

# An Introduction to Unsupervised Learning

A. Gilad Kusne

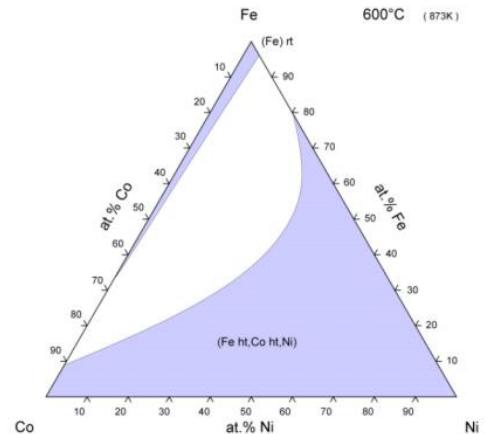
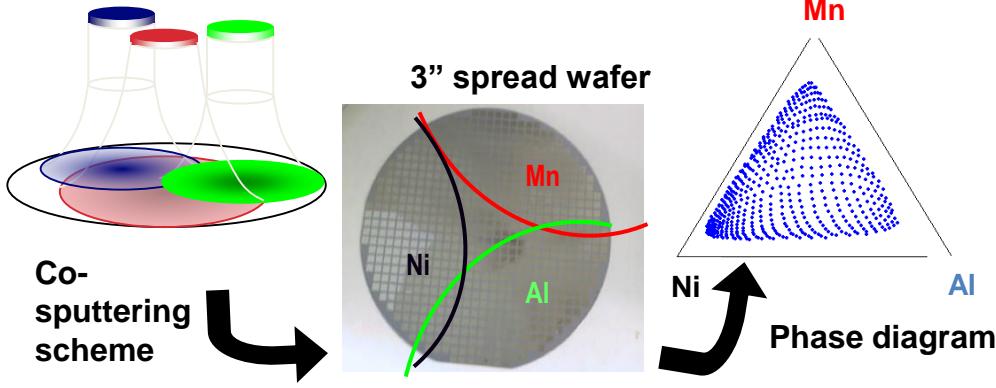
aaron.kusne@nist.gov

# Additional Information

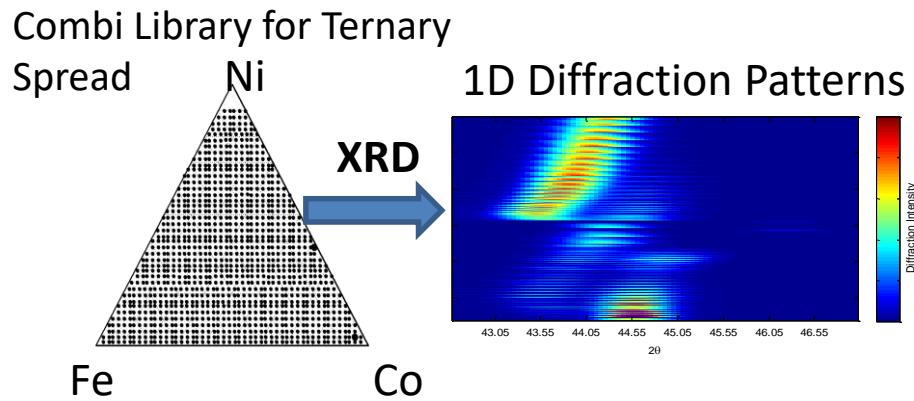
- Materials Informatics:
  - Mueller, Tim, Aaron Gilad Kusne, and Rampi Ramprasad. "**Machine Learning in Materials Science.**" *Reviews in Computational Chemistry, Volume 29*: 186-273. (2016).
- Machine Learning:
  - James, Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani. ***An introduction to statistical learning.*** 2013.
  - Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. ***The elements of statistical learning.*** 2010.
- **All Free Online!**

# Motivating Challenge: Phase Diagram Determination

Sample synthesis

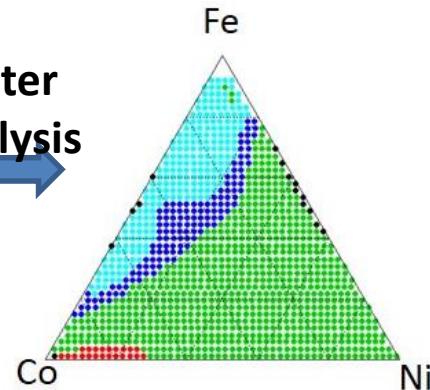


Data Analysis



Clustering Results

Cluster Analysis

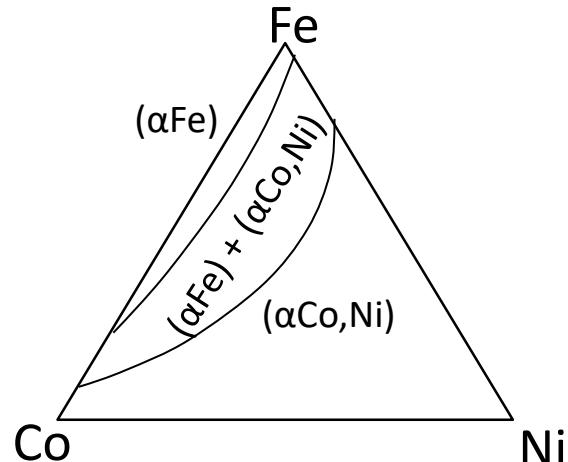


Goal: Mapping ternary structural phase diagram  
From one spread wafer.

# Using AI to Identify Structural Phase Maps

## Structural Phase Map

- Structure as a function of fabrication parameters (e.g. composition, temperature, pressure, etc.)
- Use map to predict structure of new materials.
- Structure is good predictor of important properties.

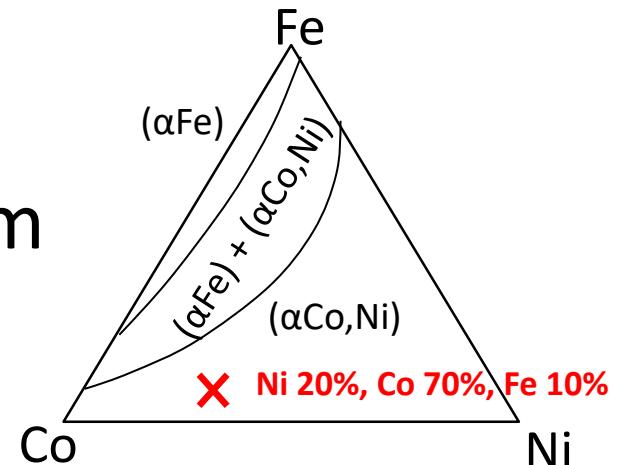


Phase Equilibria in Iron Ternary Alloys (1988) #60

# Structural Phase Mapping: Edisonian Approach

Traditional / Edisonian Approach:

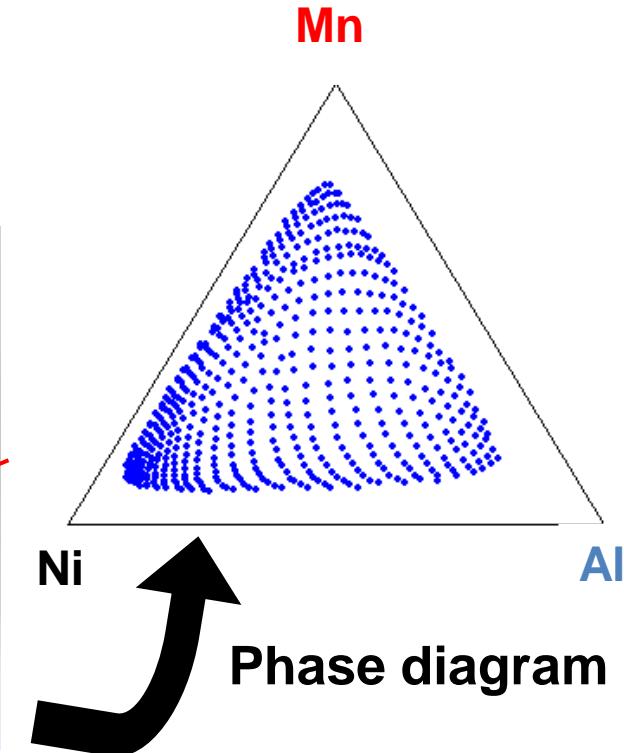
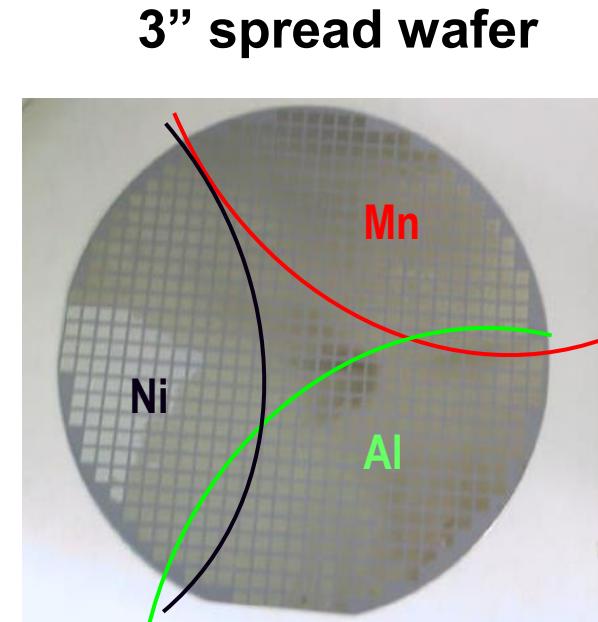
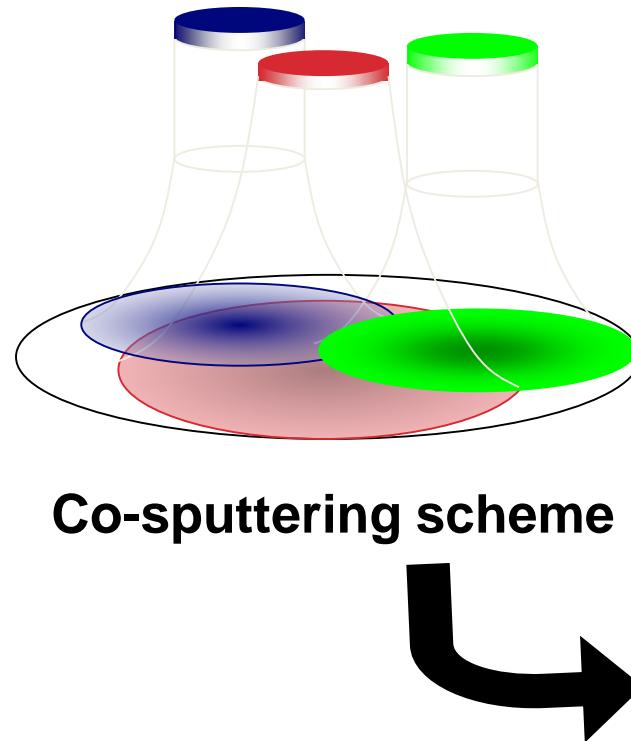
- Fabricate sample
- Measure structure
- Point placed on phase diagram
- Repeat
- This process takes years.



Phase Equilibria in Iron Ternary Alloys (1988) #60



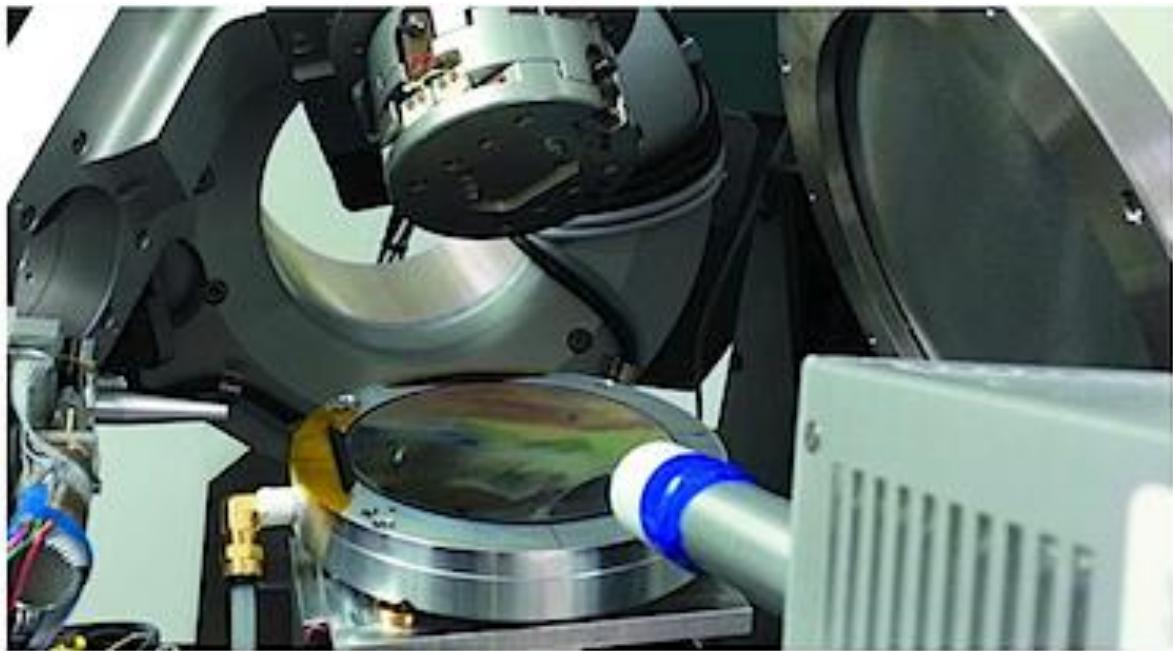
# Composition Spreads of Ternary Metallic Alloy Systems



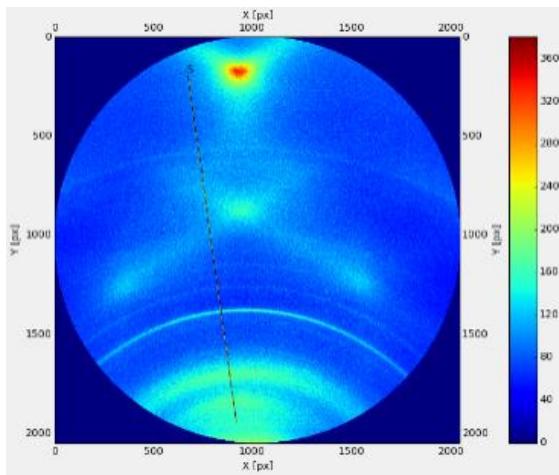
Spreads are used for mapping structural phase diagrams and functionalities (magnetic properties, electronic properties, etc.)

# HTE: Characterization

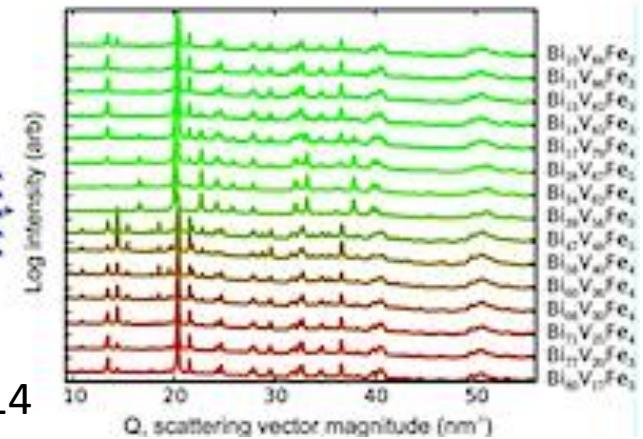
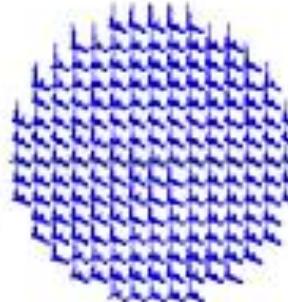
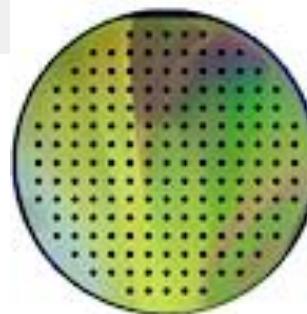
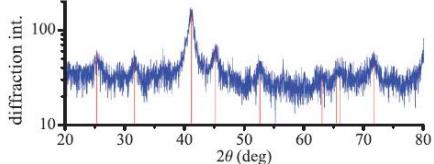
- HiTp X-ray Diffraction



Stanford Synchrotron Radiation Lightsource



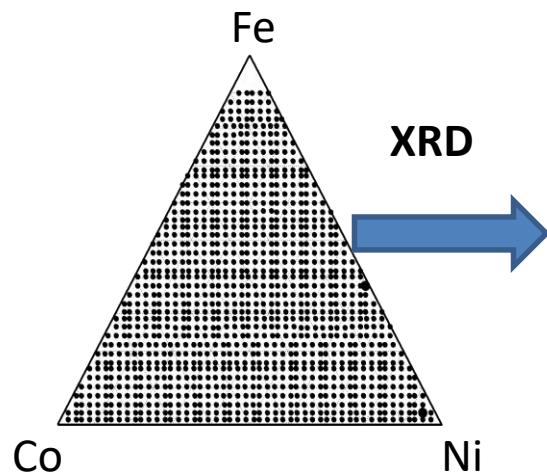
X-ray diffraction image  
converted to vector



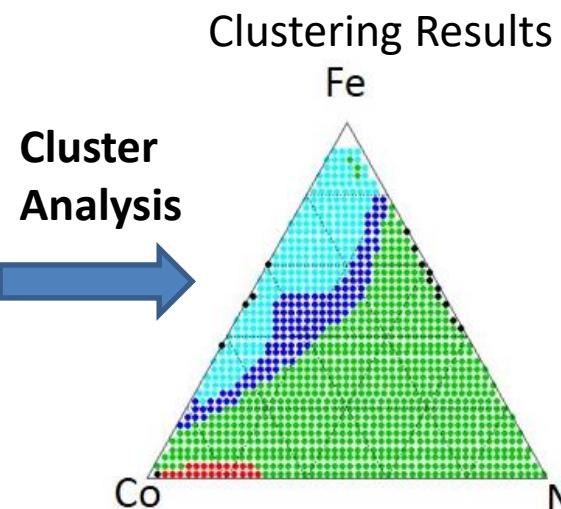
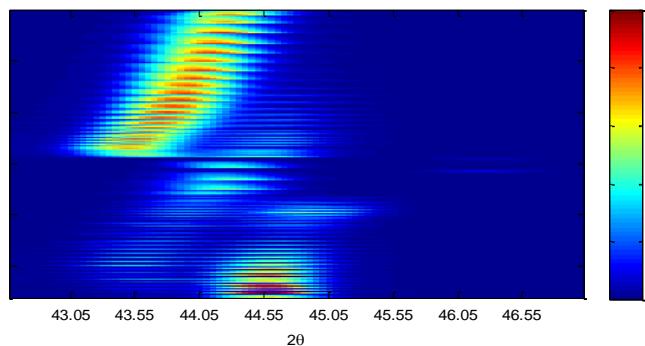
Gregoire, et al. 2014

# High Throughput Phase Mapping

Combi Library for Ternary Spread

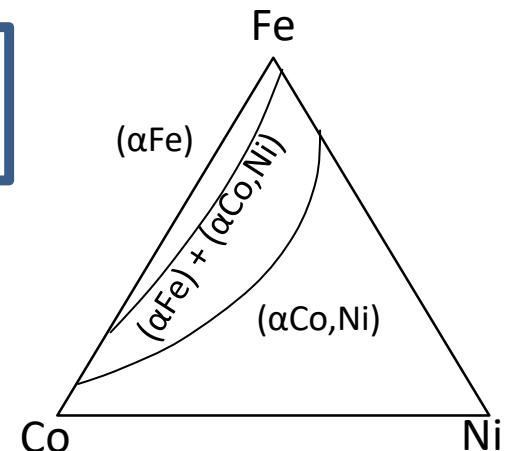


Diffraction Patterns



Goal: Mapping ternary structural phase diagram  
From one spread wafer.

Most ternary systems have not been mapped.



# Autonomous Phase Mapping

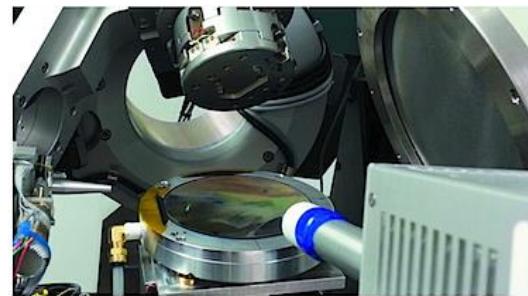
- Measurement is a time / resource sink
- For wafer of 500+ samples:
- In Lab: Takes weeks-months
- Synchrotron: Takes 5+ Hours (Every second counts)



Mn-Ni-Ge library  
535 samples

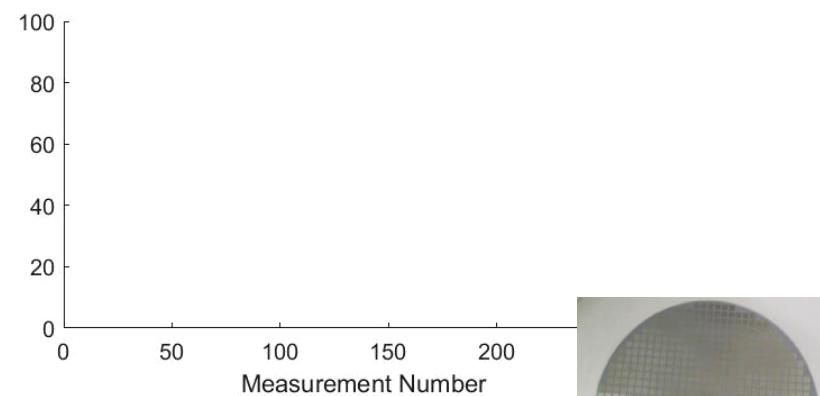
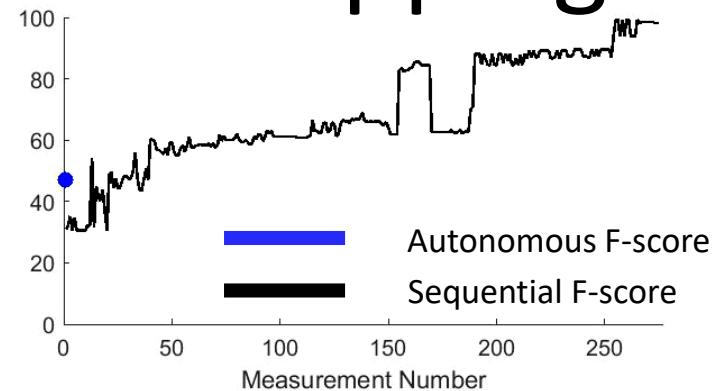
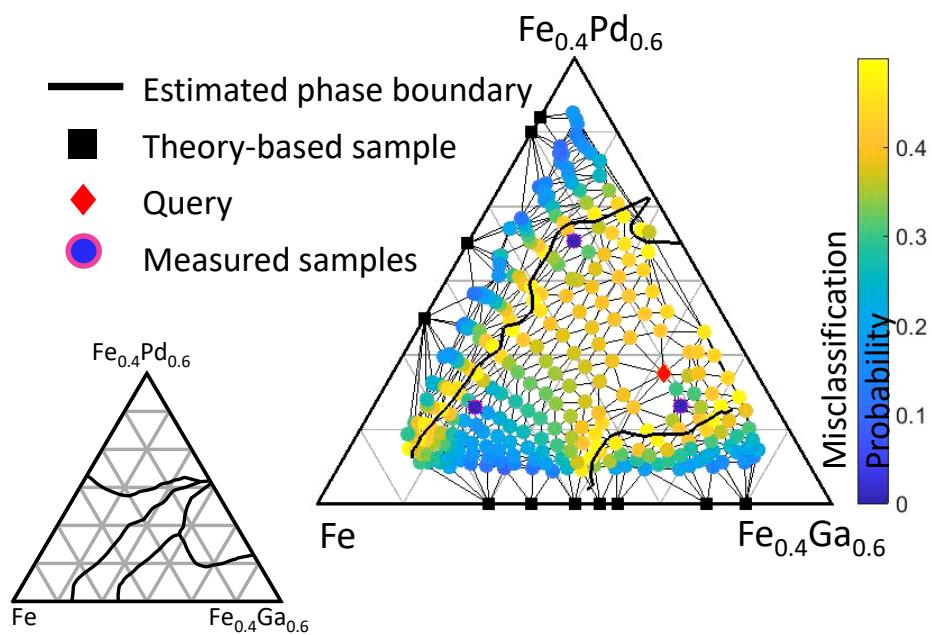


Bruker D8  
30 Minutes per sample  
2 weeks!



