

# On the Difficulties of Subspace Clustering with Missing Data

Daniel L. Pimentel-Alarcón

1<sup>st</sup> Annual Workshop on Data Sciences,  
April 17<sup>th</sup>, 2015

Joint work with Nigel Boston and Robert Nowak

# Outline

- ▶ Introduction
- ▶ What changes with missing data?
- ▶ Subspace Identifiability Problem
- ▶ Setup
- ▶ The Answer
- ▶ Application
- ▶ Conclusions



# Outline

- ▶ Introduction
- ▶ What changes with missing data?
- ▶ Subspace Identifiability Problem
- ▶ Setup
- ▶ The Answer
- ▶ Application
- ▶ Conclusions

# Introduction

We have lots of data



# Introduction

We have lots of data



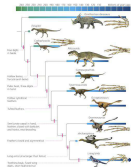
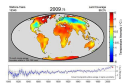
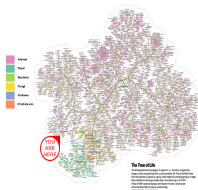
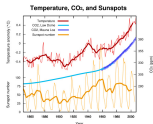
# Introduction

We have lots of data



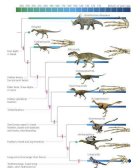
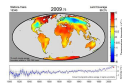
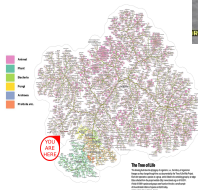
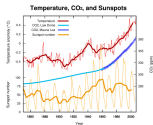
# Introduction

We have lots of data



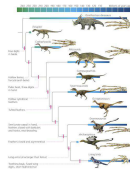
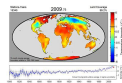
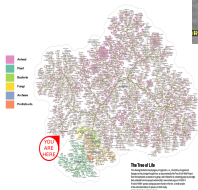
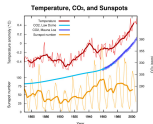
# Introduction

We have lots of data



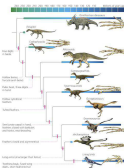
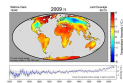
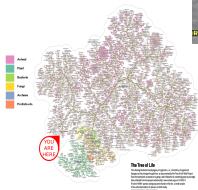
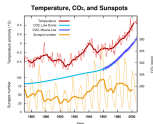
# Introduction

We have lots of data



# Introduction

We have lots of data





# Introduction

We have lots of data



# Introduction

We have lots of data



And we want to analyze it.

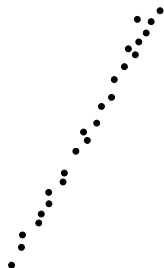
# Introduction

Linear Algebra is one of our favorite tools.

# Introduction

Linear Algebra is one of our favorite tools.

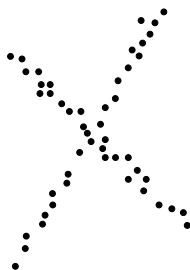
- ▶ Because data is often well-modeled by linear structures.



$$\begin{bmatrix} 1 & 2 & 1 & 3 & 2 & 1 & 3 & 1 & 2 & 2 \\ 2 & 4 & 2 & 6 & 4 & 2 & 6 & 2 & 4 & 4 \\ 3 & 6 & 3 & 9 & 6 & 3 & 9 & 3 & 6 & 6 \\ 1 & 2 & 1 & 3 & 2 & 1 & 3 & 1 & 2 & 2 \\ 2 & 4 & 2 & 6 & 4 & 2 & 6 & 2 & 4 & 4 \\ 3 & 6 & 3 & 9 & 6 & 3 & 9 & 3 & 6 & 6 \end{bmatrix}$$

# Introduction

Sometimes one subspace is not enough.

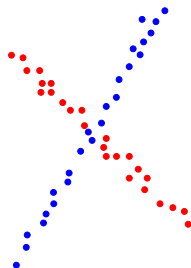


$$\begin{bmatrix} 1 & 4 & 1 & 3 & 3 & 1 & 2 & 1 & 2 & 1 \\ 2 & 4 & 2 & 6 & 3 & 2 & 2 & 2 & 4 & 1 \\ 3 & 4 & 3 & 9 & 3 & 3 & 2 & 3 & 6 & 1 \\ 1 & 8 & 1 & 3 & 6 & 1 & 4 & 1 & 2 & 2 \\ 2 & 8 & 2 & 6 & 6 & 2 & 4 & 2 & 4 & 2 \\ 3 & 8 & 3 & 9 & 6 & 3 & 4 & 3 & 6 & 2 \end{bmatrix}$$

# Introduction

Sometimes one subspace is not enough.

- Enters Subspace Clustering.

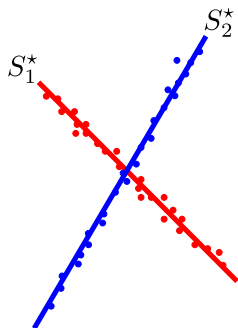


1	4	1	3	3	1	2	1	2	1
2	4	2	6	3	2	2	2	4	1
3	4	3	9	3	3	2	3	6	1
1	8	1	3	6	1	4	1	2	2
2	8	2	6	6	2	4	2	4	2
3	8	3	9	6	3	4	3	6	2

# Introduction

Sometimes one subspace is not enough.

