

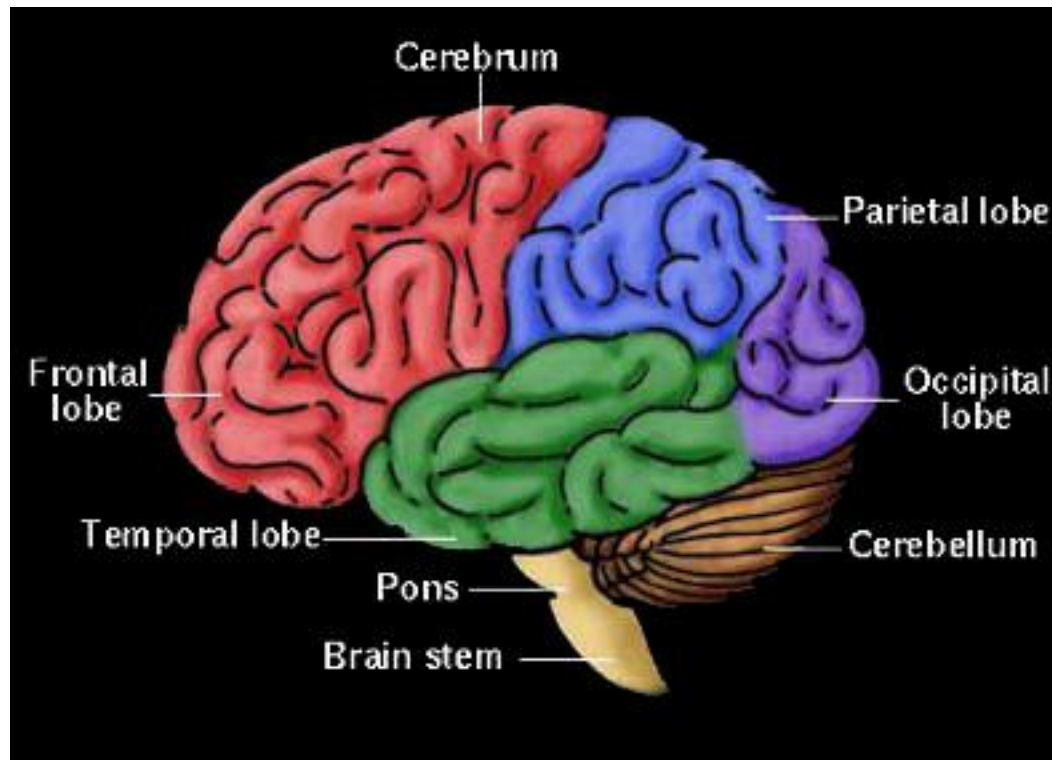
Data Mining In Modern Astronomy Sky Surveys:
*Concepts in Machine Learning,
Unsupervised Learning
& Astronomy Applications*

Ching-Wa Yip

cwyip@pha.jhu.edu; **Bloomberg 518**

Human are Great Pattern Recognizers

- Sensors: look, smell, touch, hear
- Computation: 100 billions (10^{11}) neurons



Estimated 300 million pattern recognizers
(*How to create a mind*, Kurzweil)

Machine Learning

- We want computers to perform tasks.
- It is difficult for computers to “learn” like the human do.
- We use algorithms:
 - Supervised, e.g.:
 - Classification
 - Regression
 - Unsupervised, e.g.:
 - Density Estimation
 - Clustering
 - Dimension Reduction

Machine Learning

- We want computers to perform tasks.
- It is difficult for computers to “learn” like the human do.
- We use algorithms:
 - Supervised, e.g.:
 - Classification
 - Regression
 - Unsupervised, e.g.:
 - Density Estimation
 - Clustering
 - Dimension Reduction

From Data to Information

- We don't just want data.
- We want information from the data.

Information



Database

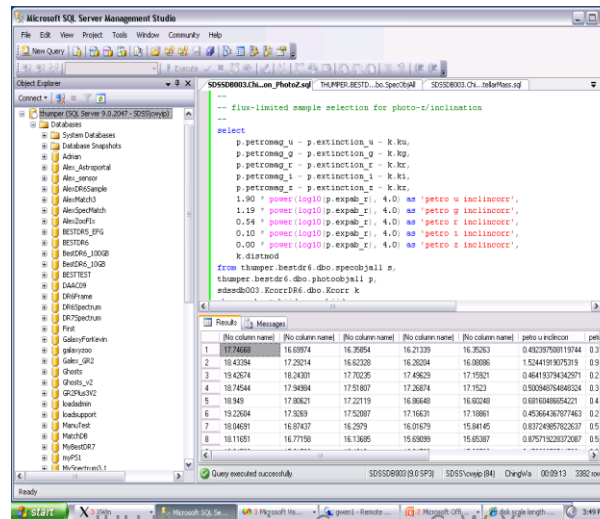


Sensors



Machine Learning

- *Not Data Mining by Human*



Expert could be Biased

(Thinking Fast and Slow, Kahneman)

challenges economic theory, according to which prices should reflect all the available information, including the weather. Ashenfelter's formula is extremely accurate—the correlation between his predictions and actual prices is above .90.

Why are experts inferior to algorithms? One reason, which Meehl suspected, is that experts try to be clever, think outside the box, and consider complex combinations of features in making their predictions. Complexity may work in the odd case, but more often than not it reduces validity. Simple combinations of features are better. Several studies have shown that human decision makers are inferior to a prediction formula even when they are given the score suggested by the formula! They feel that they can overrule the formula because they have additional information about the case, but

they are wrong more often than not. According to Meehl, there are few circumstances under which it is a good idea to substitute judgment for a formula. In a famous thought experiment, he described a formula that predicts whether a particular person will go to the movies tonight and noted that it is proper to disregard the formula if information is received that the individual broke a leg today. The name "broken-leg rule" has stuck. The point, of course, is that broken legs are very rare—as well as decisive.

Another reason for the inferiority of expert judgment is that humans are incorrigibly inconsistent in making summary judgments of complex information. When asked to evaluate the same information twice, they fre-



Expert could be Biased

(Thinking Fast and Slow, Kahneman)

challenges economic theory, according to which prices should reflect all the available information, including the weather. Ashenfelter's formula is extremely accurate—the correlation between his predictions and actual prices is above .90.

Why are experts inferior to algorithms? One reason, which Meehl suspected, is that experts try to be clever, think outside the box, and consider complex combinations of features in making their predictions. Complexity may work in the odd case, but more often than not it reduces validity. Simple combinations of features are better. Several studies have shown that human decision makers are inferior to a prediction formula even when they are given the score suggested by the formula! They feel that they can overrule the formula because they have additional information about the case, but

they are wrong more often than not. According to Meehl, there are few circumstances under which it is a good idea to substitute judgment for a formula. In a famous thought experiment, he described a formula that predicts whether a particular person will go to the movies tonight and noted that it is proper to disregard the formula if information is received that the individual broke a leg today. The name “broken-leg rule” has stuck. The point, of course, is that broken legs are very rare—as well as decisive.

Another reason for the inferiority of expert judgment is that humans are incorrigibly inconsistent in making summary judgments of complex information. When asked to evaluate the same information twice, they fre-



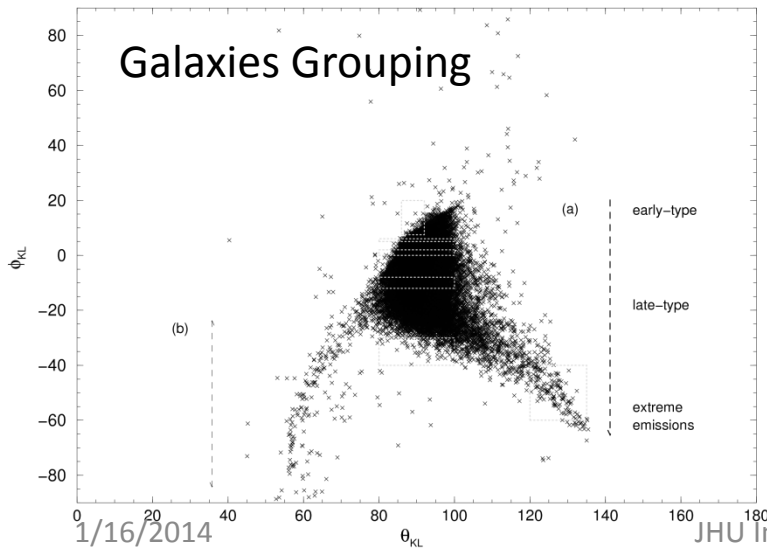
Applications of Unsupervised Learning



Consumer Clustering



Human Network Analysis



and more...

Basic Concepts in Machine Learning

- Label and Unlabeled data
- Datasets: training set and test set
- Feature space
- Distance between points
- Cost Function (or called Error Function)
- Shape of the data distribution
- Outliers

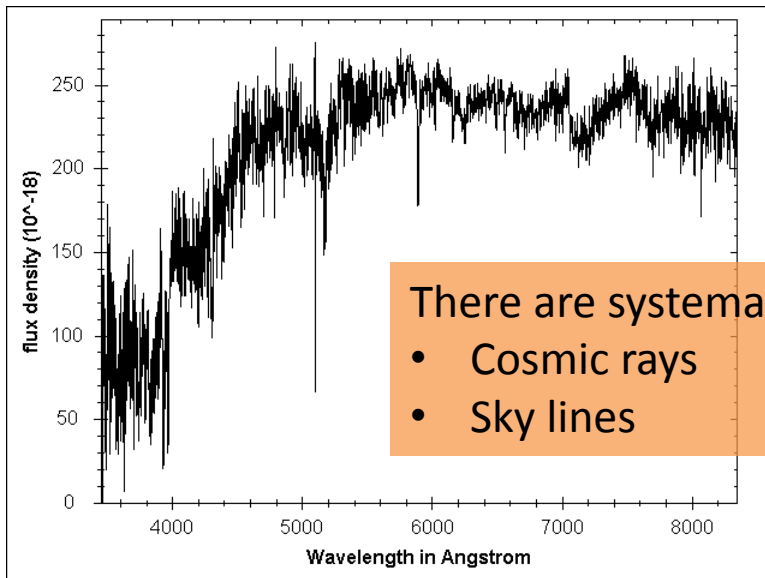
Basic Concepts in Machine Learning

- Label and Unlabeled data
- Datasets: training set and test set
- Feature space
- Distance between points
- Cost Function (or called Error Function)
- Shape of the data distribution
- Outliers

Feature Space

- The raw data may not be immediately suitable for pattern recognition.
- Feature space is spanned by the resultant data after pre-processing the raw data.
- The data usually N points (or vectors), each has M variables (or components).

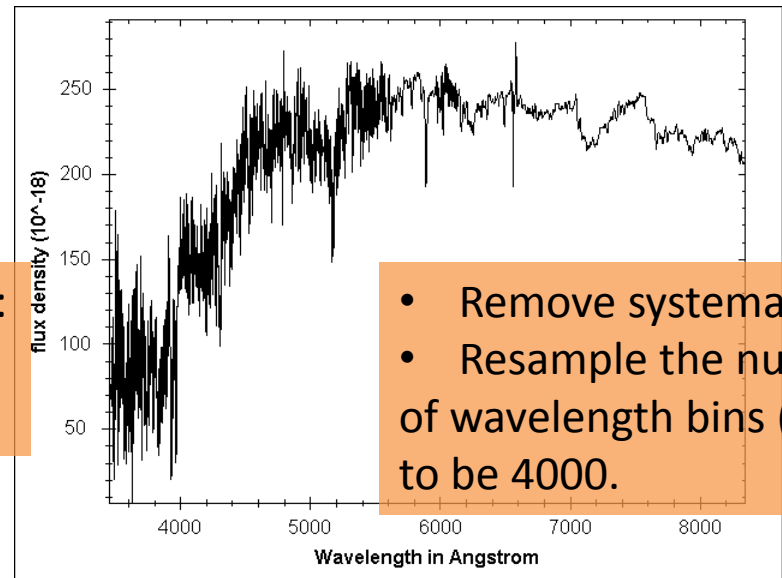
Galaxy Spectrum: Raw Data



There are systematics:

- Cosmic rays
- Sky lines

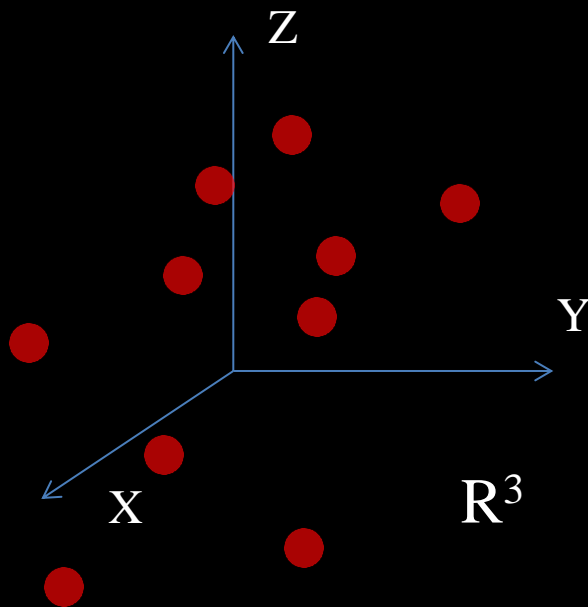
Galaxy Spectrum: After Pre-processing



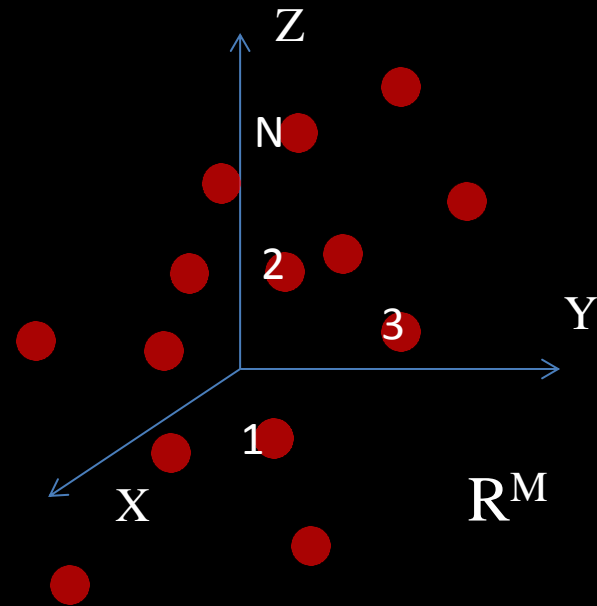
- Remove systematics.
- Resample the number of wavelength bins (M) to be 4000.