Stanford Synchrotron Radiation  
Lightsource  
30 seconds per sample  
4.5 hours

# Autonomous Phase Mapping



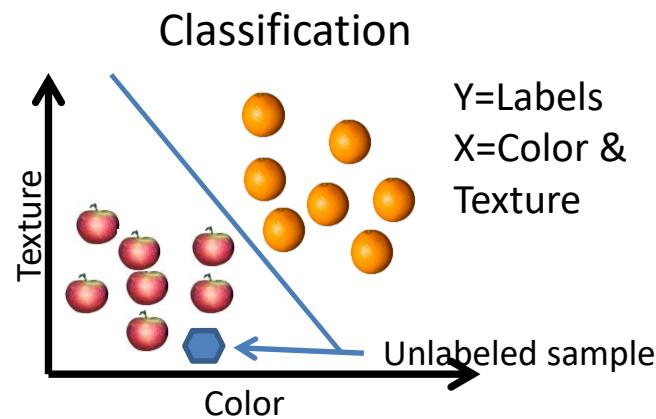
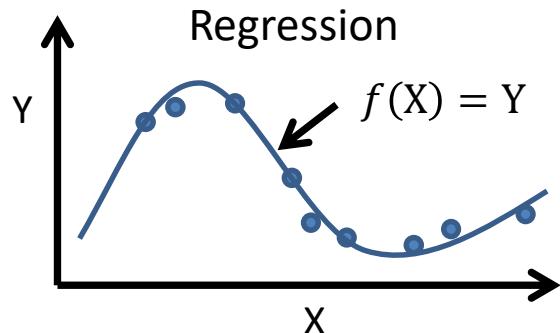
AI is controlling X-ray diffraction systems at SLAC & in the lab!



# Supervised Learning

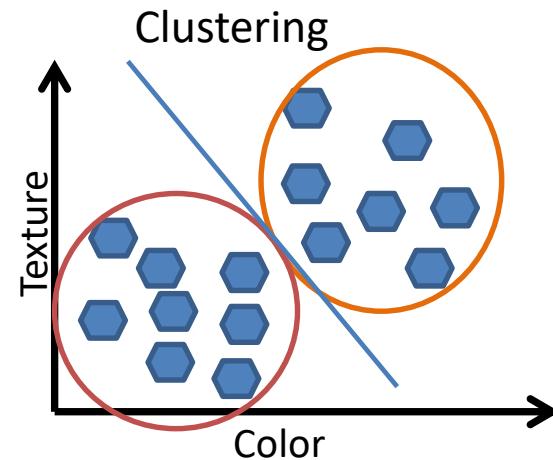
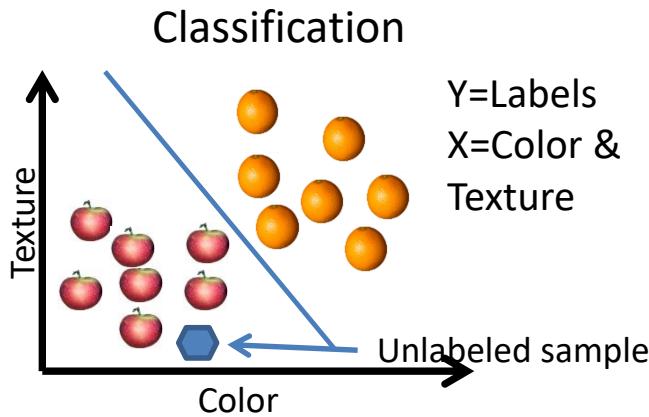
- Supervised Learning
  - Find the function  $f$  that maps the input data  $x$  to the output data  $y$ .  
 $f: x \rightarrow y$

- $y$  is continuous: Regression
- $y$  is discrete: Classification
- Conditional density  $P(f|x, y)$
- Cross-validation to check performance of supervised learning methods & determine parameters.



# Unsupervised Learning

- Unsupervised Learning
  - “Unsupervised”: We don’t have output data  $y$ .
  - Interested in relationship between the data  $x$ .
  - Learn about  $x$  from its distribution.
  - Joint marginal density  $P(x)$
  - Cross-validation for algorithm performance isn’t available.
    - Performance checked with: Heuristics & Expert analysis



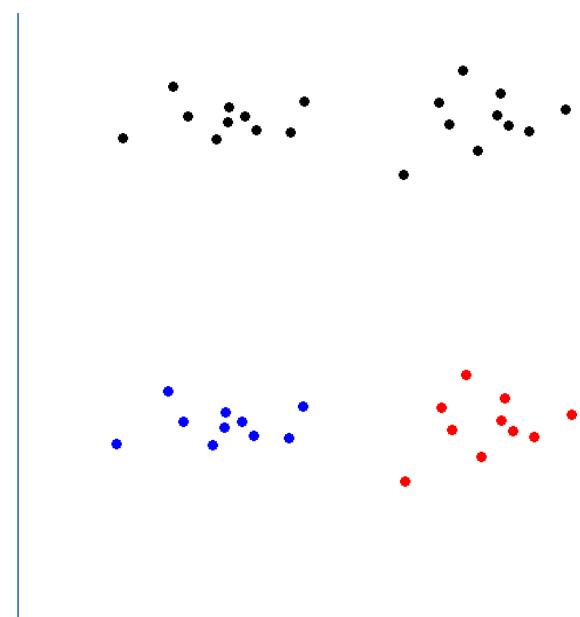
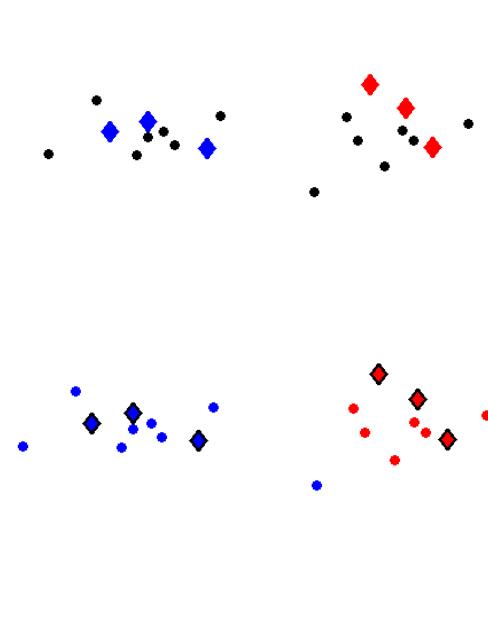
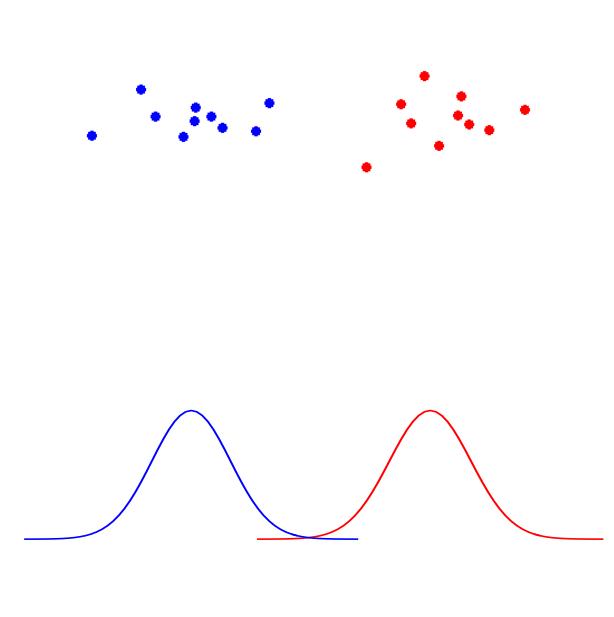
# Unsupervised Learning

- Example

Data from two Gaussian distributions, labelled by color

Classification using SVM

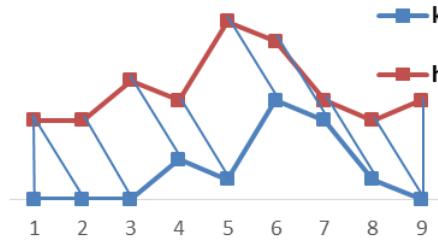
Clustering using k-means



# Unsupervised Learning

- Dissimilarity Measures

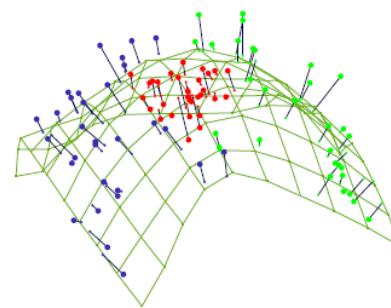
- Define difference between sample



- Spectral Unmixing

- Constituent components

$$\begin{bmatrix} \text{Spectral Data} \\ \vdots \end{bmatrix} = \begin{bmatrix} \text{Spectral Components} \\ \vdots \end{bmatrix} * \begin{bmatrix} \text{Abundance Coefficients} \\ \vdots \end{bmatrix} + \begin{bmatrix} \text{Noise} \\ \vdots \end{bmatrix}$$

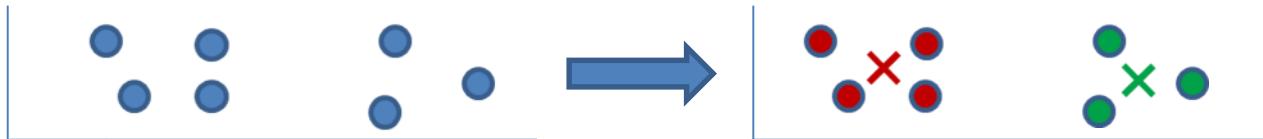


- Latent Variable Analysis

- Visualize high dimensional data
  - Smoothing, outlier detection, etc.

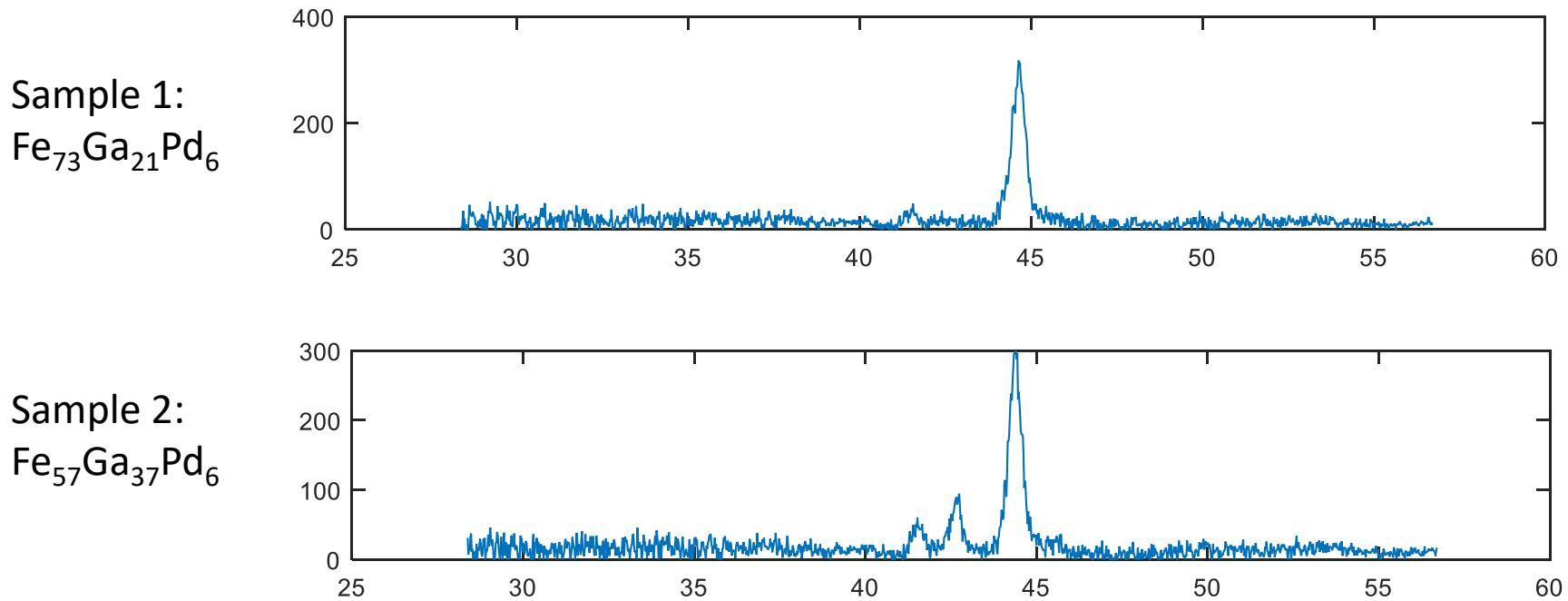
- Clustering

- Group samples into phase regions.



# Dissimilarity Measures (aka kernels)

- Quantify the dissimilarity or difference between two samples,  $d(x_1, x_2) \leftrightarrow K(x_1, x_2)$ .
- E.g. comparing sample structure via x-ray or micrographs
- "Specifying an appropriate dissimilarity measure is far more important in obtaining success with clustering than choice of clustering algorithm." – Hastie, et al.
- **Need domain knowledge!**



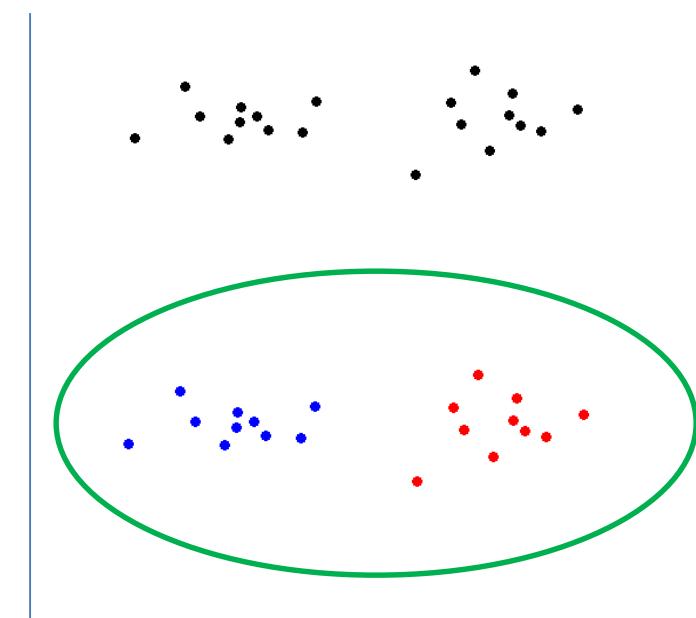
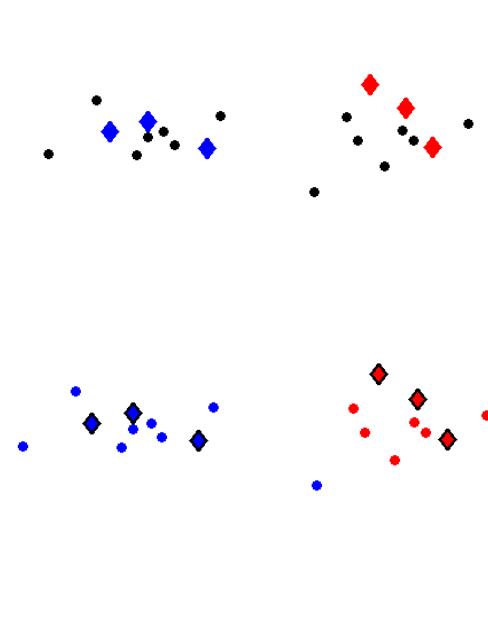
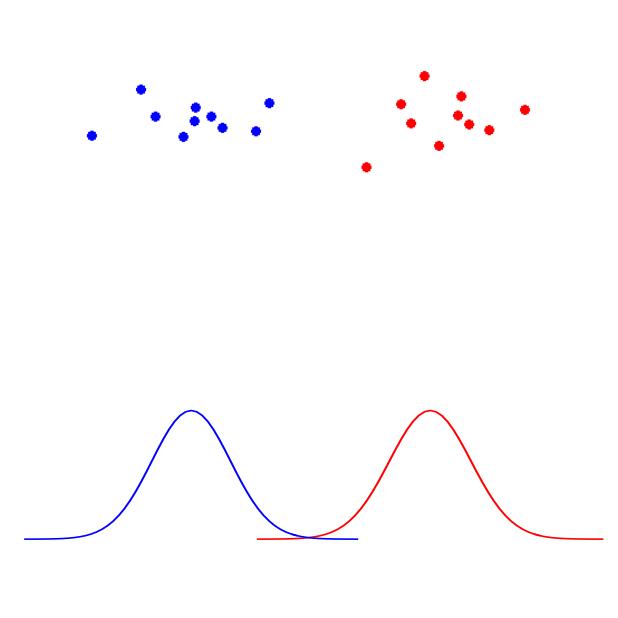
# Unsupervised Learning

- Example

Data from two Gaussian distributions, labelled by color

Classification using SVM

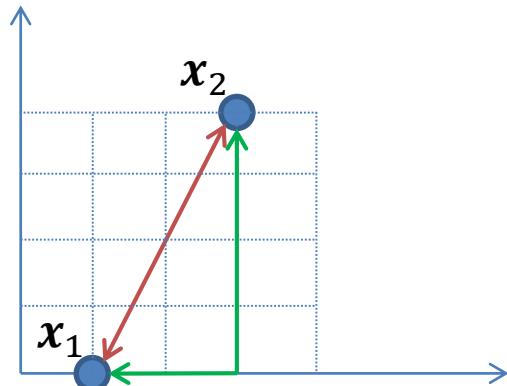
Clustering using k-means



# Dissimilarity Measures (aka kernels)

- Defining the similarity or difference between two samples,  
 $d(x_1, x_2) \leftrightarrow K(x_1, x_2)$ .
- E.g. comparing sample structure via x-ray or micrographs
- 1 – 3 dimensions:

$$d(x_1, x_2) = \sqrt{(x_{11} - x_{21})^2 + (x_{12} - x_{22})^2}$$



$$\begin{aligned}x_1 &= (x_{11}, x_{12}) \\x_2 &= (x_{21}, x_{22})\end{aligned}$$

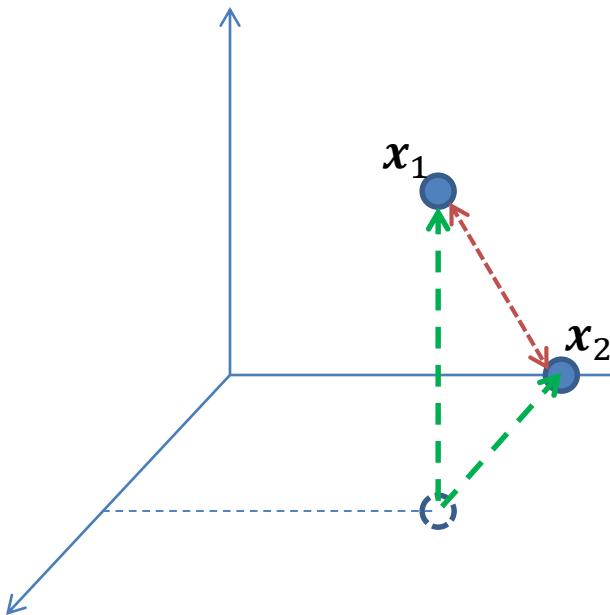
$$\text{Euclidean(L2): } d_E(x_1, x_2) = \left[ \sum_i (x_{1i} - x_{2i})^2 \right]^{1/2}$$

$$\begin{aligned}\text{Taxi-cab(L1): } d_T(x_1, x_2) &= \sum_i |x_{1i} - x_{2i}| \\&\text{City-block}\end{aligned}$$

$$\text{p-norm: } d_p(x_1, x_2) = \left[ \sum_i |x_{1i} - x_{2i}|^p \right]^{1/p}$$

# Measures (aka kernels)

- Defining the similarity or difference between two samples,  
 $d(\mathbf{x}_1, \mathbf{x}_2) \leftrightarrow K(\mathbf{x}_1, \mathbf{x}_2)$ .
  - E.g. comparing sample structure via x-ray or micrographs
  - 1 – 3 dimensions:
- $$d(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{(x_{11} - x_{21})^2 + (x_{12} - x_{22})^2}$$



$$\text{Euclidean(L2): } d_E(\mathbf{x}_1, \mathbf{x}_2) = \left[ \sum_i (x_{1i} - x_{2i})^2 \right]^{1/2}$$

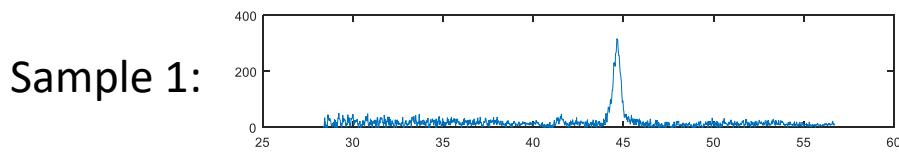
$$\text{Taxi-cab(L1): } d_T(\mathbf{x}_1, \mathbf{x}_2) = \sum_i |x_{1i} - x_{2i}|$$

$$\text{p-norm: } d_p(\mathbf{x}_1, \mathbf{x}_2) = \left[ \sum_i |x_{1i} - x_{2i}|^p \right]^{1/p}$$

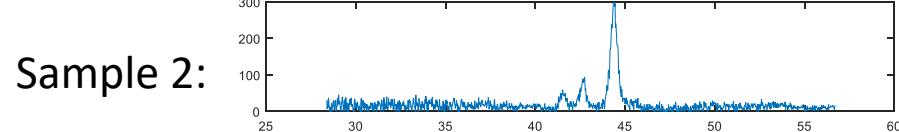
# Measures

- High dimensional data

List of intensity as a function of frequency, angle, location, etc



$$x_1 = [0 1 0 0 0 1 0 0 4 5 7 1 0 1 1 9 8 4 3 1 0 0 0 \dots]$$



$$x_2 = [0 0 0 0 0 1 0 0 0 1 4 5 7 1 0 1 1 9 8 1 0 0 0 \dots]$$

List of intensity as a function of 2 dimensional frequency, angle, location, etc



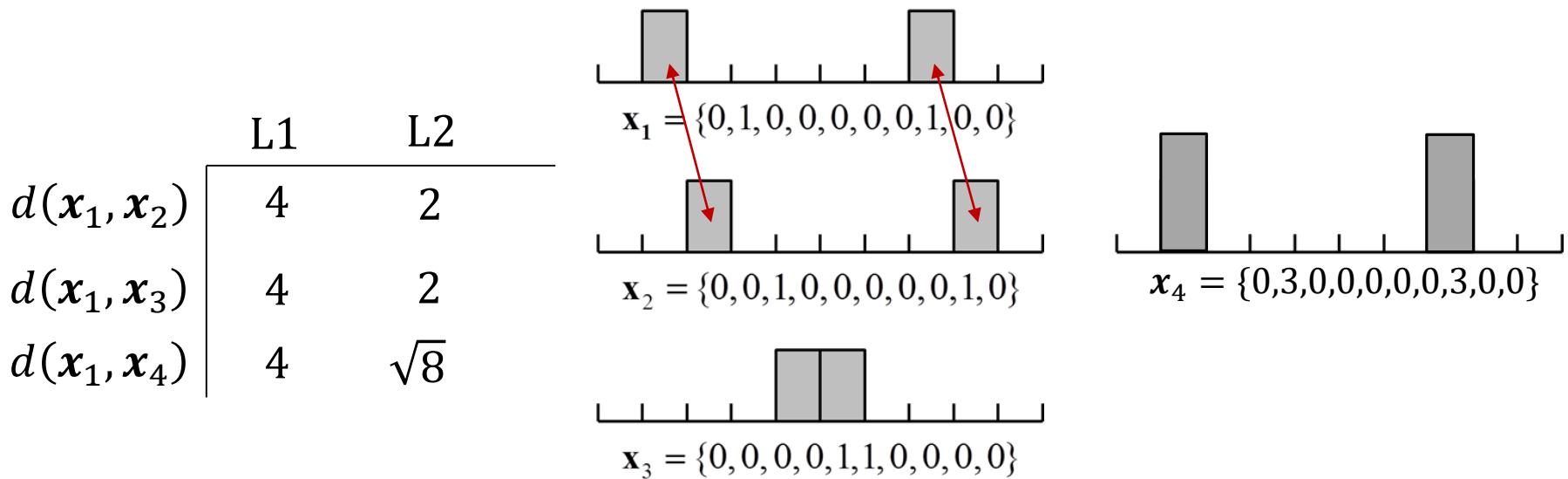
$$x_1 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 4 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$x_1 = [0 0 1 1 1 1 0 1 4 2 0 0 1 1 1 \dots]$$



# Measures

- High dimensional features.
- L1 and L2 do not preserve cross-bin features.
  - P-norm measures are bin-to-bin

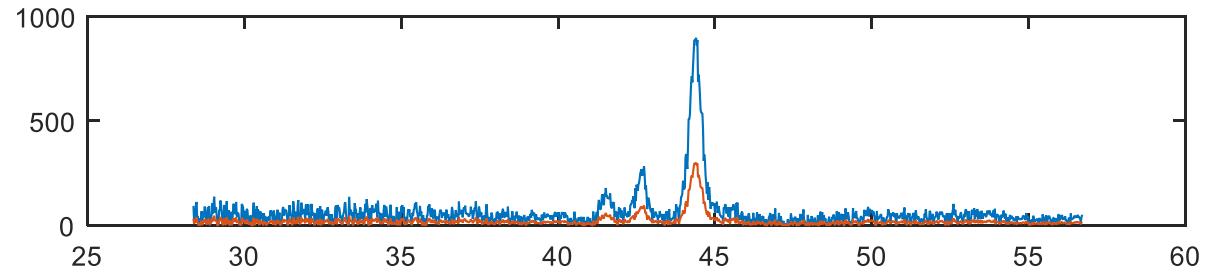


Want a dissimilarity measure that preserves features! (e.g. peaks)

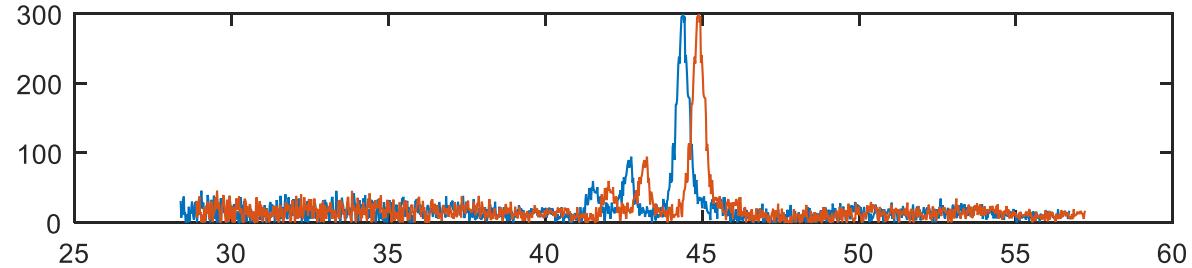
# Measures

- High dimensional & Feature preserving

Scale Invariant:  
(ignores height changes)



Translation / Shift – Invariant:  
(ignores shifts)



Want a dissimilarity measure that preserves features! (e.g. peaks)

# Measures: Background

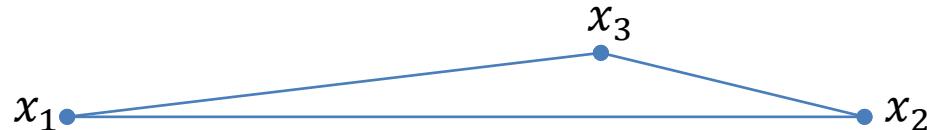
- There are many measures.
  - Types: geometric, statistical, information theoretic, feature preserving, etc.
- Metric

$$d(x_i, x_j) \geq 0, \quad \text{Non-negativity}$$

$$d(x_i, x_j) = 0 \quad \text{iff} \quad x_i = x_j, \quad \text{identity}$$

$$d(x_i, x_j) = d(x_j, x_i), \quad \text{symmetry}$$

$$d(x_i, x_k) \leq d(x_i, x_j) + d(x_j, x_k), \quad \text{Triangle inequality}$$



# Measures: Scale Invariance

- Types: geometric, statistical, information theoretic, feature preserving

Geometric: p-norm

$$d_p = \left[ \sum_i |h_i - k_i|^p \right]^{1/p}$$

$x_i$  and  $x_j$  here replaced with h and k!