- Enters Subspace Clustering.



1	4	1	3	3	1	2	1	2	1
2	4	2	6	3	2	2	2	4	1
3	4	3	9	3	3	2	3	6	1
1	8	1	3	6	1	4	1	2	2
2	8	2	6	6	2	4	2	4	2
3	8	3	9	6	3	4	3	6	2

# Introduction

That's all very nice, but... often data is missing!



# Introduction

That's all very nice, but... often data is missing!

- Example: Vision.

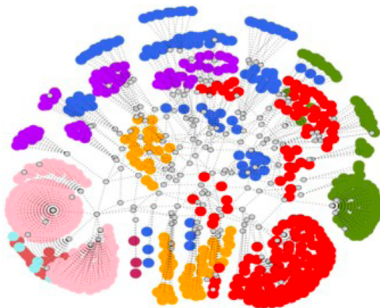


Image: Hopkins 155 Dataset

# Introduction

Often data is missing!

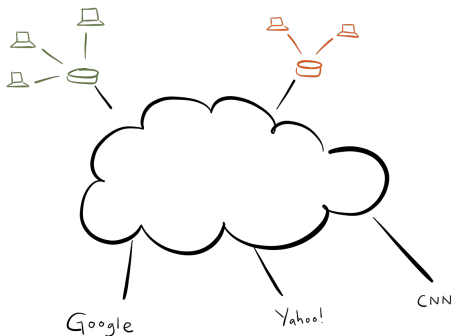
- Other example: Network topology estimation



# Introduction

Often data is missing!

- ▶ Other example: Network topology estimation



# Introduction

Often data is missing!

- ▶ Other example: Network topology estimation

$$\text{monitors} \left\{ \begin{bmatrix} 1 & \cdot & \cdot & 3 & \cdot & 3 & \cdot & 1 & 2 & \cdot \\ 2 & \cdot & 2 & \cdot & \cdot & 6 & \cdot & \cdot & 4 & \cdot \\ \cdot & \cdot & 3 & \cdot & \cdot & 9 & \cdot & 3 & 6 & \cdot \\ 1 & \cdot & 1 & 3 & 6 & \cdot & 4 & 1 & 2 & 2 \\ \cdot & 8 & \cdot & \cdot & 6 & \cdot & 4 & \cdot & \cdot & \cdot \\ \cdot & 8 & \cdot & \cdot & \cdot & \cdot & 4 & \cdot & \cdot & 2 \end{bmatrix} \right.$$

IP's

# Introduction

Often data is missing!

- ▶ Other example: Network topology estimation

$$\text{monitors} \left\{ \begin{bmatrix} 1 & \cdot & \cdot & 3 & \cdot & 3 & \cdot & 1 & 2 & \cdot \\ 2 & \cdot & 2 & \cdot & \cdot & 6 & \cdot & \cdot & 4 & \cdot \\ \cdot & \cdot & 3 & \cdot & \cdot & 9 & \cdot & 3 & 6 & \cdot \\ 1 & \cdot & 1 & 3 & 6 & \cdot & 4 & 1 & 2 & 2 \\ \cdot & 8 & \cdot & \cdot & 6 & \cdot & 4 & \cdot & \cdot & \cdot \\ \cdot & 8 & \cdot & \cdot & \cdot & \cdot & 4 & \cdot & \cdot & 2 \end{bmatrix} \right.$$

IP's

- ▶ We still want to analyze these datasets.

# Introduction

Subspace Clustering with Missing Data poses a new problem:

# Introduction

Subspace Clustering with Missing Data poses a new problem:

- ▶ **False** subspaces.

# Introduction

Subspace Clustering with Missing Data poses a new problem:

- **False** subspaces.

$$\begin{bmatrix} 1 & \cdot & \cdot & 3 & \cdot & 3 & \cdot & 1 & 2 & \cdot \\ 2 & \cdot & 2 & \cdot & \cdot & 6 & \cdot & \cdot & 4 & \cdot \\ \cdot & \cdot & 3 & \cdot & \cdot & 9 & \cdot & 3 & 6 & \cdot \\ 1 & \cdot & 1 & 3 & 6 & \cdot & 4 & 1 & 2 & 2 \\ \cdot & 8 & \cdot & \cdot & 6 & \cdot & 4 & \cdot & \cdot & \cdot \\ \cdot & 8 & \cdot & \cdot & \cdot & \cdot & 4 & \cdot & \cdot & 2 \end{bmatrix} \subset \text{span} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$



# Introduction

- ▶ We want to know how to identify **false** subspaces!

