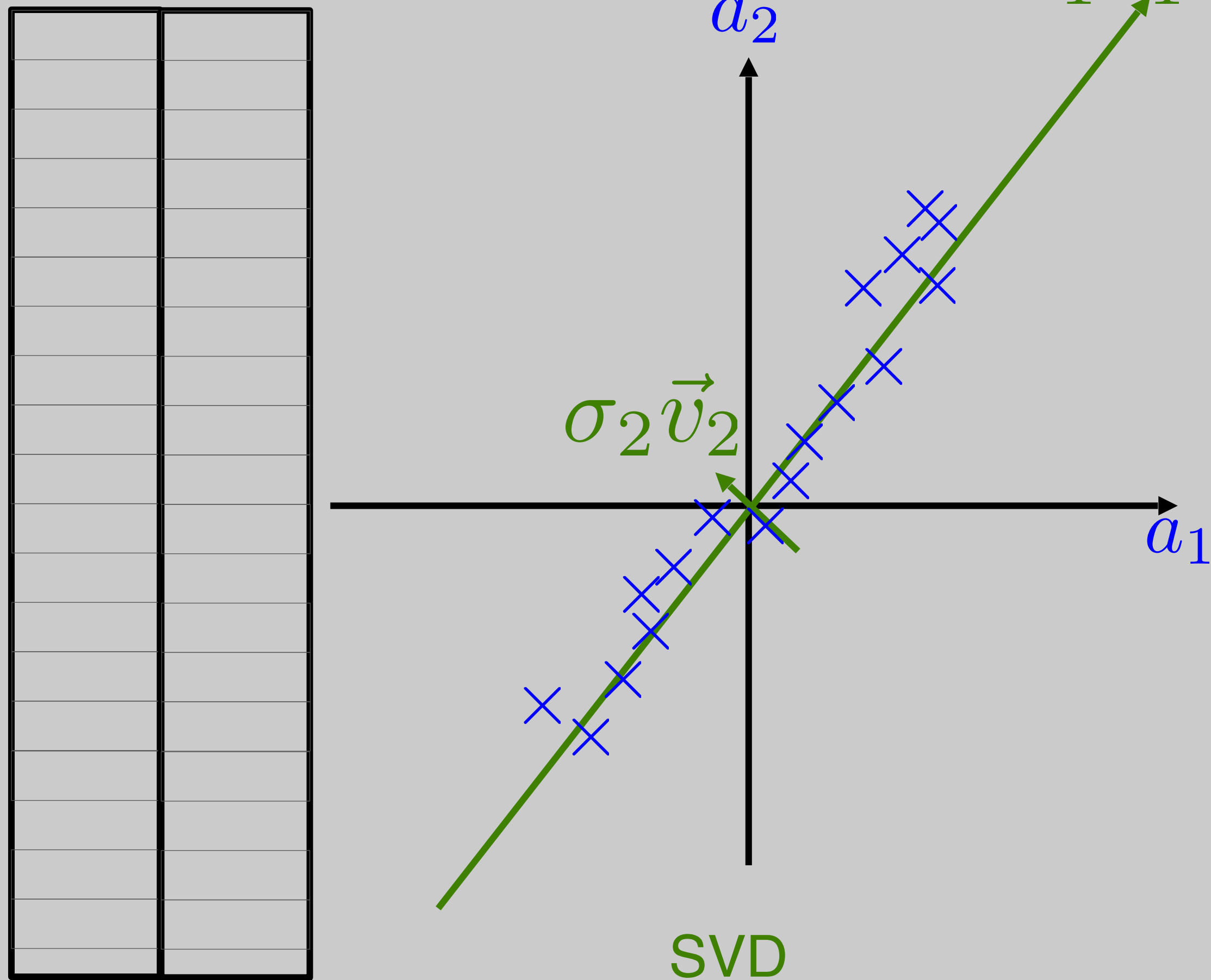




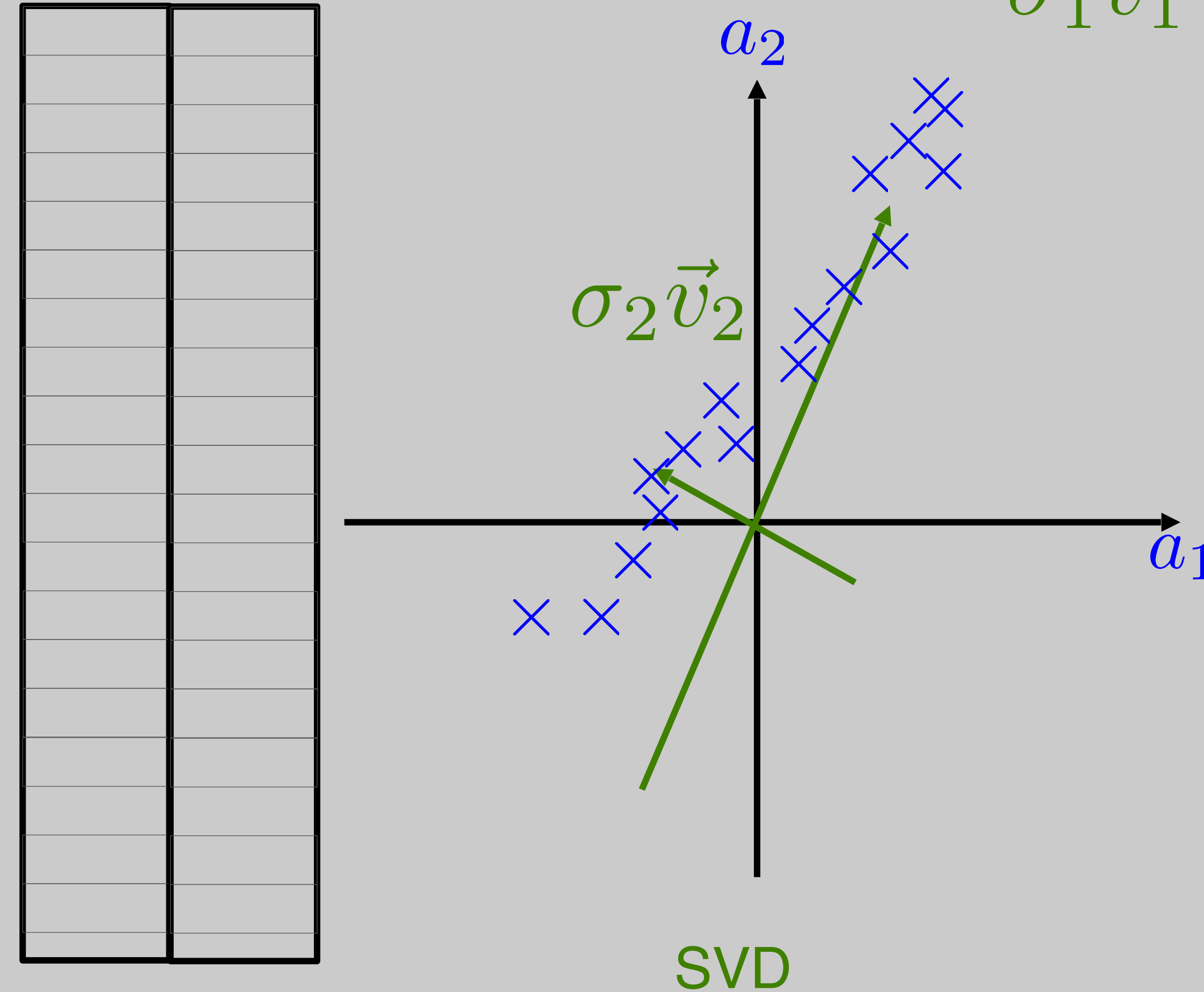
# Example -- PCA

Consider data s.t.

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



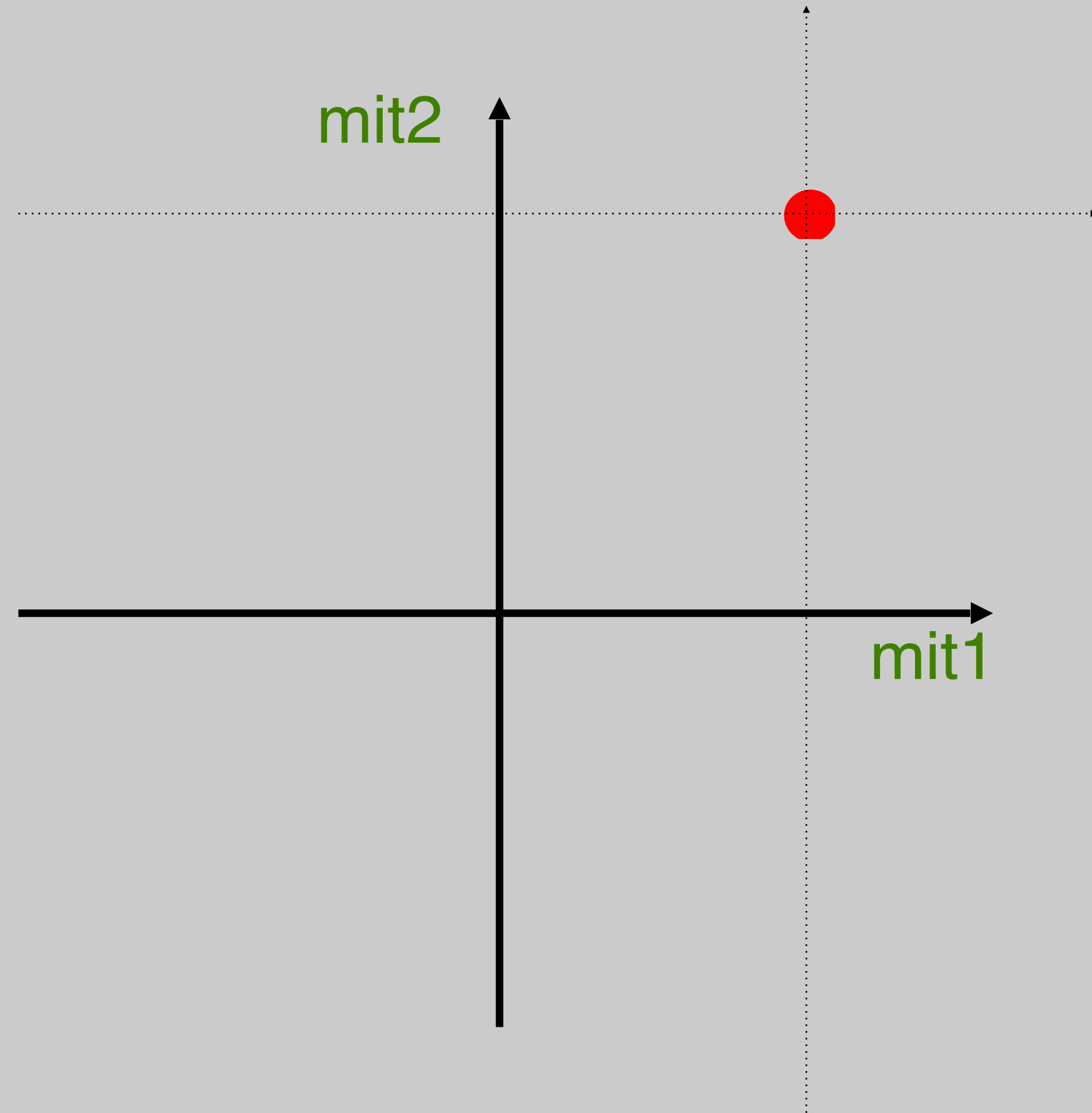
$$\vec{a}_1 \quad \vec{a}_2 \approx 3\vec{a}_1 + 1$$