Geometric: Cosine

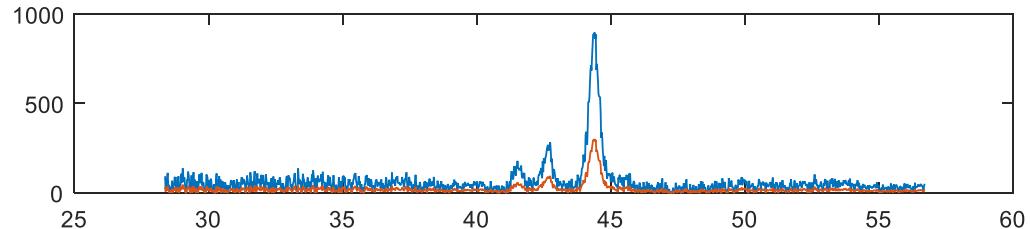
$$d = 1 - \frac{\sum_{i=1}^n (h_i \cdot k_i)}{\left( \sum_{i=1}^n h_i^2 \right)^{1/2} \left( \sum_{i=1}^n k_i^2 \right)^{1/2}}$$

Index is bin number!

Statistics: 1 - Pearson product-moment correlation

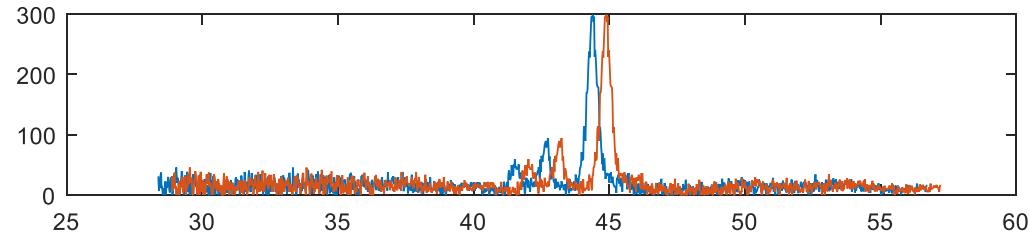
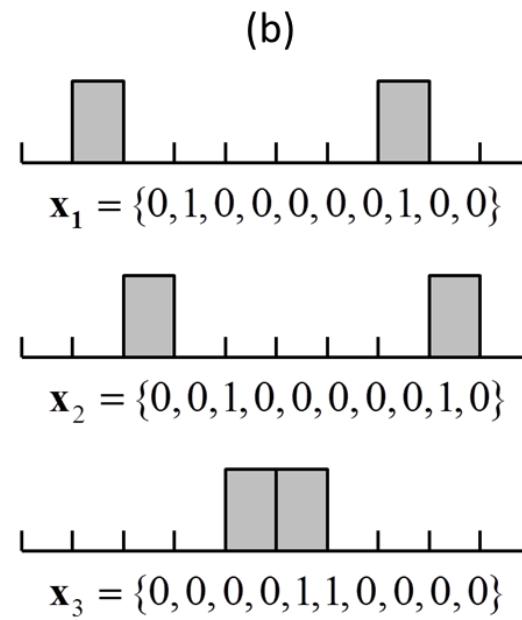
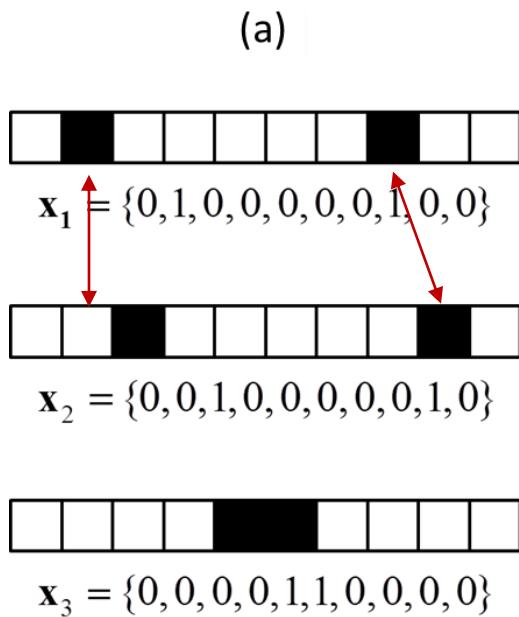
$$d = 1 - \frac{\sum_{i=1}^n (h_i - \bar{h})(k_i - \bar{k})}{\left( \sum_{i=1}^n (h_i - \bar{h}) \right)^{1/2} \left( \sum_{i=1}^n (k_i - \bar{k}) \right)^{1/2}}$$

**Cosine and Pearson metrics are scale invariant!**



# Measures: Translation Invariance

- Types: geometric, statistical, information theoretic, feature preserving

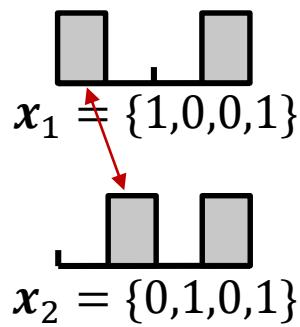


# Measures: Translation Invariance

- Types: geometric, statistical, information theoretic, feature preserving

Cross-bin Measure

$$d_A(x_1, x_2) = \sqrt{(x_1 - x_2)A(x_1 - x_2)^T}$$



Euclidean!

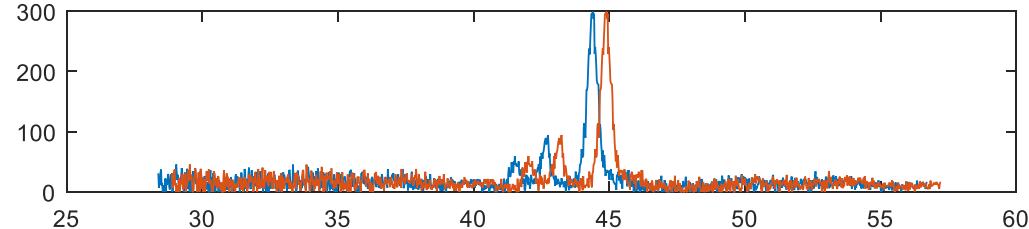
$$A_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$A_2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

$$d_{A_1}(x_1, x_2) = \sqrt{2}$$

$$d_{A_2}(x_1, x_2) = 0$$

Useful when features shift  
AND shift magnitude is known!

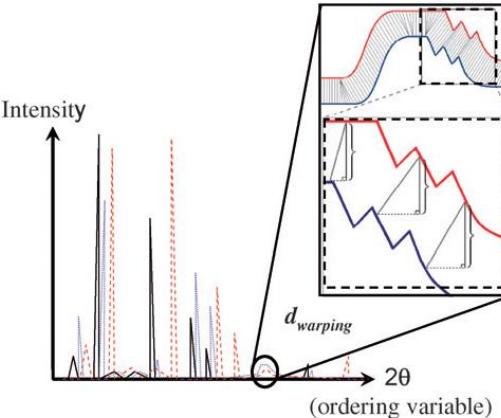
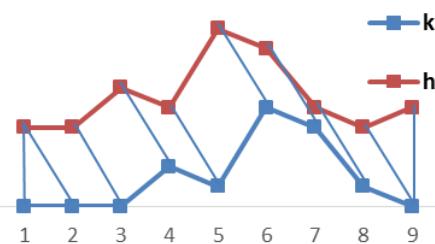


# Measures

- Types: geometric, statistical, information theoretic, feature preserving

Dynamic Time Warping

	1	2	3	4	5	6	7	8	9
h	0	0	2	1	5	4	1	0	1
k	0	0	0	2	1	5	4	1	0



Baumes, et al. 2008

Useful when features shift  
AND shift magnitude is UNknown!

# Measures

- Types: geometric, statistical, information theoretic, feature preserving

Information theory: Entropy  
Kullback-Leibler divergence

$$d(h, k) = \sum_i h_i \ln \left( \frac{h_i}{k_i} \right)$$

Information theory: Entropy  
Jensen-Shannon divergence

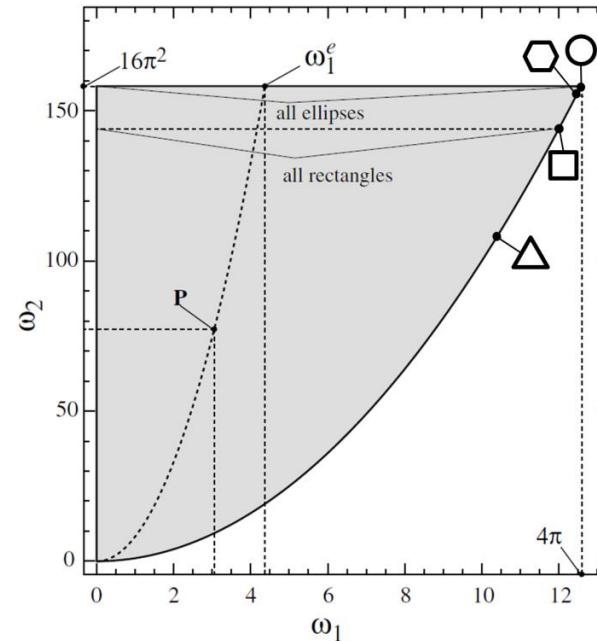
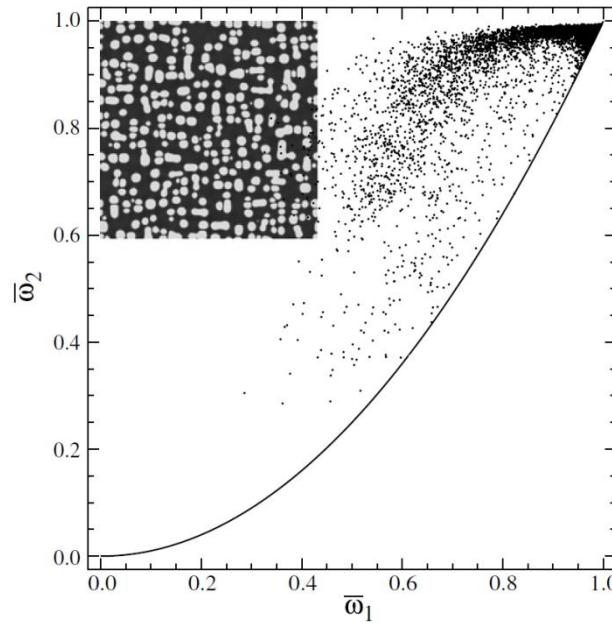
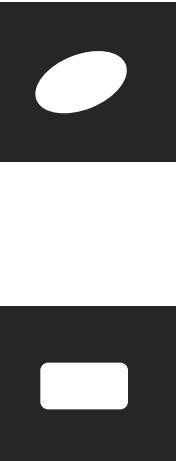
$$d(h, k) = \frac{1}{2} \left( \sum_i h_i \ln \frac{h_i}{m_i} + \sum_i k_i \ln \frac{k_i}{m_i} \right)$$

$$m_i = (h_i + k_i)/2$$

For comparing 2 probability density functions  
Quantifies mutual information

# Measures

- Types: geometric, statistical, information theoretic, feature preserving
- Metric for similarity in shape, combine:
  - Features: Moment Invariants (Sample Representation)
  - Metric: Euclidean

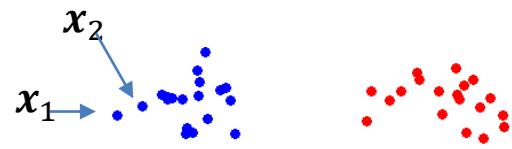


JP MacSleyné, et al.  
(2008) *Acta  
Materialia*, **56**

# Measures

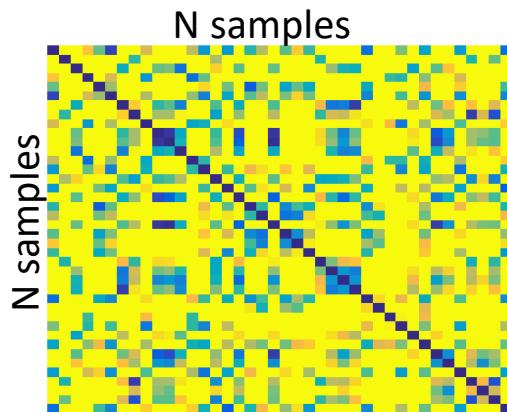
- Dissimilarity Matrix – Visualization

Data from two Gaussian distributions, labelled by color

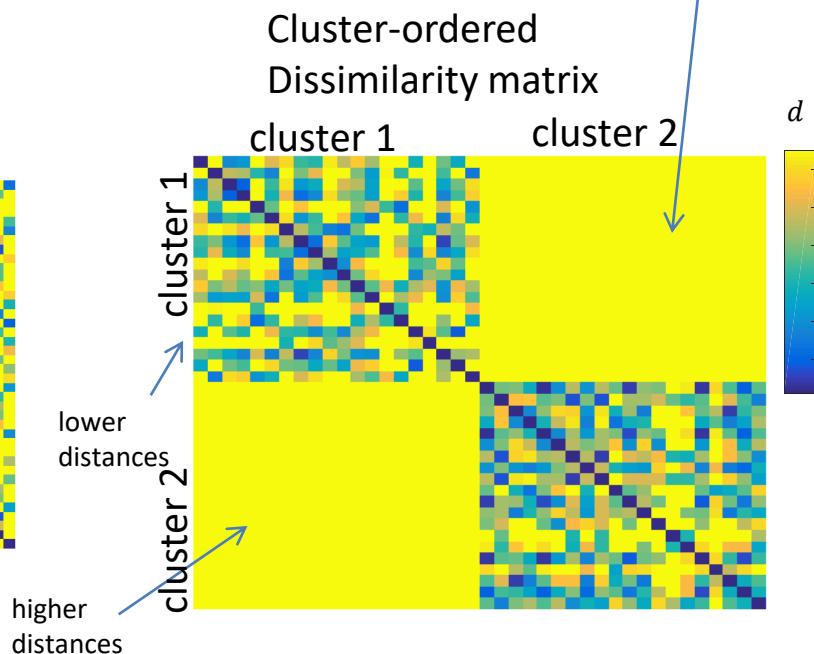


$$\begin{matrix} & x_1 & x_2 & x_3 & x_4 \end{matrix}$$
$$\begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix} \left[ \begin{array}{c} \\ \\ \\ \end{array} \right]$$

Unordered Dissimilarity matrix  $d(x_i, x_j)$



Cluster-ordered Dissimilarity matrix

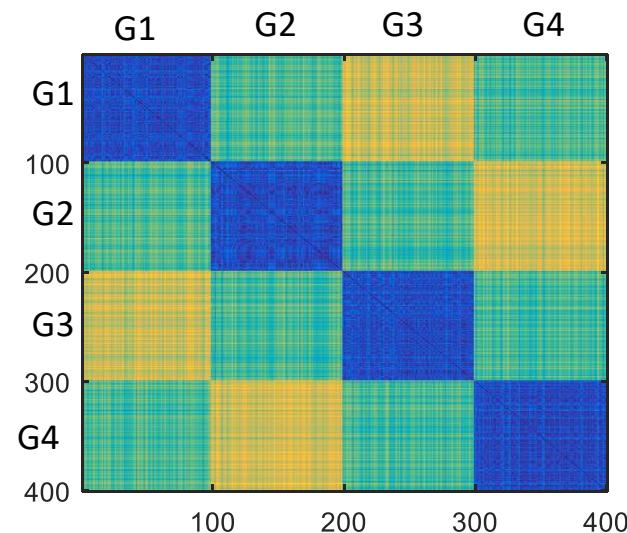
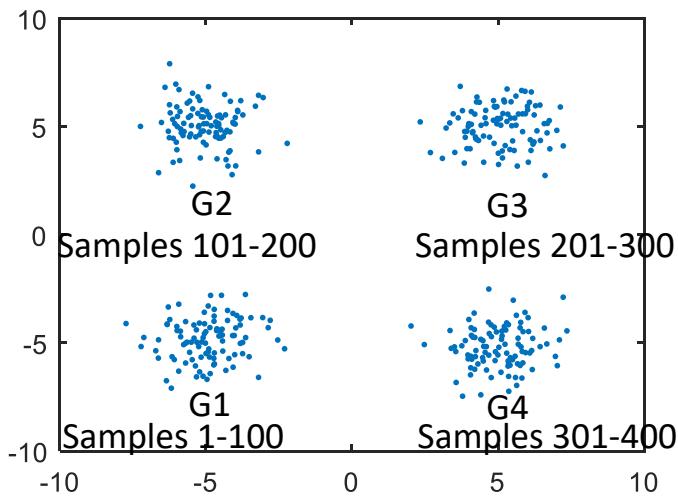


Heat Map!

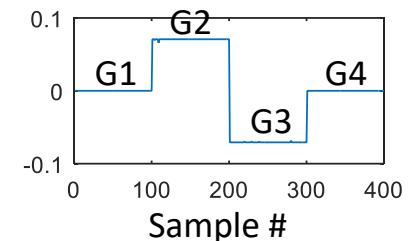
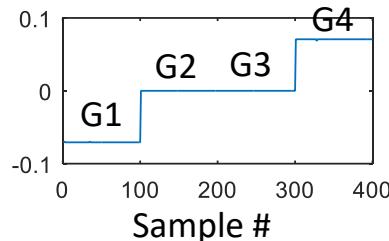
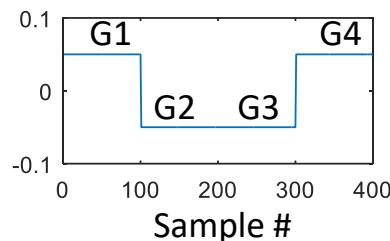
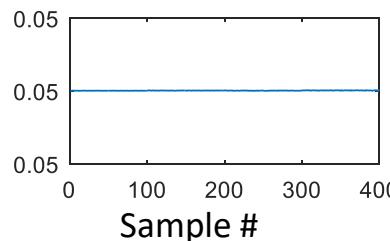
Dissimilarity (/similarity) matrix is at the core of Unsupervised Learning

# Dissimilarity Matrix

- Using the dissimilarity matrix, can form the graph Laplacian.