Intuition to High-Dimensional Data: N points, M components

$N = 10, M = 3$



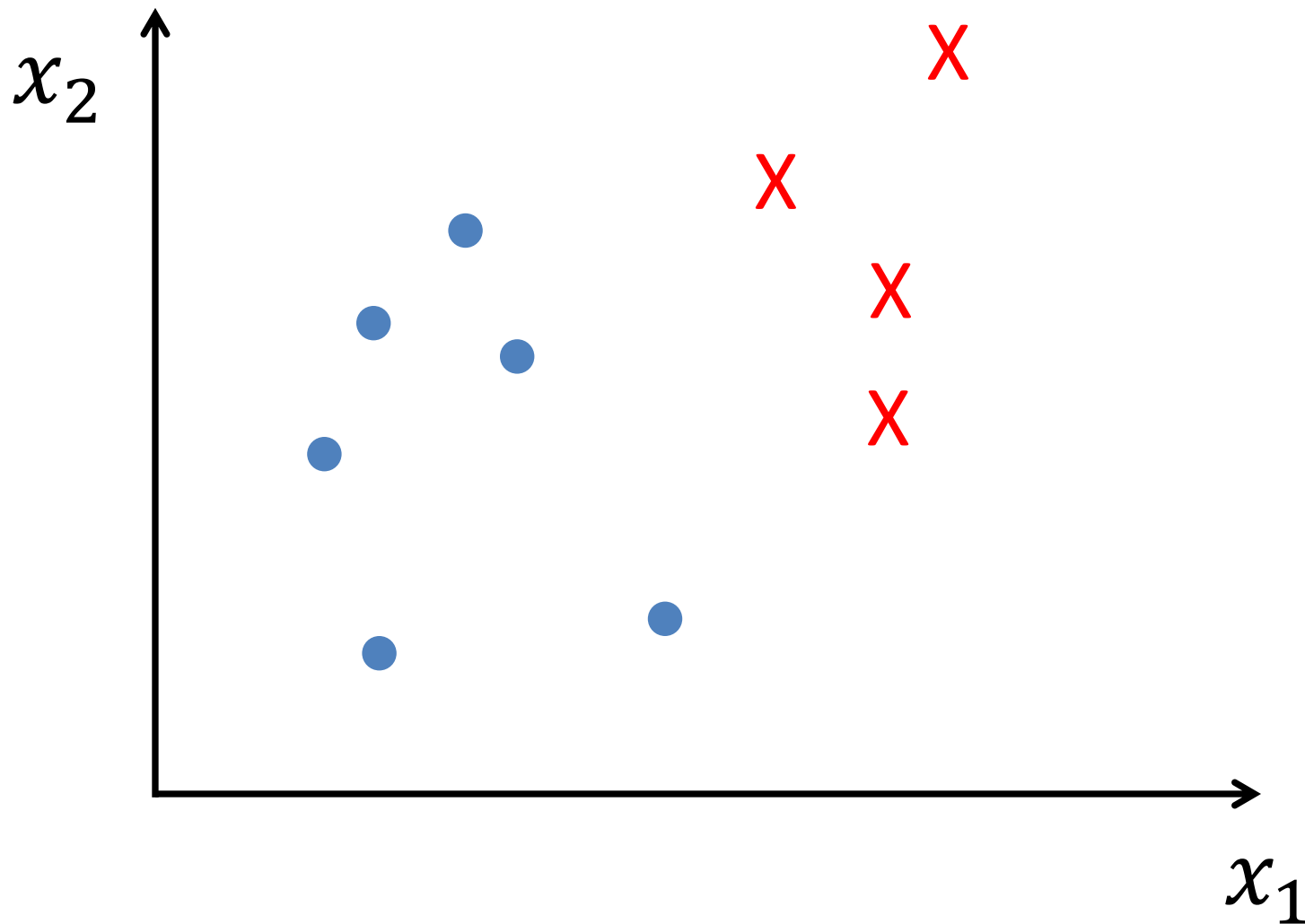
N, M



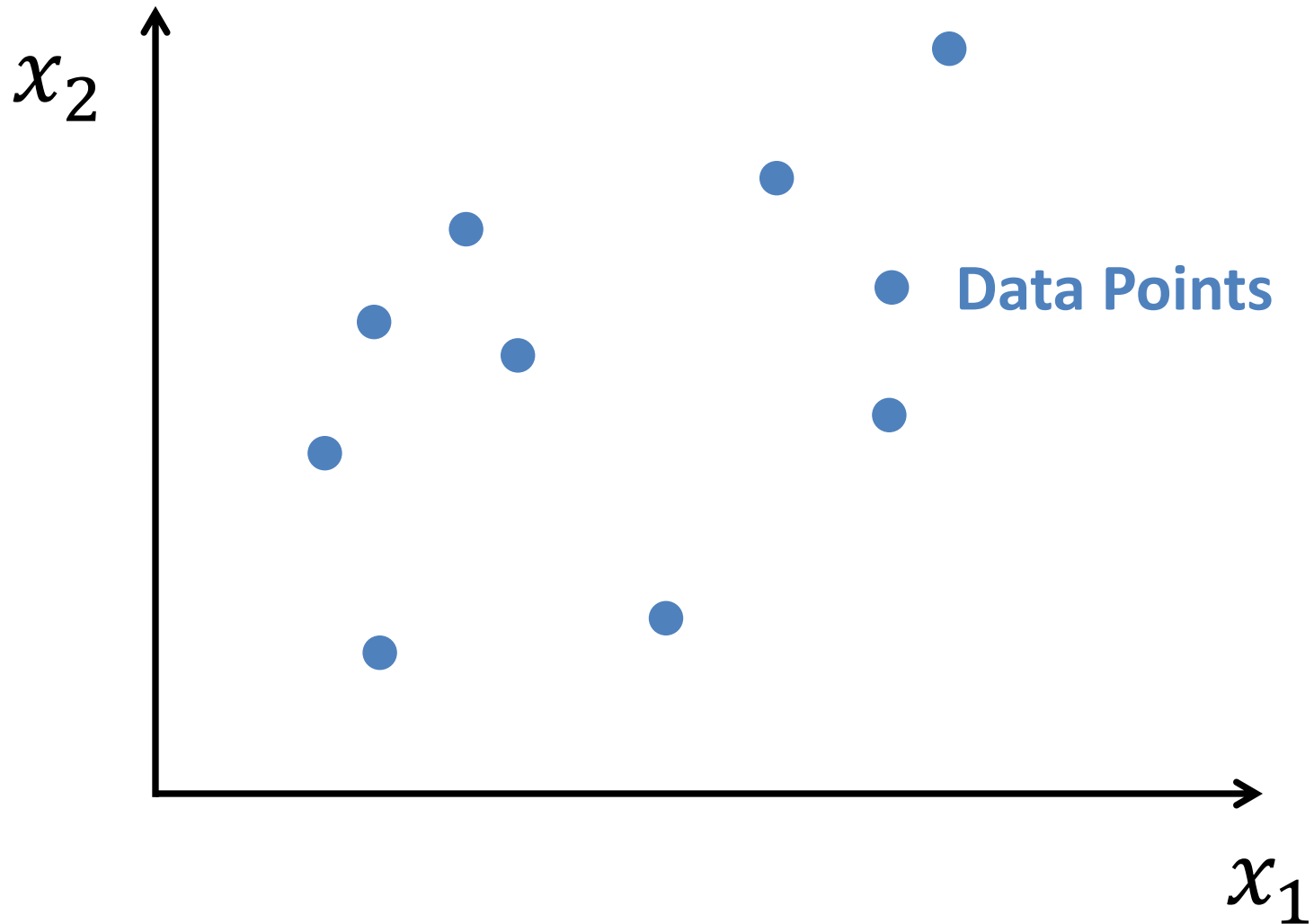
Unlabeled vs. Labeled Data

- Labeled data has an extra “label” compared with the unlabeled.
- Data labeling can be expensive (e.g., human expert).

Labeled Data



Unlabeled Data



Galaxy Images are Unlabeled Data (well, before Labeling)



(SDSS images of
CALIFA DR1;
Husemann et al. 2012,
Sanchez et al. 2012,
Walcher et al. in prep.)

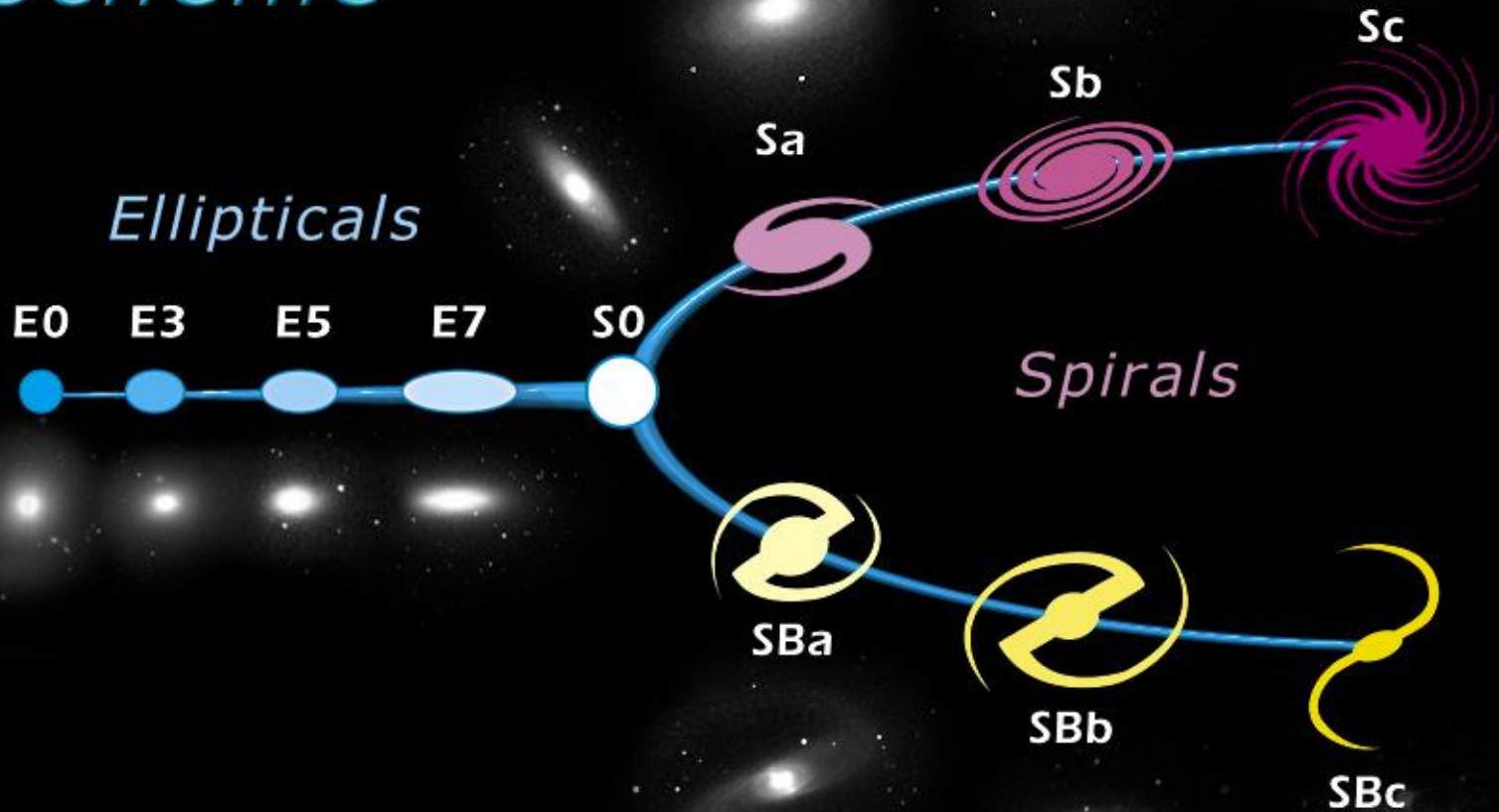
Galaxy Zoo Project

- Web users classify galaxy morphologies.
- Use majority vote to decide on the answer.