# Example -- PCA

Consider midterm data

	mit1	mit2
students		

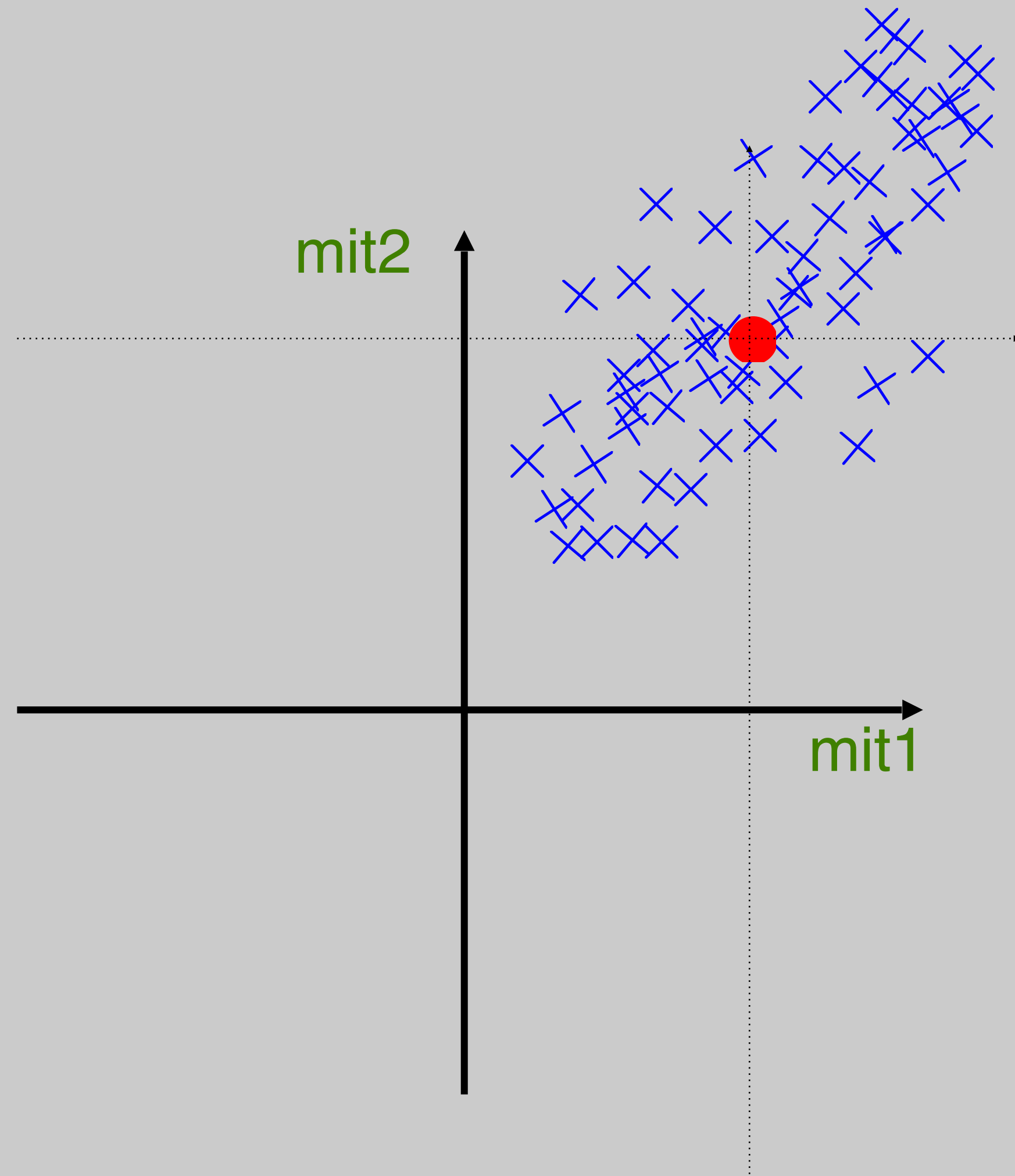




# Example -- PCA

Consider midterm data

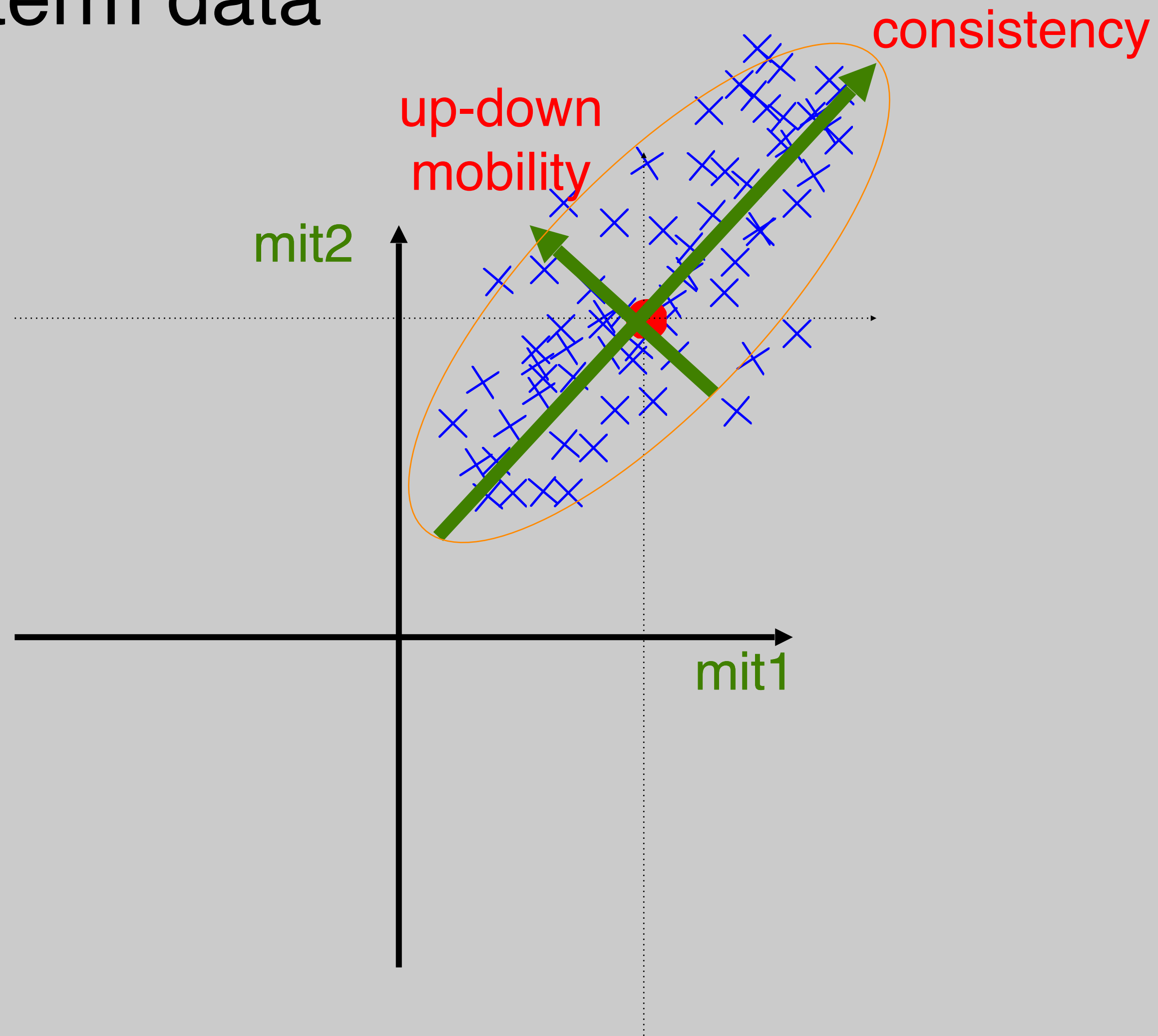
	mit1	mit2
students		



# Example -- PCA

Consider miterm data

	mit1	mit2
students		





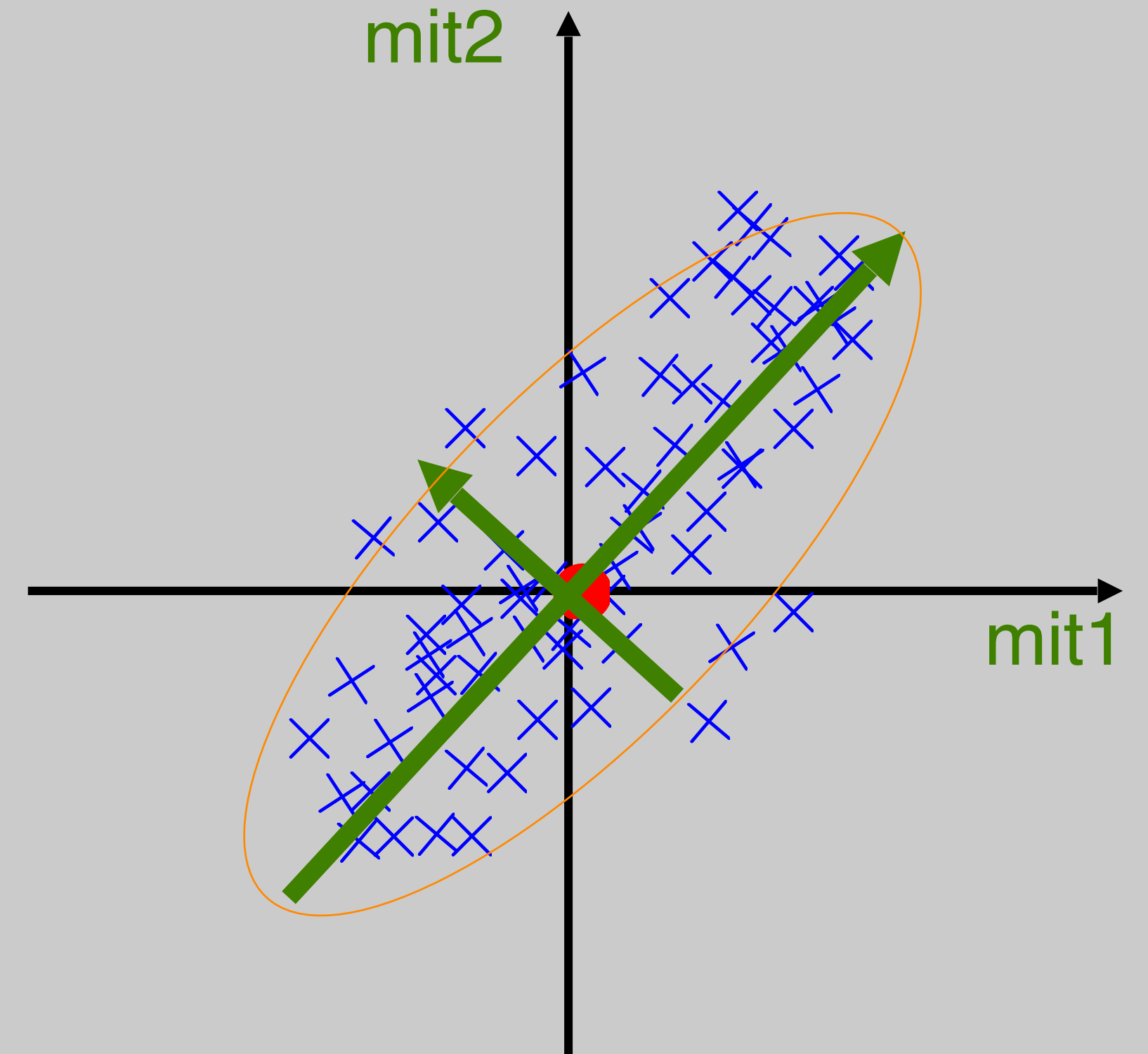
# PCA Procedure

Remove averages from column of A

From  $A^T A$ , find  $\sigma_i$ ,  $\vec{v}_i$

$\vec{v}_i$  are principal components!

	mit1	mit2
students		



# $A^T A$ as sample covariance matrix

---

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

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

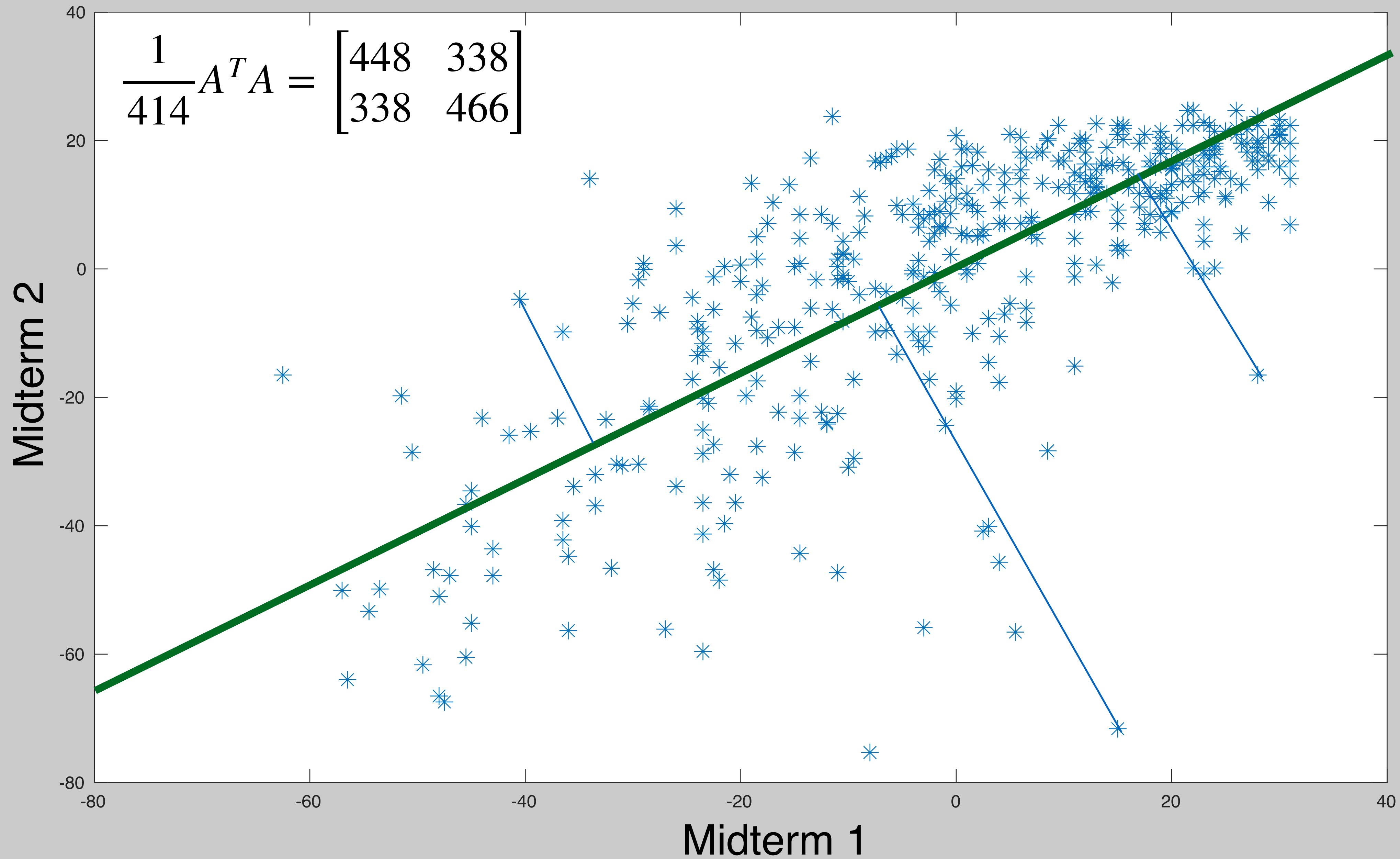
$$= \vec{a}^T \vec{a} - 2N a_\mu^2 + N a_\mu^2 = \vec{a}^T \vec{a} - N a_\mu^2$$

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

Sample Variance!