First 4 Eigenvectors of graph Laplacian

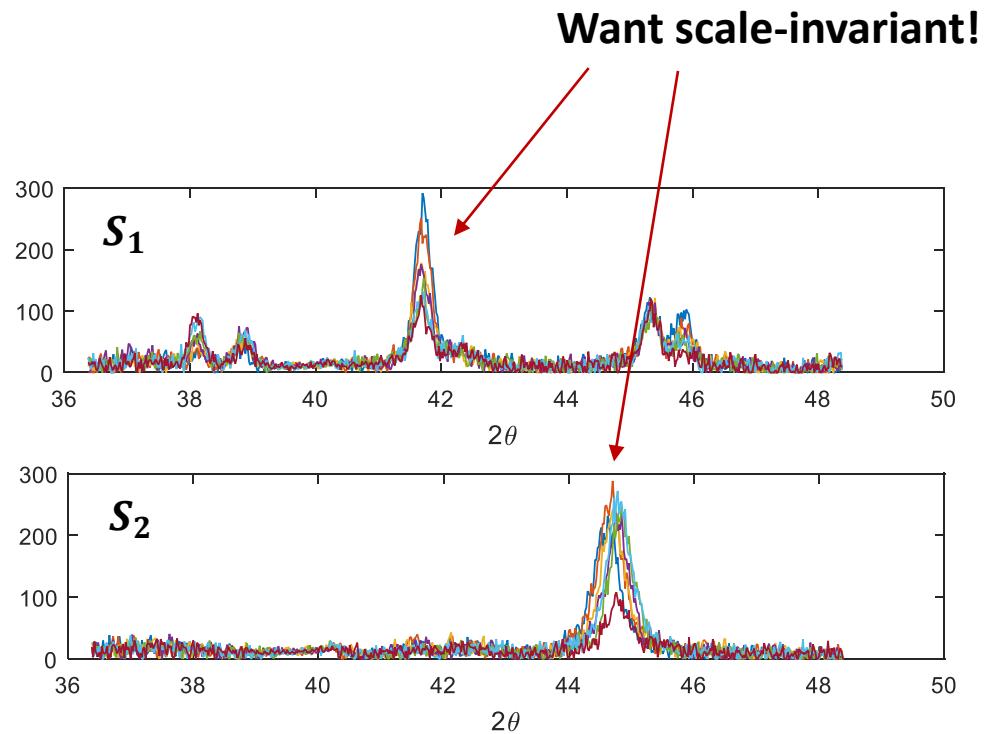
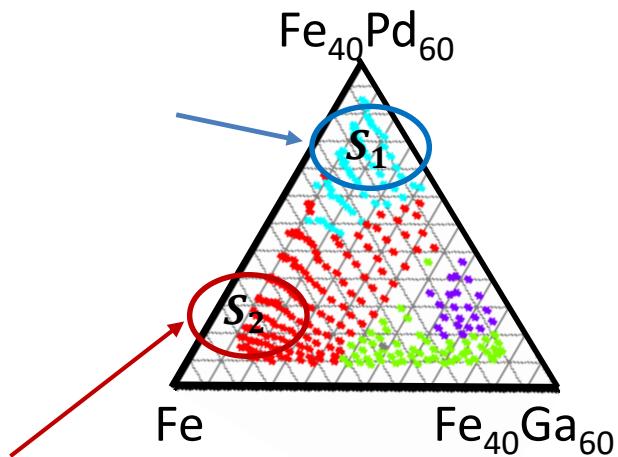


# Next: Jupyter for Dissimilarities

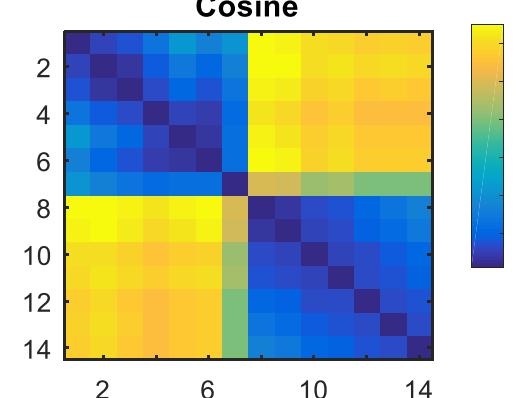
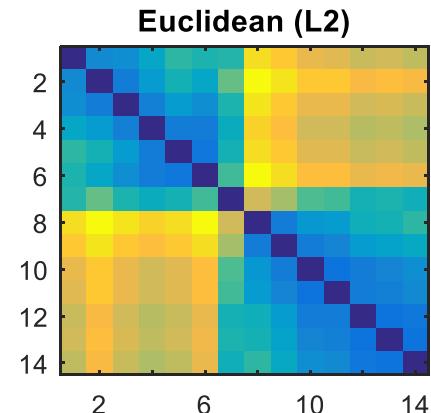
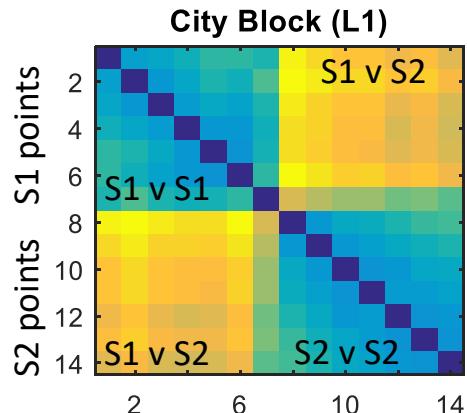
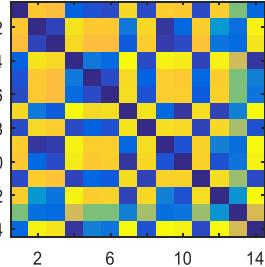
- First, we will install Variable Browser for Jupyter
- How to get the variable explorer to show up:
- Close Jupyter
- Start ‘anaconda prompt’
- Type (type one line, press enter, repeat):
  - pip install jupyter\_contrib\_nbextensions
  - jupyter contrib nbextension install --user
  - jupyter nbextension enable varInspector/main
- Open Jupyter

# Jupyter: Measures

- Example: X-ray Diffraction

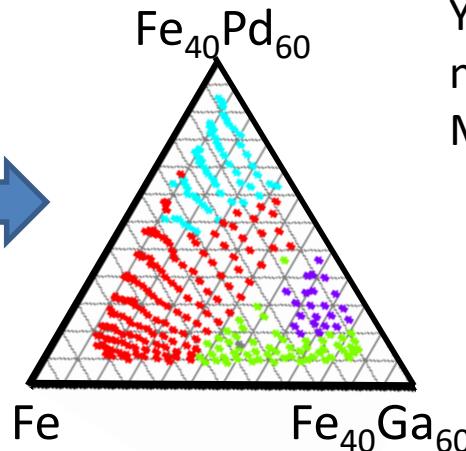
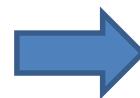
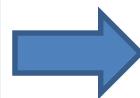
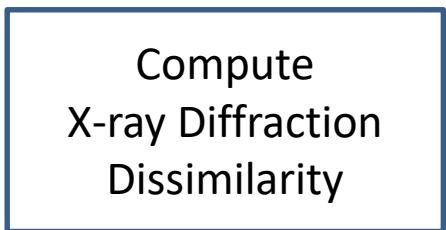
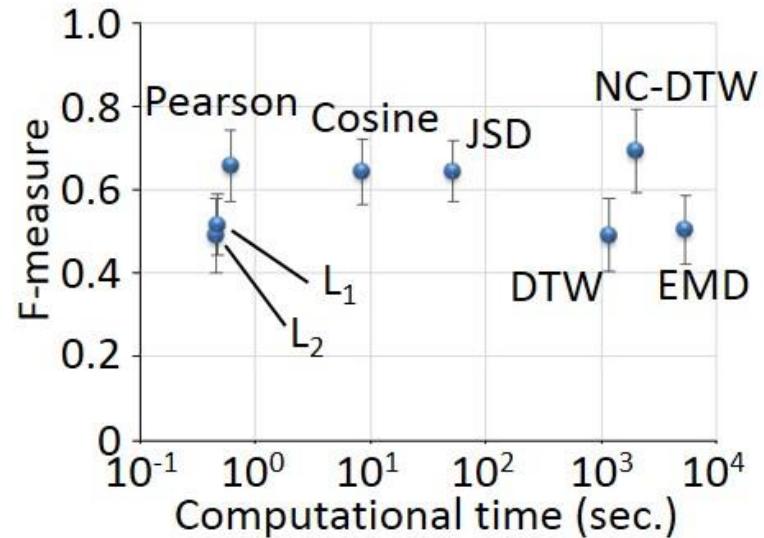


Randomly ordered



# Measures

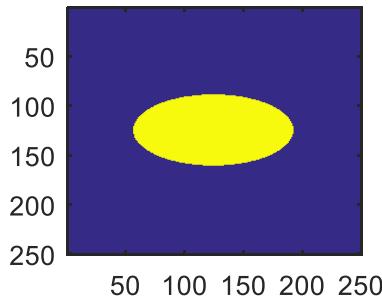
- Impact of measure on phase diagram determination
  - Accuracy
  - Computational cost



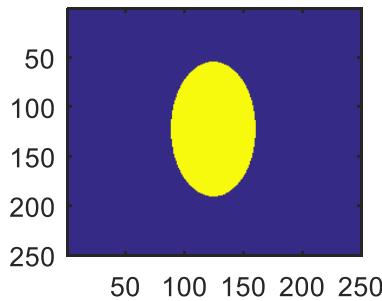
Yuma, et al. (2017)  
npj Computational Materials

Must take into account accuracy vs computation time

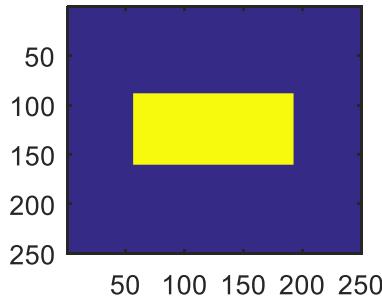
# Jupyter: Moment Invariants & Dissimilarity



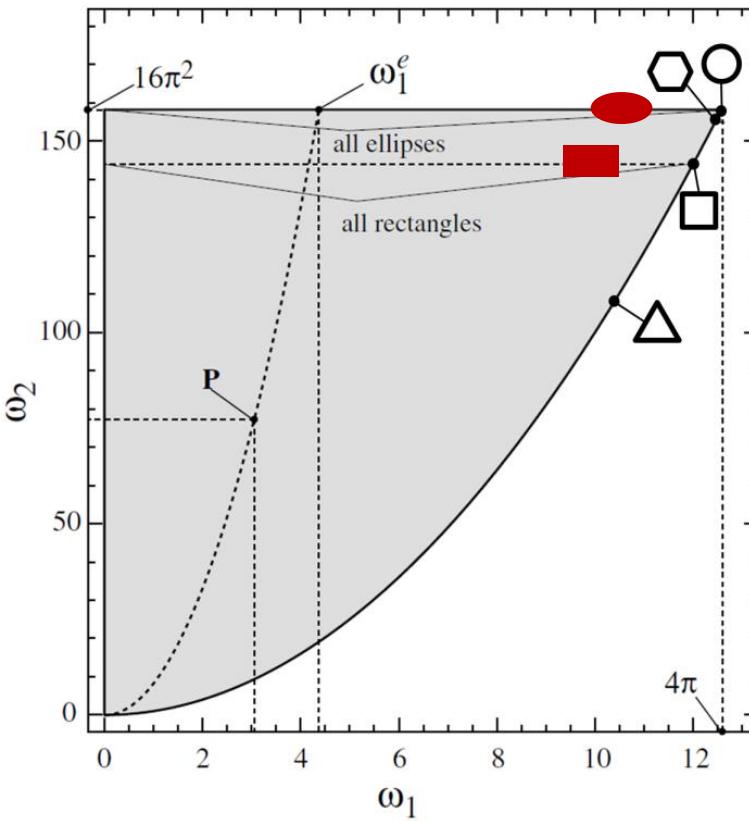
$$\begin{aligned}\omega_1 &= 10.4 \\ \omega_2 &= 157.8\end{aligned}$$



$$\begin{aligned}\omega_1 &= 10.4 \\ \omega_2 &= 157.8\end{aligned}$$



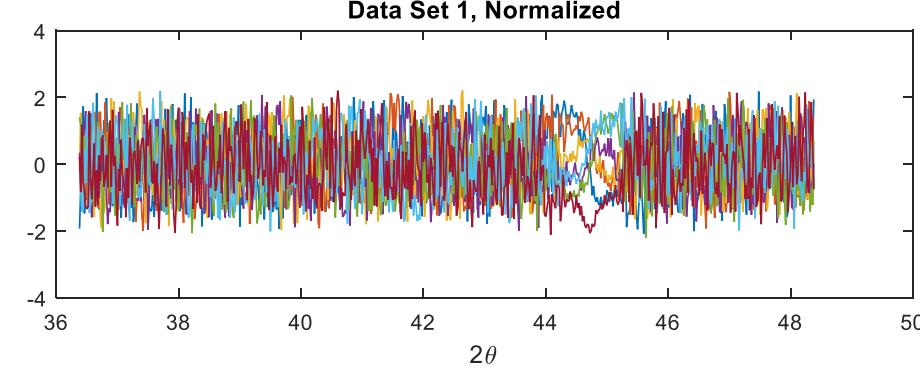
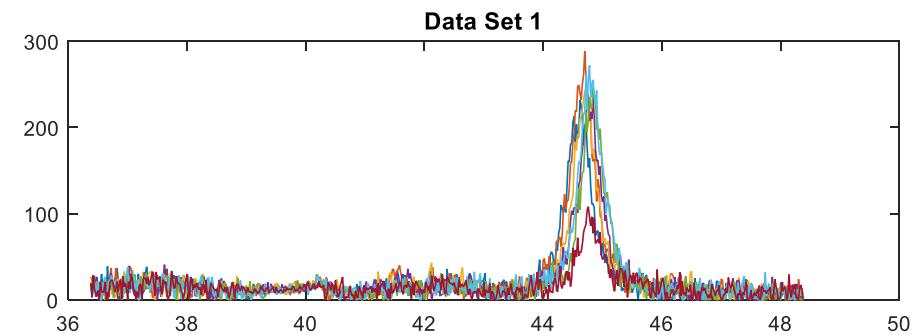
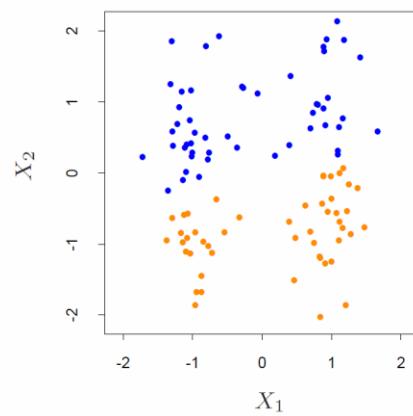
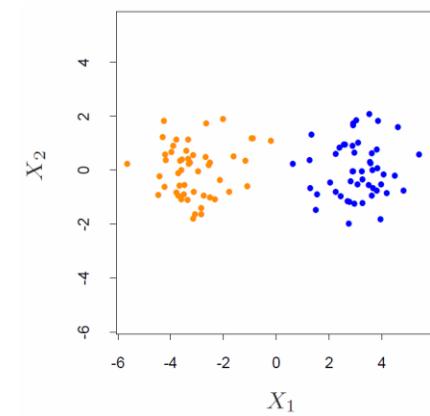
$$\begin{aligned}\omega_1 &= 9.9 \\ \omega_2 &= 143.9\end{aligned}$$



Difference between Shapes

# Data Standardization: Issues

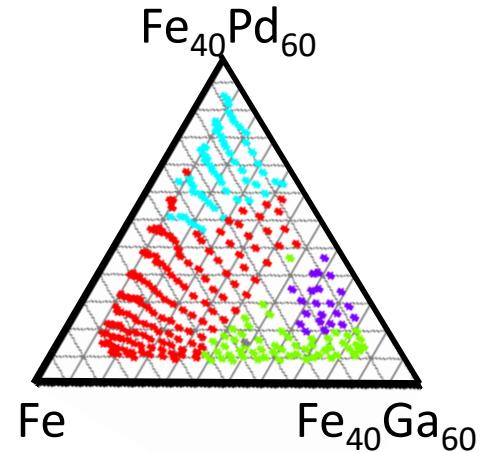
- For unsupervised learning, the goal is to identify distribution of  $X$
- Normalizing data may warp this distribution and make analysis difficult.



Questions?

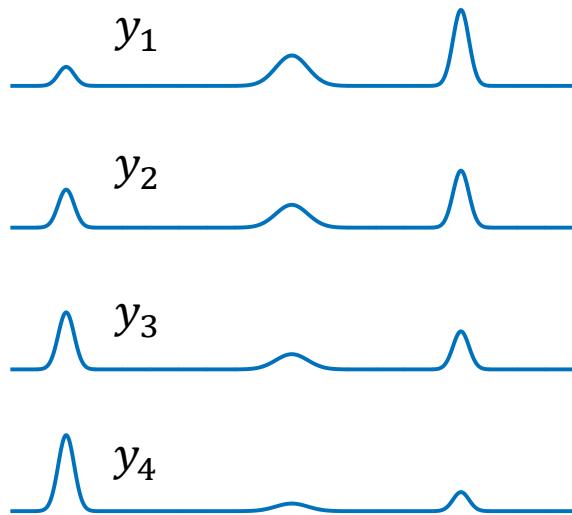
# Spectral Unmixing / Matrix Factorization

# Spectral Unmixing

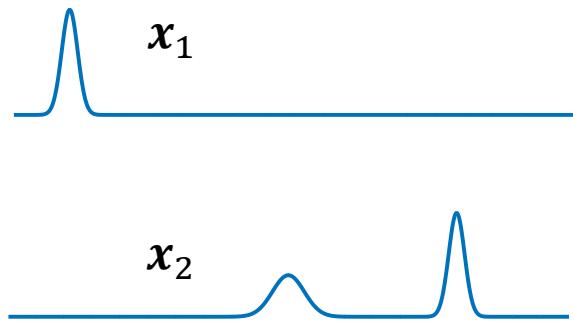


- Assume that spectra composed of linear mixture of underlying spectra.

$$Y = \beta X + \epsilon$$



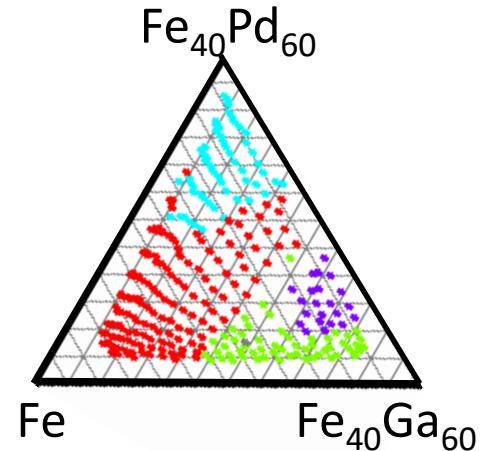
$$= \begin{bmatrix} .2 & .8 \\ .4 & .6 \\ .6 & .4 \\ .8 & .2 \end{bmatrix} *$$



**This method assumes linear mixing!**

# Spectral Unmixing

- E.g. finding constituent components in XRD



$$Y = \beta X + \epsilon$$

Abundances      Endmembers

$$\begin{bmatrix} \leftarrow & y_1 & \rightarrow \\ \vdots & & \vdots \\ \leftarrow & y_N & \rightarrow \end{bmatrix} = \begin{bmatrix} \beta_{11} & \cdots & \beta_{1L} \\ \vdots & \ddots & \vdots \\ \beta_{N1} & \cdots & \beta_{NL} \end{bmatrix} \begin{bmatrix} \leftarrow & x_1 & \rightarrow \\ \vdots & & \vdots \\ \leftarrow & x_L & \rightarrow \end{bmatrix} + \begin{bmatrix} \epsilon_{11} & \cdots & \epsilon_{1L} \\ \vdots & \ddots & \vdots \\ \epsilon_{N1} & \cdots & \epsilon_{NL} \end{bmatrix}$$

$\mathbf{y}$        $\boldsymbol{\beta}$        $\mathbf{x}$        $\boldsymbol{\epsilon}$

$$\begin{bmatrix} \text{Sample XRD} \\ \vdots \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 \\ 0.4 & 0.6 \\ 0.7 & 0.3 \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} \text{Constituent phases} \\ \vdots \end{bmatrix} + \begin{bmatrix} \text{Noise} \\ \vdots \end{bmatrix}$$

Mixture amounts

**This method assumes linear mixing!**

# Spectral Unmixing

To solve  $Y = \beta X + \epsilon$ :

$$\{\hat{X}, \hat{\beta}\} = \min_{X, \beta} \|Y - \beta X\|_2^2$$

**Least Squares!**



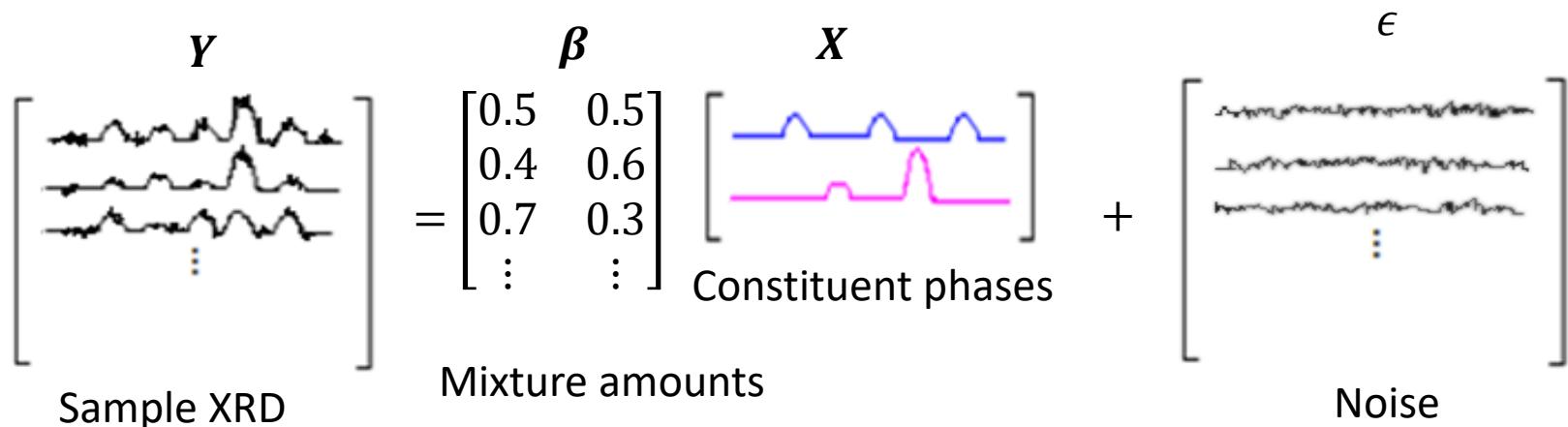
(Solved using similar tools  
as supervised learning)

- Supervised Learning:  $X$  is known
    - $x_i$  – independent variables,  $y_j$  - response / dependent variables
    - Solving for  $\beta$
  - Unsupervised: few or no  $x_i$  are known -> solve for  $X$  and  $\beta$
  - Initialize the values of  $\beta$  or  $X$
  - Iterate until convergence
  - Issue: Large solution space
    - Resolving may not provide the same results (random number seeds)
- > Constraints

# Spectral Unmixing: Non-negative

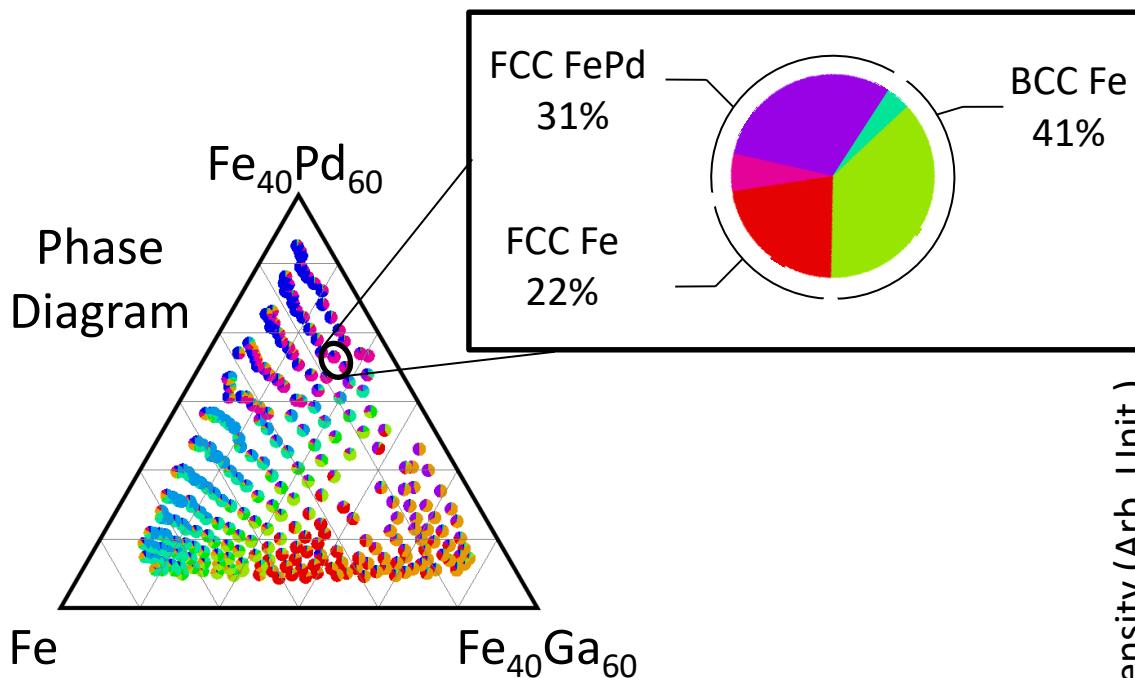
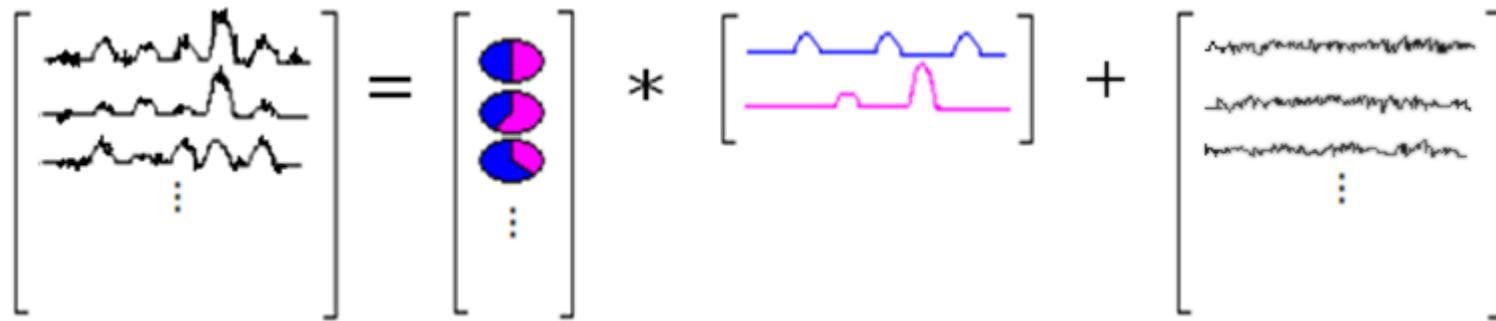
- Many data sets are non-negative ( $\geq 0$ )
  - X-ray, Raman, SEM, AFM, etc.

