# The Information-Theoretic Requirements of Subspace Clustering with Missing Data

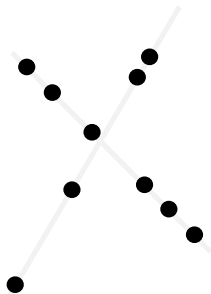**Daniel L. Pimentel-Alarcón**
Robert Nowak
*University of Wisconsin-Madison*

ICML 2016

# Subspace Clustering

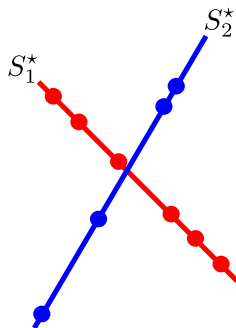- **We are given:** Columns in a Union of Subspaces.

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

# Subspace Clustering

- **We are given:** Columns in a Union of Subspaces.
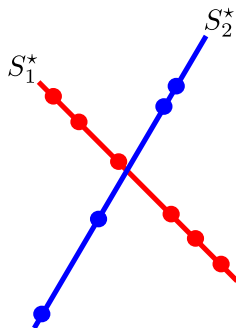- **Goal:** Cluster the columns or find the subspaces.

$$\begin{bmatrix} 1 & 4 & 1 & 3 & 3 & 1 & 2 & 1 & 2 & 1 \\ 2 & 4 & 2 & 6 & 3 & 2 & 2 & 2 & 4 & 1 \\ 3 & 4 & 3 & 9 & 3 & 3 & 2 & 3 & 6 & 1 \\ 1 & 8 & 1 & 3 & 6 & 1 & 4 & 1 & 2 & 2 \\ 2 & 8 & 2 & 6 & 6 & 2 & 4 & 2 & 4 & 2 \\ 3 & 8 & 3 & 9 & 6 & 3 & 4 & 3 & 6 & 2 \end{bmatrix}$$
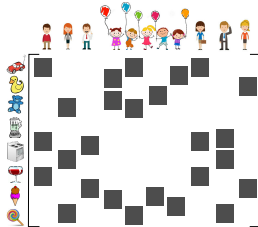
# Subspace Clustering with Missing Data (SCMD)

- **We are given:** Incomplete columns in a Union of Subspaces.
- **Goal:** Cluster the columns or find the subspaces.

# This arises in many Applications

In principle,
what do we need to succeed?

In principle,

what do we need to succeed?

To find out, let us look back at the full-data case.

# Full-data case

- Suppose I have unlimited computational power
- Say I want to identify $r$-dimensional subspaces from **complete** columns:

$$\mathbf{X} = \begin{bmatrix} 3 & 1 & 3 & 2 & 4 & 5 & 7 & 1 & 8 & 5 \\ 3 & 3 & 1 & 2 & 4 & 5 & 5 & 1 & 8 & 7 \\ 1 & 3 & 2 & 3 & 5 & 4 & 7 & 5 & 5 & 8 \\ 2 & 1 & 2 & 3 & 3 & 5 & 5 & 4 & 7 & 4 \end{bmatrix}$$

# Full-data case

Key Idea:

- $r$ columns define a *candidate* subspace.



$$\begin{bmatrix} x_1 & x_2 \\ y_1 & y_2 \\ z_1 & z_2 \end{bmatrix}$$

# Full-data case

Key Idea:

- $r$ columns define a *candidate* subspace.

# Full-data case

Key Idea:

- $r$ columns define a *candidate* subspace.
- I can certify this subspace with an $(r+1)^{\text{th}}$ column.



$$\begin{bmatrix} x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \\ z_1 & z_2 & z_3 \end{bmatrix}$$

# Full-data case

Key Idea:

- $r$ columns define a *candidate* subspace.
- I can certify this subspace with an $(r+1)^{\text{th}}$ column.



$$\begin{bmatrix} x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \\ z_1 & z_2 & z_3 \end{bmatrix}$$

# Full-data case

I can try all combinations of $r + 1$ columns (here $r = 2$).

$$\begin{bmatrix} 3 & 1 & 3 \\ 3 & 3 & 1 \\ 1 & 3 & 2 \\ 2 & 1 & 2 \end{bmatrix}$$

Linearly independent
$\Updownarrow$
Columns come from
different subspaces.

$$\begin{bmatrix} 3 & 2 & 5 \\ 3 & 2 & 5 \\ 1 & 3 & 4 \\ 2 & 3 & 5 \end{bmatrix}$$

Linearly dependent
$\Updownarrow$
Columns come from the same
subspace.