# Example midterm

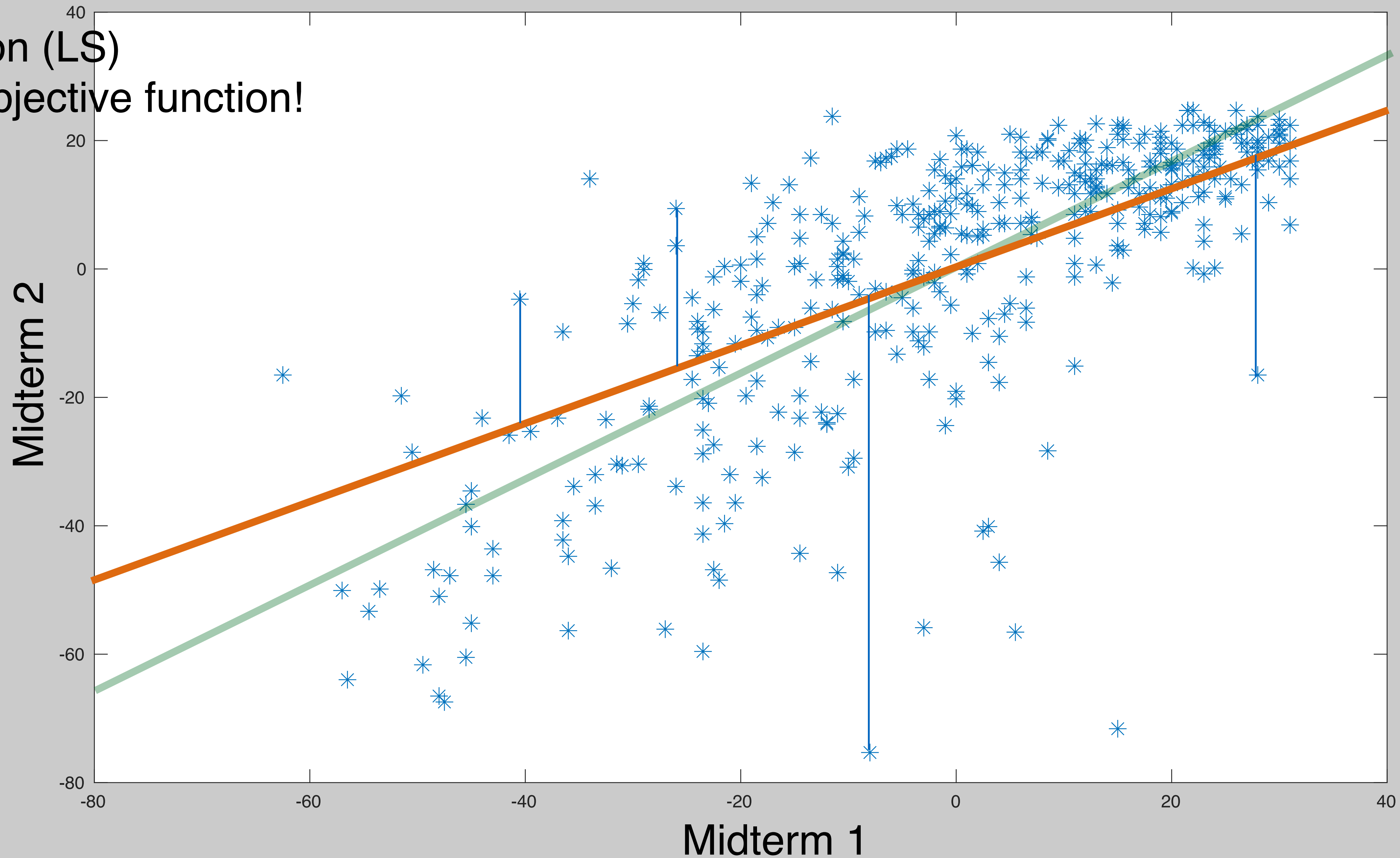
---





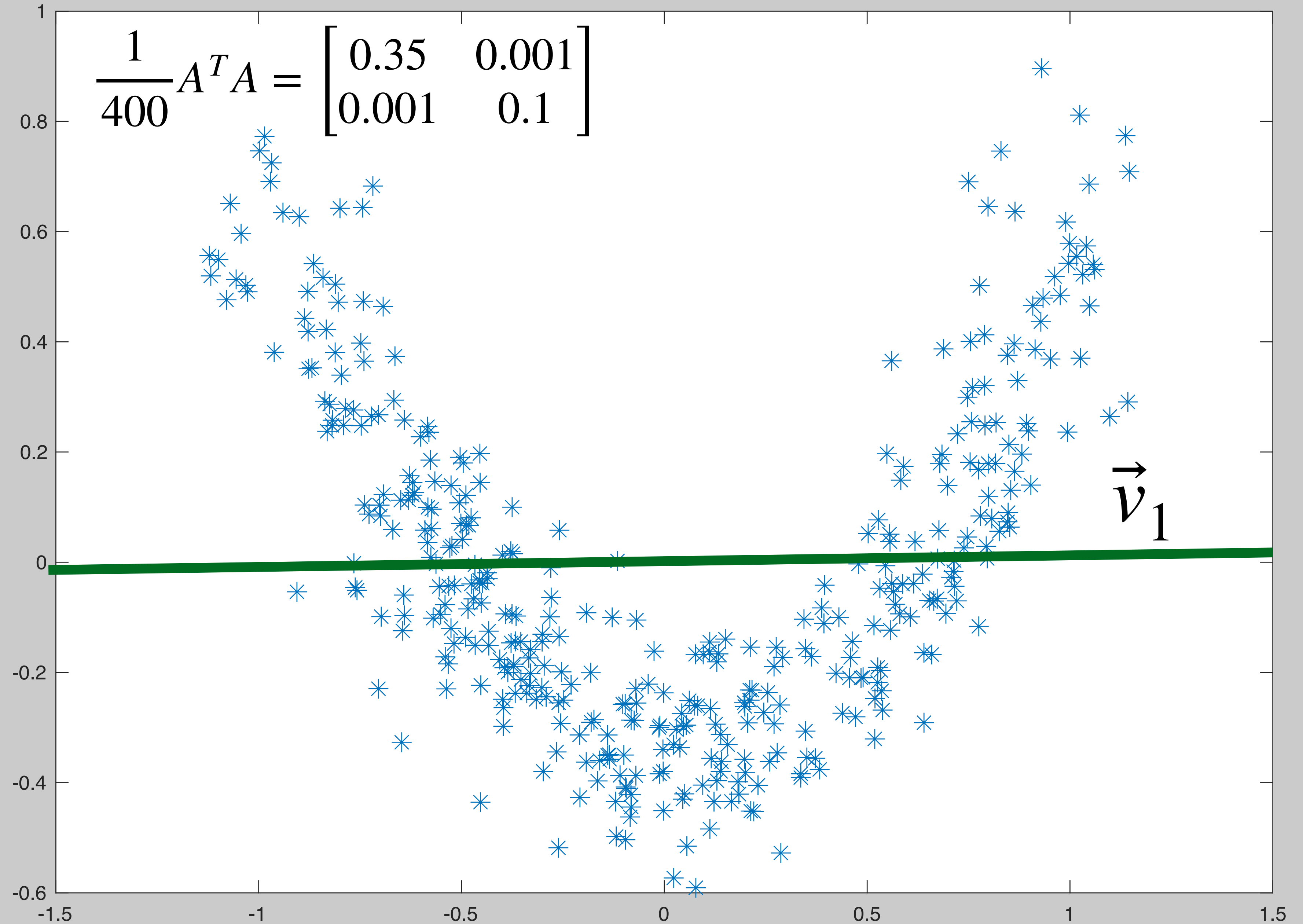
# Example midterm

Linear regression (LS)  
Not the same objective function!



# PCA captures linear correlations

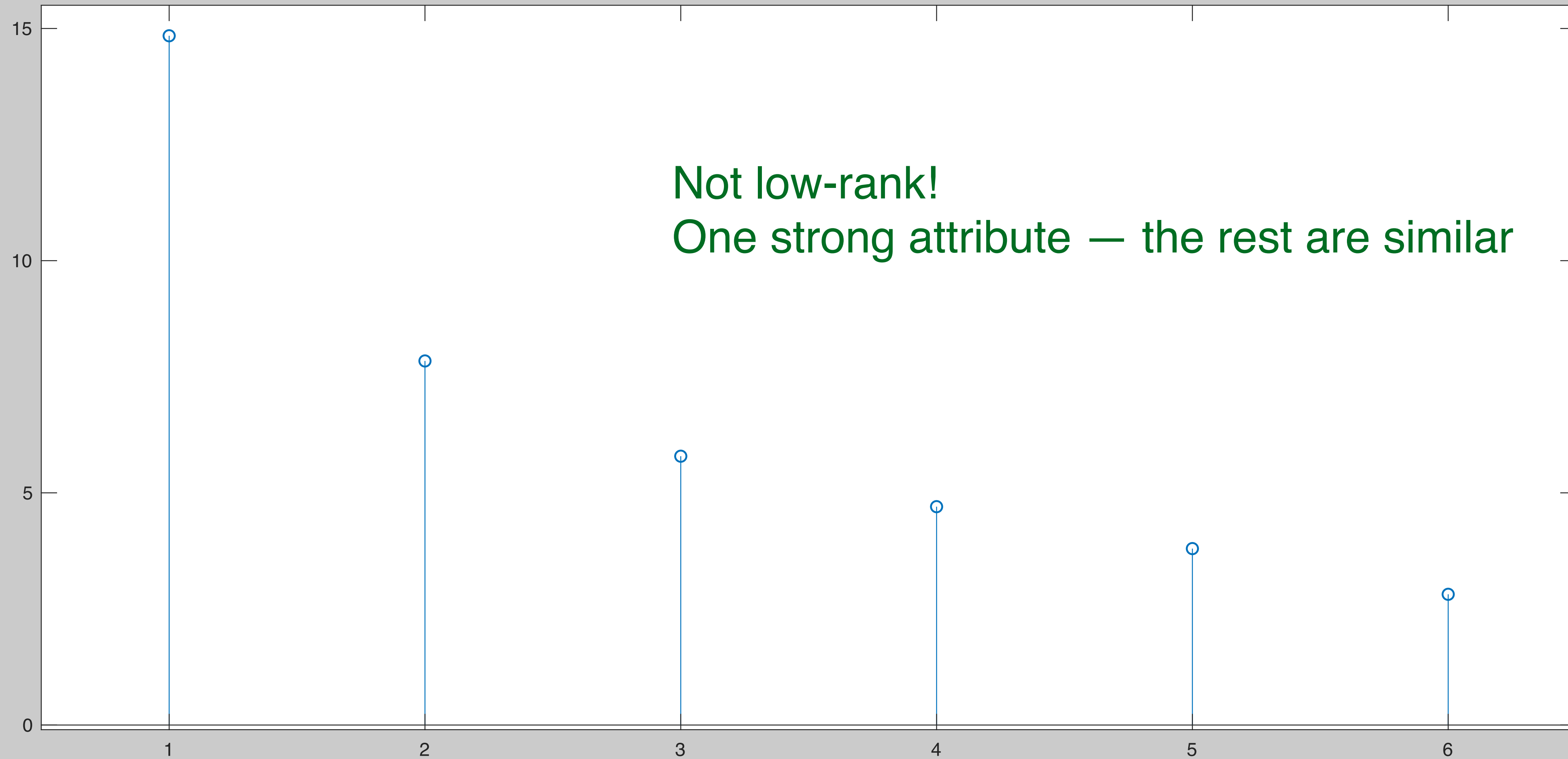
Linear subspace does not capture the low dimensionality of a one dimensional manifold