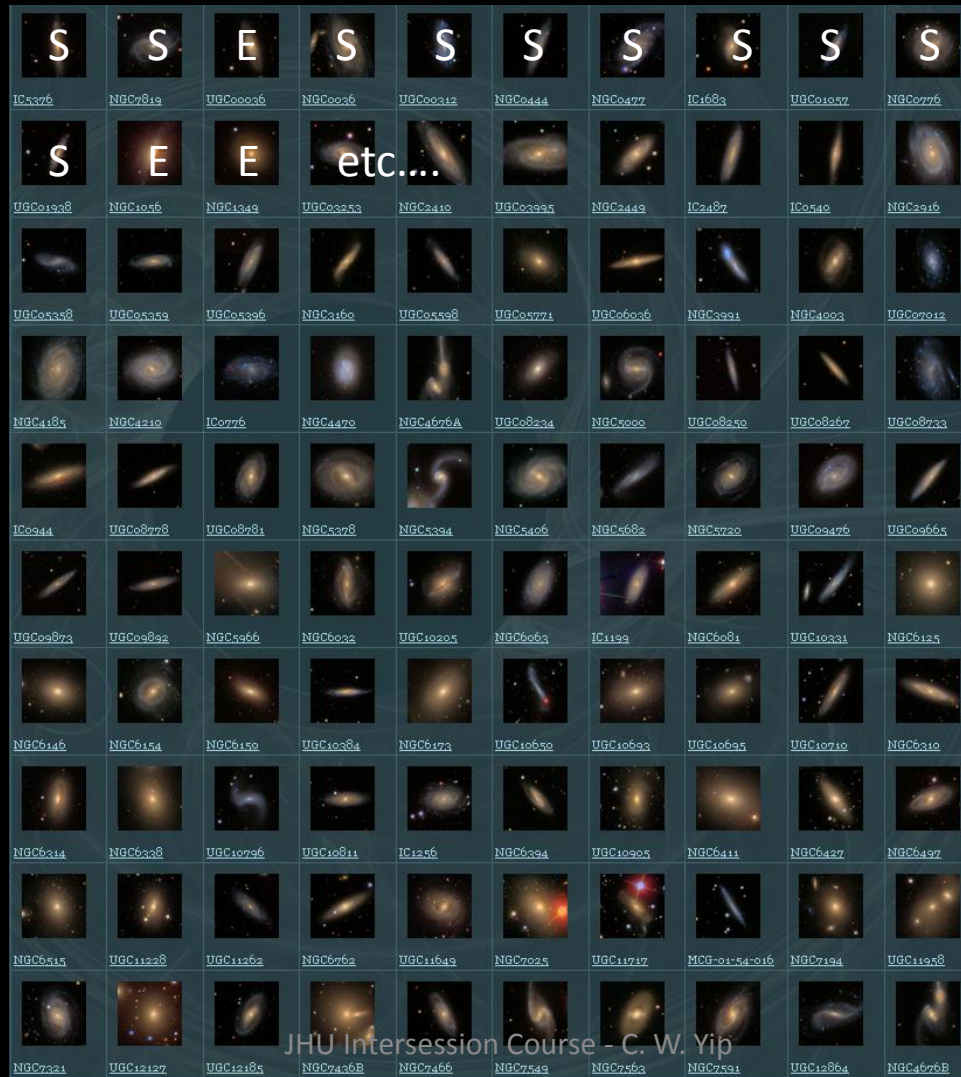
The screenshot shows the Galaxy Zoo.org website. At the top, the logo "GALAXY ZOO.org" is displayed with a stylized orange and red galaxy icon. Below the logo is a navigation menu with links: Welcome, Home, The Science, How to Take Part, Galaxy Analysis, Forum, Press, Blog, FAQ, Links, and Contact Us. There are also "Login" and "Register" buttons. The main content area features a large image of a galaxy. Below the image, there is a message to "Dear Galaxy Zoo users," thanking them for their contributions and providing information about the project's progress and future plans. On the right side, there is a login form with fields for "User Name" and "Password," a "Remember me next time" checkbox, and "Log In" and "Register" buttons. At the bottom of the page, there is a footer with copyright information and links to "The Team," "Privacy Statement," and "Copyright Notice."

Edwin Hubble's Classification Scheme



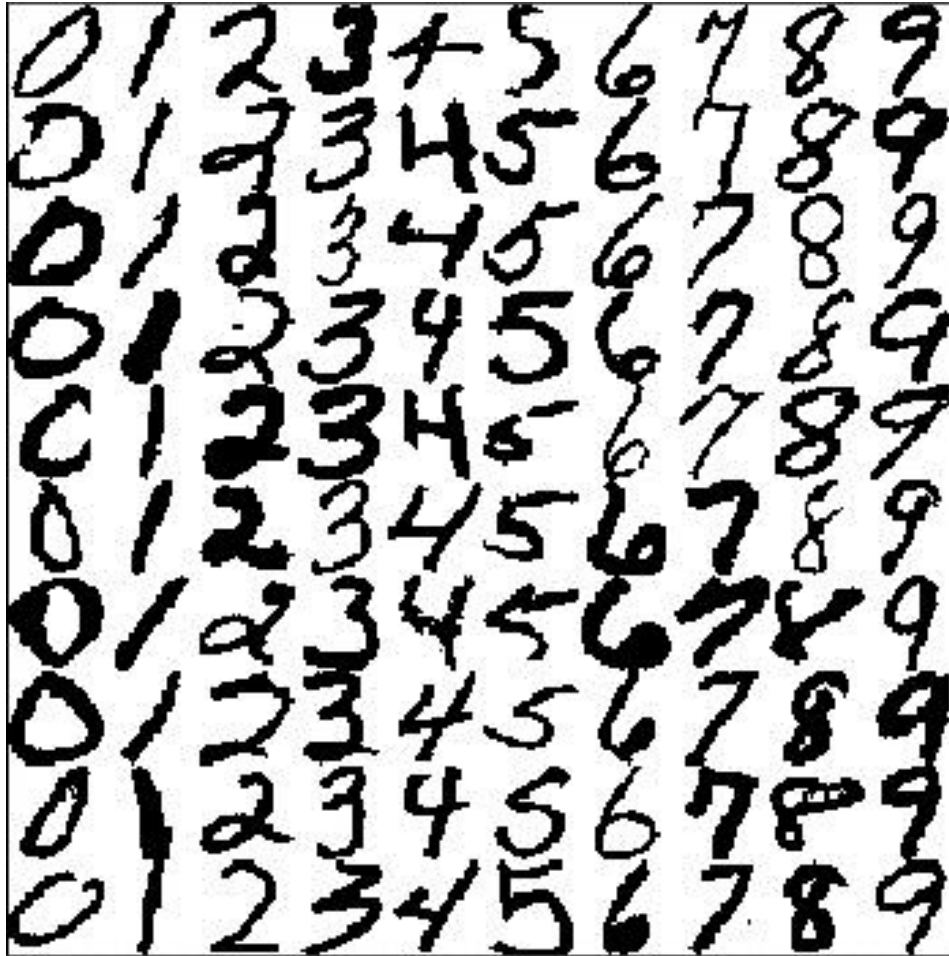
E0, E2, S0, ... etc. are the labels.

Galaxy Images are Unlabeled Data (well, before Labeling)



(SDSS images of
CALIFA DR1;
Husemann et al. 2012,
Sanchez et al. 2012,
Walcher et al. in prep.)

Digit Recognition



1/16/2014 Labels:

0 1 2 3 4 5 6 7 8 9

Distance between 2 Points in Space

- The set of data points spans a space.
- We can measure the distance between 2 points in such space.
 - E.g., Chi-sq measures the sum of distance squared between a point and model.
- If the points are close to each other:
 - They are “**neighbors**” with similar features.
- The best definition of distance (to yield concepts like “very close” and “far away”) is data-dependent.

Cost Function

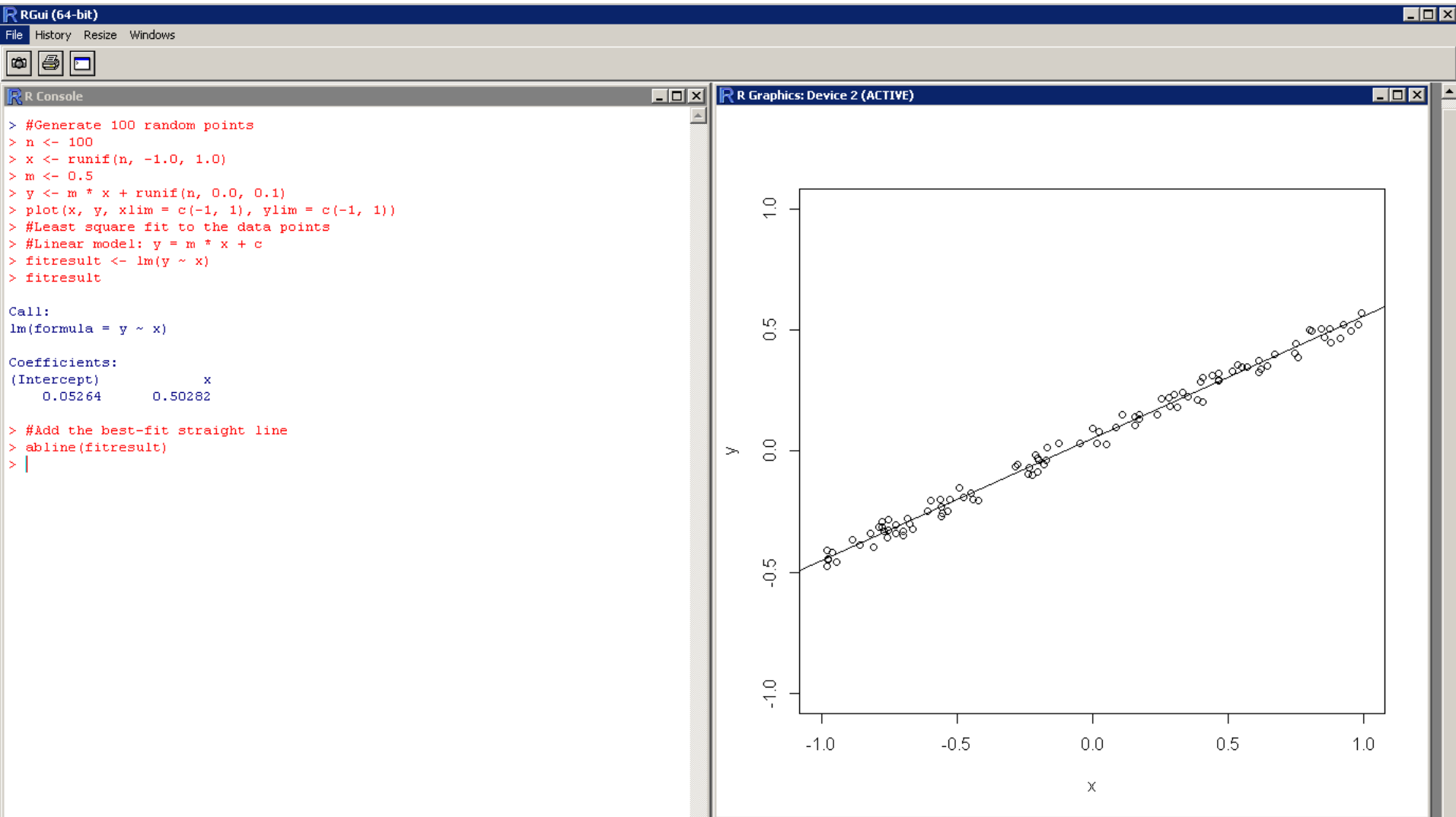
- In many machine learning algorithms, the idea is to find the model parameters θ which **minimize the cost function $J(\theta)$** :