$$Y = \beta X + \epsilon \quad y_{ij}, x_{ij}, \beta_{ij} \in \mathbb{R}^+$$



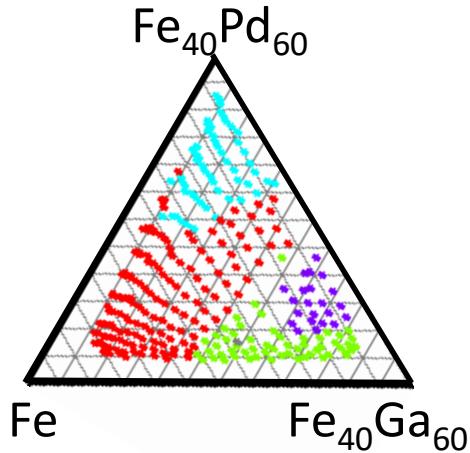
# Challenge: Phase Identification

## Non-Negative Matrix Factorization

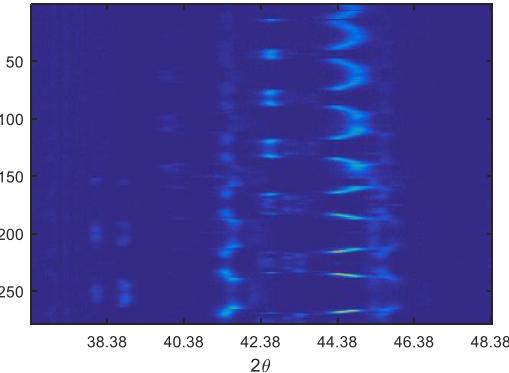


# Jupyter: Non-Negative Matrix Factorization

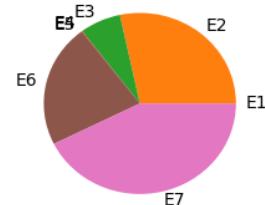
- Demo



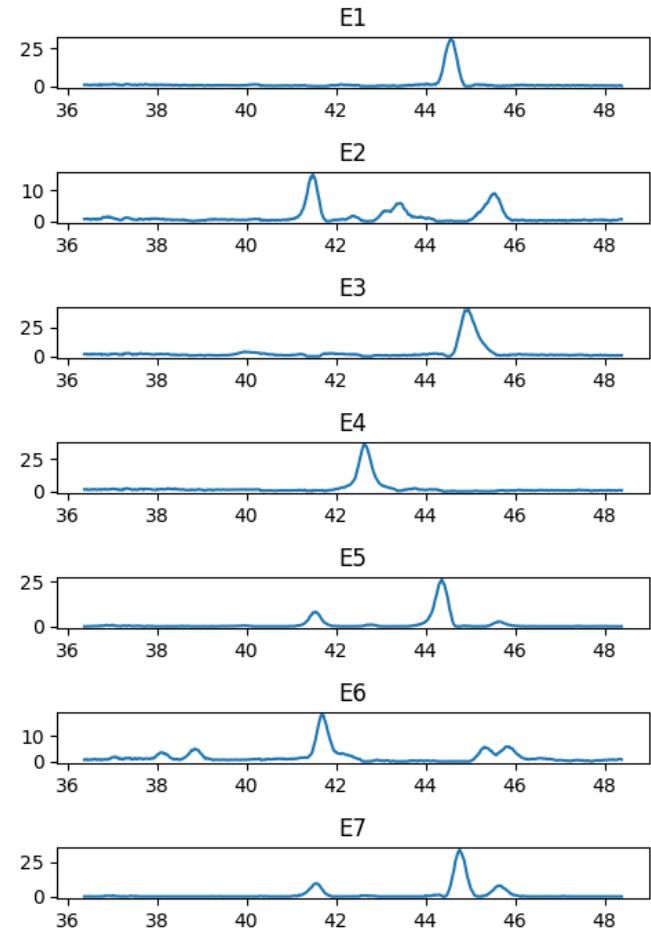
X-ray Diffraction



Abundances



Constituent phases



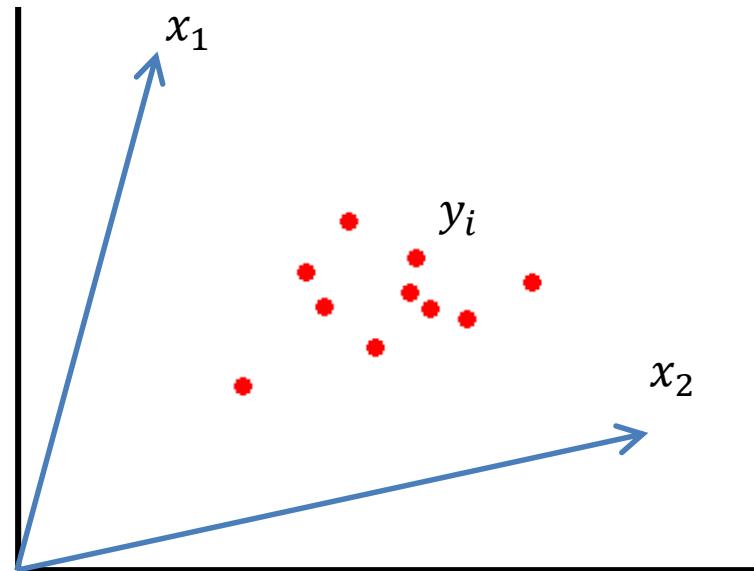
# Spectral Unmixing: Constraints

- Constraints
  - Minimize Endmember Volume

$$J = \left\| \mathbf{y}_k - \sum_{l=1}^L \mathbf{x}_l \beta_{lk} \right\|_2^2 + \alpha \sum_{k=1}^{L-1} \sum_{l=k+1}^L \|\mathbf{x}_l - \mathbf{x}_k\|_2^2$$

- Ensure that constituent phases look like the x-ray diffraction patterns of  $\mathbf{y}_k$

+ Others



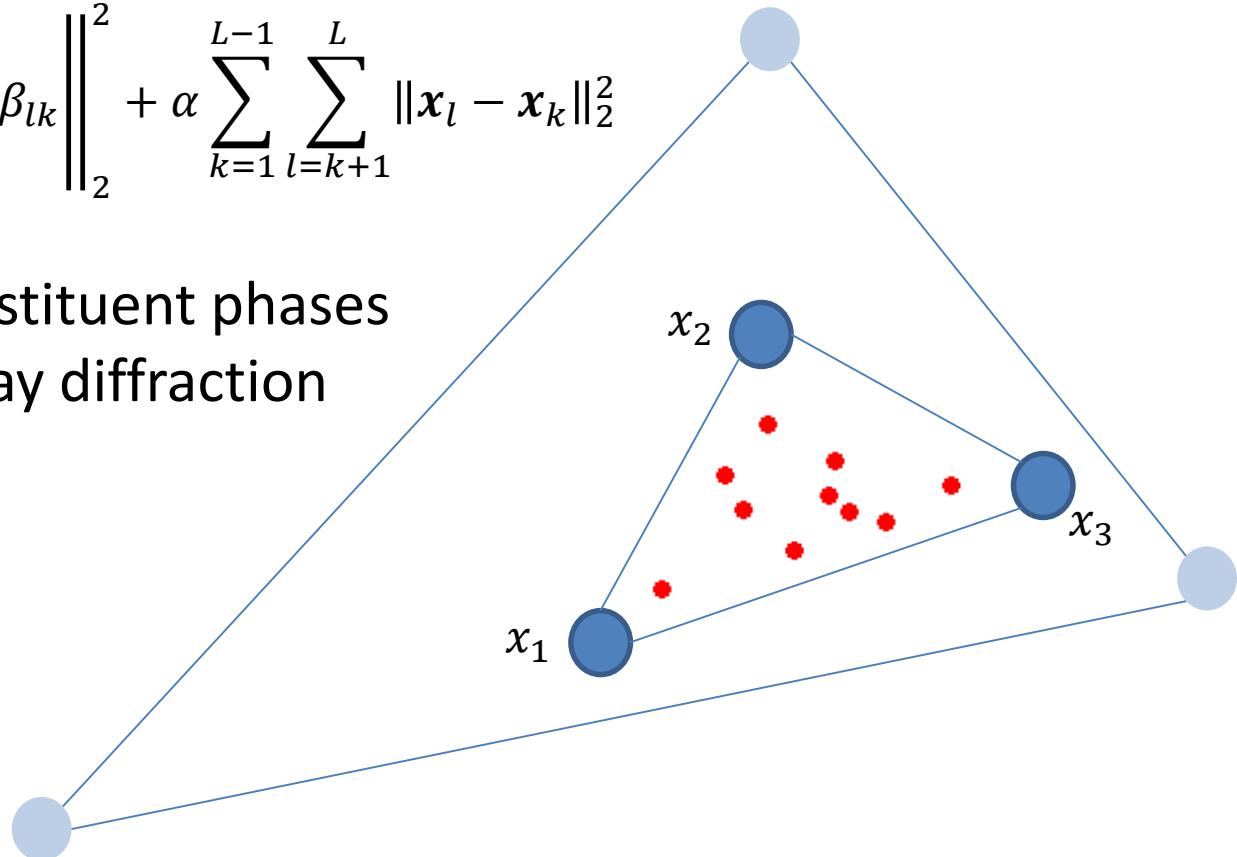
# Spectral Unmixing: Constraints

- Constraints
  - Minimize Endmember Volume

$$J = \left\| \mathbf{y}_k - \sum_{l=1}^L \mathbf{x}_l \beta_{lk} \right\|_2^2 + \alpha \sum_{k=1}^{L-1} \sum_{l=k+1}^L \|\mathbf{x}_l - \mathbf{x}_k\|_2^2$$

- Ensure that constituent phases look like the x-ray diffraction patterns of  $\mathbf{y}_k$

+ Others



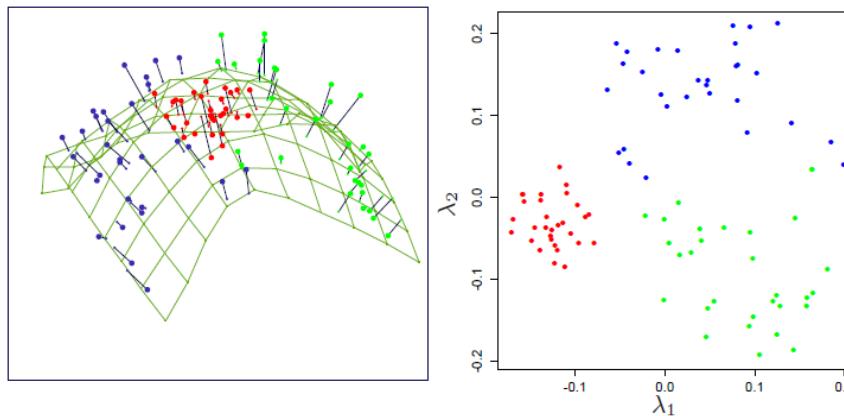
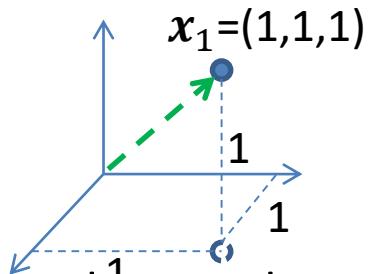
# Questions?

# Latent Variable Analysis

Visualization / Dimension Reduction  
Approximation / Remove noise  
Outlier Detection  
Identify underlying variables

# The Challenge of Data Visualization

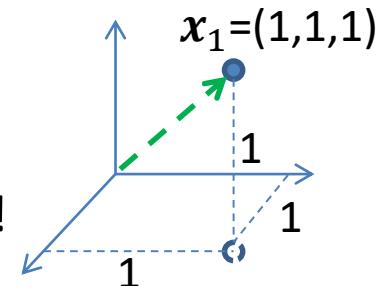
- For  $N$  data points, each of dimension  $D$ :
  - Each data point exists in  $D$  dimensional space.
  - If  $D \leq 3$  can easily visualize in native space.
  - If  $D \geq 3$ , can look for “manifold” – lower dimensional space that describes the data.



**Use LVA to visualize dimension  $\geq 3$**

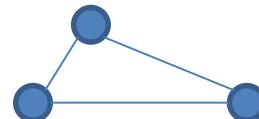
# The Challenge of Data Visualization

- Alternative: visualizing dissimilarity
  - Each data point exists in  $D$  dimensional space.
  - If  $D \leq 3$  can easily visualize in native space.
  - If  $D > 3$ , can visualize in  $N-1$  dimensional space!



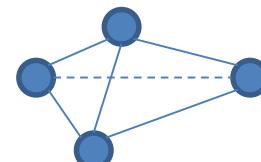
2 data points  
-> 1D

$$d(x_1, x_2)$$

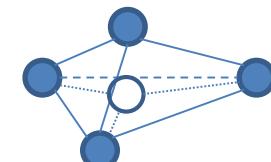


3 data points  
-> 2D

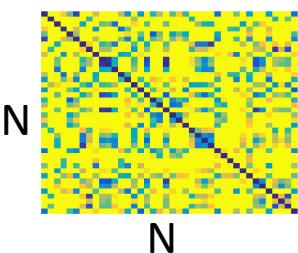
$$\begin{aligned} d(x_1, x_2) \\ d(x_1, x_3) \\ d(x_2, x_3) \end{aligned}$$



4 data points  
-> 3D



4 data points  
Flatten -> 2D

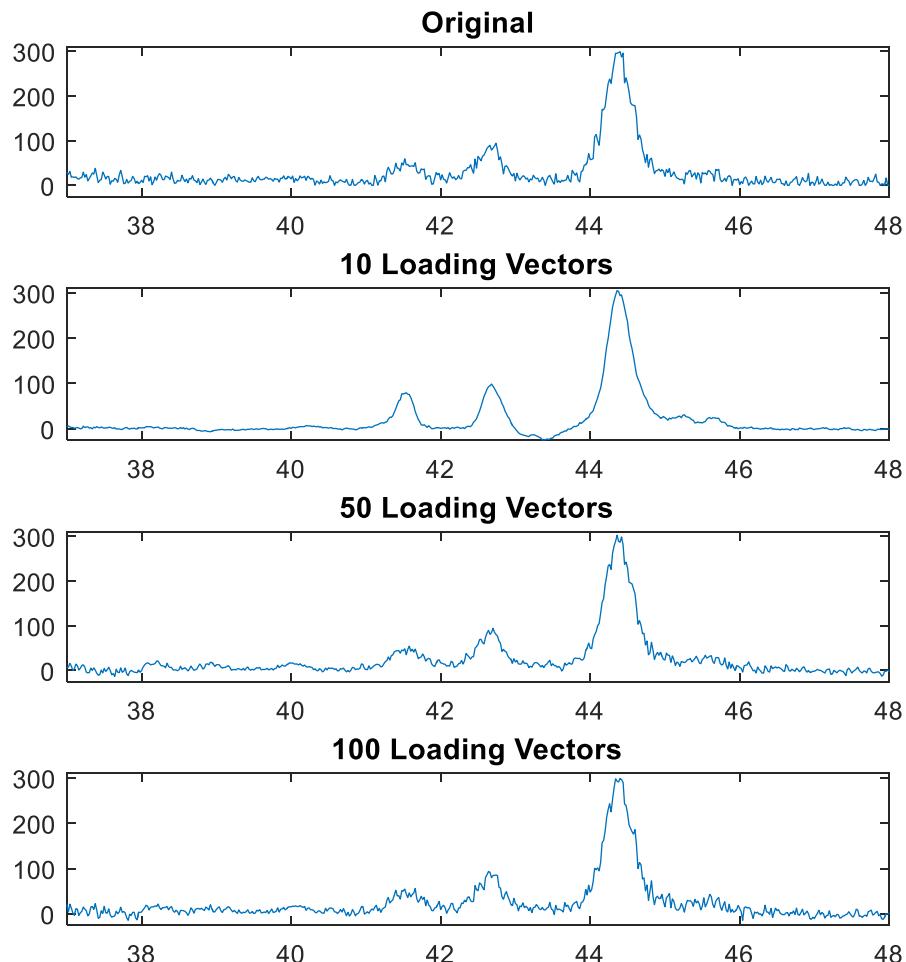


5+ data points?

Use LVA to visualize sample number  $> 4$  and/or dimension  $> 3$

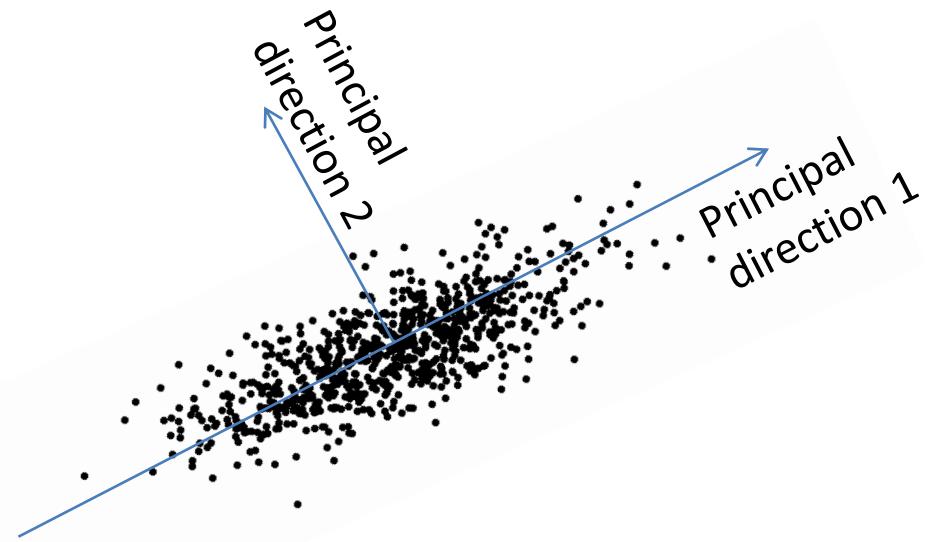
# Approximation / Smoothing

- Can use Latent Variable Analysis to smooth data.
  - Control how much high frequency signals appear in the data.
  - Approximate data



# Principal Component Analysis (PCA)

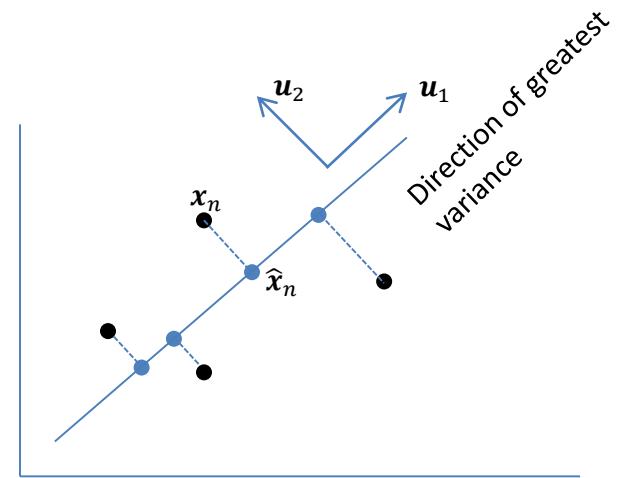
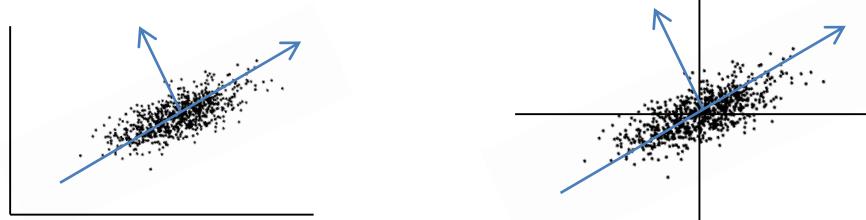
- PCA identified direction of greatest variance & variance in that direction.
- Finds orthogonal directions of greatest variance.
- Applicable in high dimensions.
- Also useful for:
  - Vis high dimensional data
  - Approximations
  - Outlier identification
  - Latent variable analysis



# Principal Component Analysis (PCA)

$$X = \begin{bmatrix} \leftarrow & x_1 & \rightarrow \\ \vdots & \ddots & \vdots \\ \leftarrow & x_N & \rightarrow \end{bmatrix} \quad X \text{ is } N \text{ samples} \times D \text{ variables}$$

**REQUIRED!**  $x_n \leftarrow (x_n - \bar{x})$  normalizing each variable



- Want to find unit vector  $\mathbf{u}_1$  that maximizes projected variance.
- Projection onto  $\mathbf{u}_1$  is  $X\mathbf{u}_1^T$
- Variance of projection is  $\lambda_1 = \frac{1}{N} (X\mathbf{u}_1^T)^T X\mathbf{u}_1^T = \frac{1}{N} \mathbf{u}_1 X^T X \mathbf{u}_1^T = \mathbf{u}_1 S \mathbf{u}_1^T / \mathbf{u}_1 \mathbf{u}_1^T$

$$S = \begin{bmatrix} \uparrow & \cdots & \uparrow \\ x_1 & \ddots & x_N \\ \downarrow & \cdots & \downarrow \end{bmatrix} \begin{bmatrix} \leftarrow & x_1 & \rightarrow \\ \vdots & \ddots & \vdots \\ \leftarrow & x_N & \rightarrow \end{bmatrix} \quad S \text{ is the Covariance Matrix, } D \times D$$

For 2D data (x,y),  $S$  is  $2 \times 2$

# Principal Component Analysis (PCA)

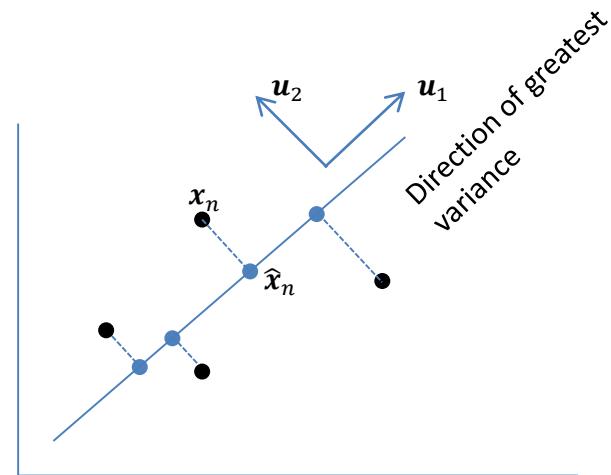
Covariance Matrix:

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

$\mathbf{X}$  is  $N \times D$   
 $\mathbf{S}$  is  $D \times D$

**REQUIRED!**  $x_n \leftarrow (x_n - \bar{x})$  normalizing each variable

- Want to find unit vector  $u_1$  that maximizes projected variance  $\lambda_1 = \mathbf{u}_1 \mathbf{S} \mathbf{u}_1^T / \mathbf{u}_1 \mathbf{u}_1^T$



Solution: Eigen analysis!