## Singular values

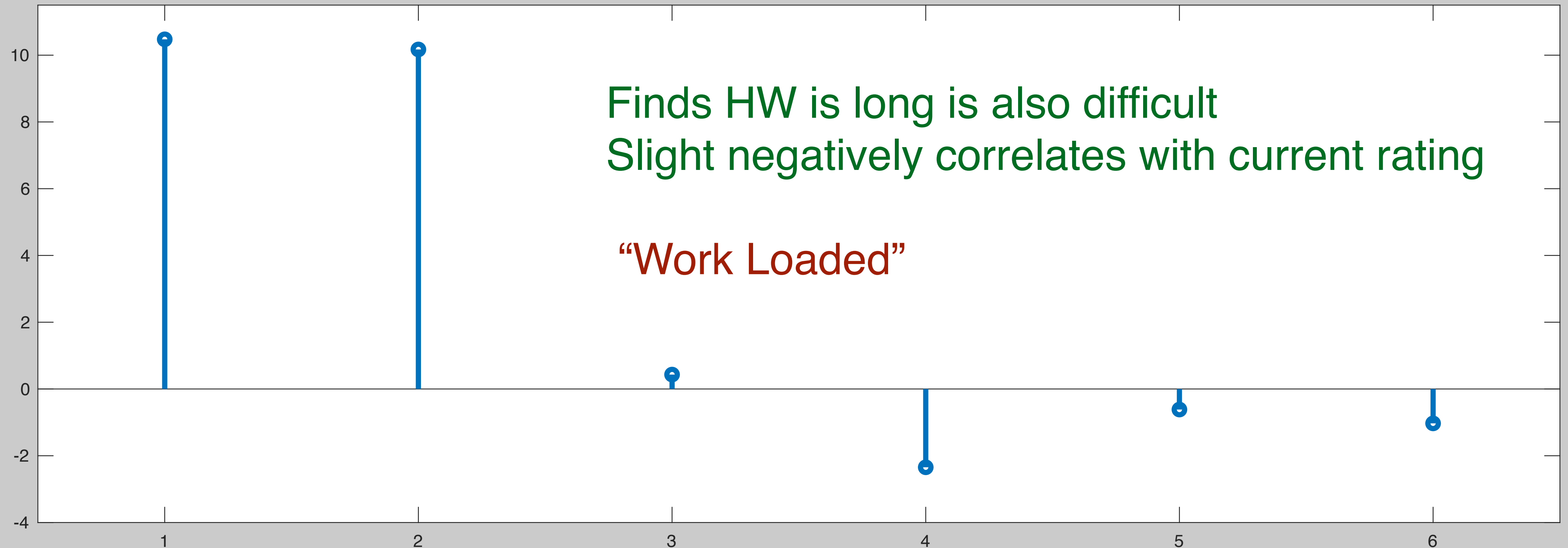
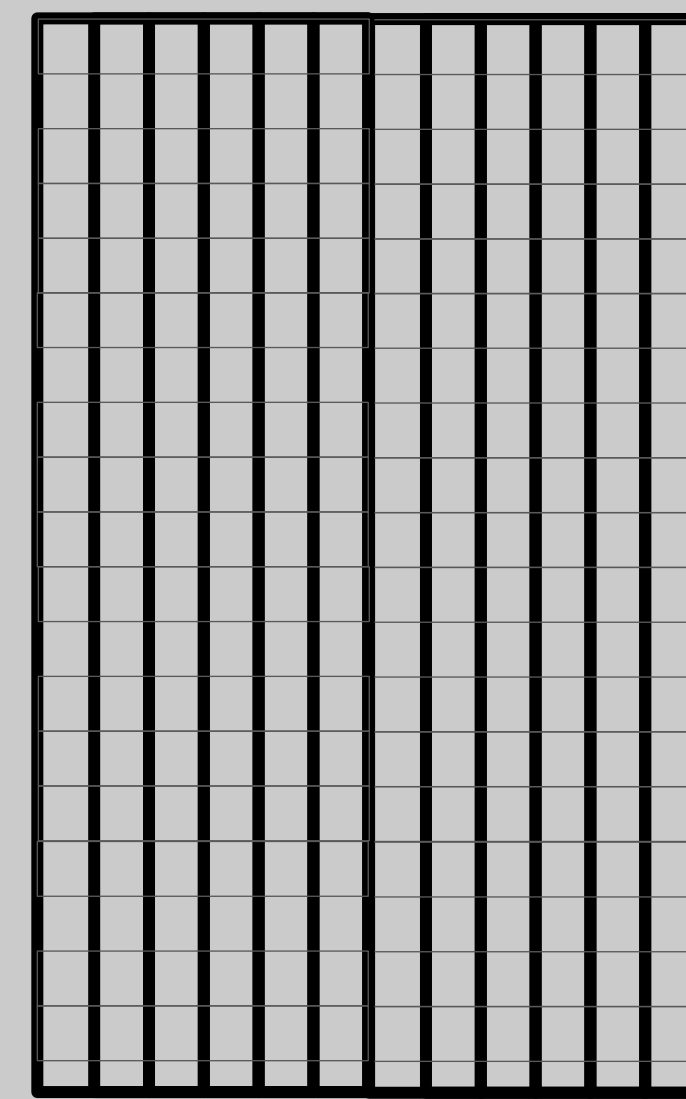


Not low-rank!  
One strong attribute — the rest are similar

# Data Science

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

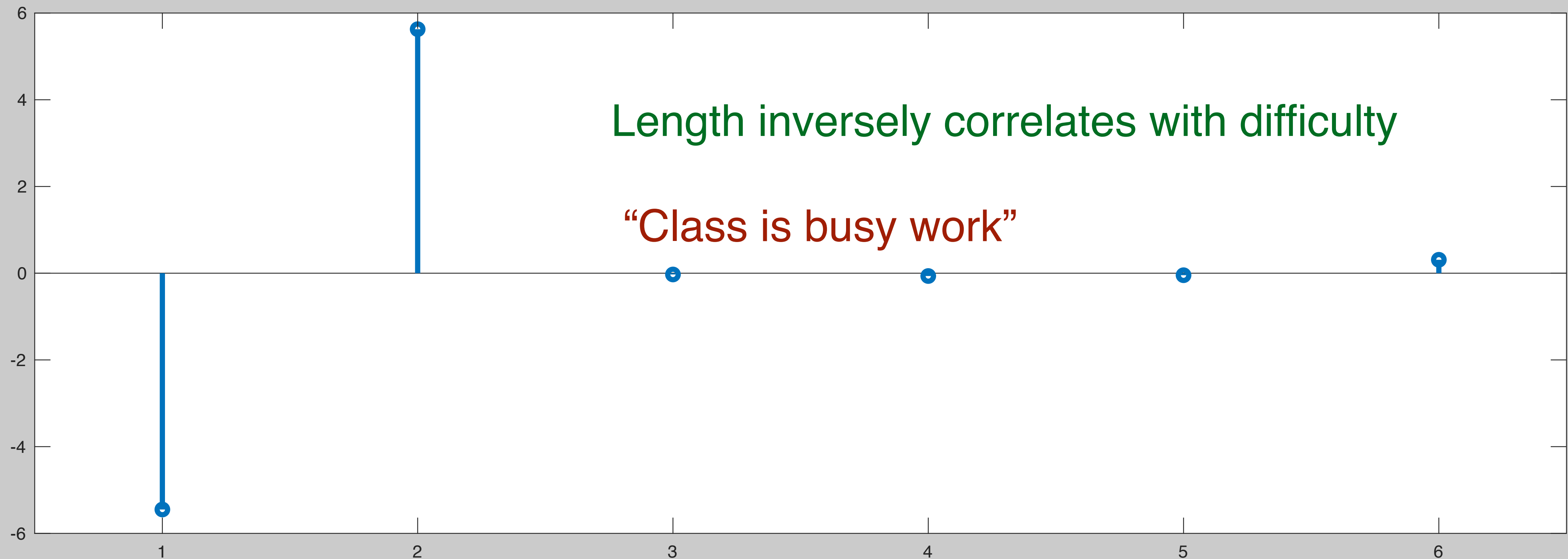
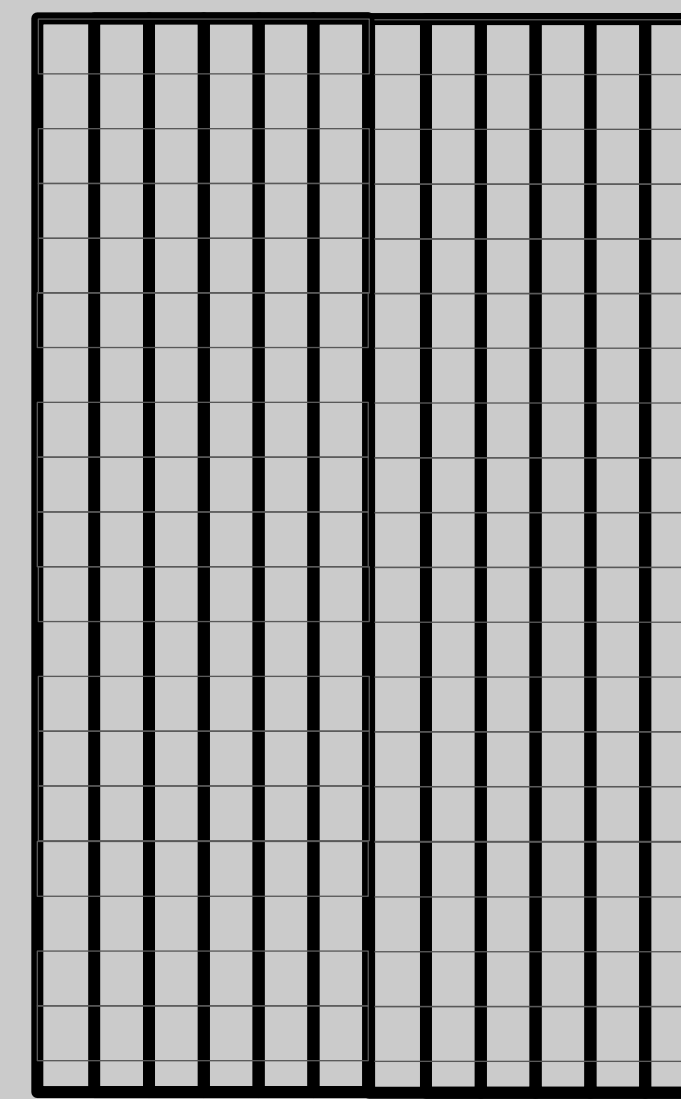
- 1) HW difficulty
- 2) HW Length
- 3) Lab hour/week
- 4) Class rating
- 5) Expectation Rating
- 6) Comfortable attending OH



# Data Science

$$A^T \vec{u}_2$$

- 1) HW difficulty
- 2) HW Length
- 3) Lab hour/week
- 4) Class rating
- 5) Expectation Rating
- 6) Comfortable attending OH

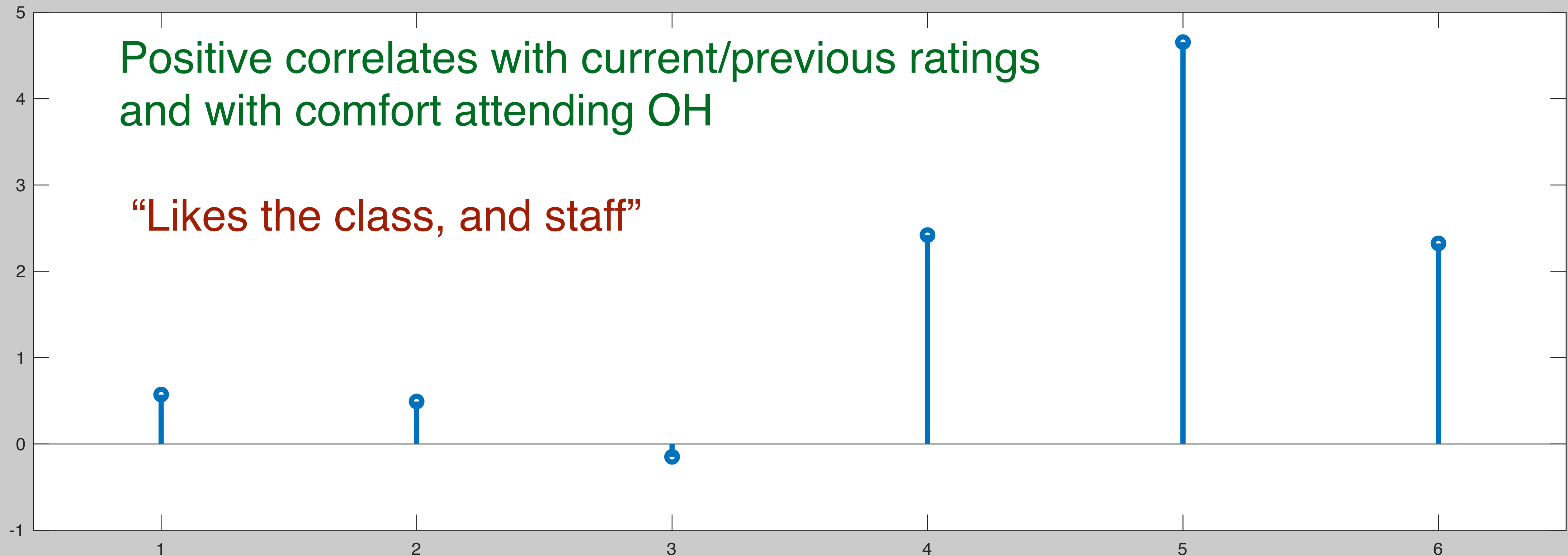
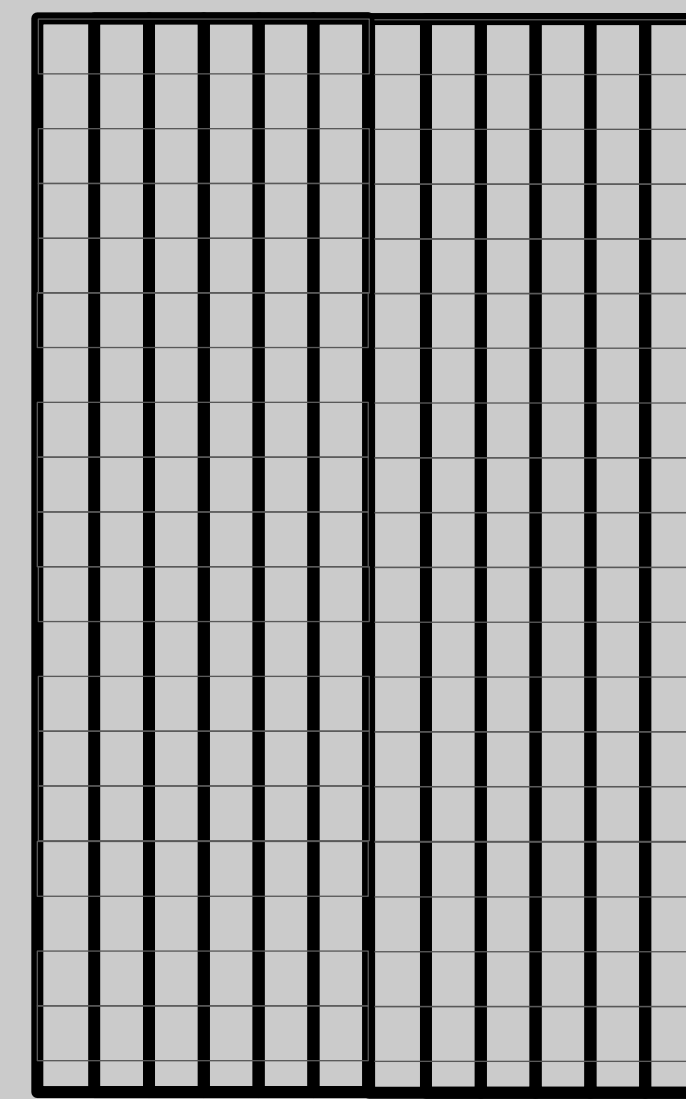