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (\text{model}(\text{data}^i) - \text{data}^i)^2$$

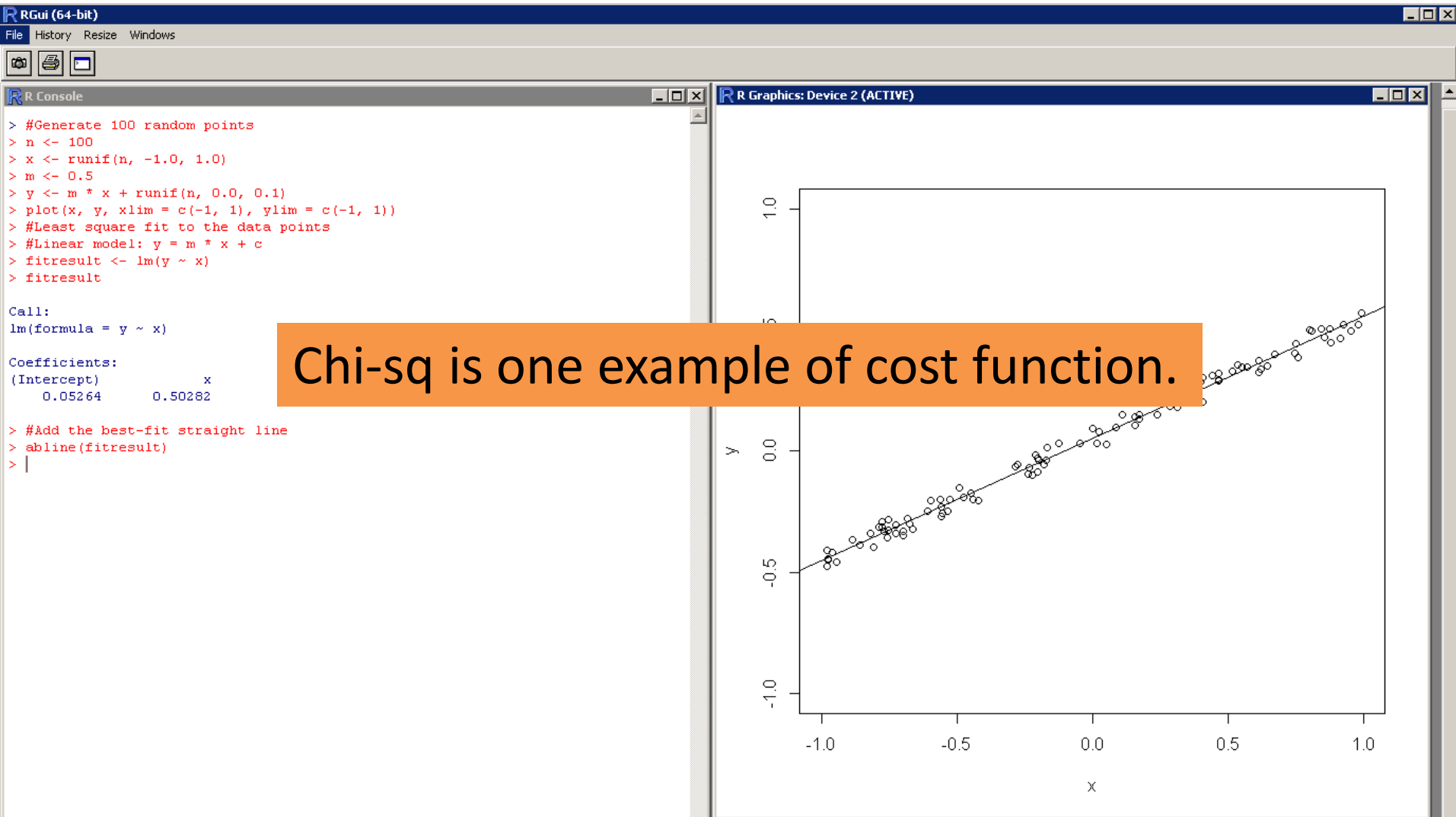
m is the size of the dataset or the training set.

- That is, we want *model* as close to *data* as possible.
- Note that *model* depends on θ .

Recall: LSQ Fitting



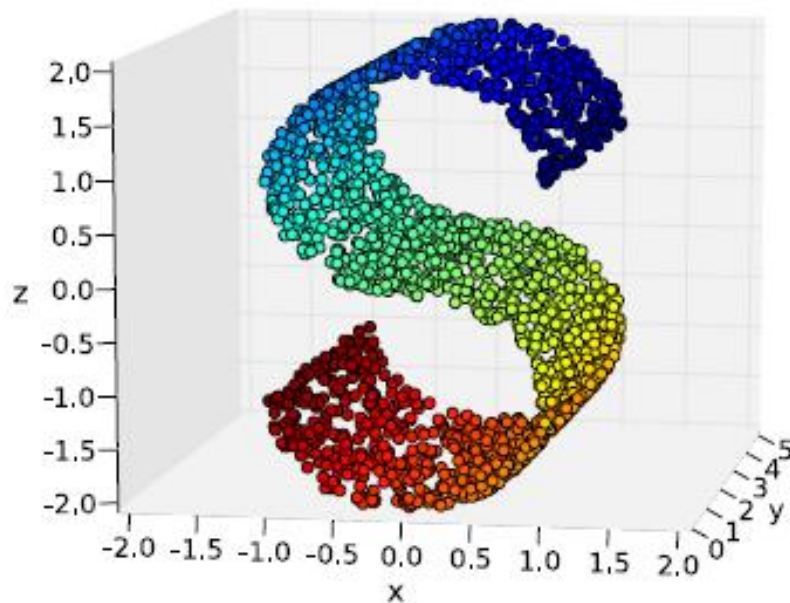
Recall: LSQ Fitting



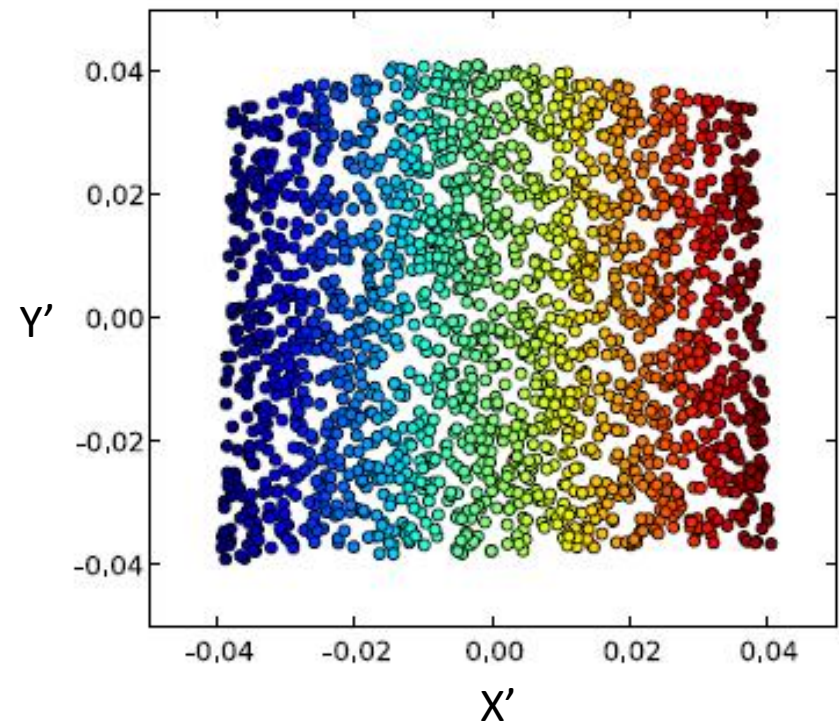
Chi-sq is one example of cost function.

Distance between Points is Non-trivial: “S Curve”

3D Space



2D Embedding Space



$$\left. \begin{array}{l} x = \sin(\theta) \\ z = \frac{\theta}{|\theta|}(\cos(\theta) - 1) \\ 0 < y < 5. \end{array} \right\} -1.5\pi \leq \theta \leq 1.5\pi$$

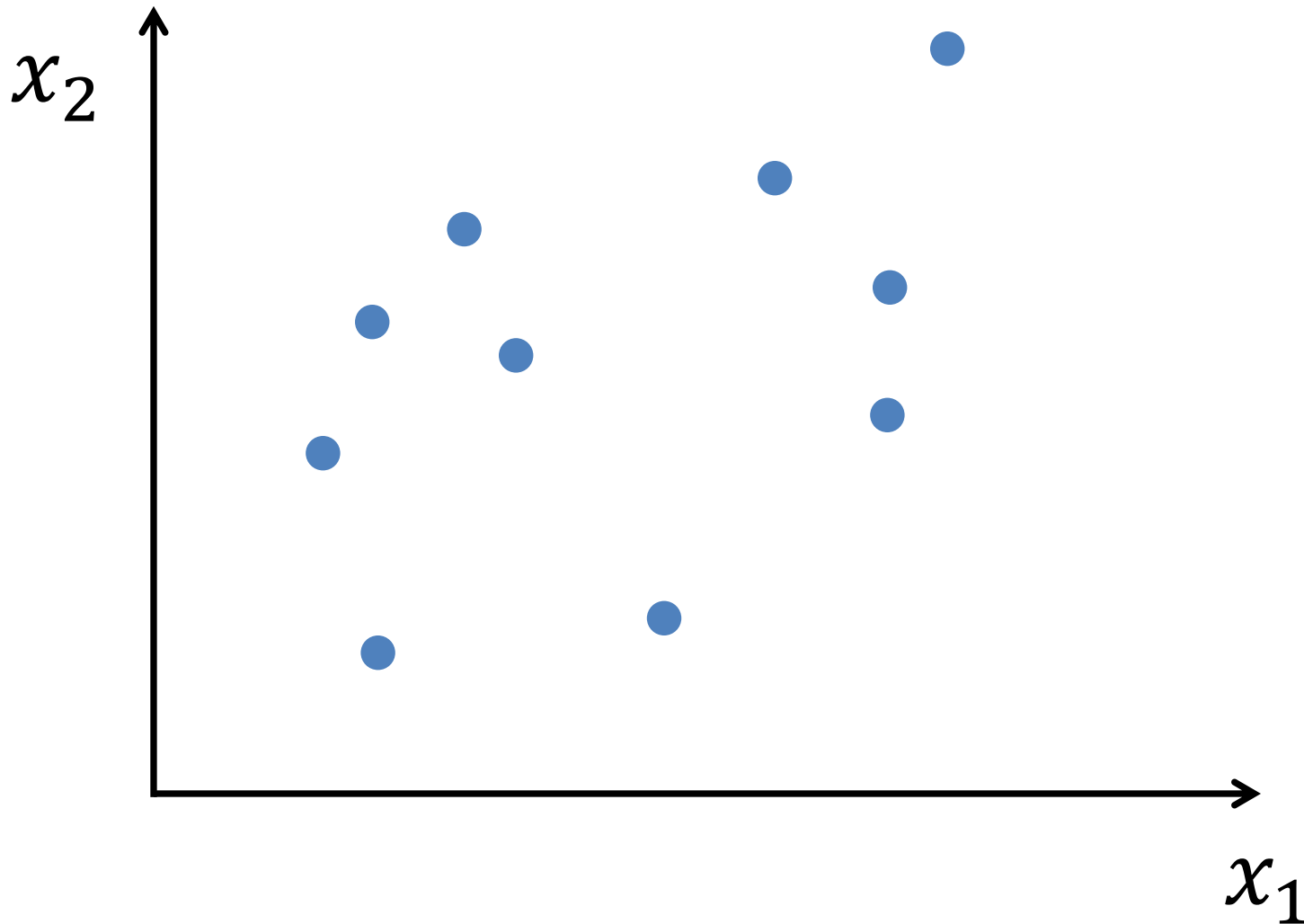
Unsupervised vs. Supervised Learning

- Unsupervised:
 - Given data $\{\mathbf{x}^1, \mathbf{x}^2, \mathbf{x}^3, \dots, \mathbf{x}^n\}$ find patterns.
 - The description of a pattern may come in the form of a function (say, $g(\mathbf{x})$).
- Supervised:
 - Given data $\{(\mathbf{x}^1, \mathbf{y}^1), (\mathbf{x}^2, \mathbf{y}^2), (\mathbf{x}^3, \mathbf{y}^3), \dots, (\mathbf{x}^n, \mathbf{y}^n)\}$ find a function such that $f(\mathbf{x}) = \mathbf{y}$.
 - \mathbf{y} are the labels.