For all orthogonal unit vectors  $u_i$ , maximize projected variance  $u_i^T \mathbf{S} u_i$

$$\mathbf{S} \mathbf{u}_1^T = \lambda_1 \mathbf{u}_1^T$$

$$\mathbf{S} \mathbf{u}_i^T = \lambda_i \mathbf{u}_i^T$$

$\lambda_i$  Variances,

$\mathbf{u}_i$  Principal directions /  
Loading vectors,

$x_n \mathbf{u}_i^T$  Principal components:  
Projections onto principal directions

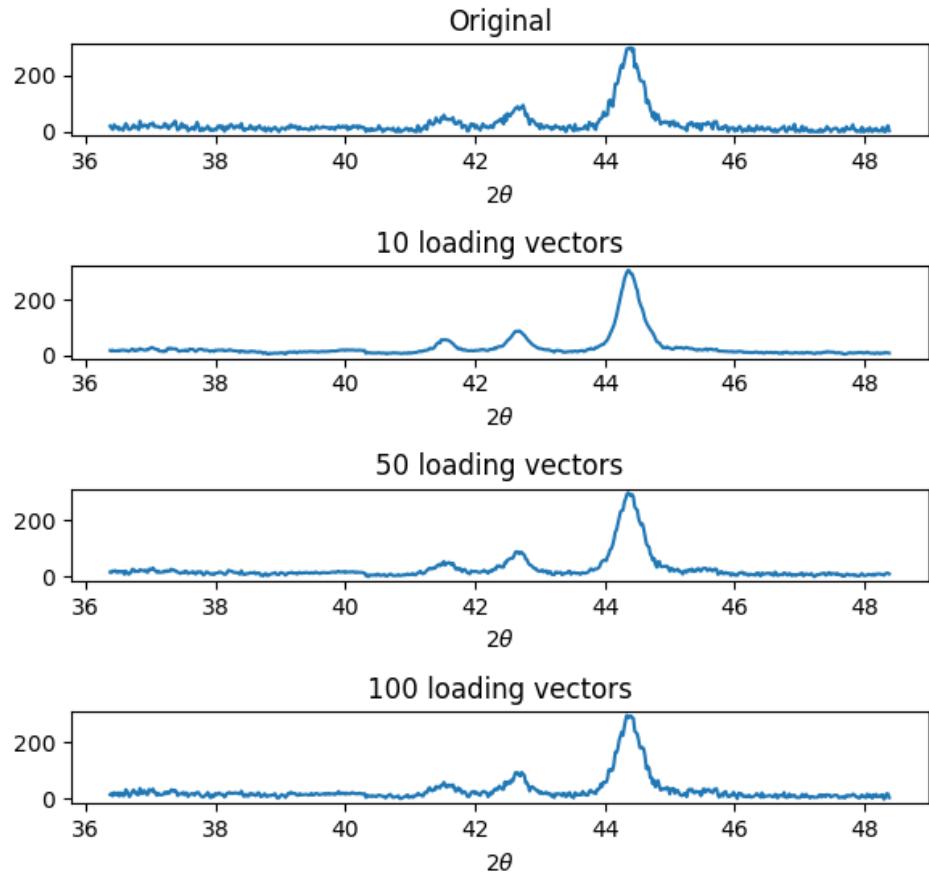
# Approximation / Smoothing

- The data is fully described by the entire set of loading vectors.

$$x_n = \sum_{i=1}^D (x_n u_i^T) u_i$$

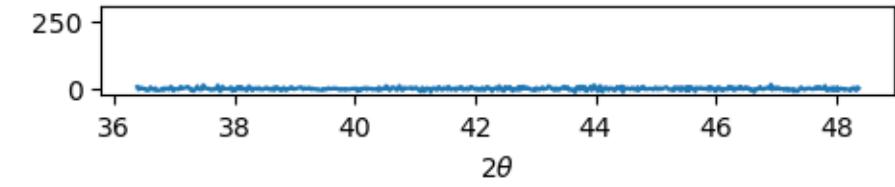
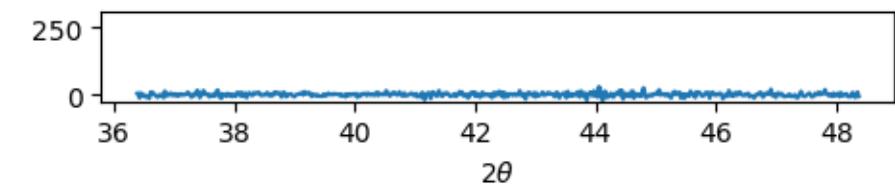
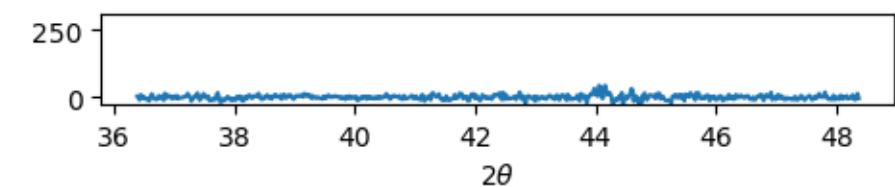
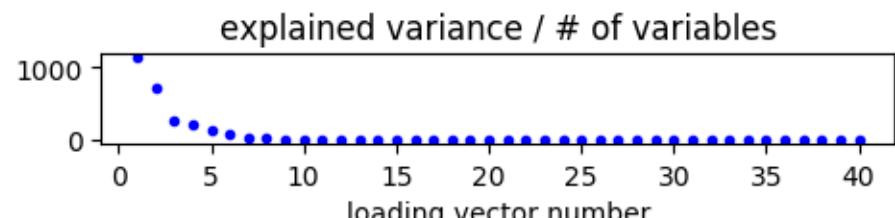
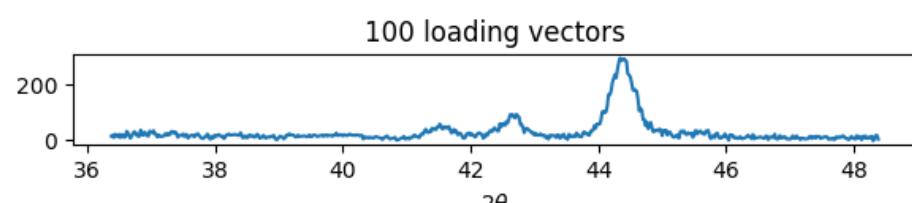
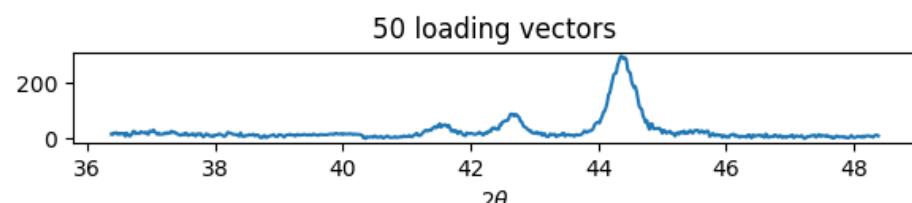
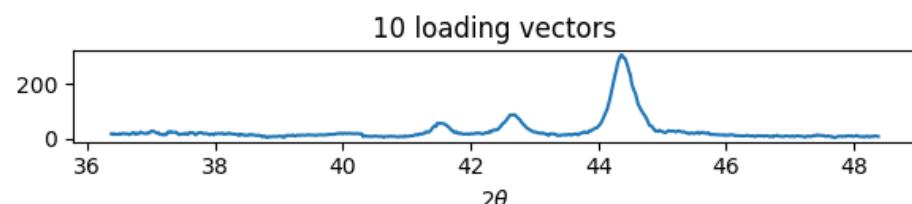
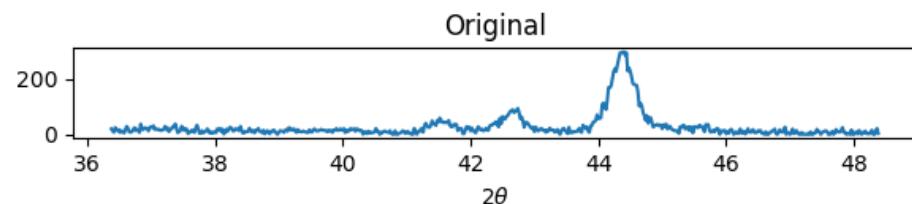
- Using only the first M loading vectors (with the largest variances) will approximate the data.

$$\hat{x}_n = \sum_{i=1}^M (x_n u_i^T) u_i, M < D$$



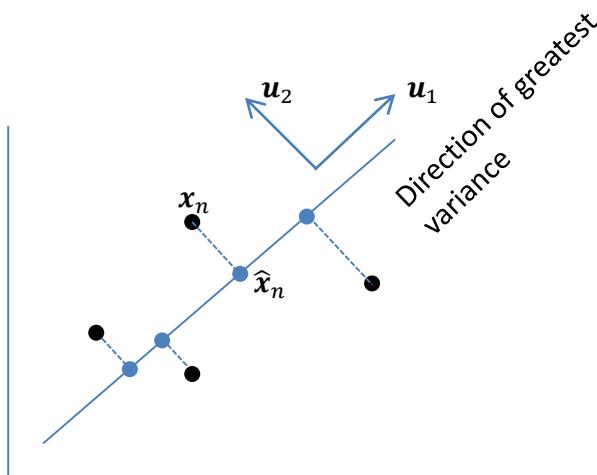
# Jupyter: Approximation / Smoothing

- Mapping onto the remaining loading vectors



# Data Visualization

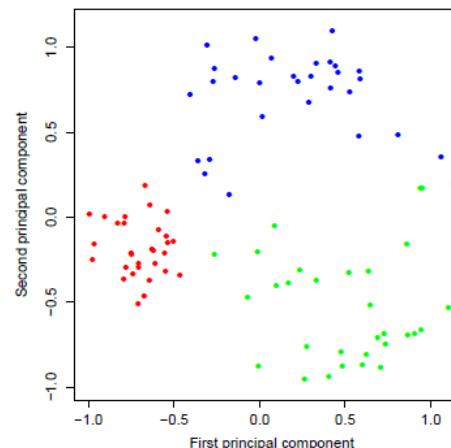
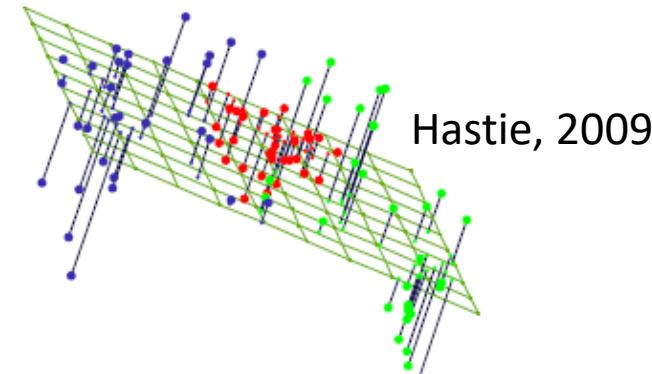
- The PCA approximation can also be used to reduce data dimensionality to simplify visualization.
- With  $M = 2$ , the data points are mapped onto the first 2 loading vectors.



$$\hat{x}_n = (x_n \ u_1^T) u_1$$

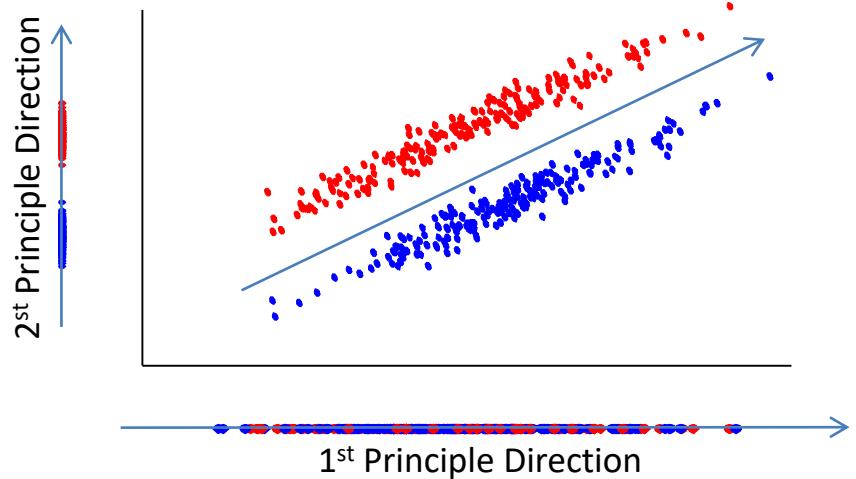
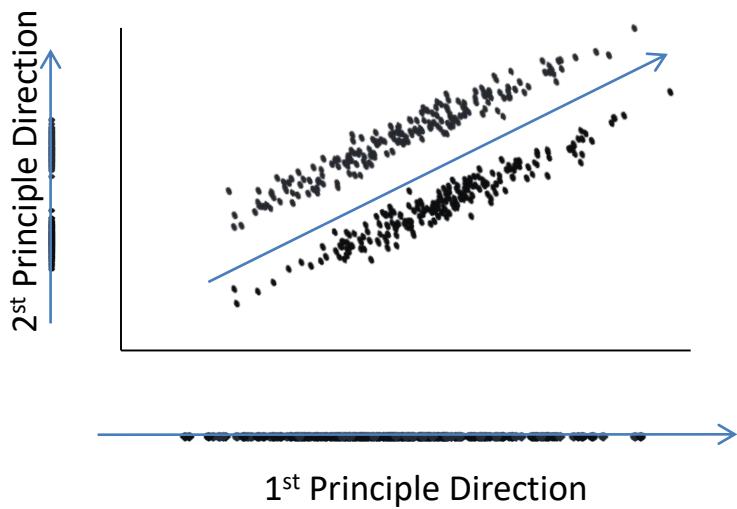
$u_1$

$\hat{x}_n$



$$\hat{x}_n = \sum_{i=1}^2 (x_n \ u_i^T) u_i$$

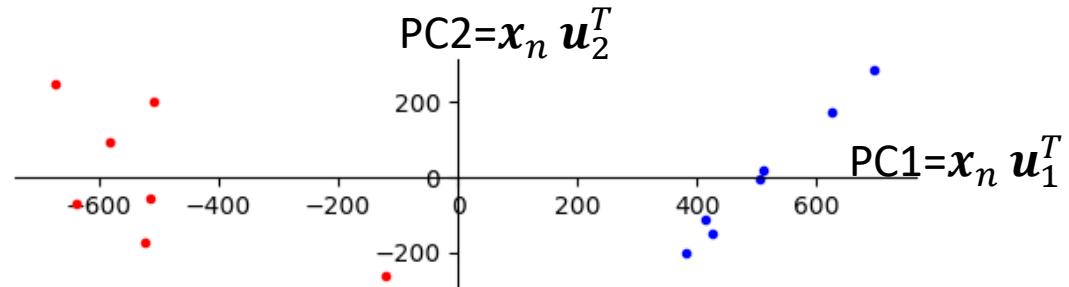
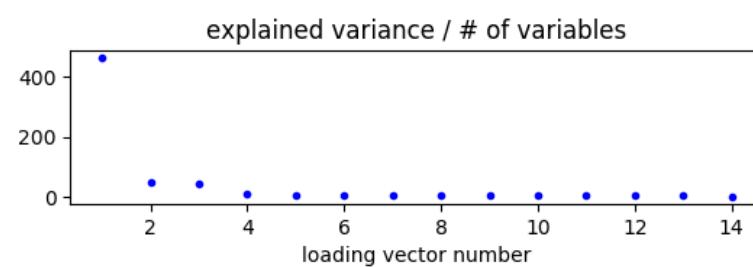
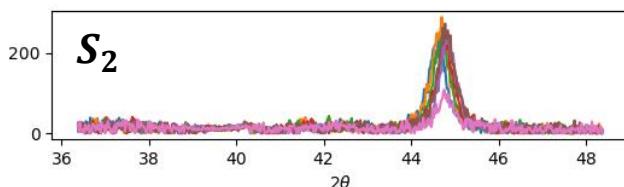
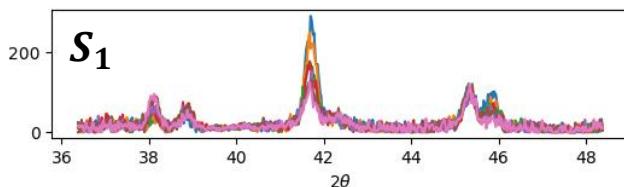
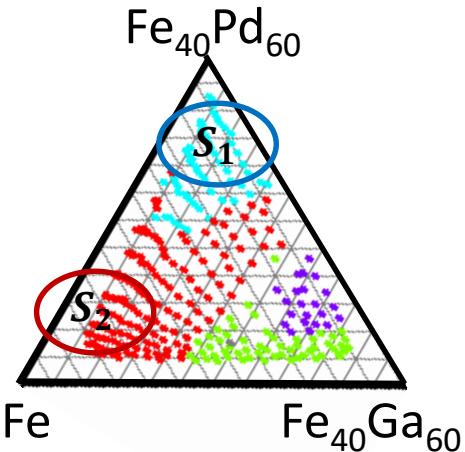
# Example of when Not to use PCA



# Jupyter: Data Visualization

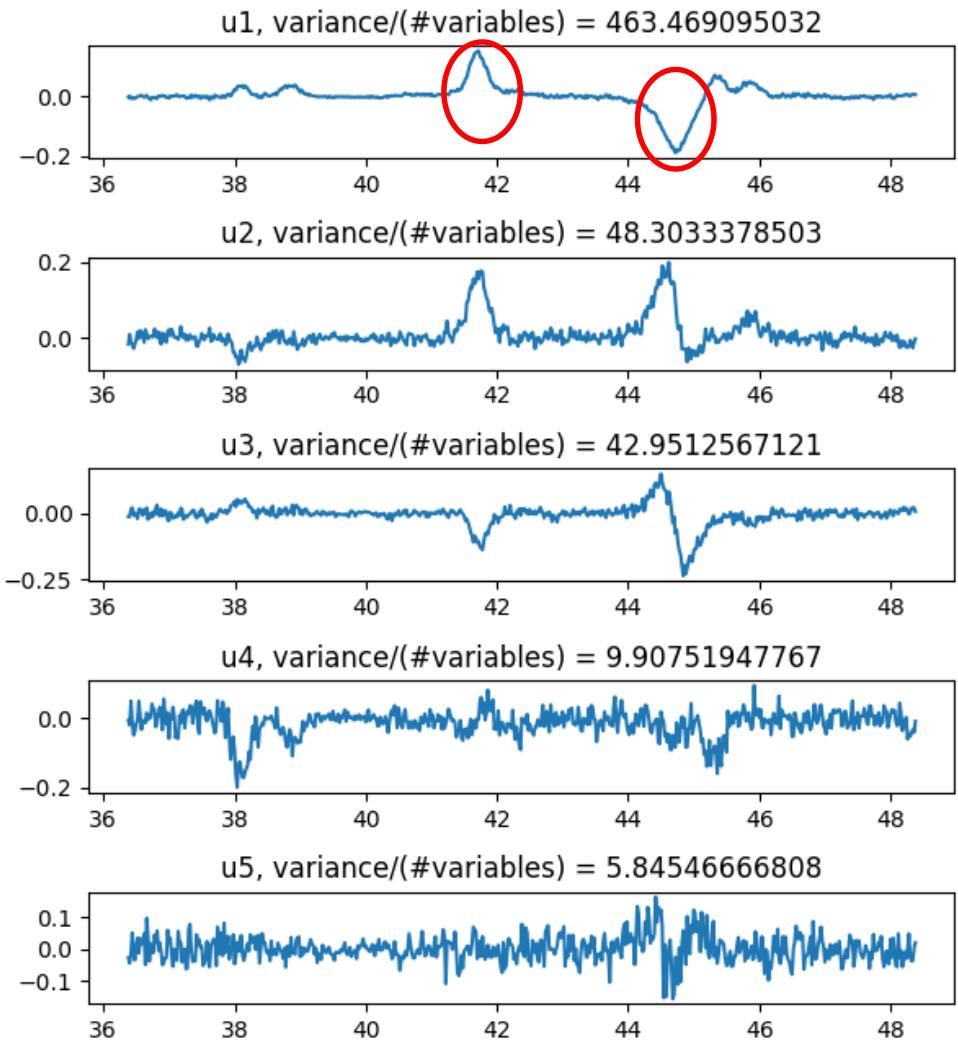
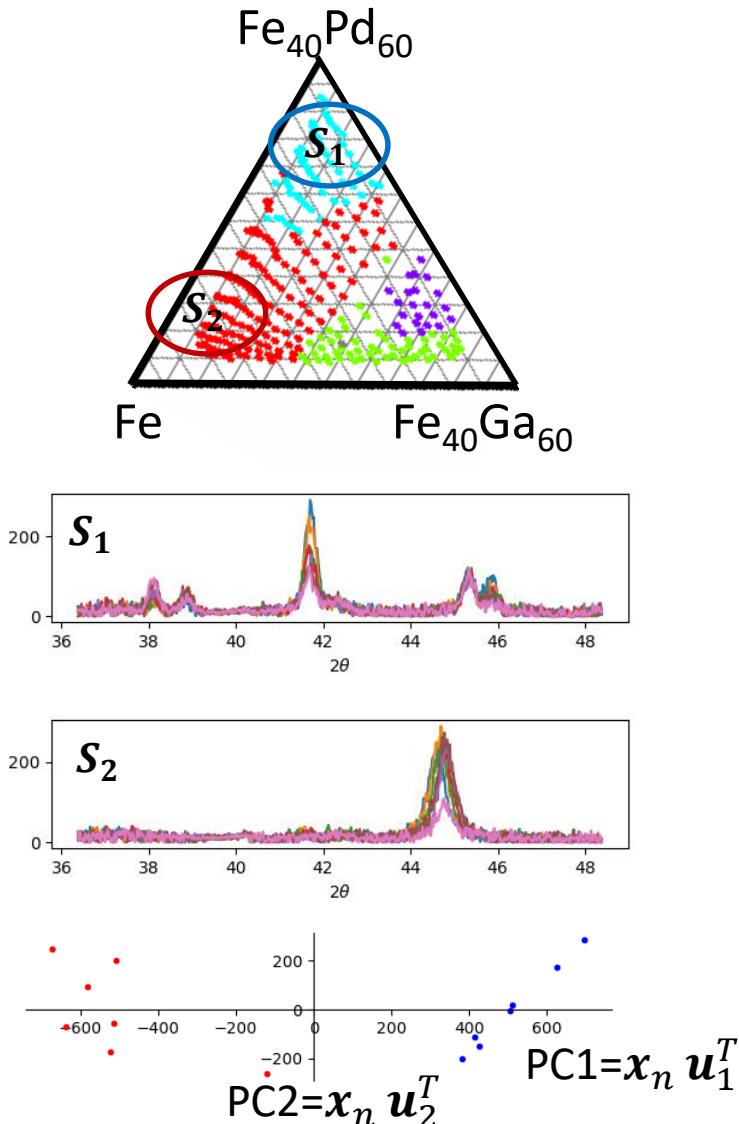
- Mapping X-ray diffraction to 2D.
- $X$  is 14 samples with 601 variables (intensity values)

$$\hat{x}_n = \sum_{i=1}^2 (x_n \mathbf{u}_i^T) \mathbf{u}_i$$



X-ray diffraction projections using PC1 and PC2  
Most the variance described by 1 loading vector!

# Jupyter: Visualize Loading Vectors



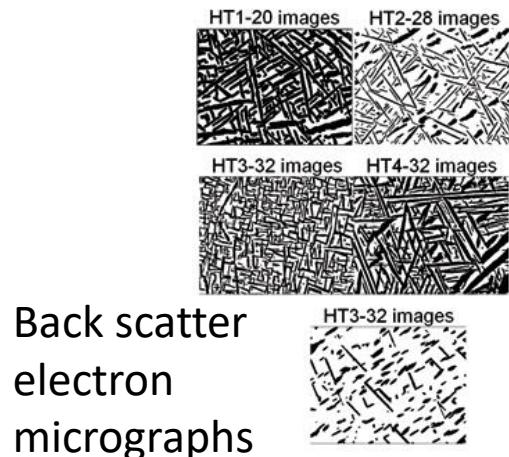
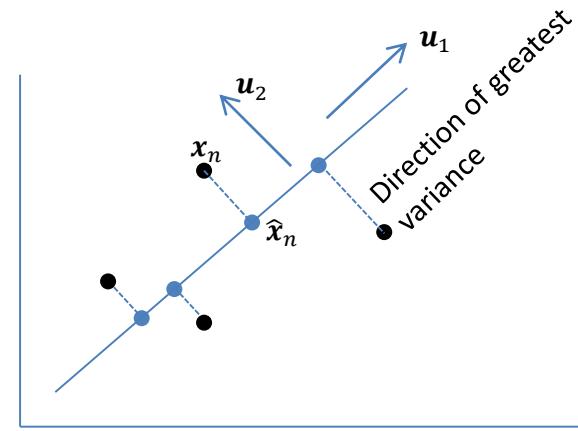
# Principle Component Analysis (PCA)

- Dimension reduction and approximation

- Select  $M < D$

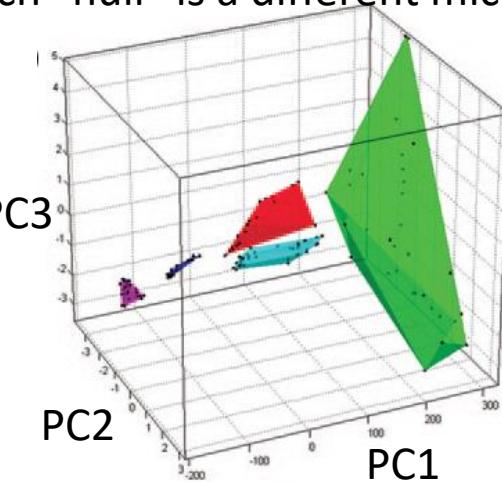
$$\hat{x}_n = \sum_{i=1}^M (x_n u_i^T) u_i$$

- $x_n$  of  $D$  dimensions  $\rightarrow \hat{x}_n$  of  $M$  dimensions
- Selecting  $M = 1-3$ , creates visualization of data.



