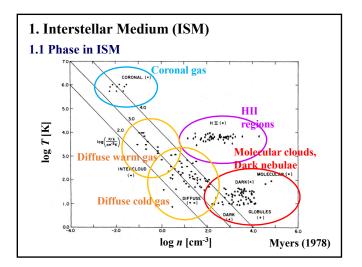
Application of High-dimensional Statistical Analysis to Astrophysical Problems

Tsutomu T. TAKEUCHI

- 1. Division of Particle and Astrophysical Science, Nagoya University, Japan
- 2. The Research Center for Statistical Machine Learning, the Institute of Statistical Mathematics

Symposium on Theories, Methodologies and Applications for Large Complex Data, Tsukuba, Japan 30-31 Oct., 2025



Collaborators

Kazuyoshi YATA (矢田 和書), Makoto AOSHIMA(青嶋 誠) Institute of Mathematics, University of Tsukuba, Japan

Kento EGASHIRA (江頭 健斗), Aki ISHII (石井 晶) Department of Information Sciences, Tokyo University of Science, Japan

Nanase HARADA (原田 ななせ), Kouichiro NAKANISHI (中西 康一郎) National Astronomical Observatory of Japan

Kohji YOSHIKAWA (吉川 耕司) Center for Computational Sciences, University of Tsukuba, Japan

Yu OGANE (大金 有羽), Ryusei R. KANO (加納龍生), Hai-Xia MA (馬 海震), Sena A. MATSUI (松井 瀬奈) Division of Particle and Astrophysical Science, Nagoya University, Japan

Suchetha COORAY (クレスチェータ) Kavli Institute Particle Astrophysics and Cosmology, Stanford University, USA

Hiroma OKUBO (大久保 宏真)

School of Science and Engineering, University of Tsukuba, Japan

Kotaro KOHNO (河野 孝太郎) Institute of Astronomy, The University of Tokyo, Japan

1.2 ISM phases and star formation

ISM has various phases

- Plasma (ionized diffuse phase)
- Neutral gas (mainly neutral hydrogen HI)
- 3. Molecular gas (mainly molecular hydrogen H₂)

Since gas must become dense enough to form stars, star formation occurs in molecular clouds. Namely,

Atomic gas ⇒ Molecular gas ⇒ Stars

0. Background: What Are Galaxies?

A galaxy is a huge agglomeration of stars, interstellar medium (ISM: gas+dust), and dark matter (DM), a complex system with a complicated interaction between each component.



Spatial scales

Spatial scales of galaxies and star formation (SF) are some orders of magnitude different:

Galaxies ~ kpc

Star formation ~ a few pc (for molecular clouds)

However, global properties of galaxies and SF activity are mysteriously correlated in various aspects!

⇒ Meso-scale physics to connect the scales of a galaxy and SF should be explored.

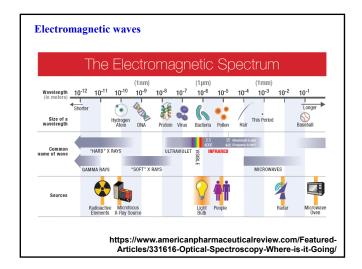
Kennicutt-Schmidt (K-S) law

Stars form in molecular cores.

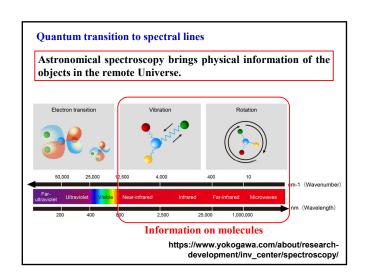
⇒ It is natural to suppose a relation between the star formation rate (SFR) and gas density. Schmidt (1959) proposed a relation

SFR $\propto \rho^n$.

- *i.* n = 1 Density controls star formation.
- *ii.* n = 2 Collision-like process plays a role for star formation
- ⇒ It is crucial to explore the properties of molecular clouds in star forming galaxies!

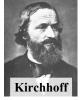


1.3 What does spectroscopy tell us? Spectroscopy https://www.atascientific.com.au/spectrometry/



1.3 What does spectroscopy tell us?

Kirchhoff and Bunsen showed that the emission lines corresponds to the absorption lines, which can be used for the identification of elements. Kirchhoff immediately pointed out that this leads to a very important application in astronomy.





http://www.hao.ucar.edu/public/education/sp/images/kirchhoff.html https://en.wikipedia.org/wiki/Robert_Bunsen

2. High-Dimensional Statistical Analysis

2.1 General situation in astrophysics

Classical statistical analysis

Sample size: *n*Data dimension: *d*

The following condition is implicitly assumed

n >> d

But this is not the case for many cases in scientific researches. Astronomers and astrophysicists have ever simply given up when they face such type of problem.