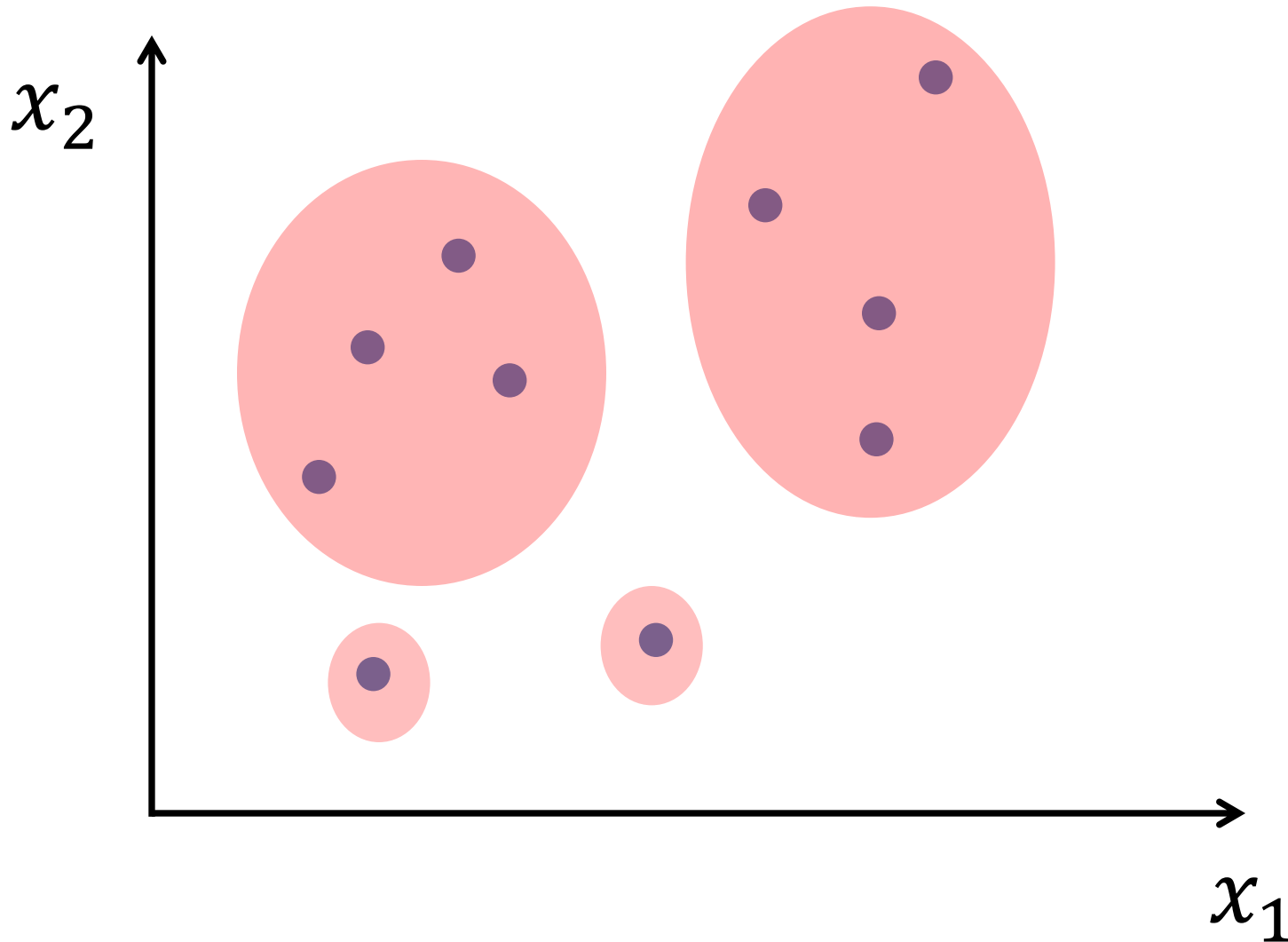
Some Areas in Unsupervised Learning

- Density Estimation
 - Kernel Density Estimation
 - Mixture of Gaussians
- Clustering
 - K Nearest Neighbor
- Dimension Reduction
 - Principal Component Analysis (Linear Technique)
 - Locally Linear Embedding (Non-Linear Technique)

Unsupervised Learning: Find Patterns in the Unlabeled Data



Unsupervised Learning: Find Patterns in the Unlabeled Data



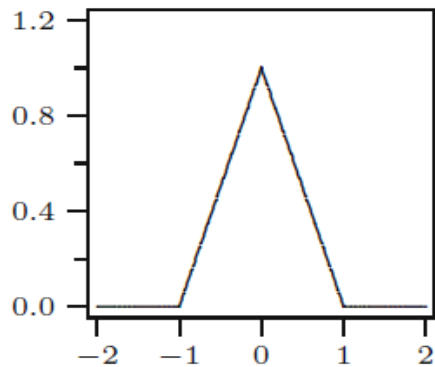
Basis Function

- Basis function allows us to parameterize the data in some handy or/and meaningful ways.
- The decomposition of a data point into basis function is usually quick.
- Examples:
 - Gaussian functions
 - Kernel functions
 - Eigenfunctions (such as “Eigenspectra” in Astronomical Spectroscopy)

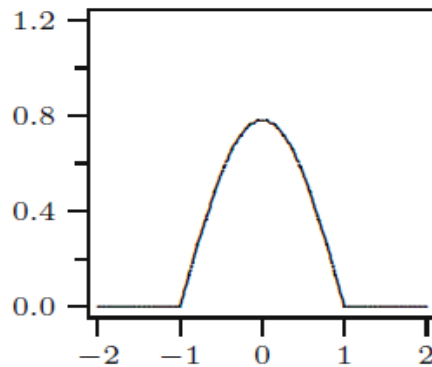
Density Estimation: Kernel Density Estimation

- Approximate a distribution by the sum of kernels (K's) of various bandwidths (h).
- Properties of kernel:
 - *K is a probability density function, $K(u) > 0$ and $\int_{-\infty}^{\infty} K(u)du = 1$.*
 - *K is symmetric around zero: $K(-u) = K(u)$*
 - *(Not Always) $K(u) = 0$ for $|u| > 1$*

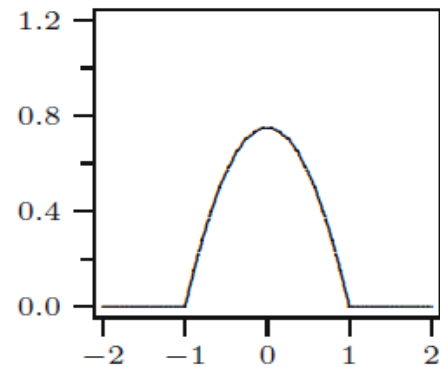
Example Kernel Functions $K(u)$



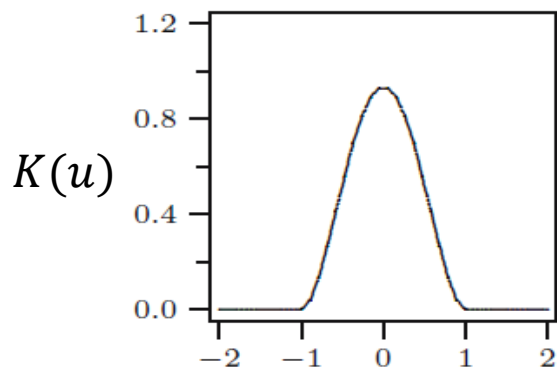
Triangular kernel



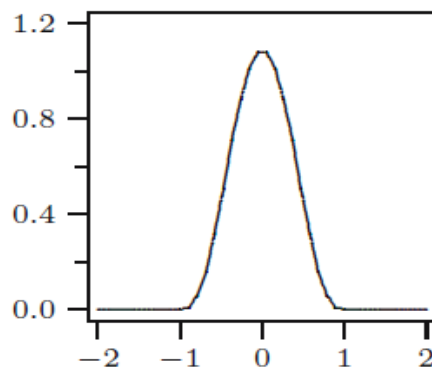
Cosine kernel



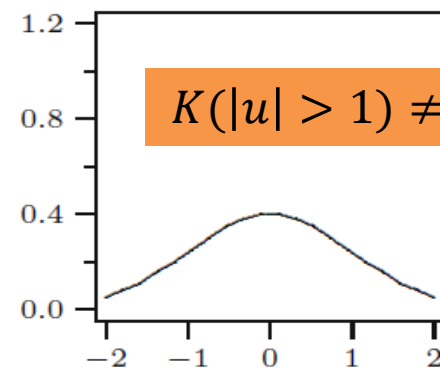
Epanechnikov kernel



Biweight kernel



Triweight kernel

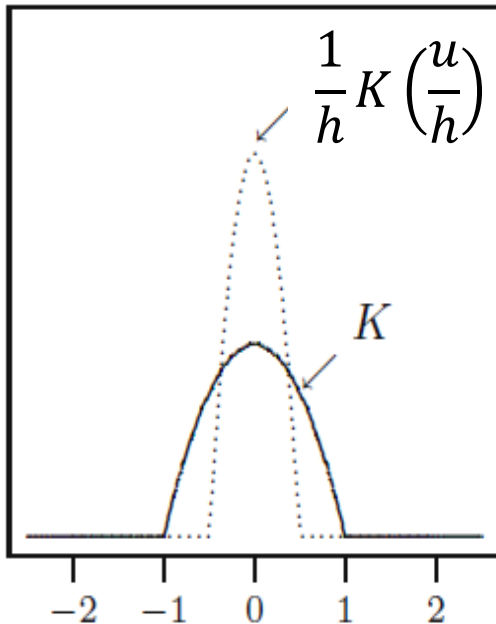


Normal kernel

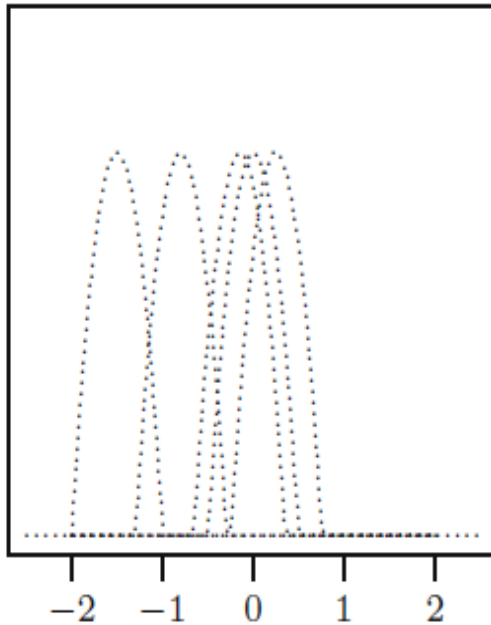
u

(Dekking et al. 2005)

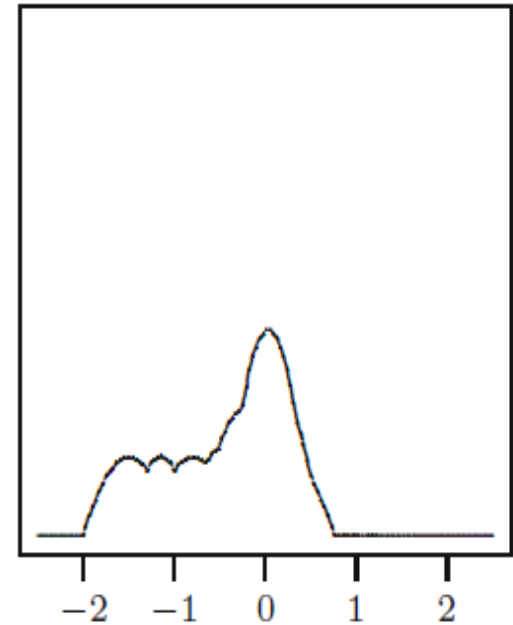
Constructing the Kernel Density Estimate