# Introduction

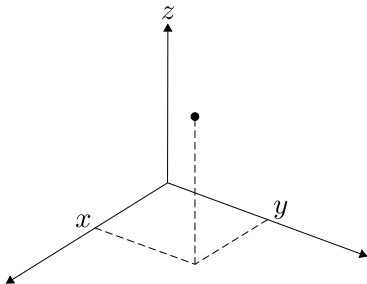
- ▶ We want to know how to identify **false** subspaces!
- ▶ We need to understand how things change when data is missing.

# Outline

- ▶ Introduction ✓
- ▶ What changes with missing data?
- ▶ Subspace Identifiability Problem
- ▶ Setup
- ▶ The Answer
- ▶ Application
- ▶ Conclusions

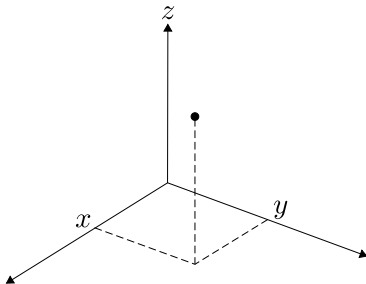
# What changes with missing data?

Say I give you one datapoint.



# What changes with missing data?

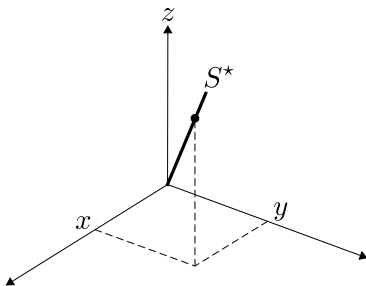
Say I give you one datapoint.



And I tell you it lies in a 1-dimensional subspace  $S^*$ .

## What changes with missing data?

Then you can uniquely identify  $S^*$ .



# What changes with missing data?

But what if data is missing?

# What changes with missing data?

But what if data is missing?

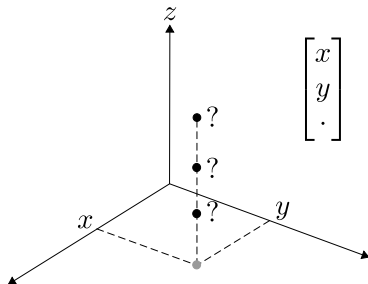
- ▶ Say I give you a point *without* the  $z$  coordinate.



# What changes with missing data?

But what if data is missing?

- Say I give you a point *without* the  $z$  coordinate.

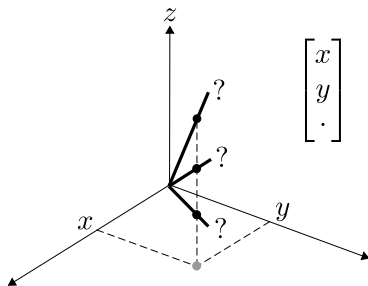


## What changes with missing data?

Then we cannot uniquely identify  $S^*$ .

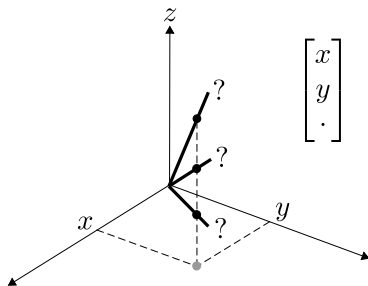
## What changes with missing data?

Then we cannot uniquely identify  $S^*$ .



## What changes with missing data?

Then we cannot uniquely identify  $S^*$ .



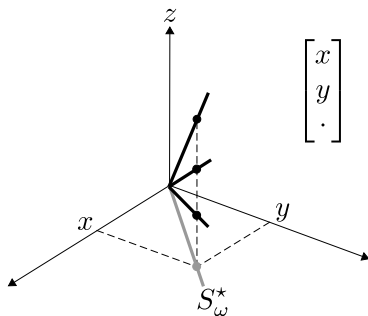
There are infinitely many *false* subspaces.

## What changes with missing data?

Nevertheless, all those *false* subspaces must satisfy one very important condition!

## What changes with missing data?

Nevertheless, all those *false* subspaces must satisfy one very important condition!



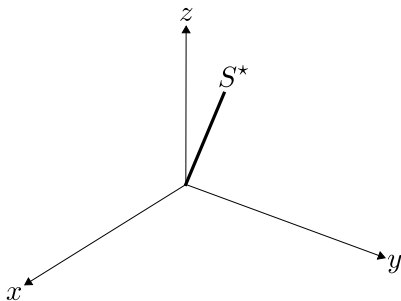
They must have the same canonical projection as  $S^*$ .

# Outline

- ▶ Introduction ✓
- ▶ What changes with missing data? ✓
- ▶ Subspace Identifiability Problem
- ▶ Setup
- ▶ The Answer
- ▶ Application
- ▶ Conclusions

# Subspace Identifiability Problem

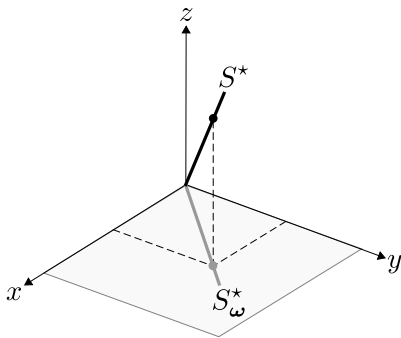
$S^* := r$ -dimensional subspace of  $\mathbb{R}^d$ ,  $r < d$ .