2. High-Dimensional Statistical Analysis

2.1 General situation in astrophysics

High-dimensional low-sample size (HDLSS) data analysis

Sample size: n Data dimension: d

For the HDLSS data, the condition is

$$n \ll d$$

This condition is often found in e.g., genomic analysis, medical analysis, etc.

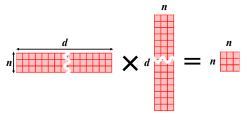
In astrophysics, for example, 2-dim spectral map such as integral field spectroscopy has this property.

2.2 Geometric Representation

Dual representation of sample covariance matrix

When we draw a set of n samples from the parent population (d > n), $\vec{x}_1, \ldots, \vec{x}_n$.

Consider a dual sample covariance matrix $(n \times n)$, $\tilde{S}_D = \frac{1}{n} \tilde{X}^T \tilde{X}$



This can be handled much more easily!

2.2 Unusual behavior of high-dimensional data

For high-dimensional data, classical limit theorems do not work. If we wrongly assume them, we would be lead to a wrong conclusion.

Simplest example: for the sample mean

$$\bar{\vec{x}} = \frac{1}{n} \sum_{i=1}^{n} \vec{x}_i$$

1. as $d/n \rightarrow 0$

$$\|\bar{\vec{x}} - \vec{\mu}\| \stackrel{P}{\rightarrow} \bar{0}$$

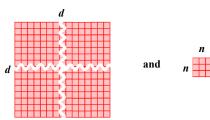
2. as $d/n \rightarrow \infty$

$$\|\bar{\vec{x}} - \vec{\mu}\| \stackrel{P}{\to} \infty$$

This striking property is referred to as the strong inconsistency.

Eigenvalues of the dual covariance matrix

When we draw a set of n samples from the parent population (d > n), $\vec{x}_1, \dots, \vec{x}_n$.



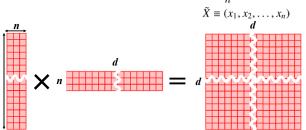
share the first n eigenvalues, i.e., the same important statistical information!

2.2 Geometric Representation

Dual representation of sample covariance matrix

When we draw a set of n samples from the parent population (d > n), $\vec{\chi}_1, \dots, \vec{\chi}_n$.

The sample covariance matrix $(d \times d)$ is $\tilde{S} = \frac{1}{n} \tilde{X} \tilde{X}^{\top}$,

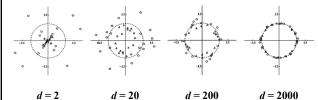


Note that this is a tremendously huge matrix!

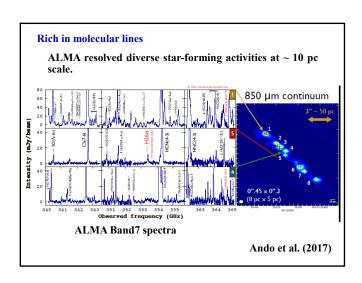
Unusual behavior of high-dimensional data: details

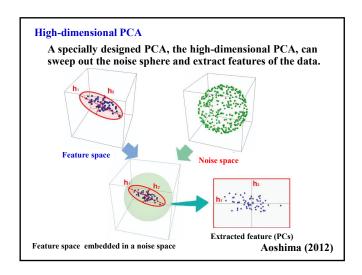
We can visualize the behavior of high-dimensional data vectors with dual representation. We omit all the mathematical details and jump onto the result.

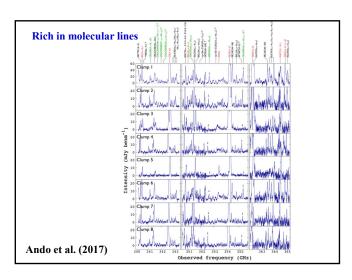
The population has a similar property with Gaussian ⇒ The eigenvectors concentrate on a sphere!!