# Data Science

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

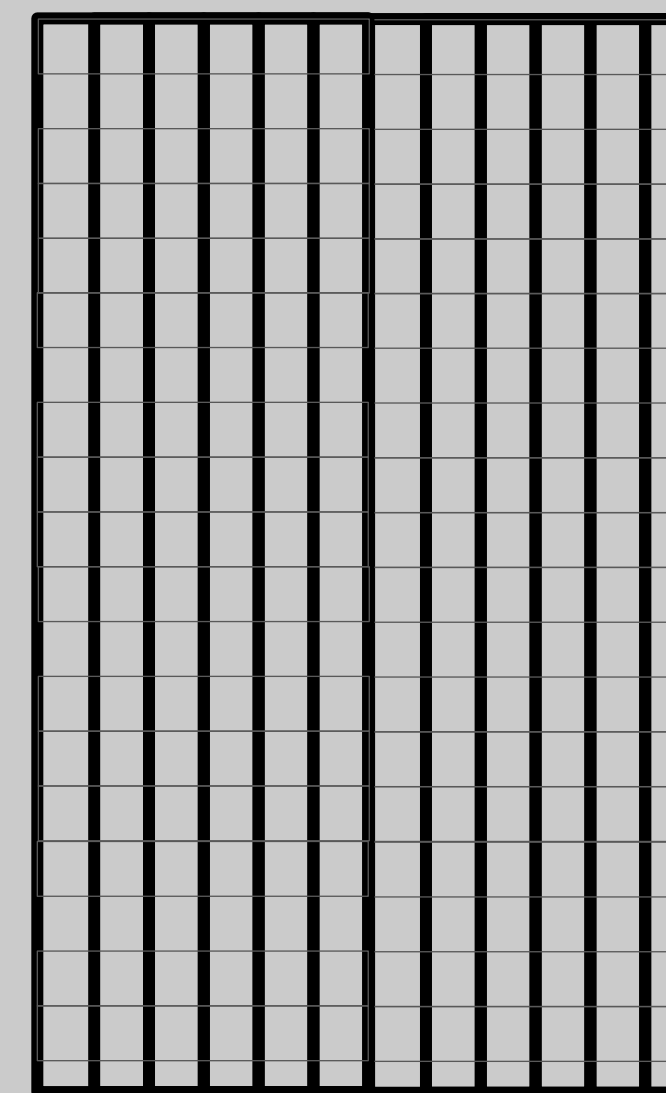
- 1) HW difficulty
- 2) HW Length
- 3) Lab hour/week
- 4) Class rating
- 5) Expectation Rating
- 6) Comfortable attending OH



# Data Science

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

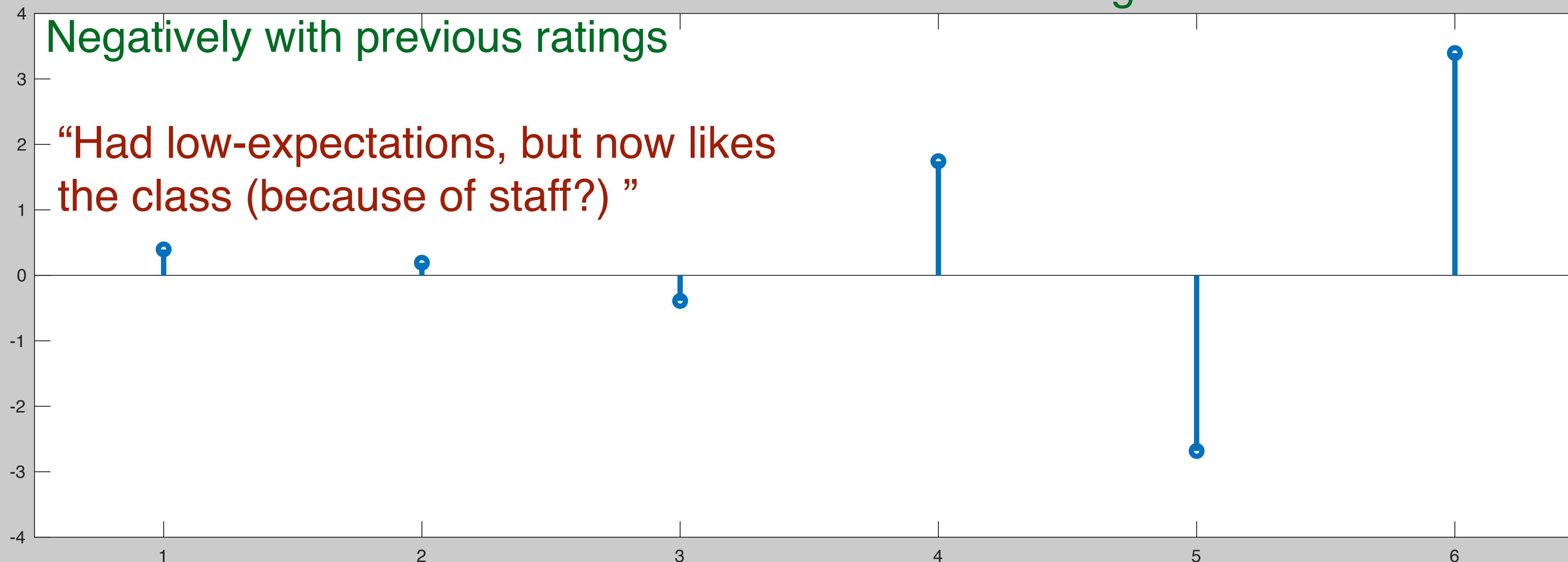
- 1) HW difficulty
- 2) HW Length
- 3) Lab hour/week
- 4) Class rating
- 5) Expectation Rating
- 6) Comfortable attending OH



Positive correlates with current & with comfort attending OH

Negatively with previous ratings

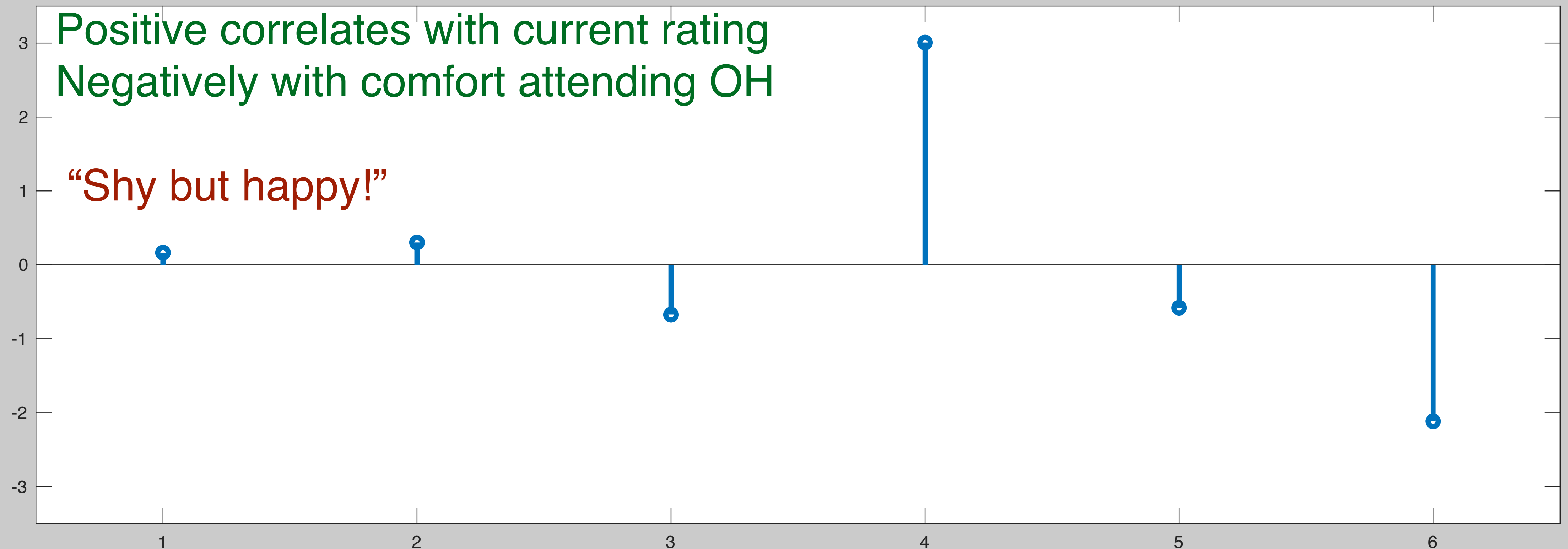
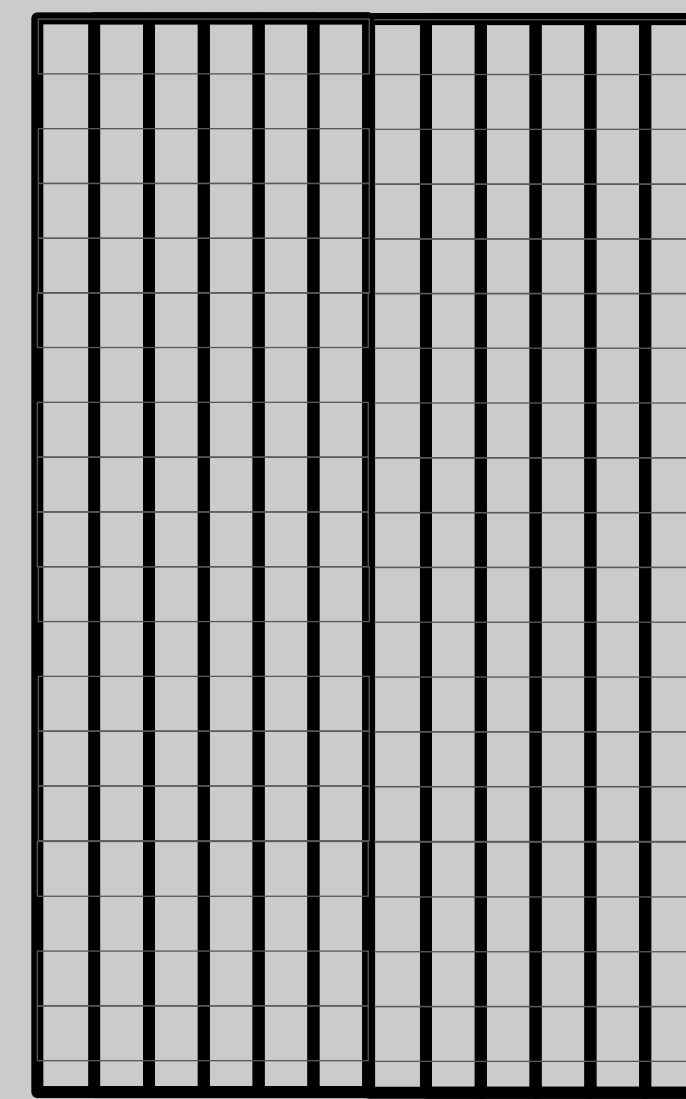
“Had low-expectations, but now likes the class (because of staff?) ”



# Data Science

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

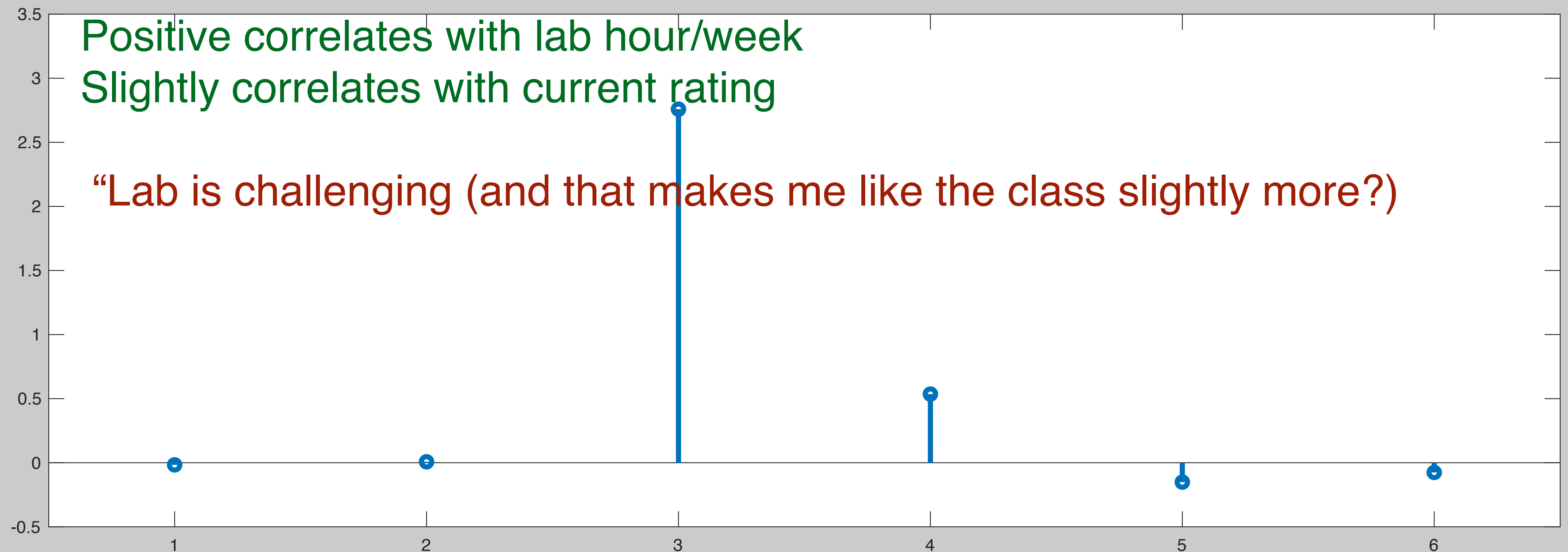
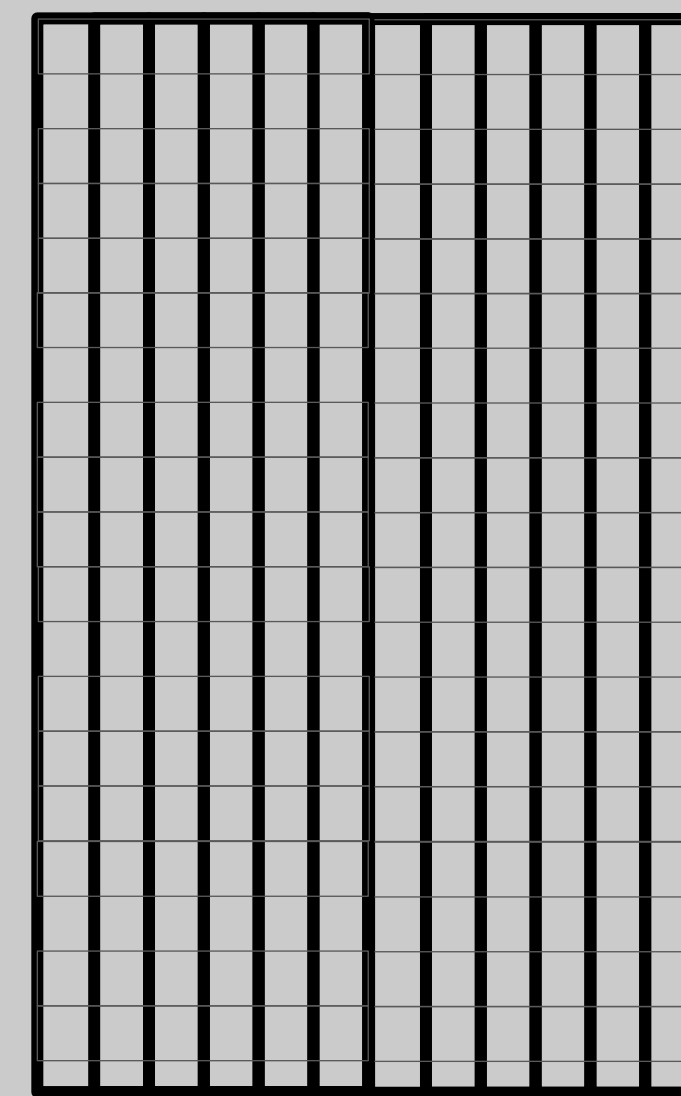
- 1) HW difficulty
- 2) HW Length
- 3) Lab hour/week
- 4) Class rating
- 5) Expectation Rating
- 6) Comfortable attending OH