# Subspace Identifiability Problem

$S_{\omega}^* :=$  Projection of  $S^*$  onto a canonical subspace.

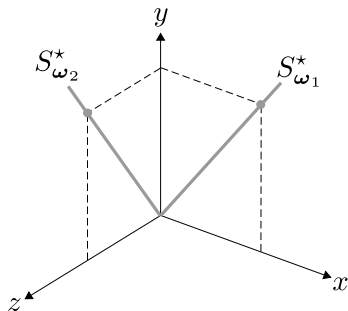


# Subspace Identifiability Problem

Suppose I don't tell you  $S^*$ ...

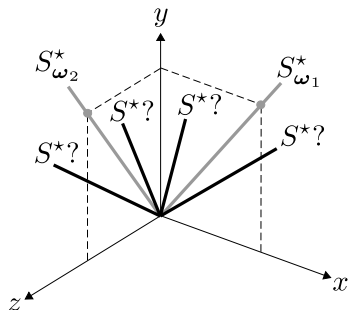
# Subspace Identifiability Problem

Suppose I don't tell you  $S^*$ ... but I give you a set of projections of  $S^*$  onto some canonical subspaces.



# Subspace Identifiability Problem

Suppose I don't tell you  $S^*$ ...but I give you a set of projections of  $S^*$  onto some canonical subspaces.

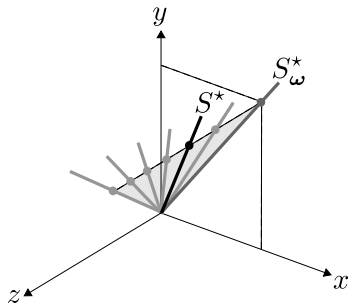


Can you uniquely determine  $S^*$  from this set of projections?

# Subspace Identifiability Problem

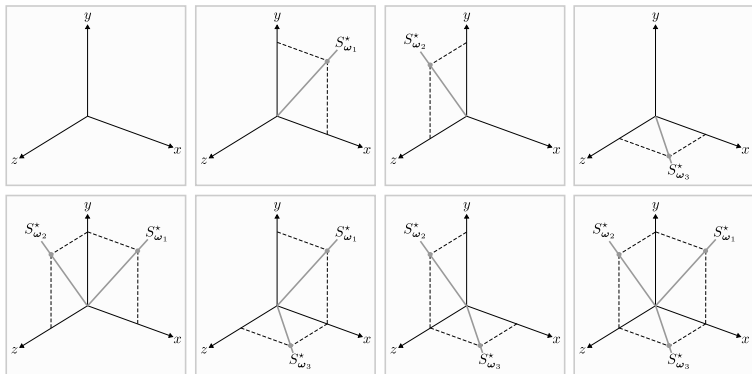
Is this even possible?

- There might be many subspaces that agree with the projections.



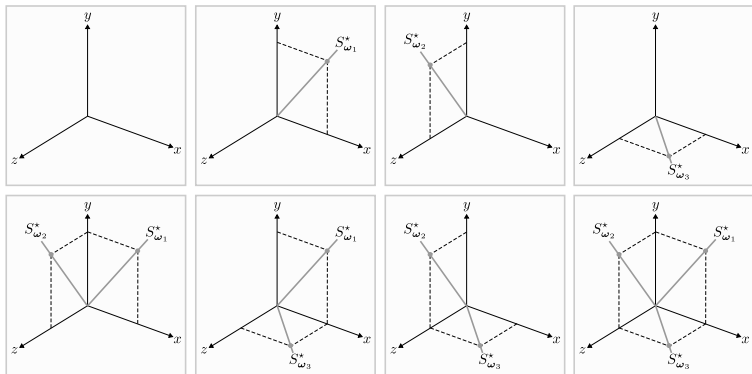
# Subspace Identifiability Problem

Well... it depends on which set of projections I give you.



# Subspace Identifiability Problem

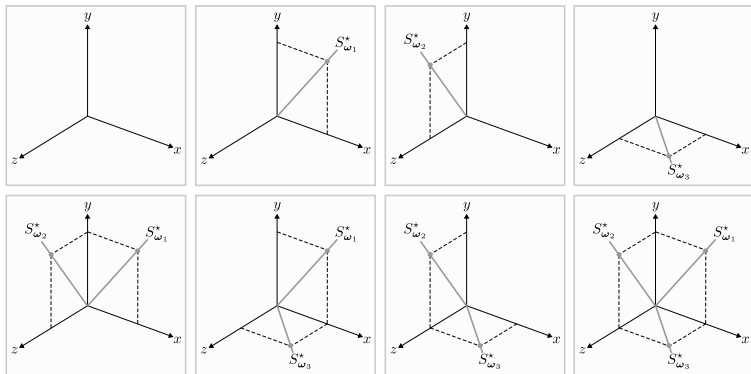
Well... it depends on which set of projections I give you.