Yata & Aoshima (2012)

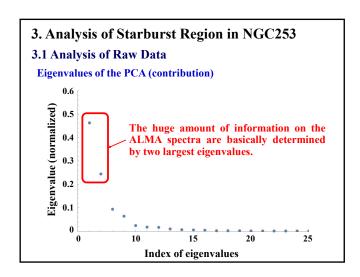


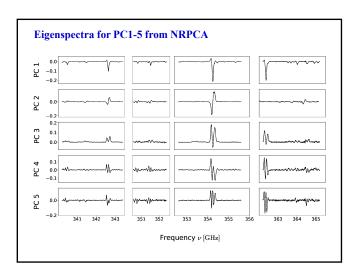


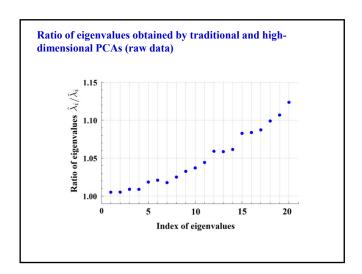


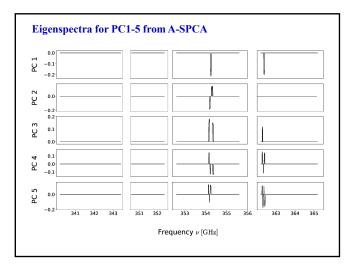
2.3 Actual data: ALMA data cube of NGC253 NGC 253: prototypal starburst Disk (out to > 20 kpc) Starburst nucleus (~ 1kpc) 2MASS (JHK) Jarrett et al. (2003)

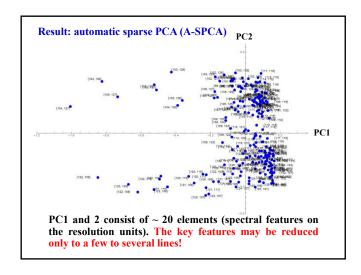
2.4 Structure of the Data Data: Ando et al. (2017) ~ spatial dimension 231 × spectral dimension 2248 ⇒ A case with n = 231 and d = 2248 (n << d) Problems from astrophysical side • Too much information on spectra. • Too large variety of spectral lines compared to n. We apply the high-dimensional statistical analysis to the ALMA spectral mapping data of NGC253.