# 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

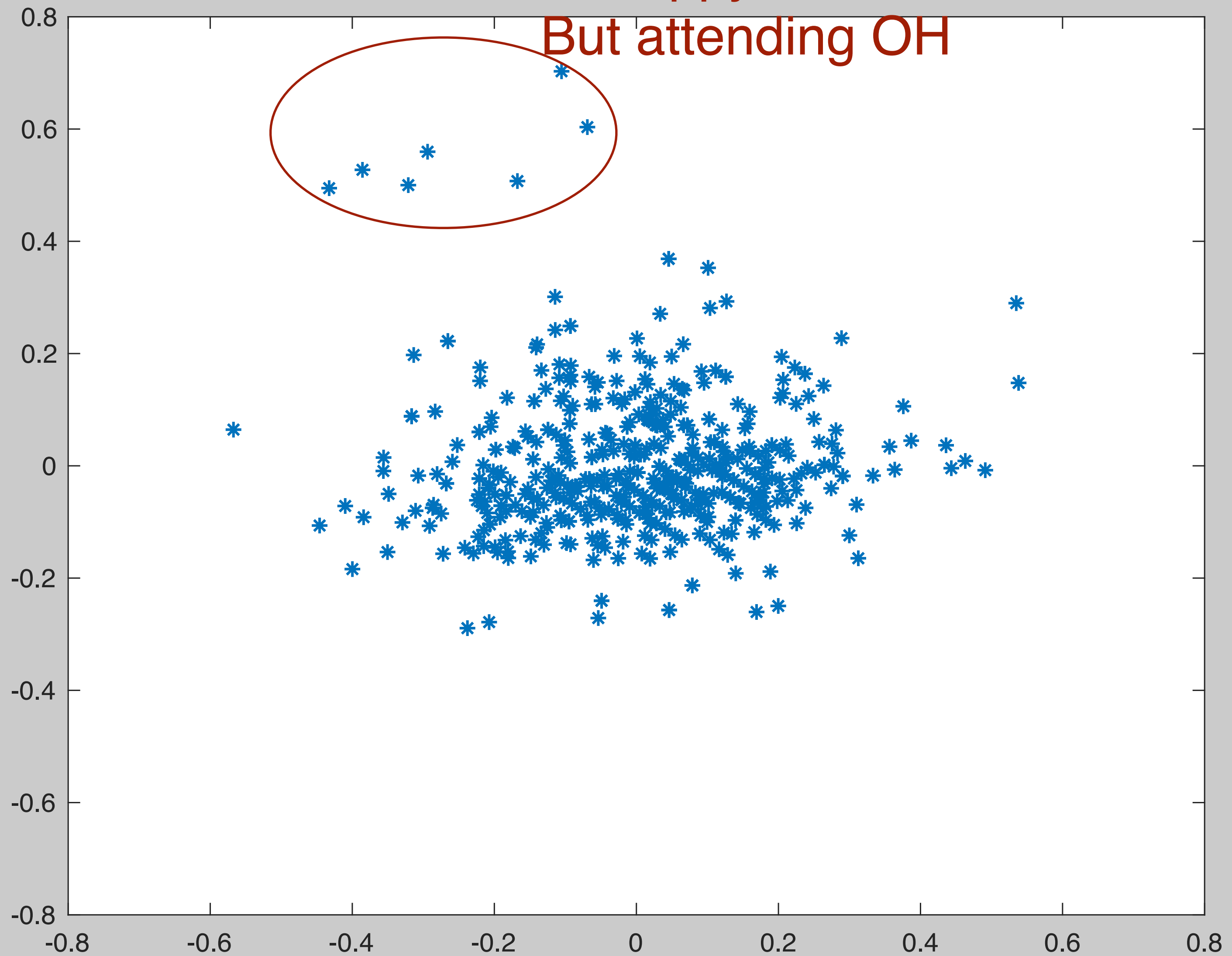




# Data Science

Lab is challenging

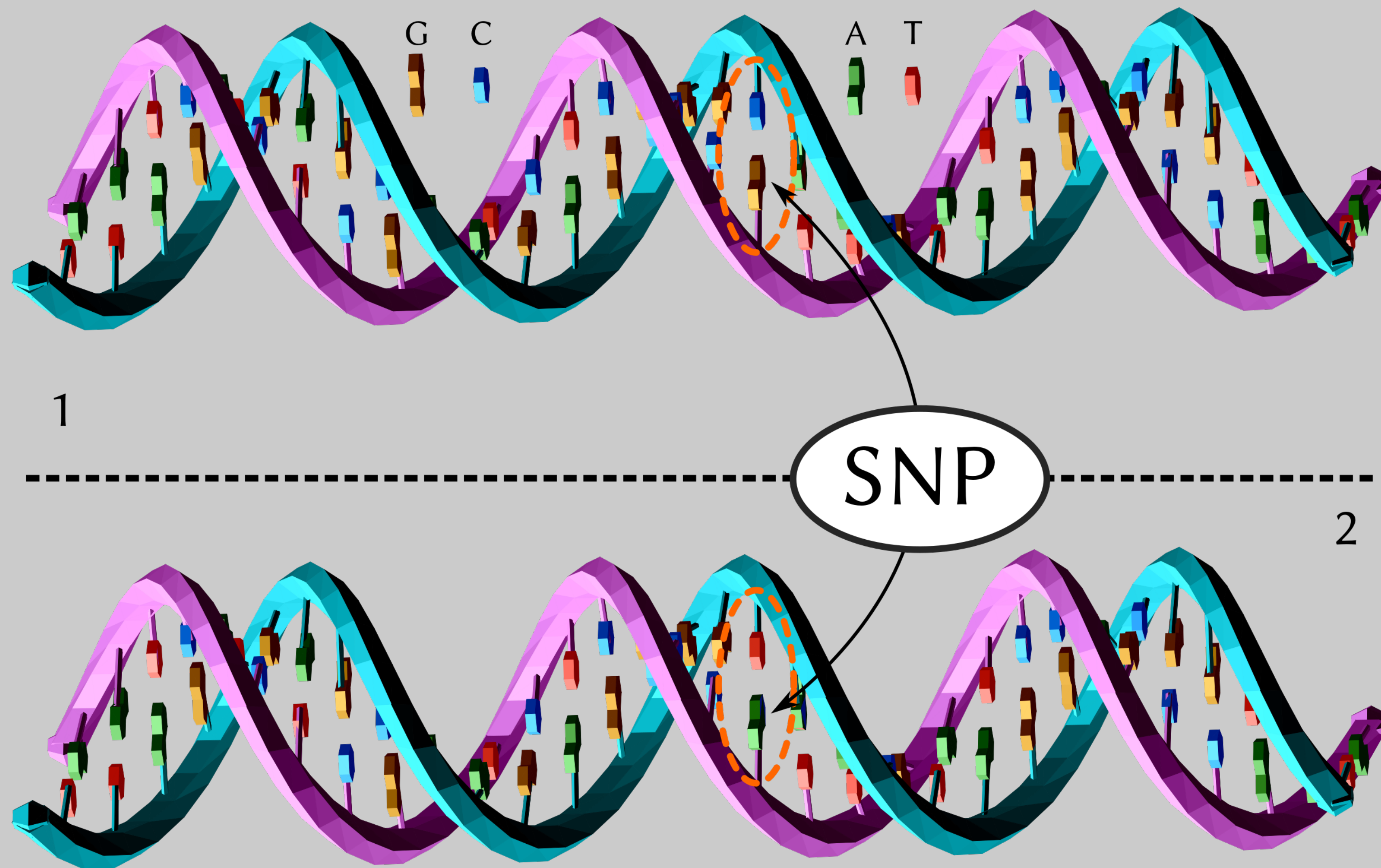
PC2  
 $A\vec{v}_6$



Lab is very challenging  
Unhappy  
But attending OH

$A\vec{v}_5$  PC5 happy, no OH

# Genetic Data



# PCA in Genetics Reveals Geography

Genes mirror geography within Europe  
*Nature* **456**, 98-101 (6 November 2008)

Study:

Characterized genetic variations in 3,000 Europeans from 36 Countries

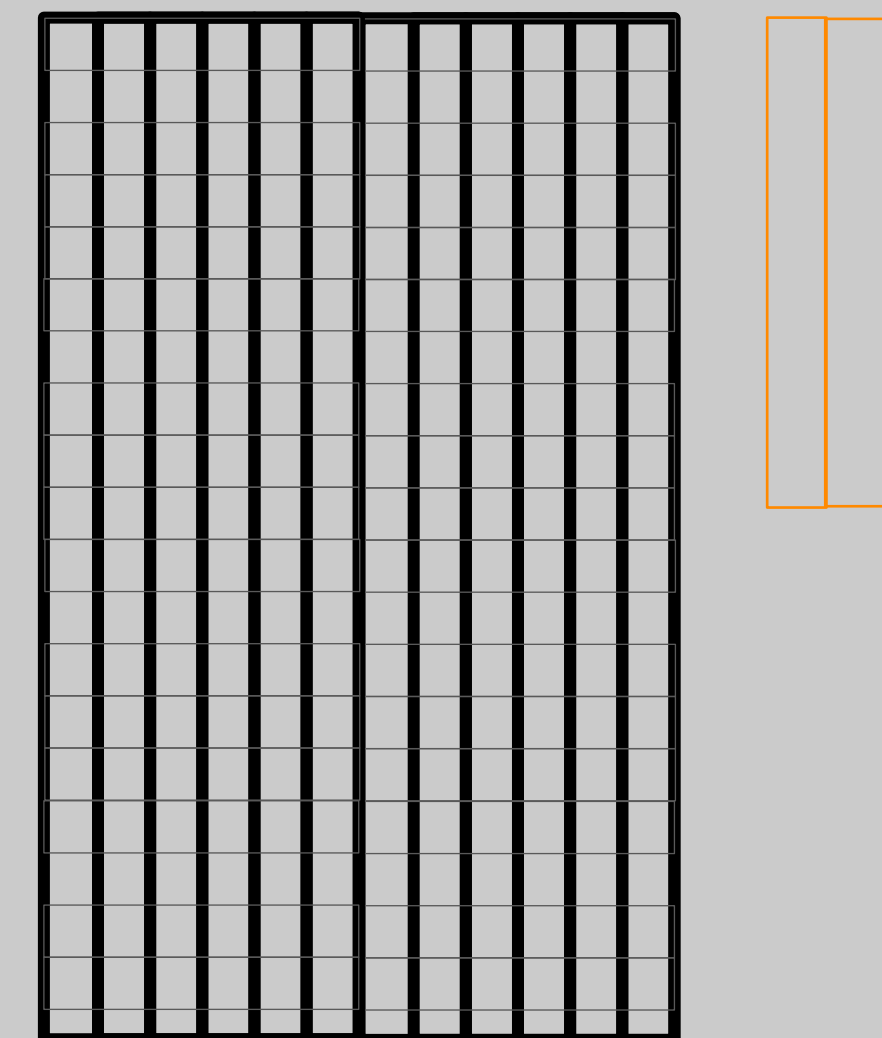
Built a matrix of 200K SNPs (single nucleotide polymorphisms)