Kernel and scaled kernel



Shifted kernel

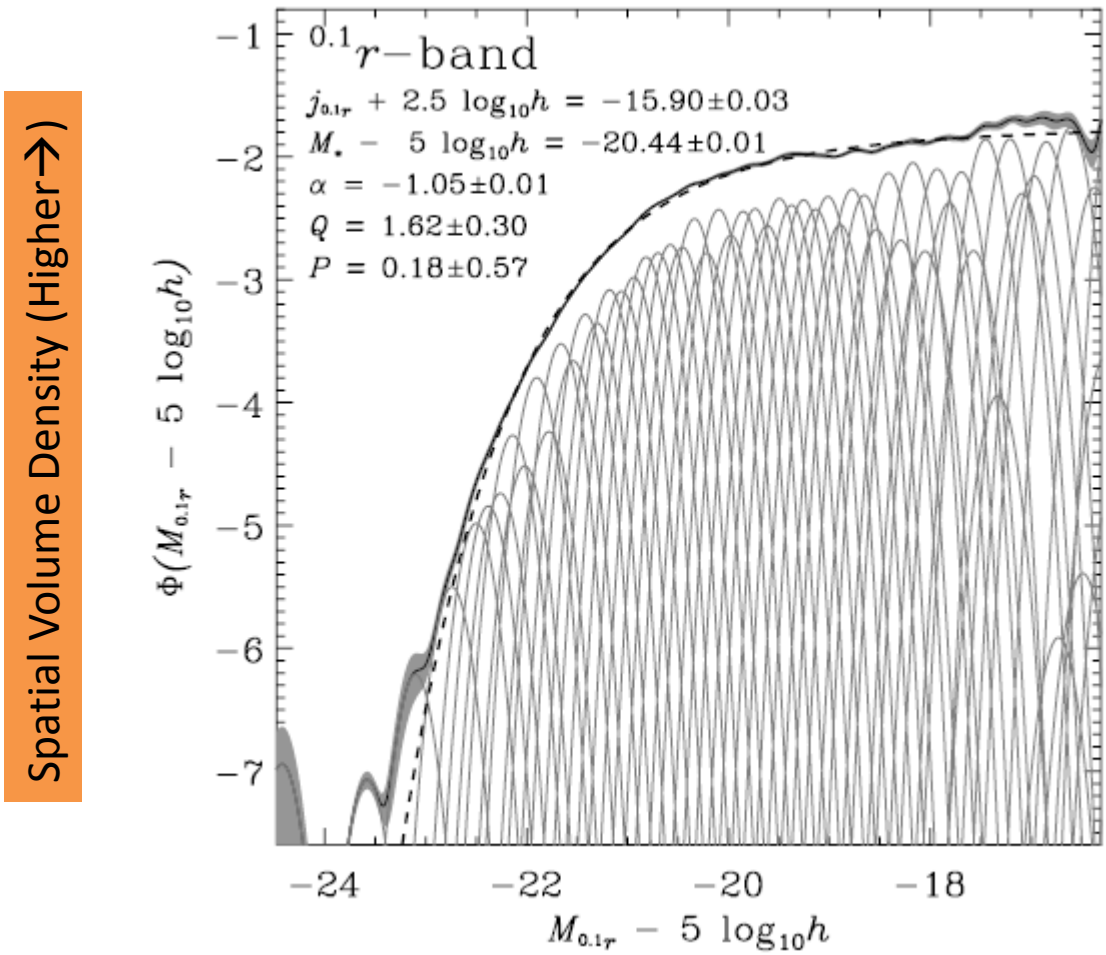


Kernel density estimate

$$f = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{u - x^i}{h}\right)$$

Note $\int_{-\infty}^{\infty} f(u) du = 1$

Luminosity Function of Nearby SDSS Galaxies



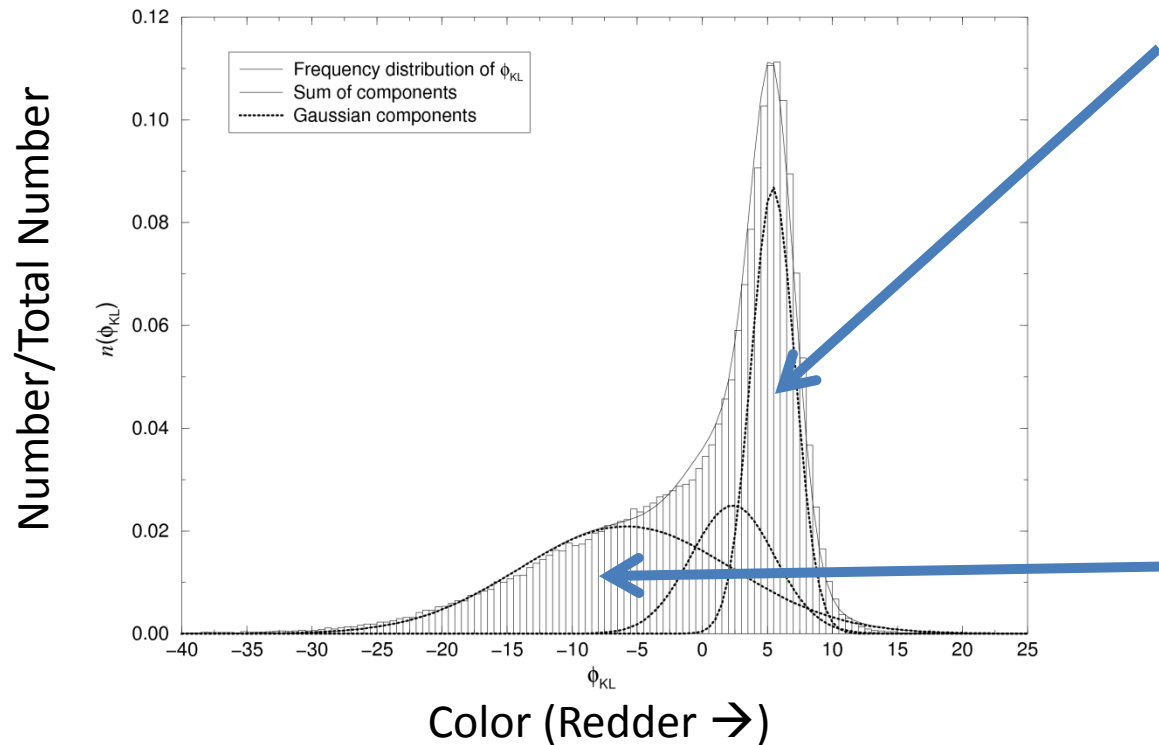
Spatial Volume Density (Higher →)

Notice the height of the kernels actually vary, slightly different from the discussed kernel density estimation.

(Blanton et al. 2003)

Magnitude of Galaxy (Dimmer →)

Color Distribution of Nearby Galaxies



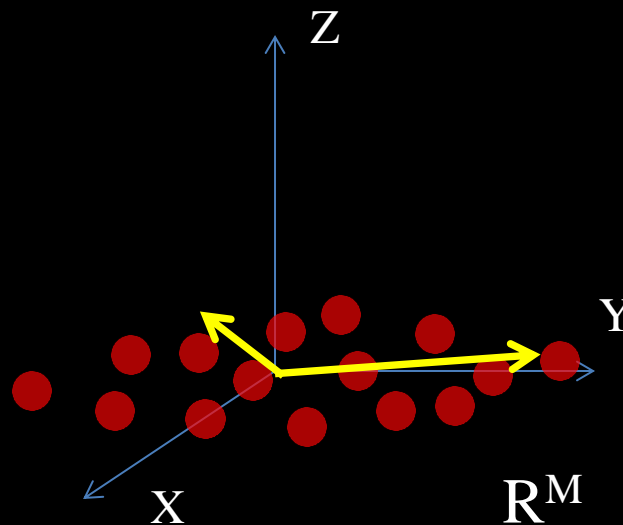
(Yip, Connolly, Szalay, et al. 2004)

Principal Component Analysis

- Perhaps the most used technique to parameterize high-dimensional data.
- Best for linear data distribution.
- Many applications in both academia and industry: dimension reduction, data parameterization, classification problems, image decomposition, audio signal separation, ..., etc.

High-Dimensional Data may lie in Lower-Dimensional Space (or Manifold)

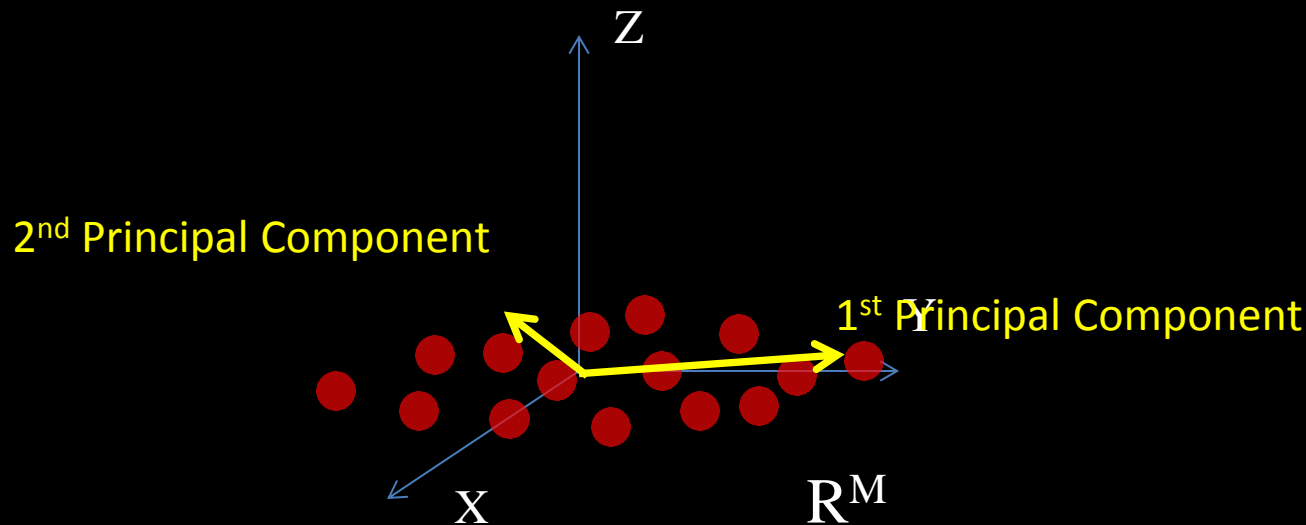
N points, M components



PCA finds the orthogonal directions in the data space which encapsulate maximum sample variances.

High-Dimensional Data may lie in Lower-Dimensional Space (or Manifold)

N points, M components

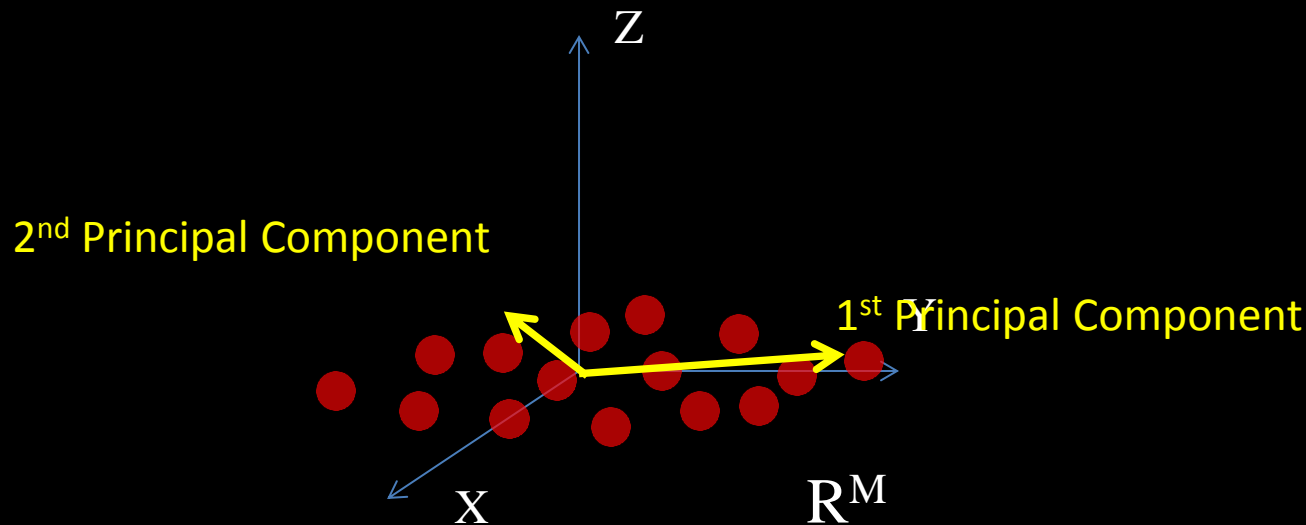


PCA finds the orthogonal directions in the data space which encapsulate maximum sample variances.

A principal component is also called **eigenvector**, or **eigenfunction**.

High-Dimensional Data may lie in Lower-Dimensional Space (or Manifold)

N points, M components



We reduce the dimension of the problem from M to 2.

PCA finds the orthogonal directions in the data space which encapsulate maximum sample variances.

A principal component is also called **eigenvector**, or **eigenfunction**.

Properties of PCA

- PCA decomposes the data into a set of **eigenfunctions**.
- The eigenfunctions are orthogonal (perpendicular) to each other:
 - The dot product between two eigenvectors is zero.
- The set of eigenfunctions **maximize the sample variance** of the data. Therefore, the data can be decomposed into a handful of eigenfunctions (for linear data distribution).
- For non-linear data, PCA may fail (i.e., we need many orders of eigenfunctions).
- In galaxy spectroscopy, the basis functions are called “eigenspectra” (Connolly et al. 1995).