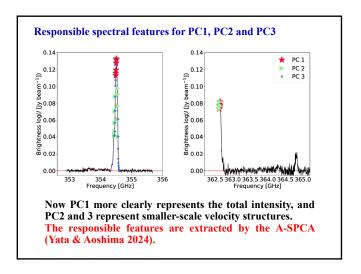


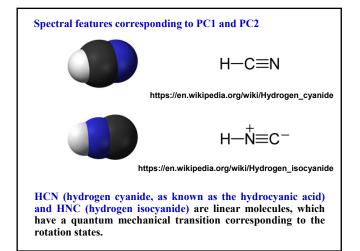


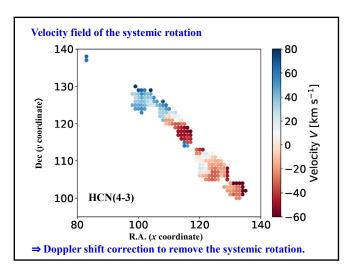


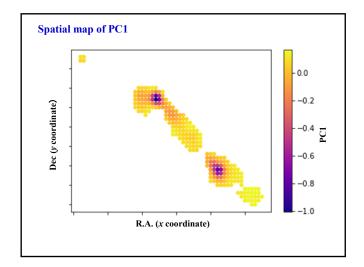


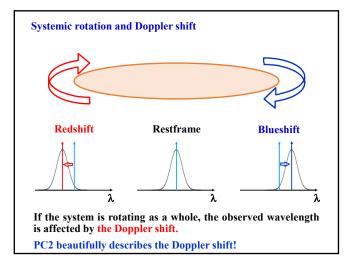


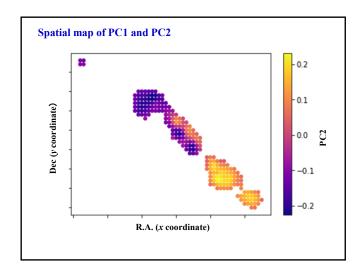


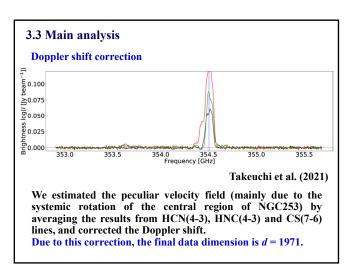


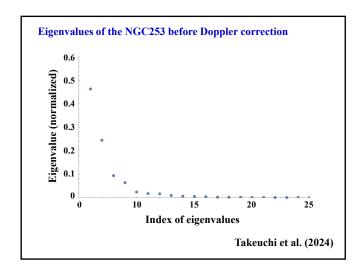


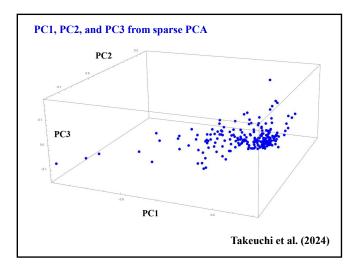


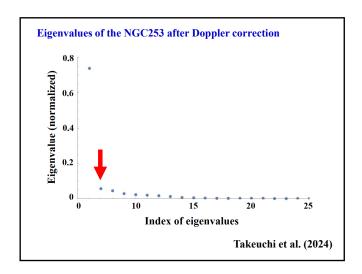


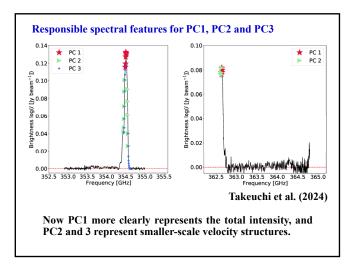


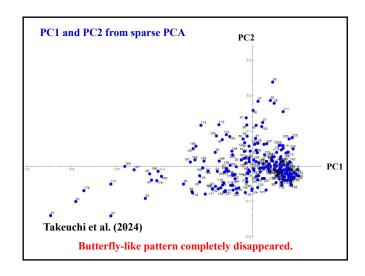


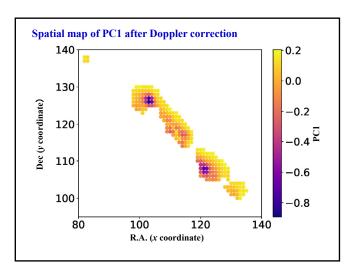


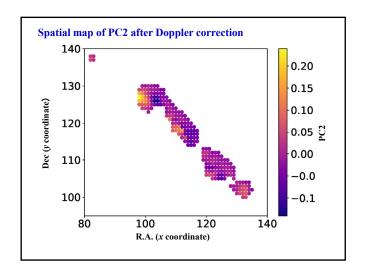










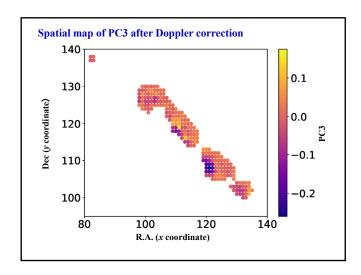


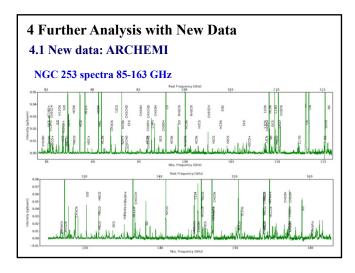
What do we see from the Doppler-corrected map?

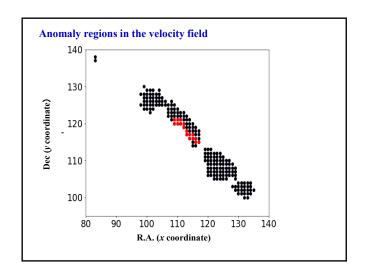
NGC253

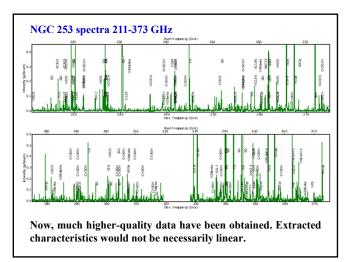
- Pure starburst: SFR in the central molecular zone is 2 M_{\odot} yr 1 (Rieke et al. 1980; Keto et al. 1999)
- Intense outflow (Matsubayashi et al. 2009; Bolatto et al. 2013)

Indeed the outflow phenomenon is mainly delineated by PC3.









4.2 Early result: finding erroneous data range

(Presented only on site)

5.2 Prospect and difficulty in the analysis of HI forest

Difficulty



The background quasars are very rare.