# Full-data case

We can try combinations of $r + 1$ columns until we identify all subspaces

$$S_1^\star = \mathrm{span} \begin{bmatrix} 1 & 3 \\ 3 & 1 \\ 3 & 2 \\ 1 & 2 \end{bmatrix} \qquad\qquad S_2^\star = \mathrm{span} \begin{bmatrix} 3 & 2 \\ 3 & 2 \\ 1 & 3 \\ 2 & 3 \end{bmatrix}$$

Then we can trivially cluster the columns.

$$\begin{bmatrix} 3 & 1 & 3 & 2 & 4 & 5 & 7 & 1 & 8 & 5 \\ 3 & 3 & 1 & 2 & 4 & 5 & 5 & 1 & 8 & 7 \\ 1 & 3 & 2 & 3 & 5 & 4 & 7 & 5 & 5 & 8 \\ 2 & 1 & 2 & 3 & 3 & 5 & 5 & 4 & 7 & 4 \end{bmatrix}$$

# Full-data case

Key idea:

> $r + 1$ **complete** columns *fit* in an $r$-dimensional subspace
> $\Updownarrow$
> Columns come from the same subspace.



$$\begin{bmatrix} x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \\ z_1 & z_2 & z_3 \end{bmatrix}$$

# What changes with missing data?

- We don't know where points really are!
- Say I give you a point *without* the $z$ coordinate.



$$\begin{bmatrix} x_1 \\ y_1 \\ \cdot \end{bmatrix}$$

# What changes with missing data?

- We don't know where points really are!
- Say I give you a point *without* the $z$ coordinate.
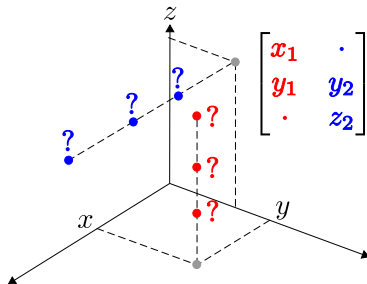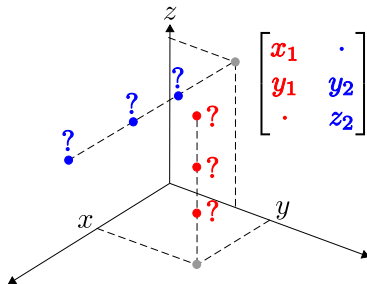- Say I give you an other point *without* the $x$ coordinate.

# What changes with missing data?

- We don't know where points really are!
- Say I give you a point *without* the $z$ coordinate.
- Say I give you an other point *without* the $x$ coordinate.



$$\begin{bmatrix} x_1 & \cdot \\ y_1 & y_2 \\ \cdot & z_2 \end{bmatrix}$$

- Is there only one subspace that agrees with these columns?

# What changes with missing data?

- ∃ **False** subspaces that can fit arbitrarily many **incomplete** columns from different subspaces.

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

# What changes with missing data?

- So even with **unlimited** computational power:
  - We could run into **false** subspaces!
  - And get a wrong clustering!
- How can we guarantee that this won't happen?
- We need to make sure that columns are observed in the right places.