PCA Eigenspectra ($e_{i\lambda}$) Representation of Galaxy Spectra

$$f_\lambda = \sum_i a_i e_{i\lambda}$$

(i runs from 1 to the number of eigenspectra)

Minimize reconstruction error with respect to a_i 's:

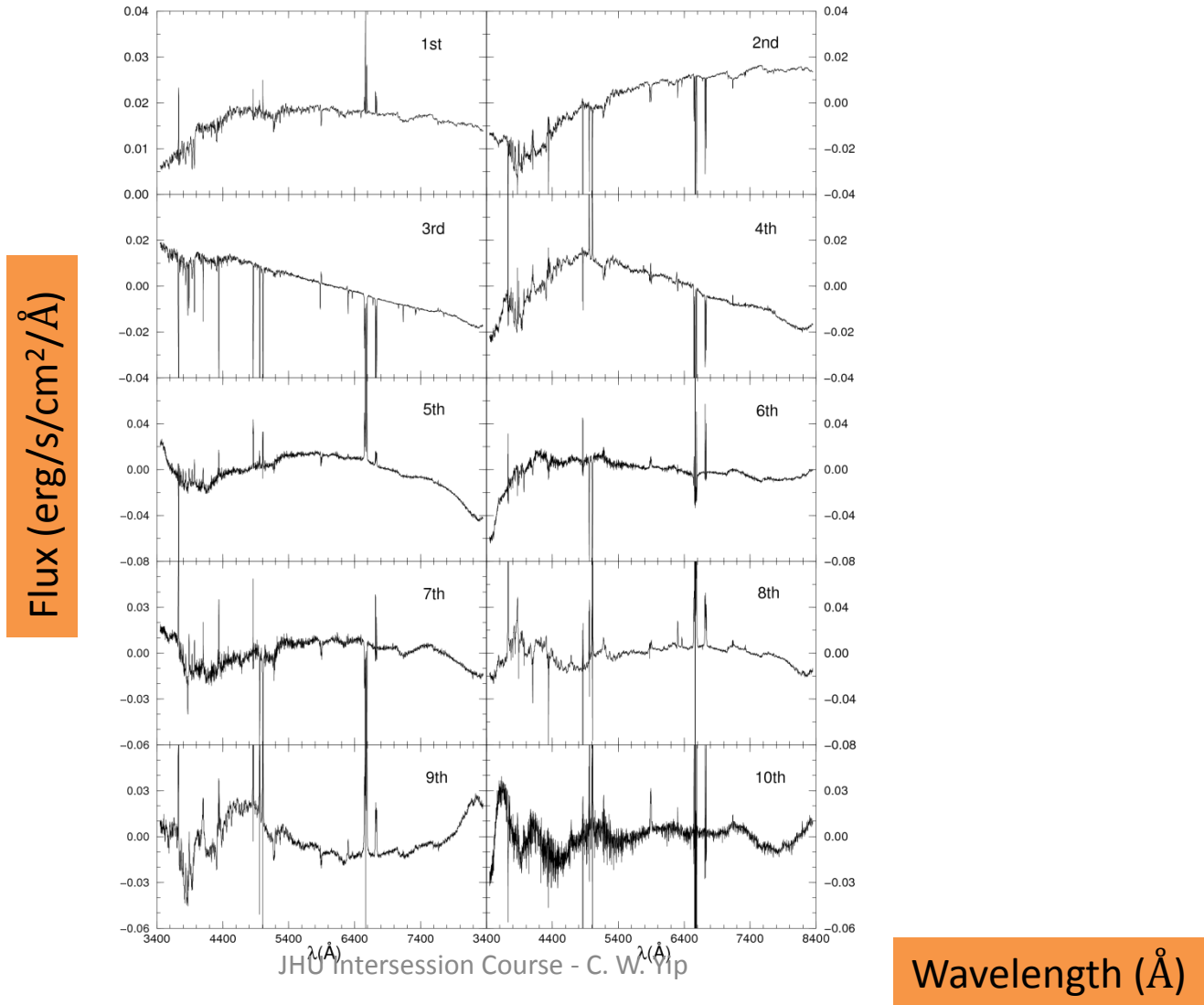
$$x^2 = \sum_\lambda w_\lambda (f_\lambda - \sum_i a_i e_{i\lambda})^2$$

Get the weights for each basis function:

$$a_i = a_i(w_\lambda, e_{i\lambda}, f_\lambda)$$

(Connolly & Szalay 99)

Eigenspectra of Nearby Galaxies (SDSS: $N = 170,000$; $M = 4000$)

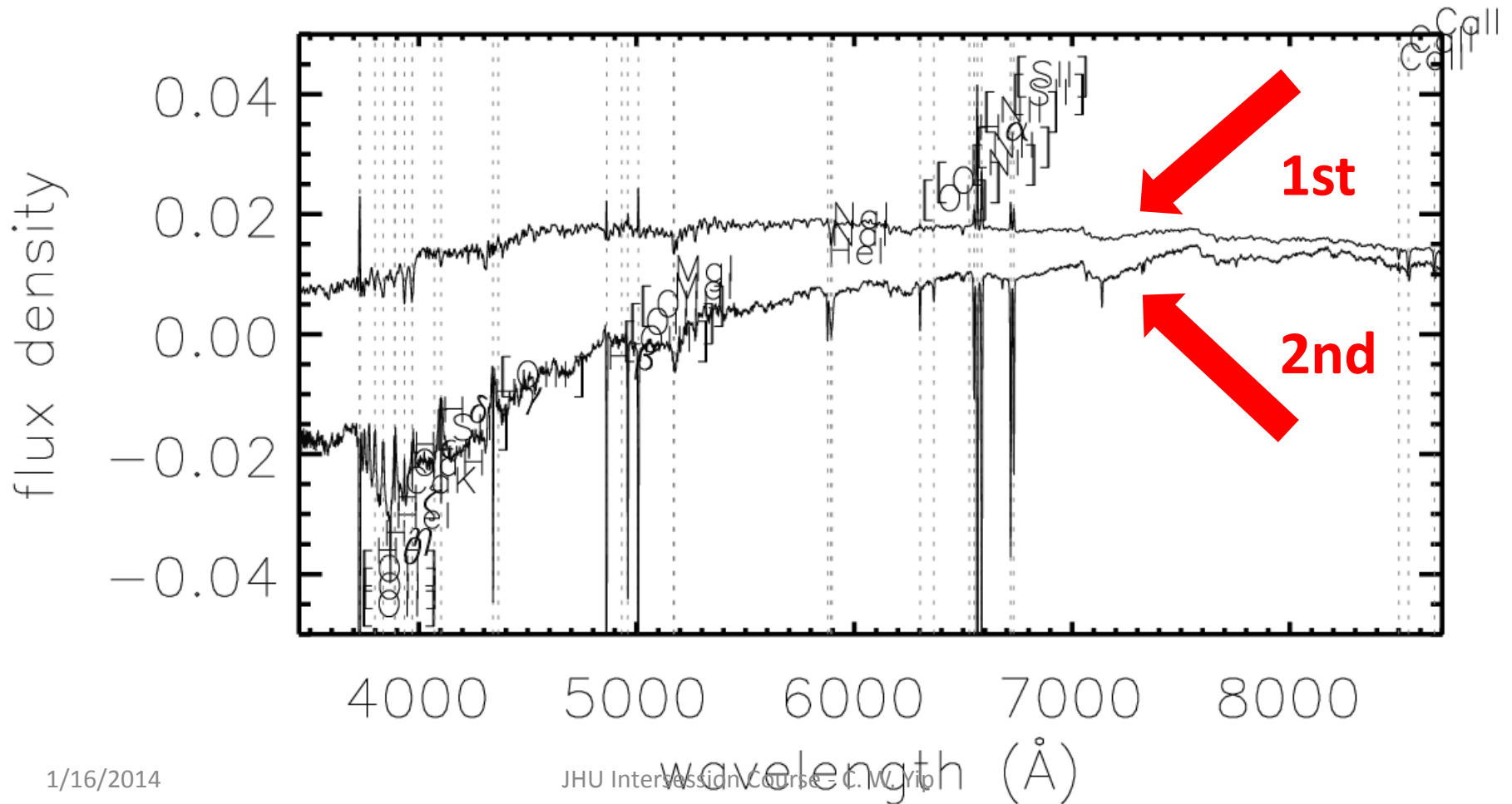


PCA as a way for Signal Separation

- The basis functions has physical meanings. Hence each order (or called “mode”) of function can be considered as a **signal**.
- In PCA, the basis functions are orthogonal. They point to different direction in the data space and are **statistically independent**.

Eigenspectra of Half a Million Galaxy Spectra

- 2nd mode: galaxy type (steepness of the spectral slope)



Eigenspectra of Half a Million Galaxy Spectra

- 3rd mode: post-starburst activities (stronger the absorptions, weaker the H α)

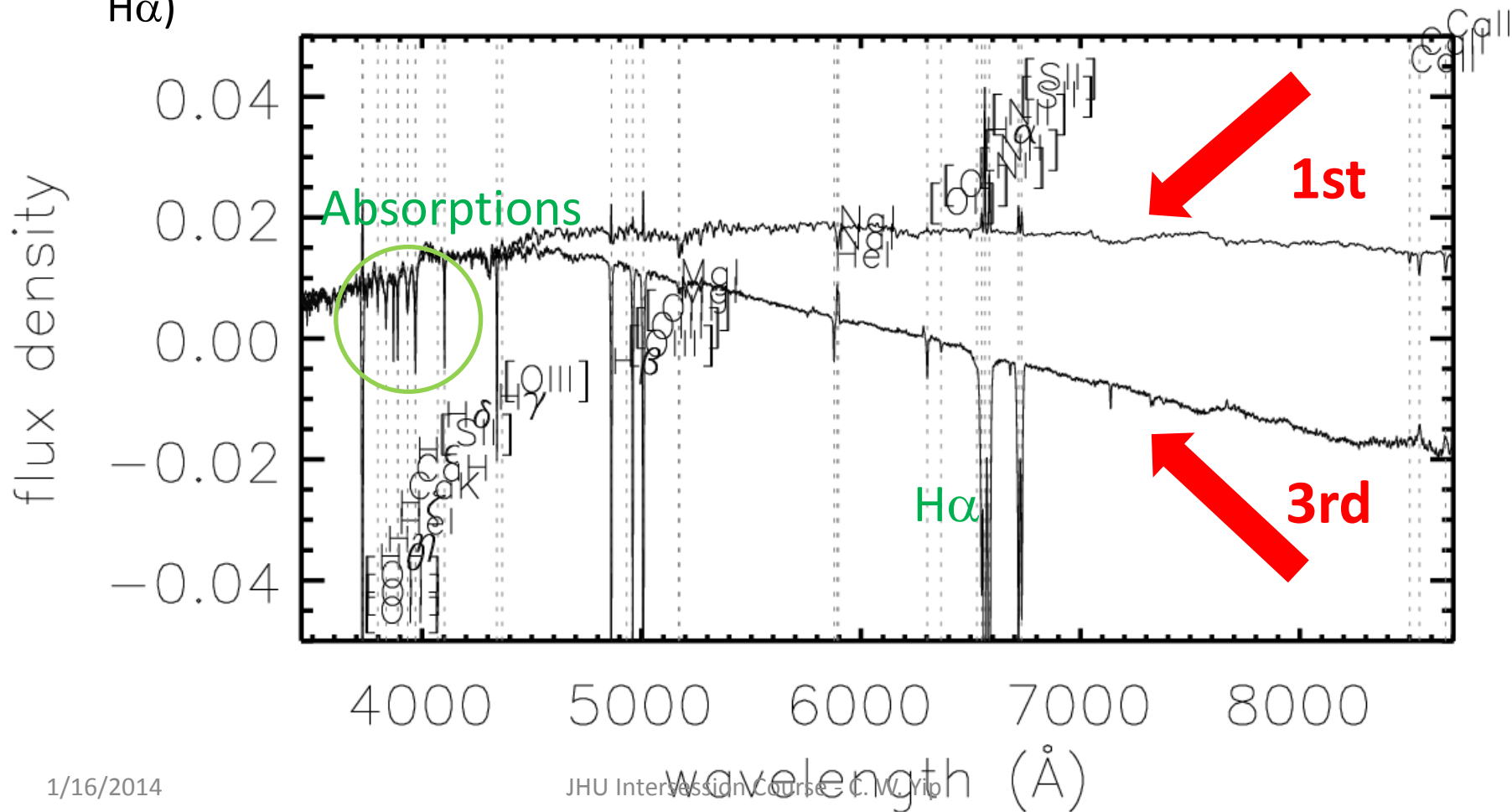


Image Reconstruction

Original

Reconstruction using 50 Eigenfaces



(Everson &
Sirovich 1995)

Matrix Representation of Data

- Many datasets are made up with N objects and M variables.
- Matrix provides a handy way to represent the data.
- An added advantage is that many algorithms can be expressed conveniently in matrix form.