Computed largest 2 principle components

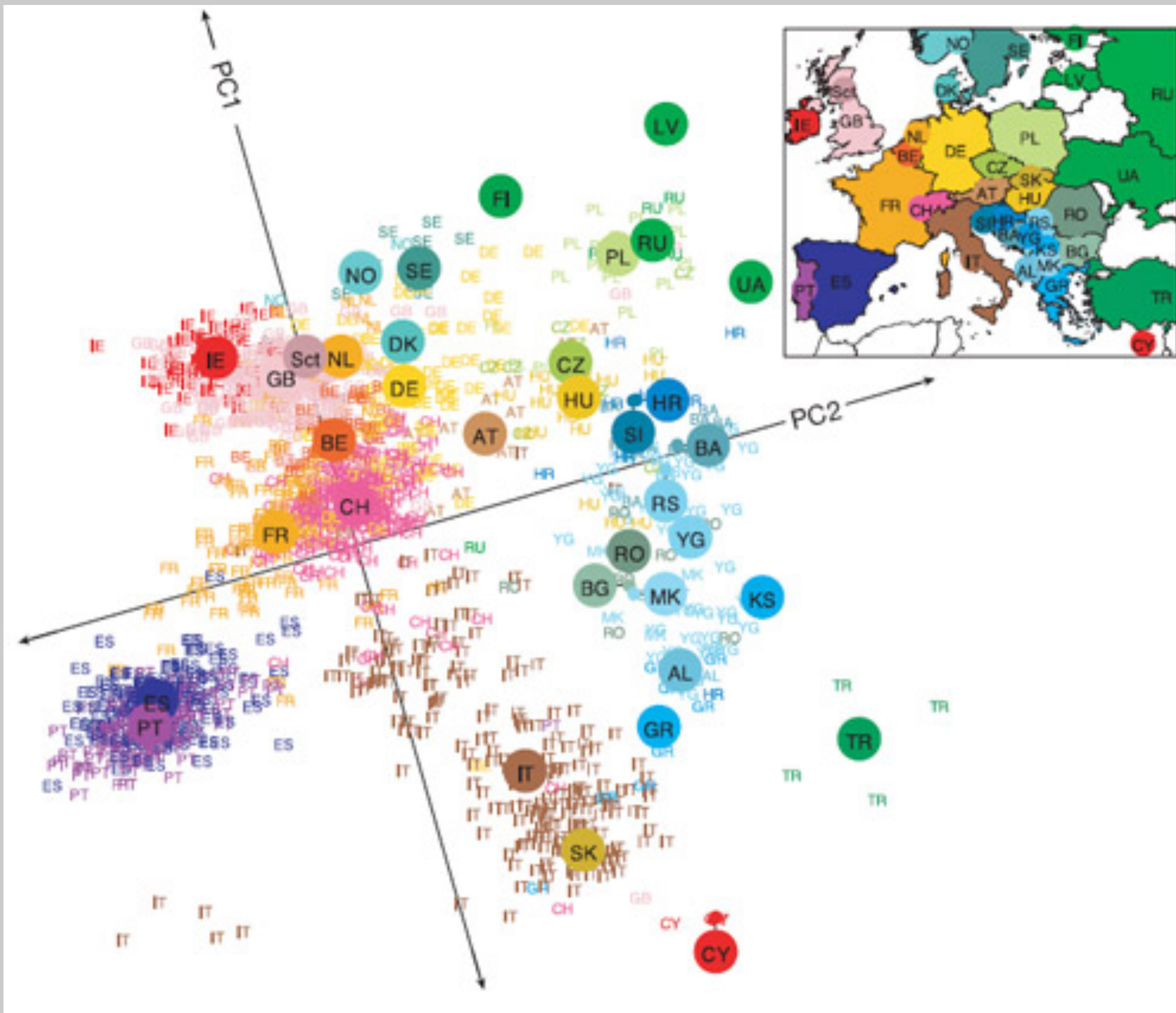
Projected subjects on 2 dimensional data

Overlayed the result on the map of Europe

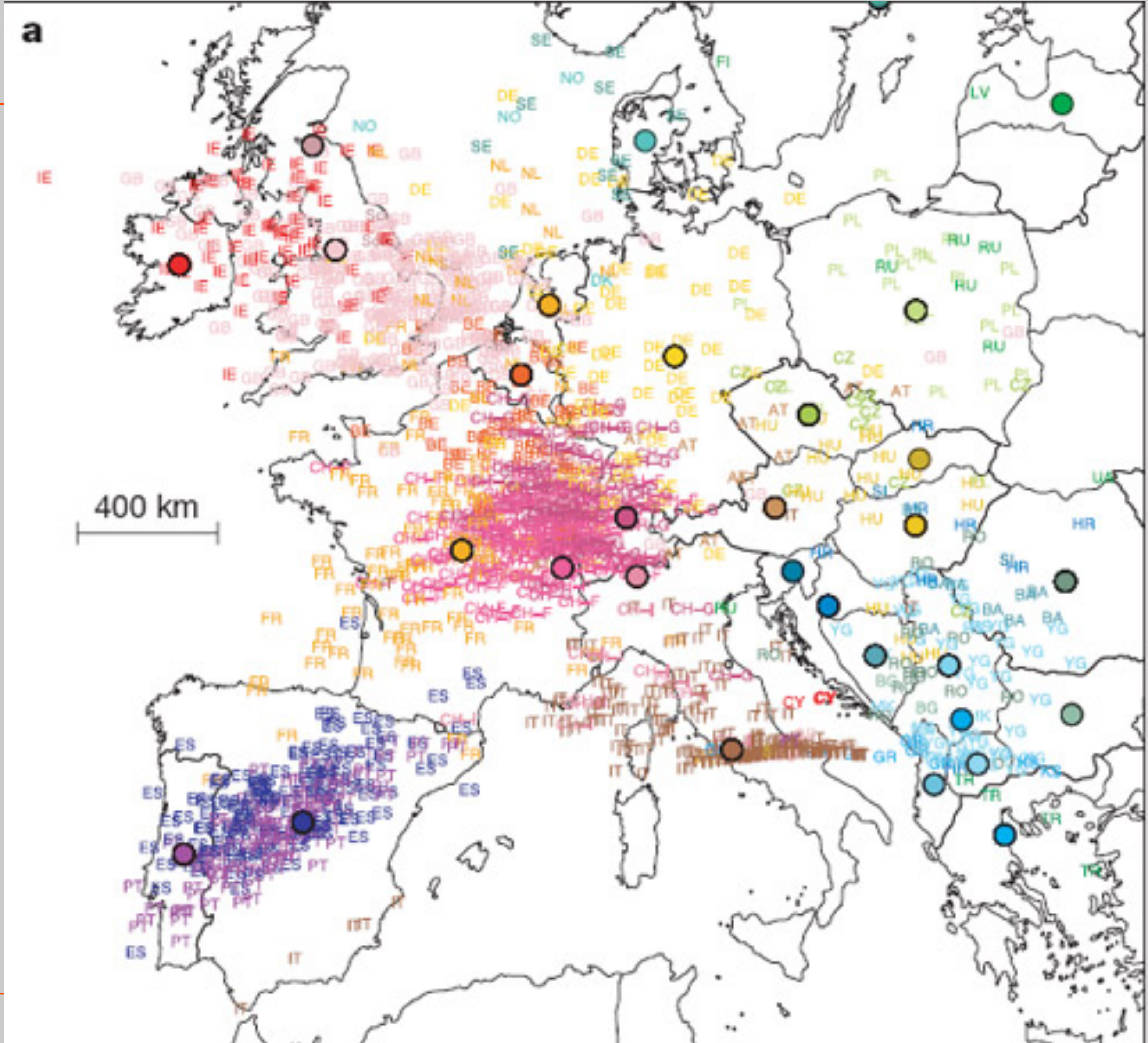
$$A\vec{v}_1 \quad A\vec{v}_2$$







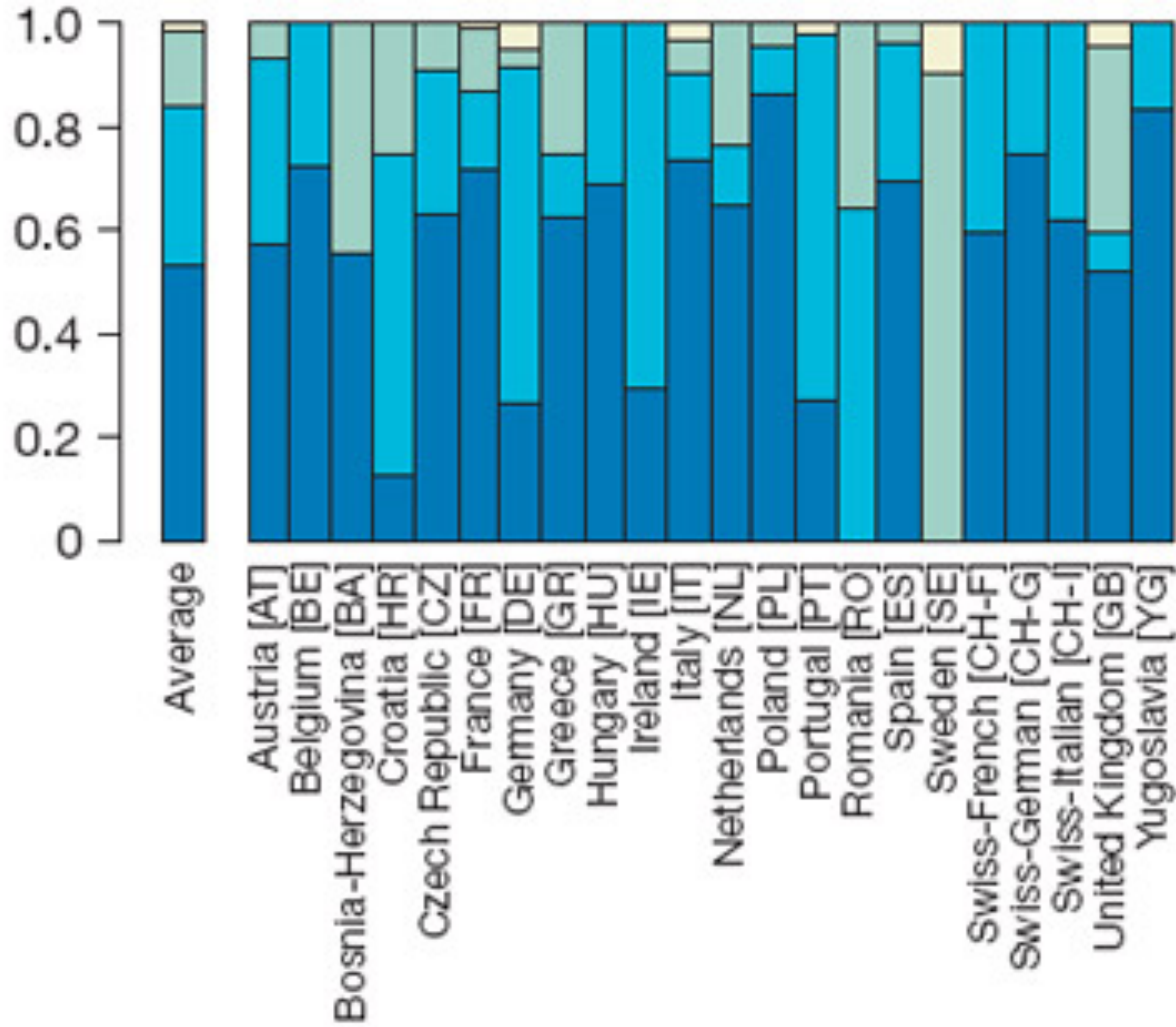






## Prediction accuracy

- 1,200–2,500 km
- 800–1,200 km
- 400–800 km
- 0–400 km



# 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/>



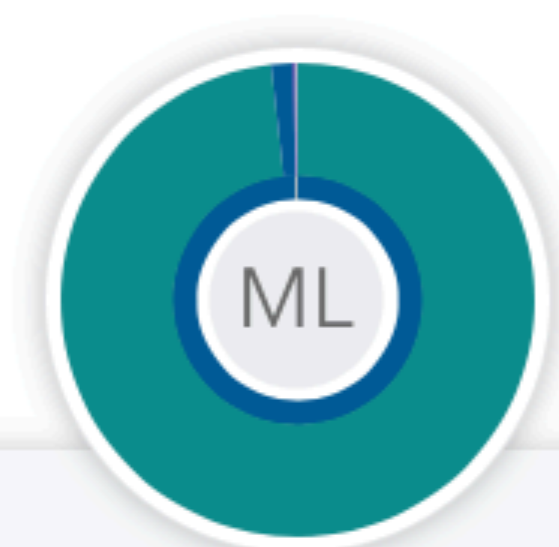
## Ancestry Composition

Your DNA tells the story of who you are and how you're connected to populations around the world. Trace your heritage through the centuries and uncover clues about where your ancestors lived and when.