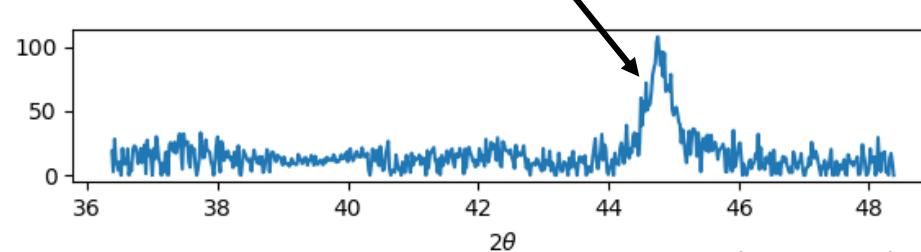
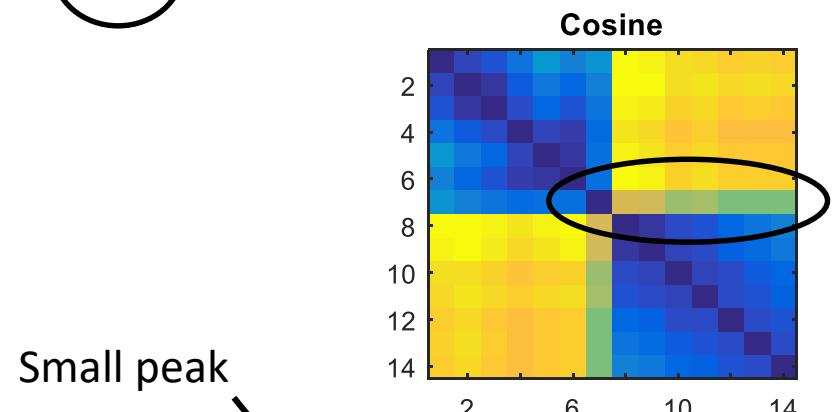
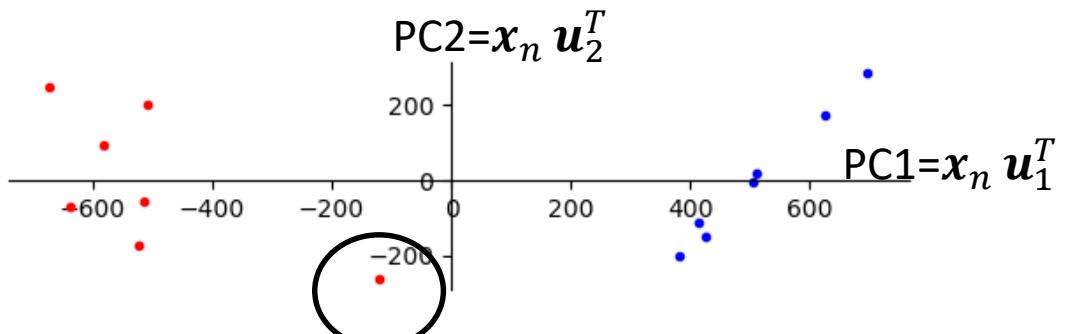
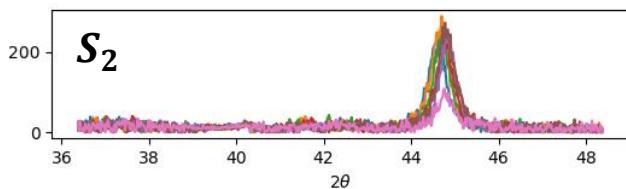
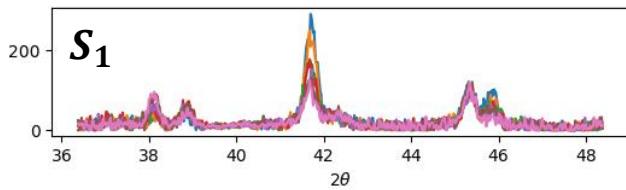
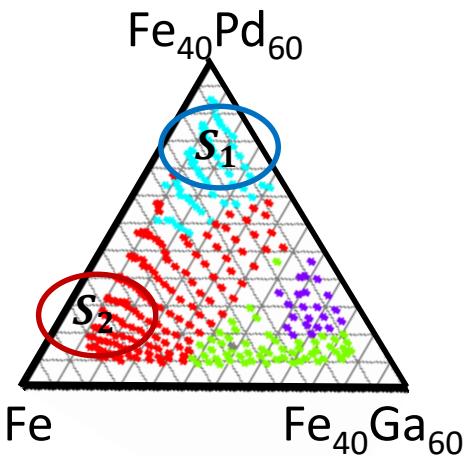
Convert to  
2-point statistics

Each “hull” is a different microstructure!



Kalidindi  
(2011)  
JOM

# Jupyter: Outlier Identification



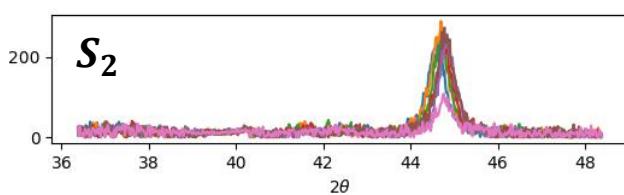
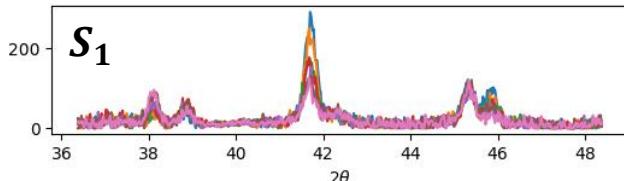
Previously sample 7, now 14

# Jupyter: Compare PCA to NMF for Endmember Determination

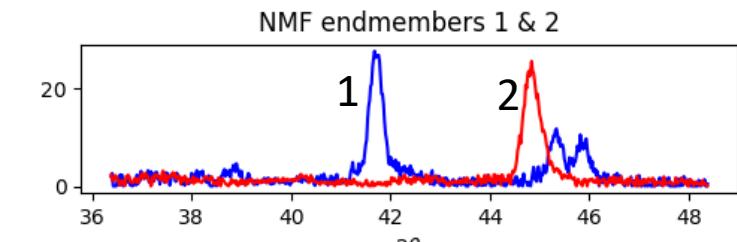
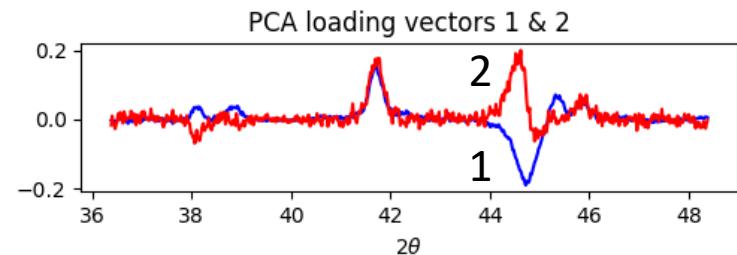
- PCA as Matrix Factorization

$$\hat{x}_n = \sum_{i=1}^M \mathbf{u}_i (x_n \mathbf{u}_i^T) \longrightarrow \mathbf{X} = \mathbf{U}\boldsymbol{\beta} + \epsilon$$

- Are  $\mathbf{u}_i$  constituent components/phases?



X-ray diffraction patterns



Prototypes!

# PCA for High Dimensional Data

- What if the number of variables D is large?

- Covariance matrix  $S$  is  $D \times D$

- Image of 1000 pixels by 1000 pixels = vector of 1E6
  - $S$  is  $1E6 \times 1E6$  (1 float = 4 bytes) or **4 Terabyte!**

- High computation and memory cost.

- Solution:

1. Pre-multiply by  $X$

$$\frac{1}{N} XX^T (\mathbf{X} \mathbf{u}_i^T) = \lambda_i (\mathbf{X} \mathbf{u}_i^T)$$

2. Define  $\mathbf{v}_i^T = \mathbf{X} \mathbf{u}_i^T$

$$\frac{1}{N} XX^T \mathbf{v}_i^T = \lambda_i \mathbf{v}_i^T$$

[ $N \times N$ ]

3. Regain  $\mathbf{u}_i$ :

$$\mathbf{u}_i^T = \frac{1}{(N\lambda_i)^{1/2}} X^T \mathbf{v}_i^T$$

N-1 eigenvalues haven't changed.  
( $S$  has  $(D-N+1)$  eigenvalues of 0)

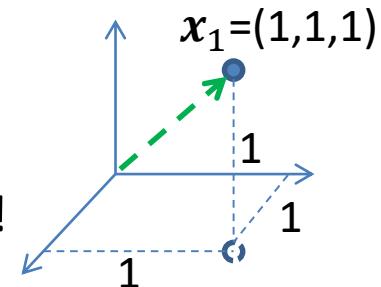
# Principle Component Analysis (PCA)

- Kernel PCA
  - Replace  $K = \mathbf{XX}^T$ ,  $K_{ij} = x_i x_j^T$  with  $K(x_i, x_j)$
  - E.g.  $K(x_i, x_j) = \exp(-\|x_i - x_j\|^2/\sigma^2)$
- Constrained / Sparse methods
  - e.g. sparse loading vectors (vectors with variables pushed to zero)

$$\min_{\theta, v} \sum_{i=1}^N \|x_i - \theta v^T x_i\|_2^2 + \lambda \|v\|_2^2 + \lambda_1 \|v\|_1$$

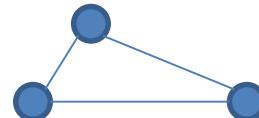
# The Challenge of Data Visualization

- Visualizing dissimilarity
  - Each data point exists in  $D$  dimensional space.
  - If  $D \leq 3$  can easily visualize in native space.
  - If  $D > 3$ , can visualize in  $N-1$  dimensional space!



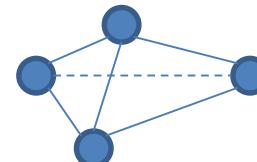
2 data points  
-> 1D

$$d(x_1, x_2)$$

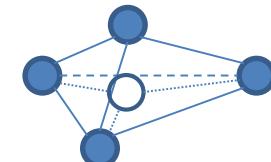


3 data points  
-> 2D

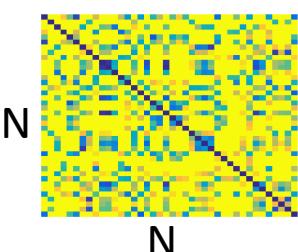
$$\begin{aligned} d(x_1, x_2) \\ d(x_1, x_3) \\ d(x_2, x_3) \end{aligned}$$



4 data points  
-> 3D



4 data points  
Flatten -> 2D

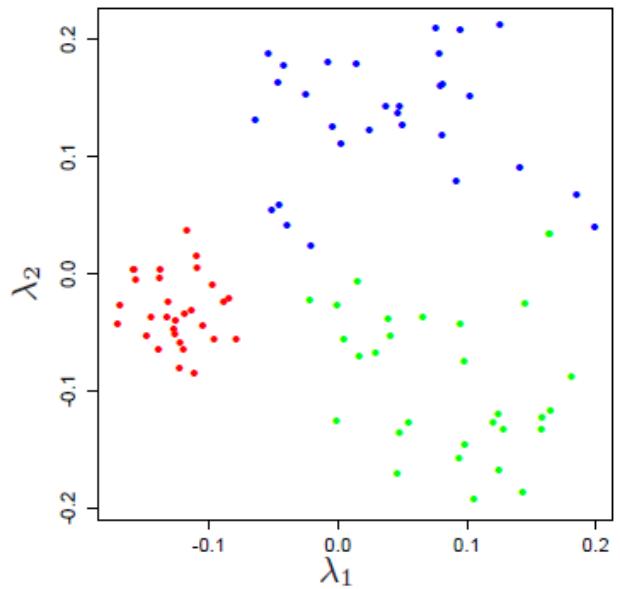
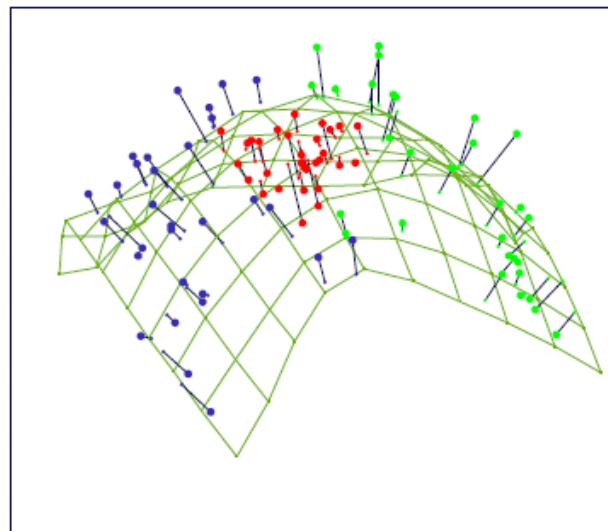
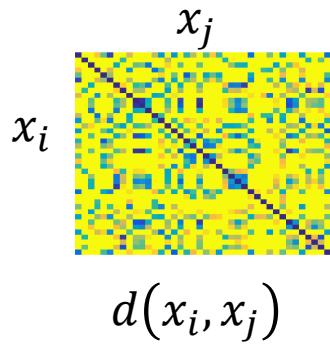


5+ data points?

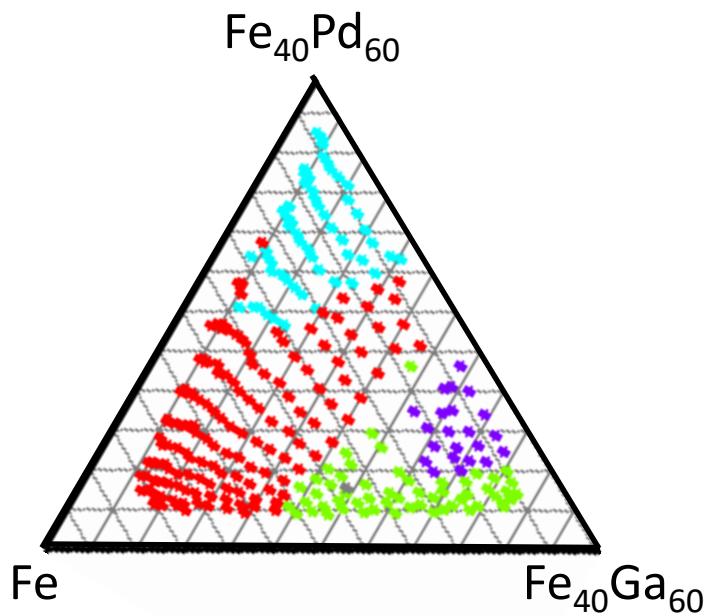
Use LVA to visualize sample number  $> 4$  and/or dimension  $> 3$

# Multidimensional Data Scaling (MDS)

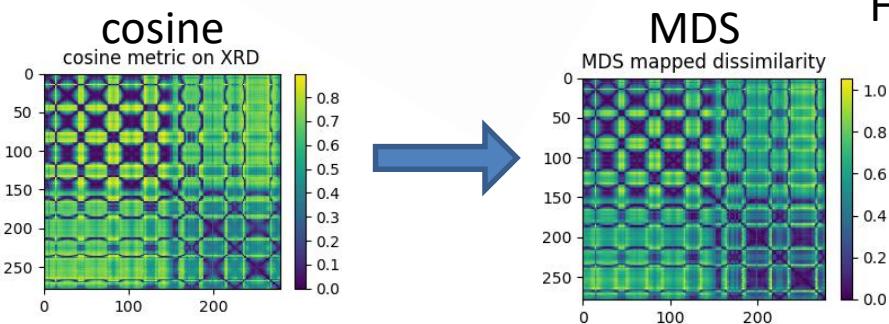
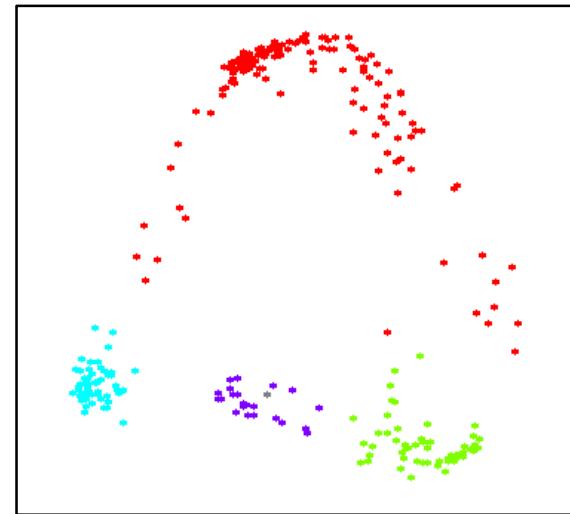
$$\text{Stress function} = \min \sum_{i,j} [d(x_i, x_j) - d(z_i, z_j)]^2$$



# Multidimensional Data Scaling (MDS)

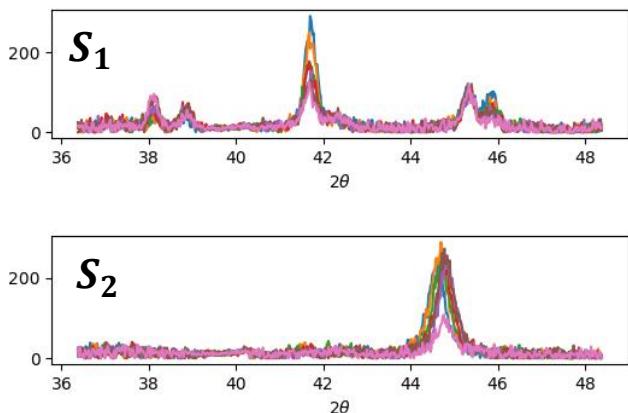
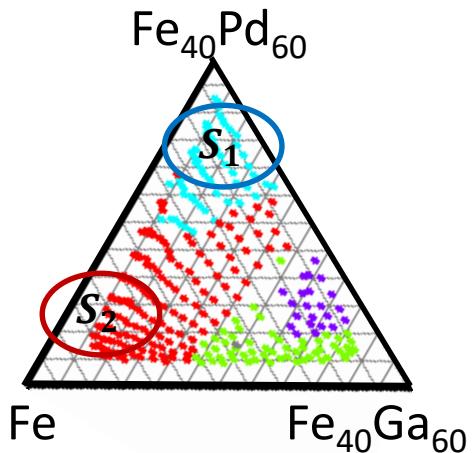


Multi-Dimensional Data Scaling

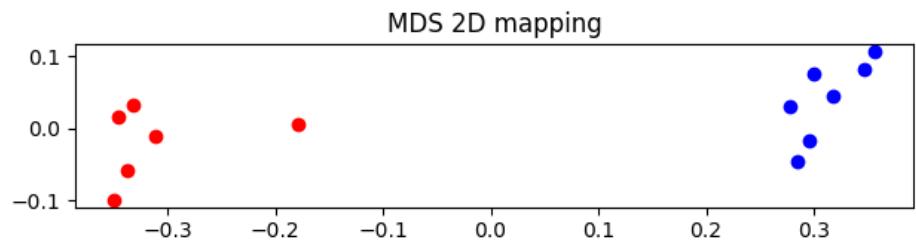


High dimensional data mapped down to 2D

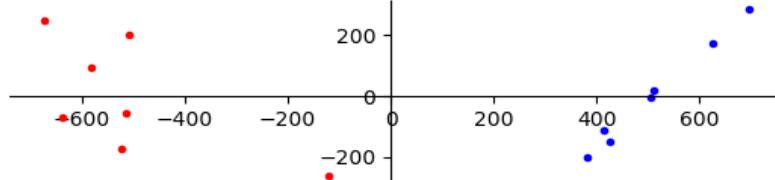
# Jupyter: Multidimensional Data Scaling (MDS)



High dimensional data mapped down to 2D

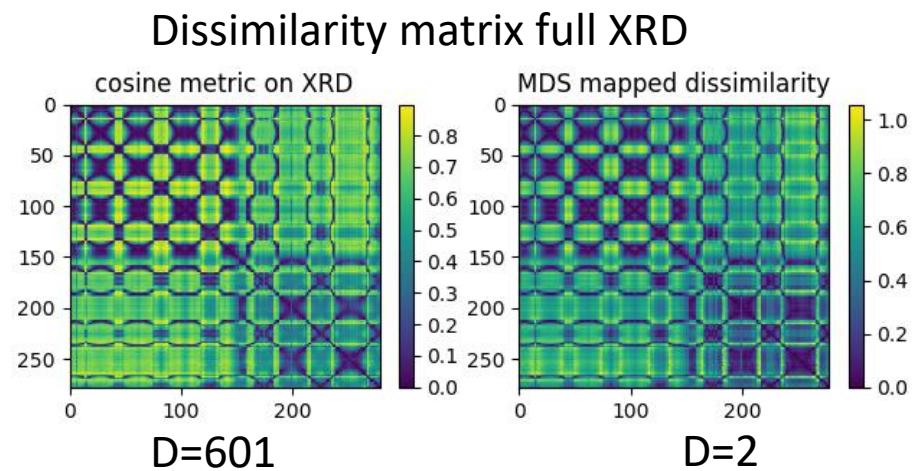
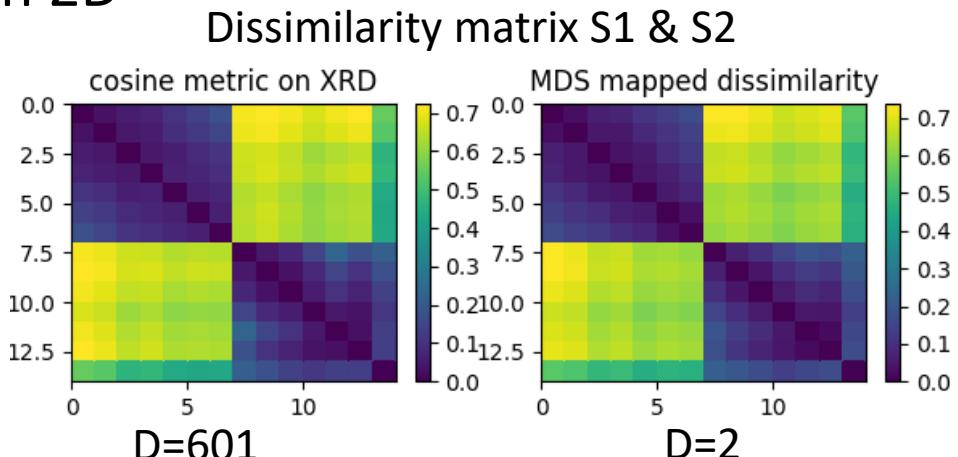
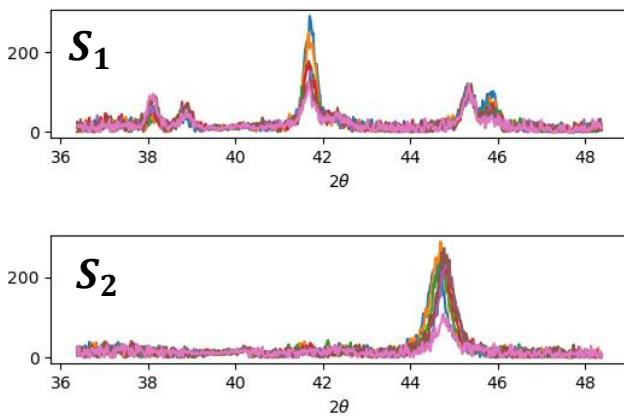
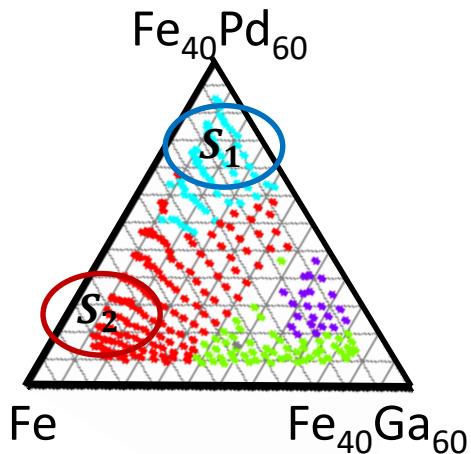


PCA 2D mapping



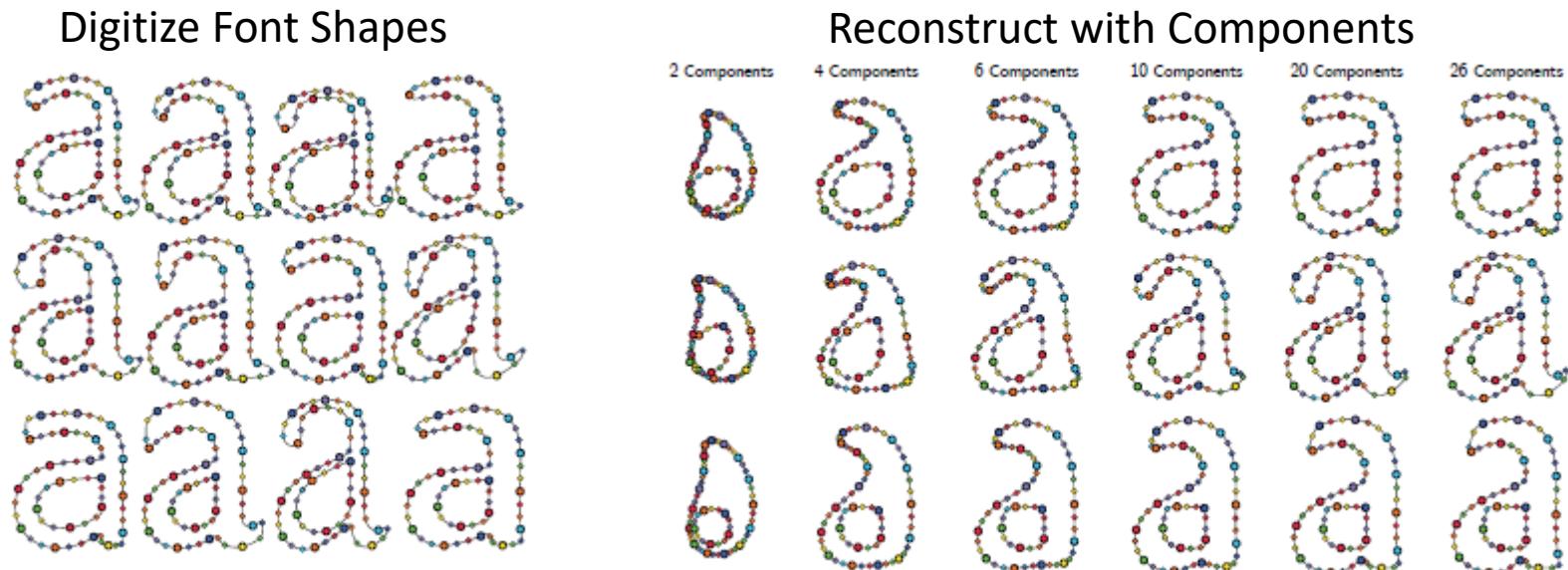
# Jupyter: Multidimensional Data Scaling (MDS)