Example:

Size $M \times N$ Data Matrix for Galaxy Spectra

$A_{ij} = \text{Flux at the } i\text{th wavelength for the } j\text{th galaxy}$

$$A = \begin{matrix} & & & & \text{Galaxy ID} \\ & & & & \\ & & & & \\ & & & & \\ \text{Wavelength} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix} \begin{matrix} A_{11} & A_{12} & A_{13} & \cdots & A_{1N} \\ A_{21} & A_{22} & A_{23} & \cdots & A_{2N} \\ \vdots & \ddots & & & \vdots \\ A_{M1} & A_{M2} & A_{M3} & \cdots & A_{MN} \end{matrix}$$

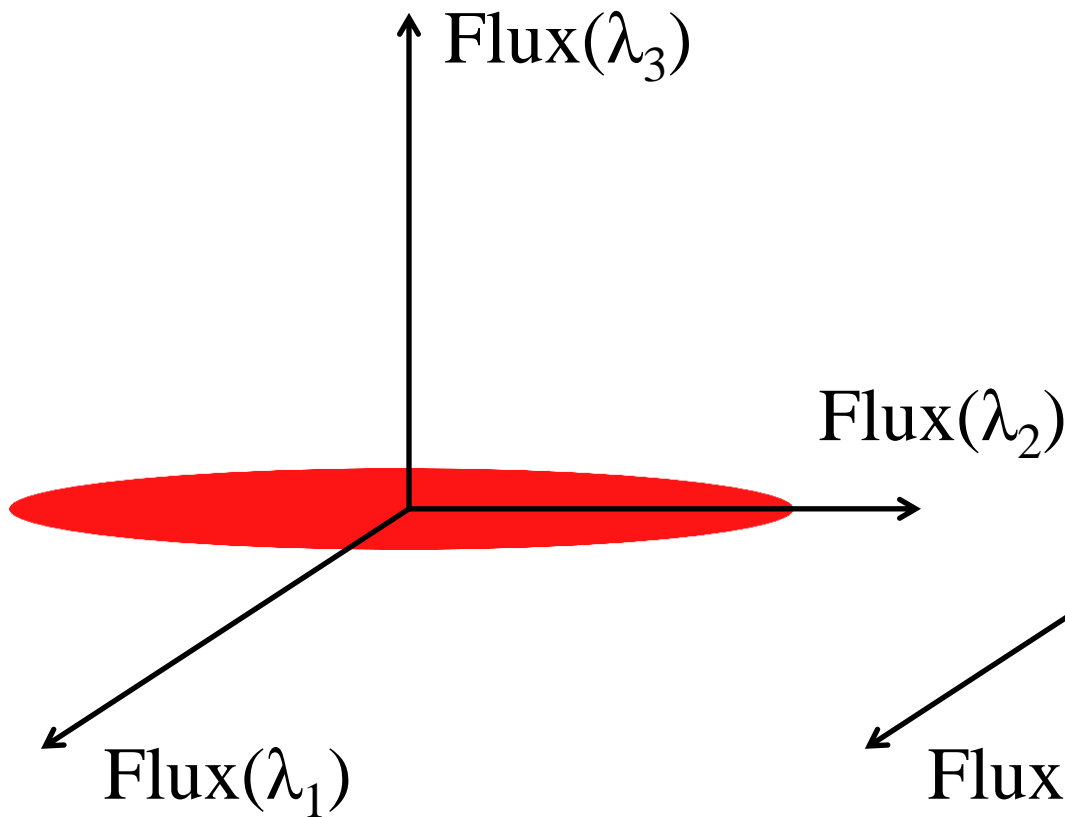
* PCA in matrix notation: see handout.

Compressing the Object Space: CUR Matrix Decomposition

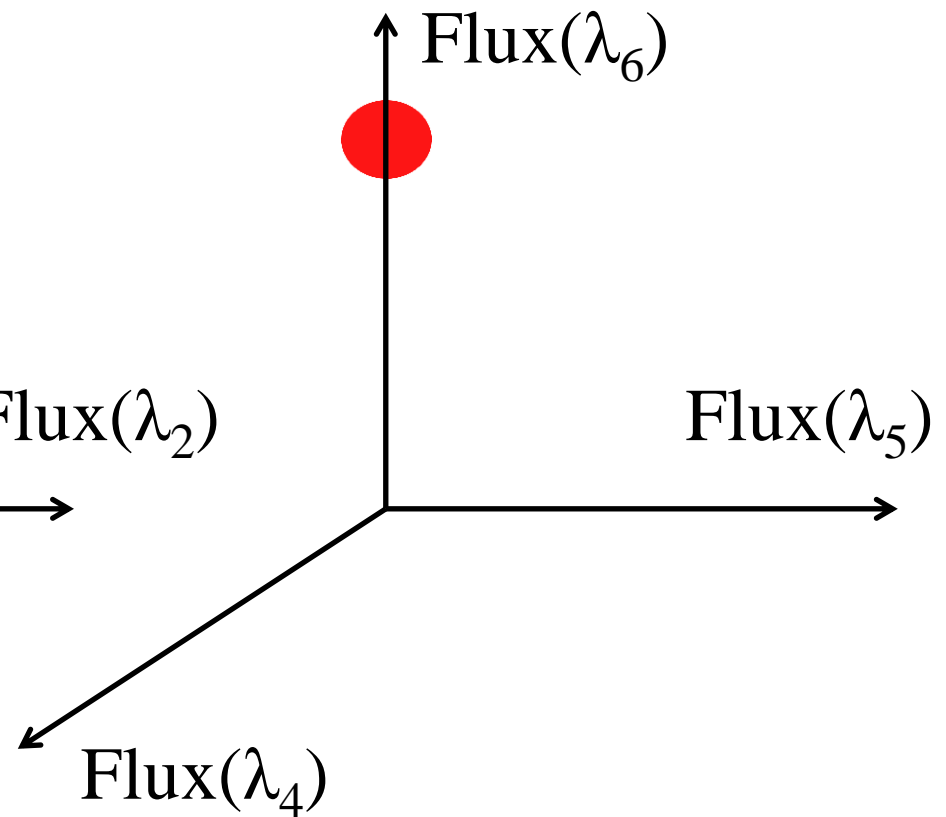
- In PCA, the number of components in the data vectors remain intact.
- CUR Matrix Decomposition provides a dramatically new way to compress big data. This approach compresses the variable space:
 - The number of components in the data vectors decreases.

Some Wavelengths are More Informative: Leverage Score per Each Variable

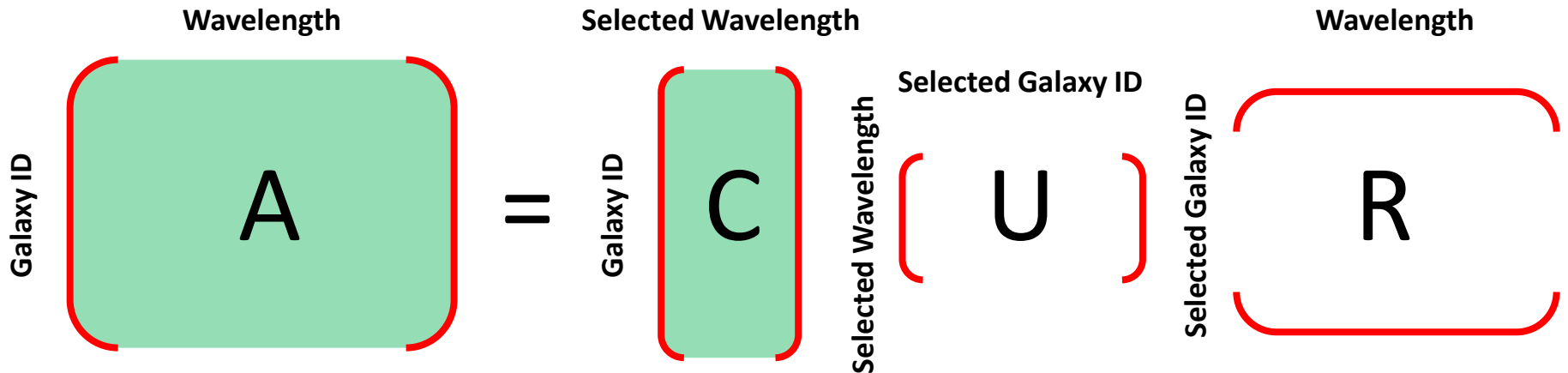
Sample Variance



Sparsity



CUR Matrix Decomposition (Mahoney & Drineas 2009)



CUR approximates data matrix:

$$\min || A - CUR ||_F$$

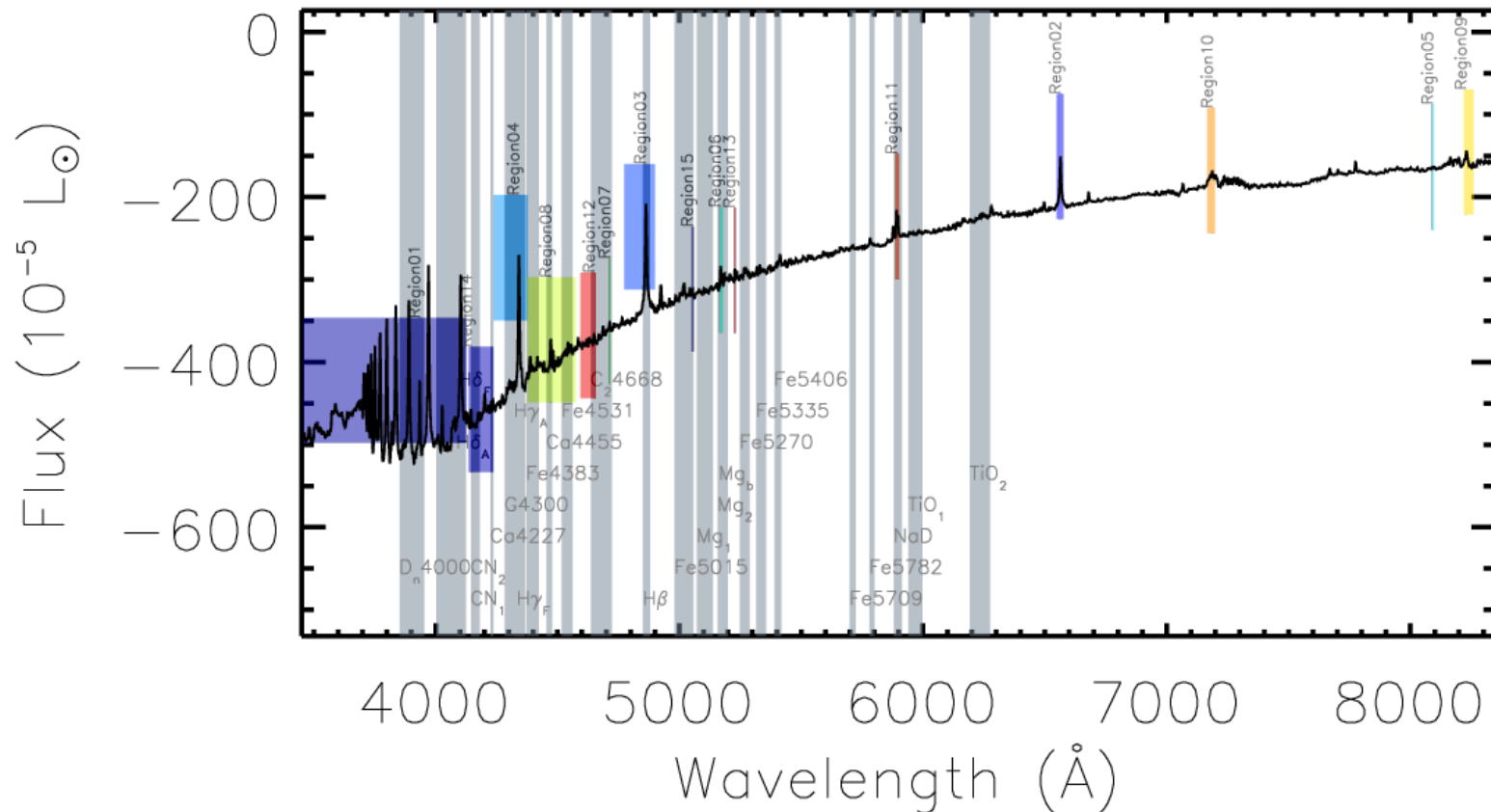
Frobenius Norm

- A matrix norm is a number for representing the amplitude of a matrix.
- The Frobenius norm is defined as follows:

$$\|A\|_F = \sqrt{\sum_{i=1}^M \sum_{j=1}^N |A_{ij}|^2}$$

- In CUR Matrix Decomposition, the matrix A is the difference between the data matrix and its approximated matrix. The Frobenius norm therefore measures the “distance” between the two matrices.

Find Important Regions in Multi-Dimensional Data: Galaxy Spectra



Discussion of Jan 7 Homework

- Total number of pixels in the CCD = 1024×1024 pixels = 1,048,576 pixels \sim 1 MPixels.
- GAIA needs $10^9/60/60/24/365$ years = 31.7 years = 31 years 8 months \sim 32 years.