[Summary](#)

[Scientific Details](#)

[Frequently Asked Questions](#)



<b>Michael Lustig</b>	<b>100%</b>
<b>European</b>	<b>99.7%</b>
● Ashkenazi Jewish	98.2% >
● Broadly European	1.5%
Trace Ancestry	0.2% v
Unassigned	0.1% v

[See all tested populations](#)



Updated: August 17, 2020 ⓘ

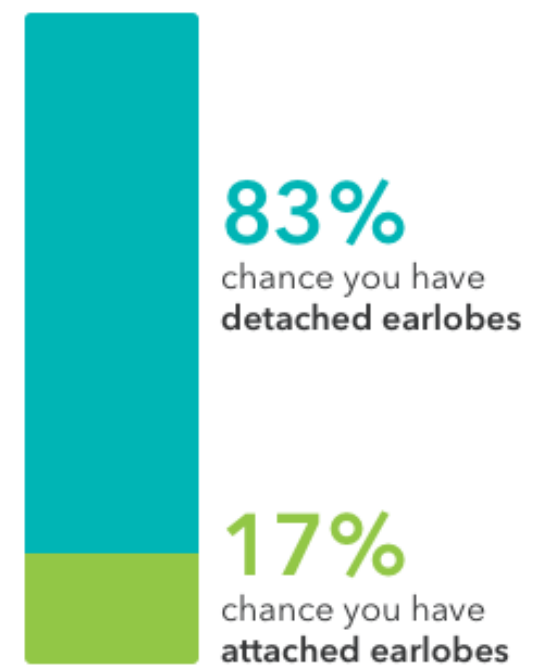


# Physical features



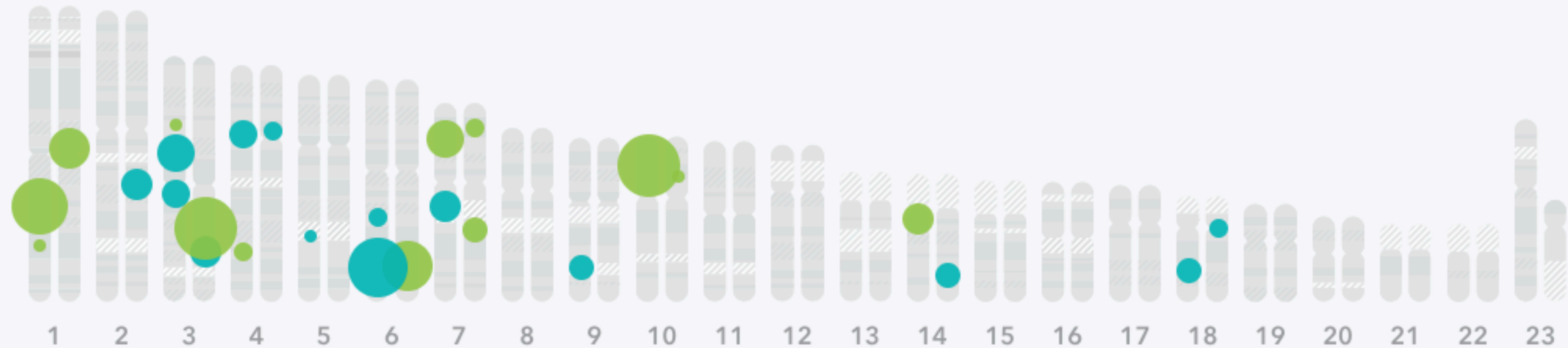
## Earlobe Type

Michael, your genetics predict



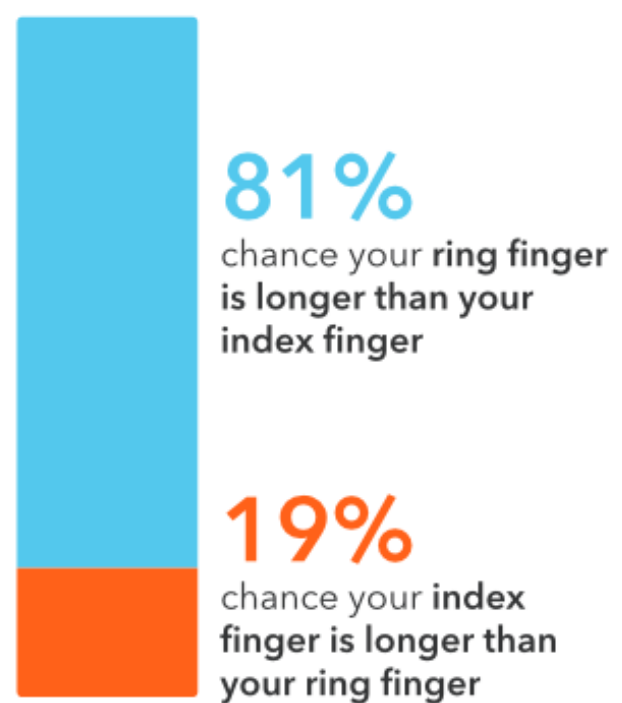
## Breakdown of your genetics

The bigger the circle, the stronger the effect your variants have on your overall chances.



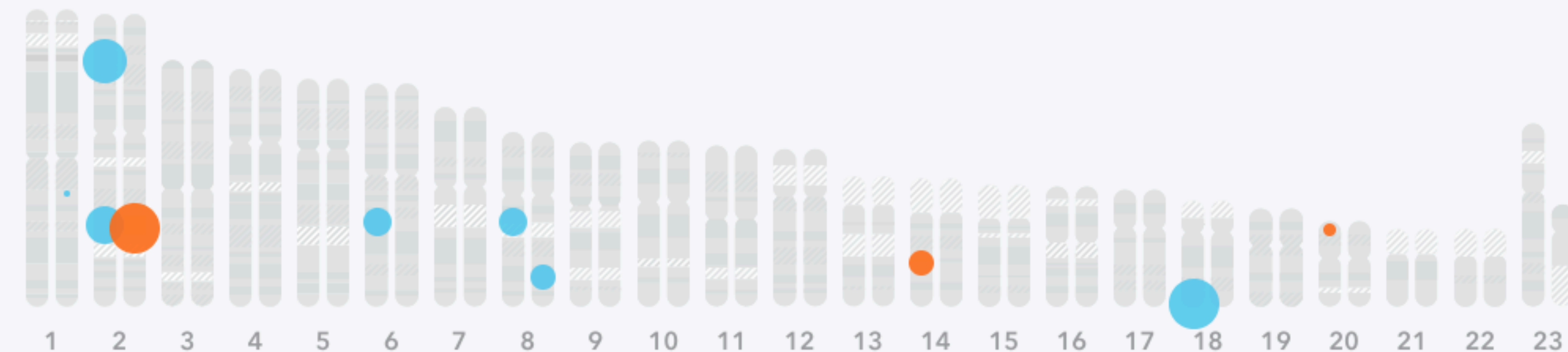
## Finger Length Ratio

Michael, your genetics predict



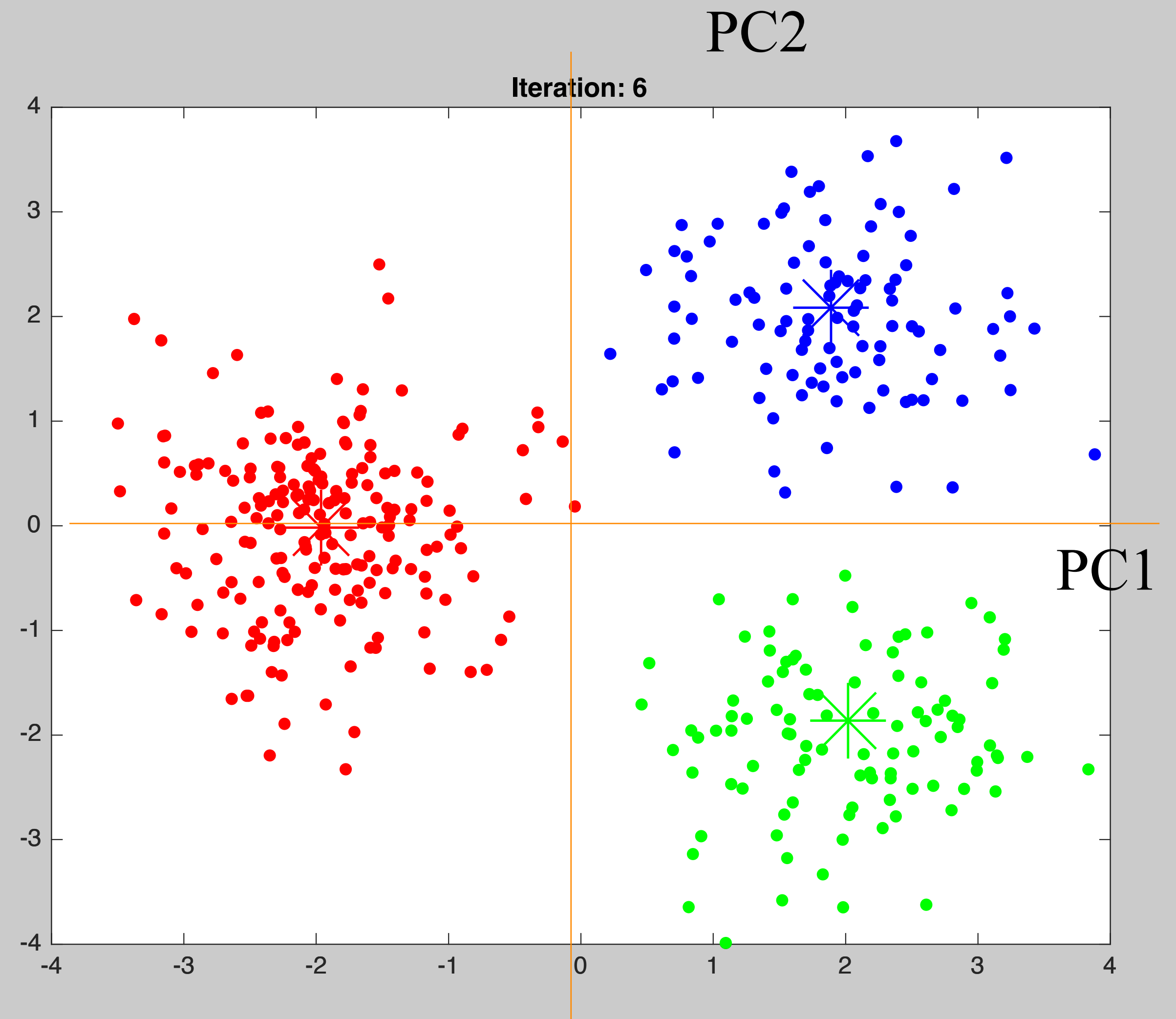
## Breakdown of your genetics

The bigger the circle, the stronger the effect your variants have on your overall chances.



# Labeled VS non labeled Classification

Word1  
Word2  
Word3  
Word4  
Word5  
Word6  
Word7  
Word8



# Labeled VS non labeled Classification

Word1

Word2

Word3

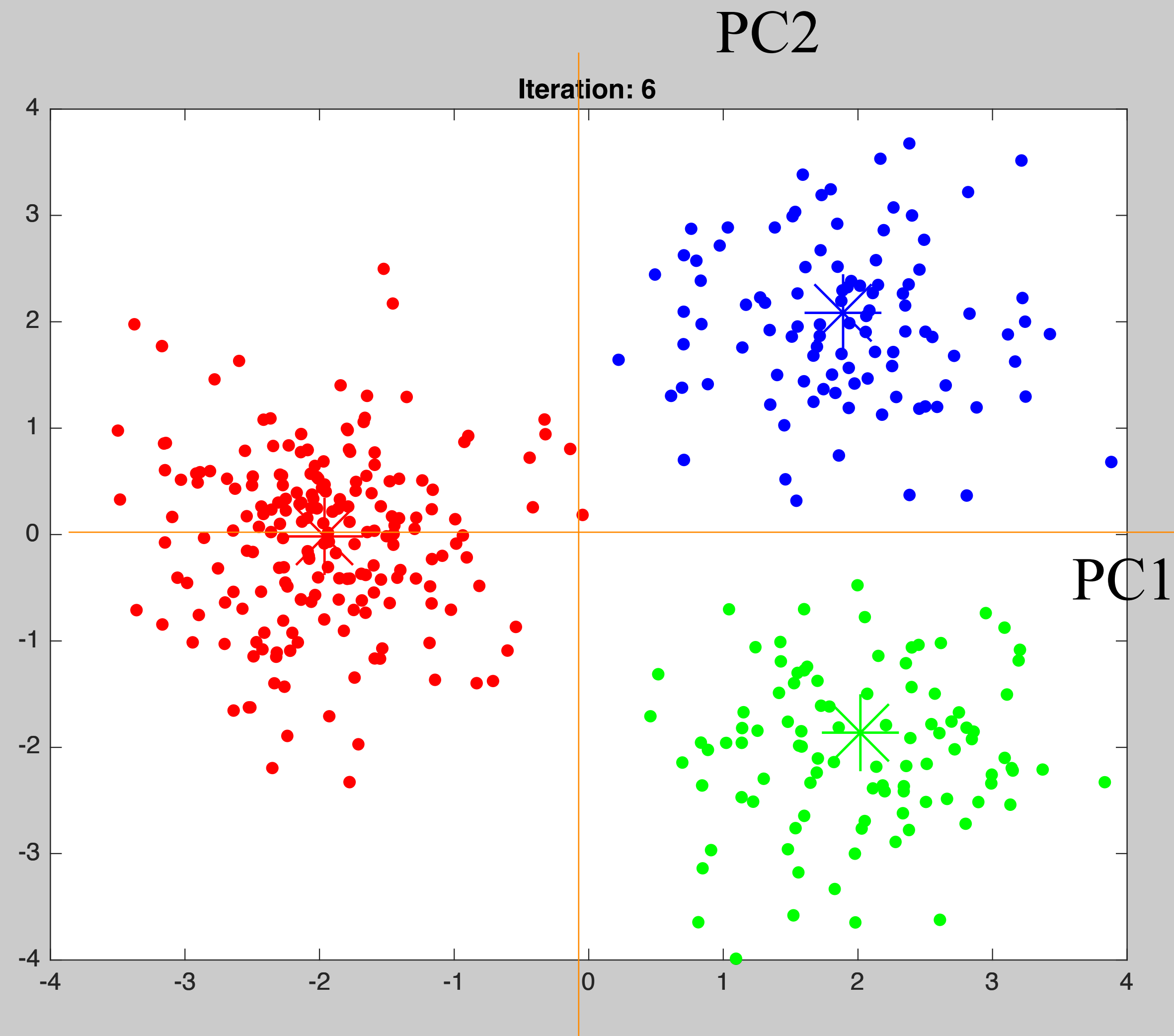
Word4

Word5

Word6

Word7

Word8



# Labeled VS non labeled Classification

---

“Banana”

“Banana”

“Banana”

“Mango”

“Mango”

“Mango”

“Chop”

“Chop”

“Chop”

PC2

PC1