Even by the next-generation radio observational facility Square Kilometre Array (SKA), only a few tens of quasars are expected., while the absorption lines are numerous.

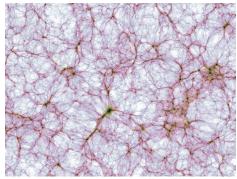
⇒ HDLSS data!

We constructed a new analysis method based on the highdimensional statistical analysis.

5. Analysis of the HI Forest 5.1 What is the HI forest? Quasar Observer v... [MHz] Ciardi et al. (2013)

5.3 Basic analysis with cosmological simulation





https://www.tng-project.org/media/

5.2 Prospect and difficulty in the analysis of HI forest

Prospect



- · The HI forest carry the information on the spatial distribution of primordial galaxies. This provides a very important clue to the formation of first galaxies.
- · The HI forest absorption line systems have evolved into galaxies at later epochs of the Universe. Their evolution might be reflected to the absorption lines.

5.3 Basic analysis with cosmological simulation

Quantification of the spatial distribution of HI gas

Absorption line frequency is described as

$$v_{\rm abs} = \frac{v_{\rm 21\,cm}}{1 + z_{\rm abs}}$$

i.e., the restframe frequency is shifted by cosmological redshift.

After some conversion of the observable, we have a data matrix of cosmic density fluctuation $\delta^{(j)}(z_i)$ $(d \times n)$ as

$$\vec{X} = \begin{pmatrix} \delta^{(1)}(z_1) & \delta^{(2)}(z_1) & \dots & \delta^{(n)}(z_1) \\ \delta^{(1)}(z_2) & \delta^{(2)}(z_2) & \dots & \delta^{(n)}(z_2) \\ \vdots & & \ddots & \vdots \\ \delta^{(1)}(z_d) & \delta^{(2)}(z_d) & \dots & \delta^{(n)}(z_d) \end{pmatrix}$$

5.3 Basic analysis with cosmological simulation

Quantification of the spatial distribution of HI gas

Two-point correlation function

$$\xi(z_1, z_2) \equiv \langle \delta(z_1)\delta(z_2) \rangle = \xi(|z_1 - z_2|)$$

From the data, we have a set of correlation functions in the form of the covariance matrix $\boldsymbol{\varXi}$

$$\Xi = \begin{pmatrix} \xi^{(1)}(0) & \xi^{(2)}(|z_1 - z_2|) & \dots & \xi^{(n)}(|z_1 - z_d|) \\ \xi^{(1)}(|z_2 - z_1|) & \xi^{(2)}(0) & \dots & \xi^{(n)}(|z_2 - z_d|) \\ \vdots & & \ddots & \vdots \\ \xi^{(1)}(|z_d - z_1|) & \xi^{(2)}(|z_d - z_2|) & \dots & \xi^{(n)}(0) \end{pmatrix}$$

This contains fundamental information on the cosmic matter density field.

6. Summary

- 1. Spectroscopic mapping and similar methods are fundamentally important to reveal the ISM physics, but the data are high-dimensional low sample size.
- 2. We applied the high-dimensional PCA on the NGC253 spectral map. ALMA mapping data are typically HDLSS in general, and in this case n = 231 and d = 2228.
- 3. NRPCA and A-SPCA can work very well to extract physical information from the HDLSS spectral map.
- 4. High-dimensional PCA also works as a powerful tool to find a part of the data with errors.

5.3 Basic analysis with cosmological simulation Quantification of the spatial distribution of HI gas High-dimensional PCA of the HI forest absorption lines Output Description of the HI forest absorption lines Cang et al. (2025)

We indeed observe the effect of the cosmic evolution of the absorption line system.

6. Summary

 High-dimensional methods are also very useful for the analysis of neutral hydrogen absorption line systems on distant QSO spectra. This is a promising method for future cosmological studies.

If you are interested in details, see
Takeuchi, T. T., et al. 2024a, ApJS, 271, 44
Takeuchi, T. T., et al. 2024b, Proceedings of the Institute
of Statistical Mathematics (統計數理), 72(2), 273
Cang, J.-S., et al. 2025, Advancing Astrophysics:
Preparing for Science with the SKAO, in press

5.3 Basic analysis with cosmological simulation Quantification of the spatial distribution of HI gas High-dimensional PCA of the HI forest absorption lines 20 PCI W PCI W

We try to quantify the evolution and specify the key factor to

characterize the evolution.