# What do I mean observed in the right places?

We say $\mathbf{X_\Omega}$ is observed in the right places if every matrix $\mathbf{X}'_{\mathbf{\Omega}'}$ formed with a *proper* subset of the columns in $\mathbf{X_\Omega}$ satisfies

$$\#RowsWithObservations(\mathbf{X}'_{\mathbf{\Omega}'}) \geq \#Columns(\mathbf{X}'_{\mathbf{\Omega}'}) + r.$$

# What do I mean observed in the right places?

We say $\mathbf{X}_{\boldsymbol{\Omega}}$ is observed in the right places if every matrix $\mathbf{X}'_{\boldsymbol{\Omega}'}$ formed with a *proper* subset of the columns in $\mathbf{X}_{\boldsymbol{\Omega}}$ satisfies

$$\#RowsWithObservations(\mathbf{X}'_{\boldsymbol{\Omega}'}) \;\geq\; \#Columns(\mathbf{X}'_{\boldsymbol{\Omega}'}) + r.$$

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

$$\underbrace{m(\boldsymbol{\Omega}')}_{3} \;\ngeq\; \underbrace{n(\boldsymbol{\Omega}')/r + r}_{4}$$

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

$$\underbrace{m(\boldsymbol{\Omega}')}_{4} \;\geq\; \underbrace{n(\boldsymbol{\Omega}')/r + r}_{4}$$

# How do we know if columns come from same subspace?

Key intuition:

$d - r + 1$ **incomplete** columns (observed in the right places)
*behave* as **one complete** column.

$$
\begin{bmatrix}
4 & 1 & \cdot \\
4 & 3 & 1 \\
\cdot & 3 & 2 \\
3 & \cdot & 2
\end{bmatrix}
\underbrace{\phantom{xxx}}_{d-r+1}
\qquad \sim \qquad
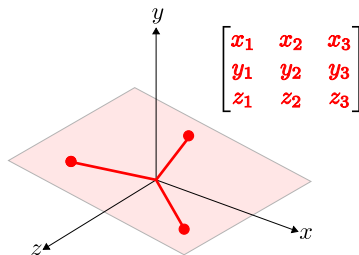\begin{bmatrix}
1 \\
3 \\
3 \\
1
\end{bmatrix}
\underbrace{\phantom{x}}_{1}
$$

# How do we know if columns come from same subspace?

Recall:

> $r + 1$ **complete** columns *fit* in an $r$-dimensional subspace
> $\Updownarrow$
> Columns come from the same subspace.

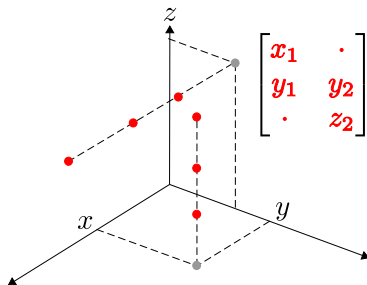# How do we know if columns come from same subspace?

Analogously:

> **Theorem (P.-A., Nowak, ICML '16)**
>
> $r + 1$ *sets of* $d - r + 1$ *incomplete columns (observed in the right places) fit in an* $r$-*dimensional subspace*
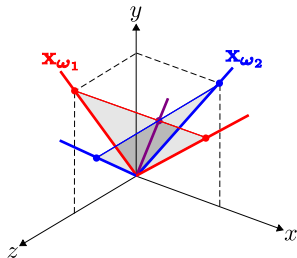>
> $\Updownarrow$
>
> *Columns come from the same subspace.*



$$\begin{bmatrix} x_1 & \cdot \\ y_1 & y_2 \\ \cdot & z_2 \end{bmatrix}$$

# Main Idea of the Proof

Each column with $r+1$ samples imposes one polynomial constraint on the subspaces that can *agree* with it.

$$\mathbf{X_\Omega} = \begin{bmatrix} \mathbf{x}_{\boldsymbol{\omega}_1} & \mathbf{x}_{\boldsymbol{\omega}_2} \\ \cdot & 2 \\ 1 & 2 \\ 1 & \cdot \end{bmatrix}$$



- A subspace $S$ *agrees* with $\mathbf{X_\Omega}$ $\iff$ $\begin{cases} f_1(S_{\boldsymbol{\omega}_1} | \mathbf{x}_{\boldsymbol{\omega}_1}) = 0 \\ f_2(S_{\boldsymbol{\omega}_2} | \mathbf{x}_{\boldsymbol{\omega}_2}) = 0 \end{cases}$ .