Can you tell which are *the good sets*?

# Subspace Identifiability Problem

Well... it depends on which set of projections I give you.



Can you tell which are *the good sets*?

This is what we focused on: which are *the good sets*.

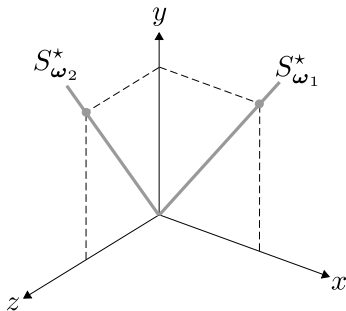


# Outline

- ▶ Introduction ✓
- ▶ What changes with missing data? ✓
- ▶ Subspace Identifiability Problem ✓
- ▶ **Setup**
- ▶ The Answer
- ▶ Application
- ▶ Conclusions

## Setup

The columns of  $\Omega$  will index the given projections.



$$\Omega = \begin{bmatrix} \omega_1 & \omega_2 \\ 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix}$$

## Setup

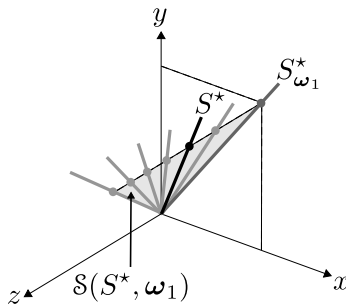
- ▶  $\text{Gr}(r, \mathbb{R}^d) :=$  Grassmannian manifold of  $r$ -dimensional subspaces in  $\mathbb{R}^d$ .

## Setup

- ▶  $\text{Gr}(r, \mathbb{R}^d) :=$  Grassmannian manifold of  $r$ -dimensional subspaces in  $\mathbb{R}^d$ .
- ▶  $\mathcal{S}(S^*, \Omega) :=$  Set of  $r$ -dimensional subspaces that agree with  $S^*$  on  $\Omega$ .

# Setup

- ▶  $\text{Gr}(r, \mathbb{R}^d) :=$  Grassmannian manifold of  $r$ -dimensional subspaces in  $\mathbb{R}^d$ .
- ▶  $\mathcal{S}(S^*, \Omega) :=$  Set of  $r$ -dimensional subspaces that agree with  $S^*$  on  $\Omega$ .

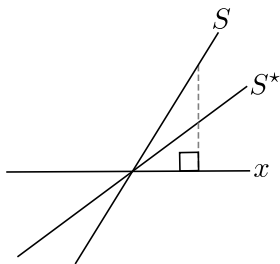


# Setup

- ▶  $S^*$  is  $r$ -dimensional.

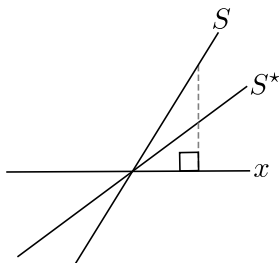
# Setup

- ▶  $S^*$  is  $r$ -dimensional.
- ▶ The projection of  $S^*$  onto  $\leq r$  canonical coordinates gives no information about  $S^*$ .



# Setup

- ▶  $S^*$  is  $r$ -dimensional.
- ▶ The projection of  $S^*$  onto  $\leq r$  canonical coordinates gives no information about  $S^*$ .



- ▶  $\Rightarrow$  Assume w.l.o.g. that all projections are onto  $r + 1$  canonical coordinates.



## Setup

- For any matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ :

$$\Omega' = \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}}_{n(\Omega') := \# \text{columns}} \left. \vphantom{\begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}} \right\} m(\Omega') := \# \text{nonzero rows}$$

## Setup

- For any matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ :

$$\Omega' = \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}}_{n(\Omega') := \# \text{columns}} \left. \vphantom{\begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}} \right\} m(\Omega') := \# \text{nonzero rows}$$

- $d - r$  projections are *necessary*, so we will assume w.l.o.g.

$$n(\Omega) = d - r.$$

# Outline

- ▶ Introduction ✓
- ▶ What changes with missing data? ✓
- ▶ Subspace Identifiability Problem ✓
- ▶ Setup ✓
- ▶ **The Answer**
- ▶ Application
- ▶ Conclusions

# The Answer

## Theorem (P.-A., Nowak, Boston, '14)

*For almost every  $S^\star$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^\star$  is the only subspace in  $\mathcal{S}(S^\star, \mathbf{\Omega})$  if and only if for every matrix  $\mathbf{\Omega}'$  formed with a subset of the columns in  $\mathbf{\Omega}$ ,*

$$m(\mathbf{\Omega}') \geq n(\mathbf{\Omega}') + r.$$

## The Answer

For **almost every**  $S^\star$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^\star$  is the only subspace in  $\mathcal{S}(S^\star, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

## The Answer

For **almost every**  $S^\star$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^\star$  is the only subspace in  $\mathcal{S}(S^\star, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

---

There is a set of measure zero of *bad* subspaces that we wouldn't identify.