- If our columns are observed in the right places, only the true subspaces will be consistent with the constraints.
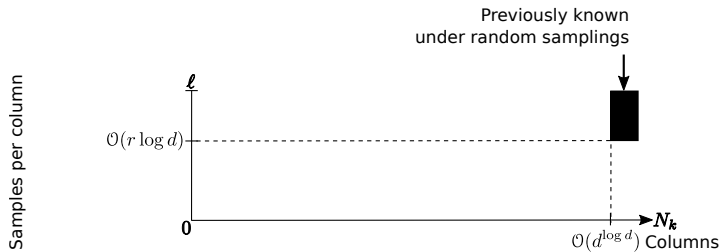
# Extend Linear Algebra Results to Missing Data

|  | Uniquely define a subspace | Certify/Discard subspaces |
|---|---|---|
| Full data | $r$ | $r+1$ |
| Missing data |  |  |

# Extend Linear Algebra Results to Missing Data

|  | Uniquely define a subspace | Certify/Discard subspaces |
|---|---|---|
| Full data | $r$ | $r + 1$ |
| Missing data | $(r + 1)(d - r)^{*}$ | $(r + 1)(d - r + 1)^{*}$ |

\* Observed in the right places.

# The Big Picture

Samples per column

Previously known
under random samplings

$\ell$

$\mathcal{O}(r \log d)$

$\mathbf{0}$

$N_k$

$\mathcal{O}(d^{\log d})$ Columns

# The Big Picture

# The Big Picture

# The Big Picture



Samples per column

$\ell = \dfrac{r(d-r)}{N_k} + r$

Previously known
under random samplings

$\ell$

$\mathcal{O}(r \log d)$
$\mathcal{O}(\max\{r, \log d\})$
$r+1$

$\mathbf{0}$

$r+1$
$r(d-r)$
$(r+1)(d-r+1)$
$\mathcal{O}(d^{\log d})$ Columns

$\boldsymbol{N_k}$

Impossible

Possible, if entries
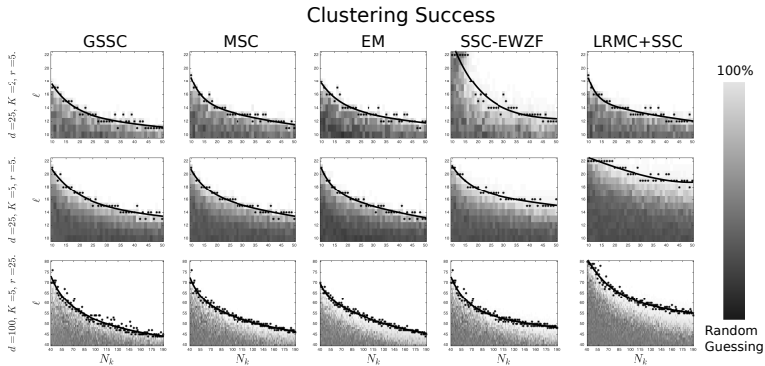are observed in
the right places

Possible, under
random samplings

# The Big Picture



Clustering Success

How am I on time?

Now we know when we
should be able to succeed.

Now we know when we
should be able to succeed.

Now the question is: How?

# Algorithms?

# Algorithms?

Computational Resources

|  | Efficient | Prohibitive |
|---|---|---|
| **Efficient** | | |
| **Prohibitive** | HRMC | Who cares. |

Number of Samples

# Algorithms?

|  | Computational Resources | |
|---|---|---|
|  | **Efficient** | **Prohibitive** |
| **Efficient** |  | $\mathcal{F} = 0$ |
| **Prohibitive** | HRMC | Who cares. |

Number of Samples

# Algorithms?



Computational Resources

|  | Efficient | Prohibitive |
|---|---|---|
| **Efficient** (Number of Samples) | EM<br>k-GROUSE<br>SSC-EWZF<br>LRMC+SSC | $\mathcal{F} = 0$ |
| **Prohibitive** (Number of Samples) | **HRMC** | Who cares. |

# Algorithms?

Computational Resources

Efficient                Prohibitive

| | Efficient | Prohibitive |
|---|---|---|
| **Efficient** | EM<br>k-GROUSE<br>SSC-EWZF<br>LRMC+SSC  **?** | $\mathcal{F} = 0$ |
| **Prohibitive** | HRMC | Who cares. |

Number of Samples

Thanks.

# What do I mean generic?