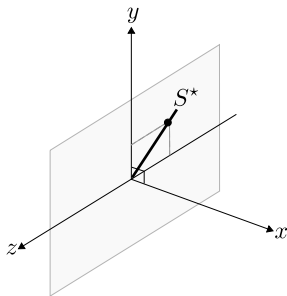
# The Answer

For **almost every**  $S^*$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^*$  is the only subspace in  $\mathcal{S}(S^*, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

---

There is a set of measure zero of *bad* subspaces that we wouldn't identify.



# The Answer

For almost every  $S^*$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^*$  is **the only** subspace in  $\mathcal{S}(S^*, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$



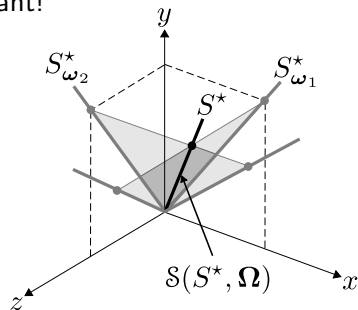
# The Answer

For almost every  $S^*$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^*$  is **the only** subspace in  $\mathcal{S}(S^*, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

---

This is what we want!



## The answer

For almost every  $S^*$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^*$  is the only subspace in  $\mathcal{S}(S^*, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

## The answer

For almost every  $S^*$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^*$  is the only subspace in  $\mathcal{S}(S^*, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

---

This is the answer!

## The answer

For almost every  $S^*$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^*$  is the only subspace in  $\mathcal{S}(S^*, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

---

This is the answer!

Every subset of  $n$  columns of  $\Omega$  has at least  $n + r$  nonzero rows.

## The answer

For almost every  $S^*$ , with respect to the uniform measure over  $\text{Gr}(r, \mathbb{R}^d)$ ,  $S^*$  is the only subspace in  $\mathcal{S}(S^*, \Omega)$  if and only if for every matrix  $\Omega'$  formed with a subset of the columns in  $\Omega$ ,

$$m(\Omega') \geq n(\Omega') + r.$$

---

This is the answer!

Every subset of  $n$  columns of  $\Omega$  has at least  $n + r$  nonzero rows.

$$\Omega = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \text{Check: } \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

# Outline

- ▶ Introduction ✓
- ▶ What changes with missing data? ✓
- ▶ Subspace Identifiability Problem ✓
- ▶ Setup ✓
- ▶ The Answer ✓
- ▶ **Application**
- ▶ Conclusions

# Application

Low-Rank Matrix Completion (LRMC)

# Application

## Low-Rank Matrix Completion (LRMC)

- ▶ Given a subset of entries in a rank  $r$  matrix, exactly recover *all* of the missing entries.

$$\mathbf{X}_{\Omega} = \begin{bmatrix} 1 & \cdot & 3 & \cdot \\ 1 & 2 & \cdot & \cdot \\ \cdot & 2 & 3 & \cdot \\ \cdot & \cdot & \cdot & 4 \\ \cdot & \cdot & \cdot & 4 \end{bmatrix} \Rightarrow \hat{\mathbf{X}} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$



# Application

## Low-Rank Matrix Completion (LRMC)

- ▶ Given a subset of entries in a rank  $r$  matrix, exactly recover *all* of the missing entries.

$$\mathbf{X}_{\Omega} = \begin{bmatrix} 1 & \cdot & 3 & \cdot \\ 1 & 2 & \cdot & \cdot \\ \cdot & 2 & 3 & \cdot \\ \cdot & \cdot & \cdot & 4 \\ \cdot & \cdot & \cdot & 4 \end{bmatrix} \Rightarrow \hat{\mathbf{X}} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

- ▶  $\sim$  Identifying the subspace spanned by the columns,  $S^*$ .

## Application

### Low-Rank Matrix Completion (LRMC)

- ▶ Given a subset of entries in a rank  $r$  matrix, exactly recover *all* of the missing entries.

$$\mathbf{X}_{\Omega} = \begin{bmatrix} 1 & \cdot & 3 & \cdot \\ 1 & 2 & \cdot & \cdot \\ \cdot & 2 & 3 & \cdot \\ \cdot & \cdot & \cdot & 4 \\ \cdot & \cdot & \cdot & 4 \end{bmatrix} \Rightarrow \hat{\mathbf{X}} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

- ▶  $\sim$  Identifying the subspace spanned by the columns,  $S^*$ . Here

$$\hat{S} = \text{span} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

## Application

How do we know we got the right completion (subspace)?

## Application

How do we know we got the right completion (subspace)?

- Maybe the real completion is:

$$\mathbf{X}_{\Omega} = \begin{bmatrix} 1 & \cdot & 3 & \cdot \\ 1 & 2 & \cdot & \cdot \\ \cdot & 2 & 3 & \cdot \\ \cdot & \cdot & \cdot & 4 \\ \cdot & \cdot & \cdot & 4 \end{bmatrix} \Rightarrow \mathbf{X} = \begin{bmatrix} 1 & 2 & 3 & 2 \\ 1 & 2 & 3 & 2 \\ 1 & 2 & 3 & 2 \\ 2 & 4 & 6 & 4 \\ 2 & 4 & 6 & 4 \end{bmatrix}$$

## Application

How do we know we got the right completion (subspace)?

- Maybe the real completion is:

$$\mathbf{X}_{\Omega} = \begin{bmatrix} 1 & \cdot & 3 & \cdot \\ 1 & 2 & \cdot & \cdot \\ \cdot & 2 & 3 & \cdot \\ \cdot & \cdot & \cdot & 4 \\ \cdot & \cdot & \cdot & 4 \end{bmatrix} \Rightarrow \mathbf{X} = \begin{bmatrix} 1 & 2 & 3 & 2 \\ 1 & 2 & 3 & 2 \\ 1 & 2 & 3 & 2 \\ 2 & 4 & 6 & 4 \\ 2 & 4 & 6 & 4 \end{bmatrix}$$

And the real subspace is

$$S^{\star} = \text{span} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 2 \end{bmatrix}$$

# Application

How do we know we got the right completion (subspace)?

# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

- ▶ Require random observed entries.



# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

- ▶ Require random observed entries.
  - ▶ May not be justified.

# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

- ▶ Require random observed entries.
  - ▶ May not be justified.
- ▶ Require incoherence

# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

- ▶ Require random observed entries.
  - ▶ May not be justified.
- ▶ Require incoherence
  - ▶ Sufficient, but not necessary condition.

# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

- ▶ Require random observed entries.
  - ▶ May not be justified.
- ▶ Require incoherence
  - ▶ Sufficient, but not necessary condition.
  - ▶ Generally unverifiable or unjustified in practice.

# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

- ▶ Require random observed entries.
  - ▶ May not be justified.
- ▶ Require incoherence
  - ▶ Sufficient, but not necessary condition.
  - ▶ Generally unverifiable or unjustified in practice.
- ▶ Work with high probability (if assumptions are met).

# Application

How do we know we got the right completion (subspace)?

Known results e.g. (Candès and Recht, '09)

- ▶ Require random observed entries.
  - ▶ May not be justified.
- ▶ Require incoherence
  - ▶ Sufficient, but not necessary condition.
  - ▶ Generally unverifiable or unjustified in practice.
- ▶ Work with high probability (if assumptions are met).

What if these assumptions are not met? How can we validate a completion?

# Application

## Corollary (P.-A., Nowak, Boston, '14)

*Let the columns of  $\mathbf{X}$  be drawn independently according to  $\nu$ , an absolutely continuous distribution with respect to the Lebesgue measure on  $S^*$ . Suppose  $\mathbf{X}_\Omega$  can be partitioned into two sets of columns,  $\mathbf{X}_{\Omega_1}$  and  $\mathbf{X}_{\Omega_2}$ , such that  $\Omega_2$  satisfies the conditions of the subspace identifiability theorem.*

*Let  $\hat{S}$  be the output of running an LRMC algorithm on  $\mathbf{X}_{\Omega_1}$ . Then for almost every  $S^*$ , and almost surely with respect to  $\nu$ ,  $\mathbf{X}_{\Omega_2}$  fits in  $\hat{S}$  if and only if  $\hat{S} = S^*$ .*

# Application

## Corollary (P.-A., Nowak, Boston, '14)

*Let the columns of  $\mathbf{X}$  be drawn independently according to  $\nu$ , an absolutely continuous distribution with respect to the Lebesgue measure on  $S^*$ . Suppose  $\mathbf{X}_\Omega$  can be partitioned into two sets of columns,  $\mathbf{X}_{\Omega_1}$  and  $\mathbf{X}_{\Omega_2}$ , such that  $\Omega_2$  satisfies the conditions of the subspace identifiability theorem.*

*Let  $\hat{S}$  be the output of running an LRMC algorithm on  $\mathbf{X}_{\Omega_1}$ . Then for almost every  $S^*$ , and almost surely with respect to  $\nu$ ,  $\mathbf{X}_{\Omega_2}$  fits in  $\hat{S}$  if and only if  $\hat{S} = S^*$ .*

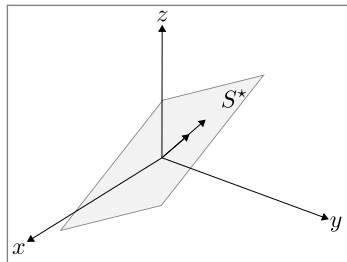
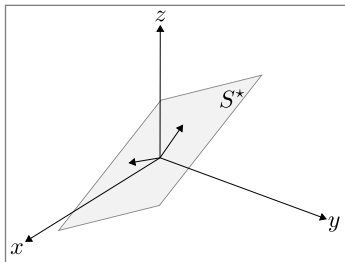


# Application

Just to make sure we have enough *useful* data

# Application

Just to make sure we have enough *useful* data



# Application

## Corollary (P.-A., Nowak, Boston, '14)

*Let the columns of  $\mathbf{X}$  be drawn independently according to  $\nu$ , an absolutely continuous distribution with respect to the Lebesgue measure on  $S^*$ . Suppose  $\mathbf{X}_\Omega$  can be partitioned into two sets of columns,  $\mathbf{X}_{\Omega_1}$  and  $\mathbf{X}_{\Omega_2}$ , such that  $\Omega_2$  satisfies the conditions of the subspace identifiability theorem.*

*Let  $\hat{S}$  be the output of running an LRMC algorithm on  $\mathbf{X}_{\Omega_1}$ . Then for almost every  $S^*$ , and almost surely with respect to  $\nu$ ,  $\mathbf{X}_{\Omega_2}$  fits in  $\hat{S}$  if and only if  $\hat{S} = S^*$ .*

# Application

## Corollary (P.-A., Nowak, Boston, '14)

*Let the columns of  $\mathbf{X}$  be drawn independently according to  $\nu$ , an absolutely continuous distribution with respect to the Lebesgue measure on  $S^*$ . Suppose  $\mathbf{X}_\Omega$  can be partitioned into two sets of columns,  $\mathbf{X}_{\Omega_1}$  and  $\mathbf{X}_{\Omega_2}$ , such that  $\Omega_2$  satisfies the conditions of the subspace identifiability theorem.*

*Let  $\hat{S}$  be the output of running an LRMC algorithm on  $\mathbf{X}_{\Omega_1}$ . Then for almost every  $S^*$ , and almost surely with respect to  $\nu$ ,  $\mathbf{X}_{\Omega_2}$  fits in  $\hat{S}$  if and only if  $\hat{S} = S^*$ .*

# Application

## Corollary (P.-A., Nowak, Boston, '14)

*Let the columns of  $\mathbf{X}$  be drawn independently according to  $\nu$ , an absolutely continuous distribution with respect to the Lebesgue measure on  $S^*$ . Suppose  $\mathbf{X}_\Omega$  can be partitioned into two sets of columns,  $\mathbf{X}_{\Omega_1}$  and  $\mathbf{X}_{\Omega_2}$ , such that  $\Omega_2$  satisfies the conditions of the subspace identifiability theorem.*

*Let  $\hat{S}$  be the output of running an LRMC algorithm on  $\mathbf{X}_{\Omega_1}$ .*

*Then for almost every  $S^*$ , and almost surely with respect to  $\nu$ ,  $\mathbf{X}_{\Omega_2}$  fits in  $\hat{S}$  if and only if  $\hat{S} = S^*$ .*

# Application

## Corollary (P.-A., Nowak, Boston, '14)

*Let the columns of  $\mathbf{X}$  be drawn independently according to  $\nu$ , an absolutely continuous distribution with respect to the Lebesgue measure on  $S^*$ . Suppose  $\mathbf{X}_\Omega$  can be partitioned into two sets of columns,  $\mathbf{X}_{\Omega_1}$  and  $\mathbf{X}_{\Omega_2}$ , such that  $\Omega_2$  satisfies the conditions of the subspace identifiability theorem.*

*Let  $\hat{S}$  be the output of running an LRMC algorithm on  $\mathbf{X}_{\Omega_1}$ .*

*Then for almost every  $S^*$ , and almost surely with respect to  $\nu$ ,  $\mathbf{X}_{\Omega_2}$  fits in  $\hat{S}$  if and only if  $\hat{S} = S^*$ .*

# Application

In contrast, our results:

# Application

In contrast, our results:

- ▶ Work for arbitrary observation schemes.



# Application

In contrast, our results:

- ▶ Work for arbitrary observation schemes.
- ▶ Work for almost every subspace.

# Application

In contrast, our results:

- ▶ Work for arbitrary observation schemes.
- ▶ Work for almost every subspace.
  - ▶ No incoherence assumption required.

# Application

In contrast, our results:

- ▶ Work for arbitrary observation schemes.
- ▶ Work for almost every subspace.
  - ▶ No incoherence assumption required.
- ▶ Hold with probability 1.

# Outline

- ▶ Introduction ✓
- ▶ What changes with missing data? ✓
- ▶ Subspace Identifiability Problem ✓
- ▶ Setup ✓
- ▶ The Answer ✓
- ▶ Application ✓
- ▶ Conclusions

# Conclusions

Now we know that:

- ▶ It is possible to uniquely identify an  $r$ -dimensional subspace  $S^*$  from its projections onto  $\Omega$ .

# Conclusions

Now we know that:

- ▶ It is possible to uniquely identify an  $r$ -dimensional subspace  $S^*$  from its projections onto  $\Omega$ .
- ▶ If and only if every subset of  $n$  columns of  $\Omega$  has at least  $n + r$  nonzero rows.

# Conclusions

Now we know that:

- ▶ It is possible to uniquely identify an  $r$ -dimensional subspace  $S^*$  from its projections onto  $\Omega$ .
- ▶ If and only if every subset of  $n$  columns of  $\Omega$  has at least  $n + r$  nonzero rows.
- ▶ Whence  $S^* = \ker \mathbf{A}^T$ .

Thanks.