- Take the information from the initial diss matrix and find the closest approximation in 2D



# Example: GP-LVM

- LVM with fonts:
- Paper: Campbell (2014) Learning a Manifold of Fonts



- [http://vecg.cs.ucl.ac.uk/Projects/projects\\_fonts/projects\\_fonts.html](http://vecg.cs.ucl.ac.uk/Projects/projects_fonts/projects_fonts.html)

# Questions?

# Clustering

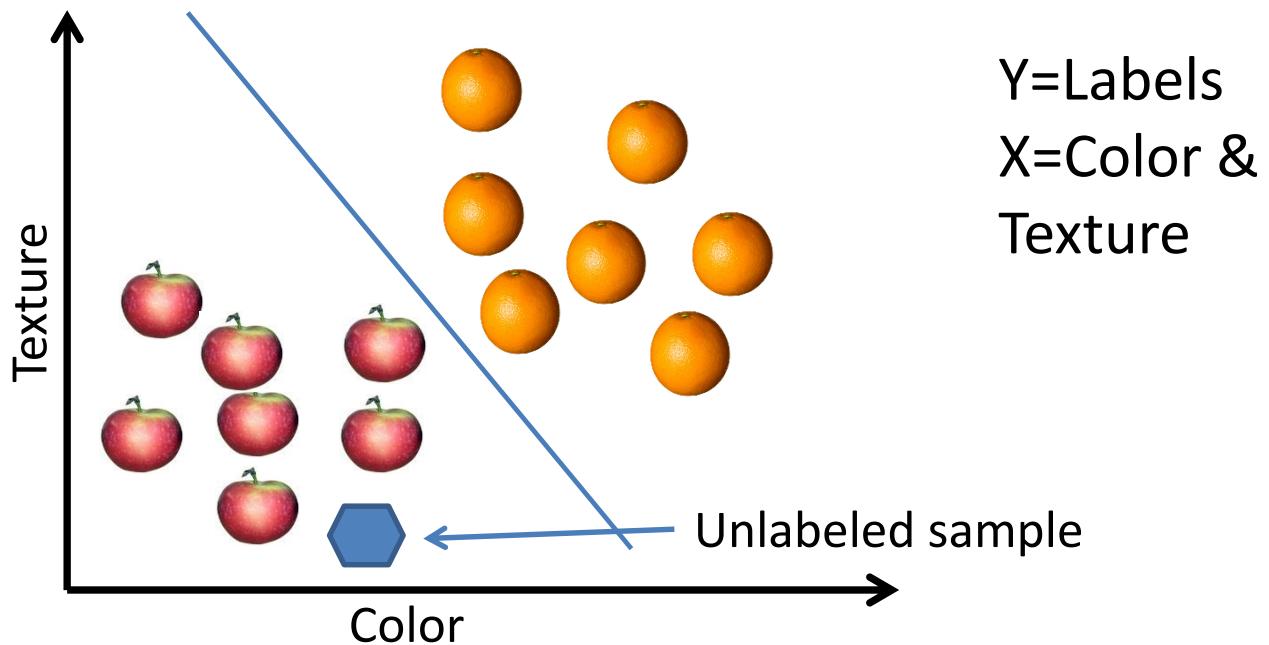
Group samples by similarity

Identify representative samples

Partition data space

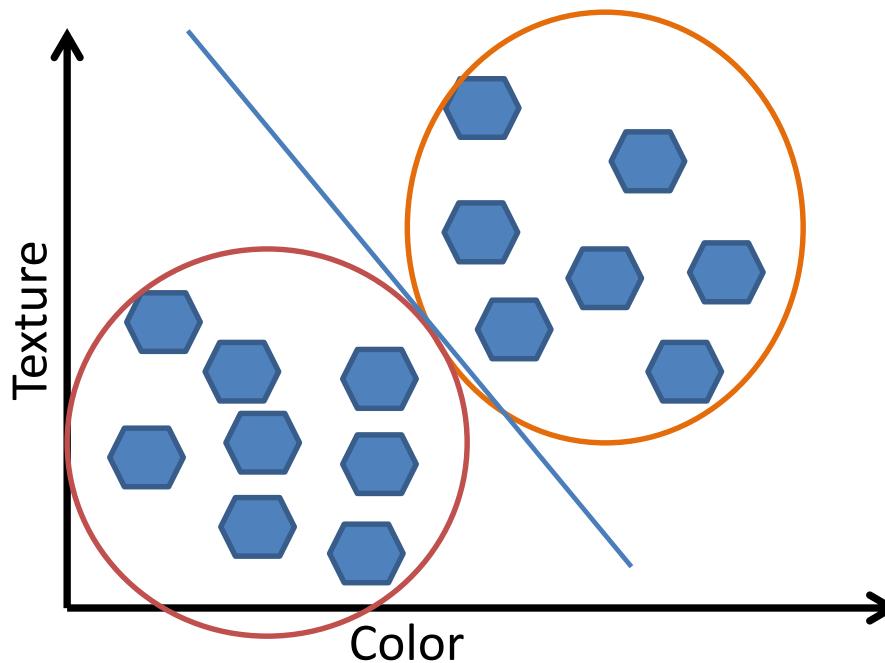
# Machine Learning

- Classification
  - Have labels for some points
  - Want a “rule” that will accurately assign labels to new points.



# Machine Learning

- Clustering
  - No labels
  - Samples in a cluster are more similar to each other than to the samples in other clusters.
- Partition the space into cluster regions.
- Identify representative data.



# Clustering

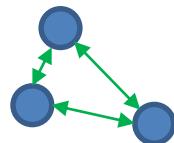
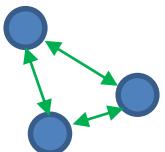
- Different types:
  - Combinatorial methods – cluster based on the data itself
  - Mixture Models
    - Assumes the data is drawn from a set of PDFs
    - Labels each sample with its parent distribution.
  - Mode-seeking
    - Assumes the data is drawn from a set of PDFs
    - Maps each sample to the mode of the parent distribution
- Some clustering methods
  - K-means
  - (Agglomerative) Hierarchical Cluster Analysis
  - Gaussian Mixture Model
  - Spectral Clustering

# Clustering

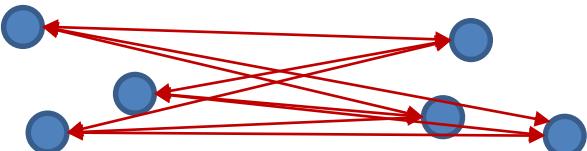
- Combinatorial Clustering Methods

- $T = \text{Sum of all distances between points} = \sum_{i,j} d(x_i, x_j)$
- $T = W(C) + B(C) \rightarrow W(C) = T - B(C)$

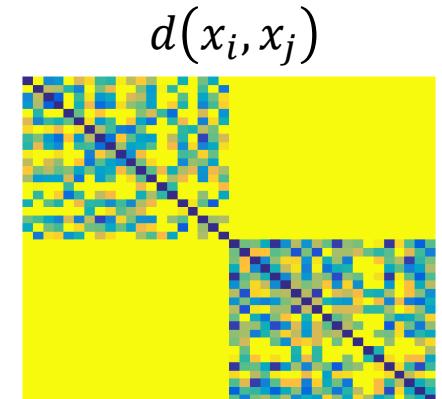
Find cluster labels  $C(i)$  so that:



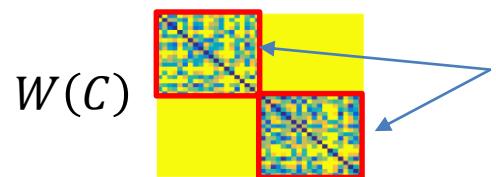
Minimize Within-Cluster scatter,  $W(C)$



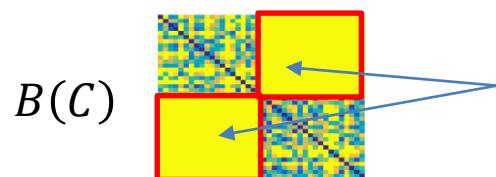
Maximize Between-cluster scatter,  $B(C)$



Dissimilarity Matrix  
for Clustered points!



Within-Cluster  
dissimilarities



Between-Cluster  
dissimilarities

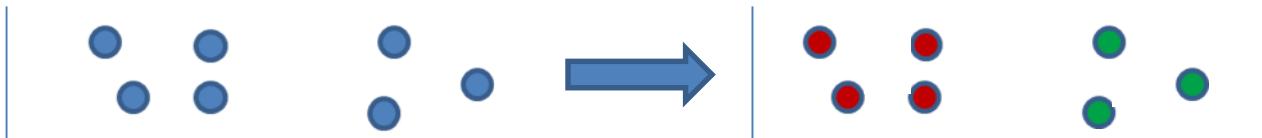
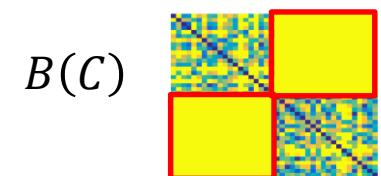
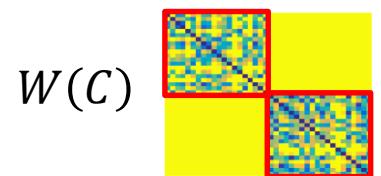
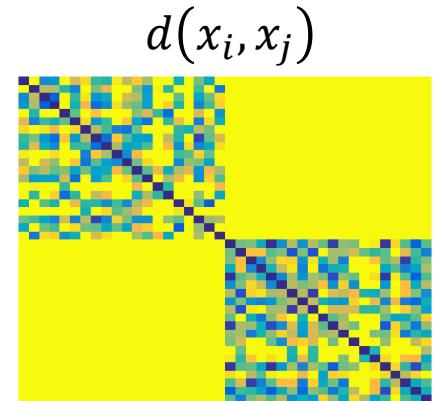
# Clustering

- Combinatorial Methods: Search through all possible assignments!

- Want to Minimize  $W(C)$  (and equivalently, maximize  $B(C)$ )
- The number of unique cluster assignments, for  $N$  samples and  $k$  clusters is:

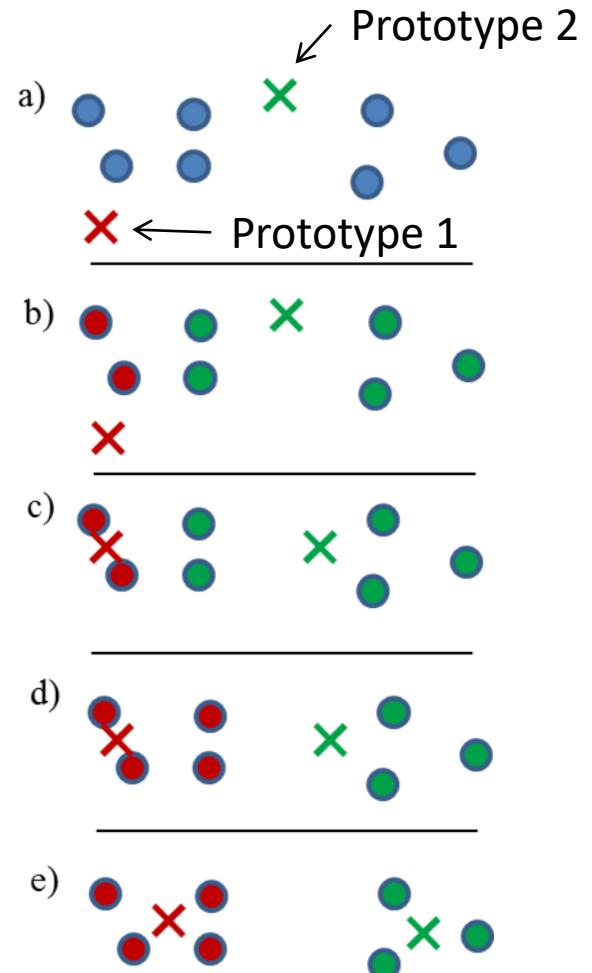
$$S(N, K) = \frac{1}{K!} \sum_{i=1}^k (-1)^{K-i} \binom{K}{i} k^N$$

- For  $N = 10$ , and  $k = 4 \rightarrow 34,105$  possibilities.
- For  $N = 19$ , and  $k = 4 \rightarrow \sim 10^{10}$
- Brute force is Not reasonable!
- Iterative greedy decent



# K-Means Algorithm

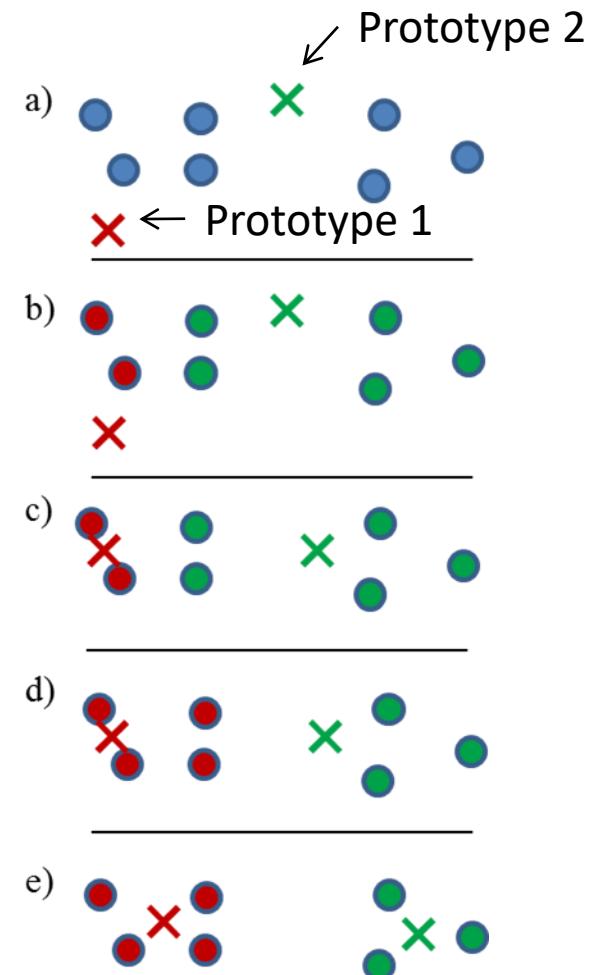
- Number of clusters  $k$  is known
  - Initialize the prototype for each cluster – often picked randomly
  - 1. Assign each point to the nearest prototype – defining clusters
  - 2. Compute the mean for each cluster.  
prototype <- mean
  - Iterate (1) and (2) until convergence.
  - Minimizes within cluster scatter
- $$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \|x_i - \bar{x}_k\|^2$$
- Prototype k
- Finds local minima. **Repeat** & choose result with lowest  $W$ .



# Combinatorial Method: K-Means Algorithm

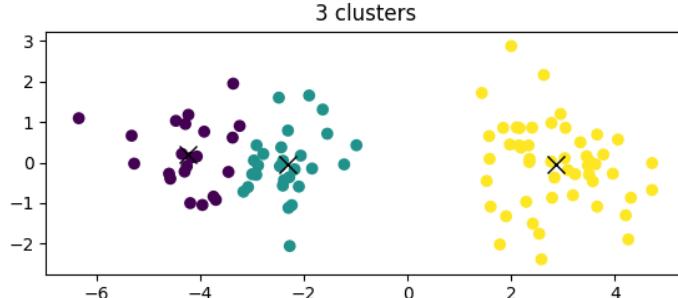
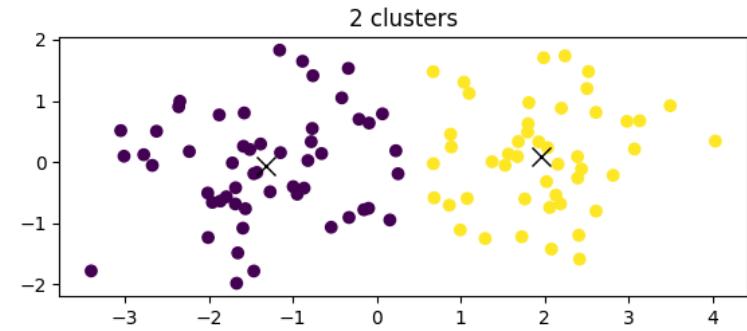
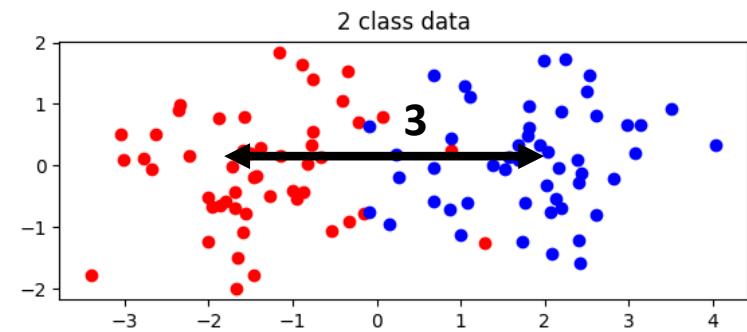
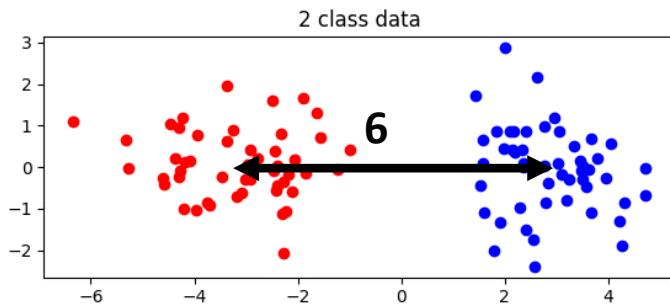
- Finds local minima. Repeat & choose result with lowest W.
- Requires Metric

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \|x_i - \bar{x}_k\|^2$$



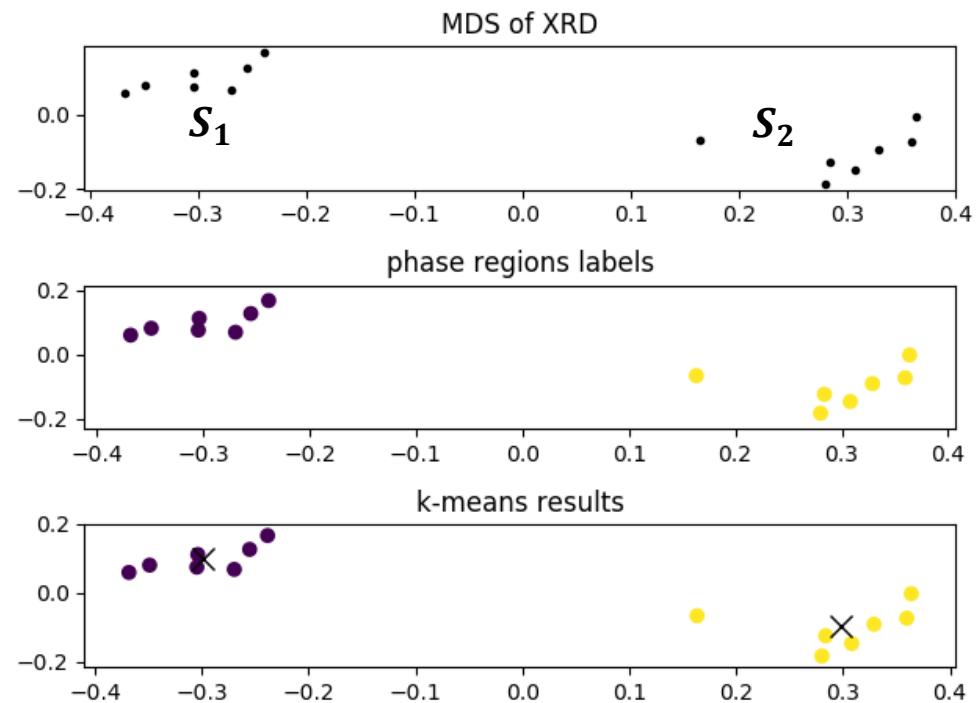
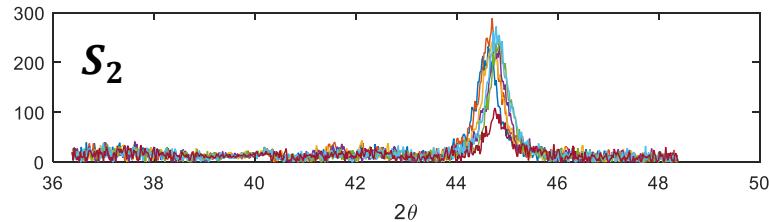
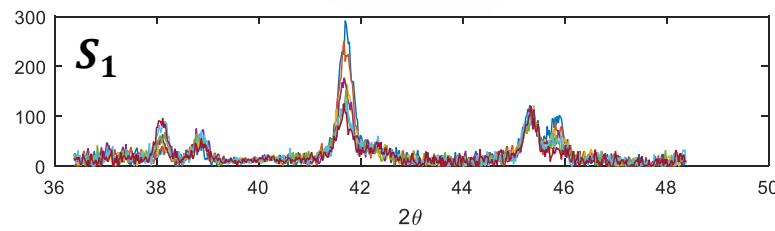
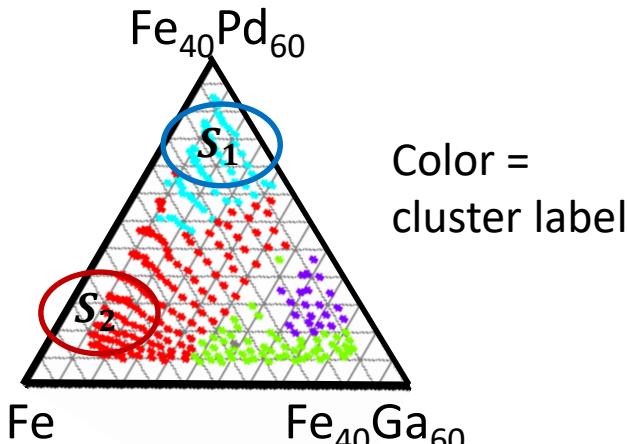
# Jupyter: K-Means

- Example, 2D Gaussian PDF



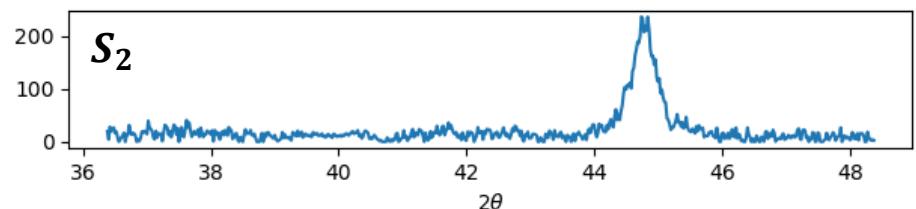
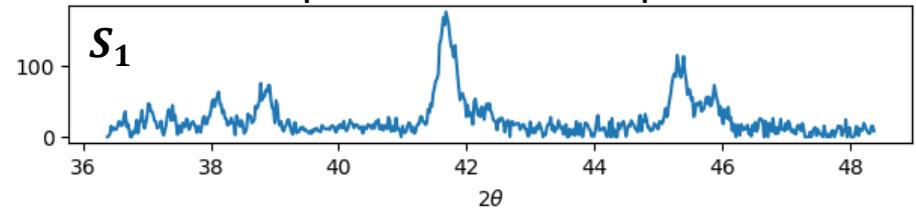
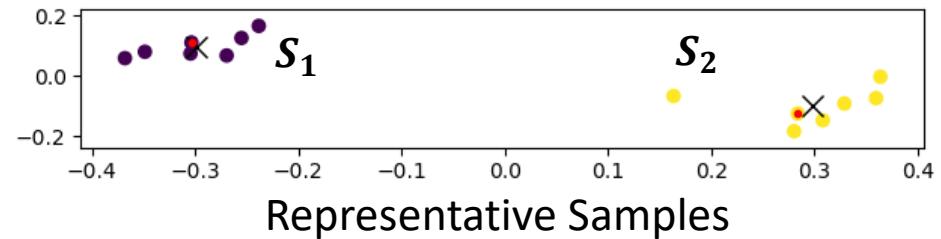
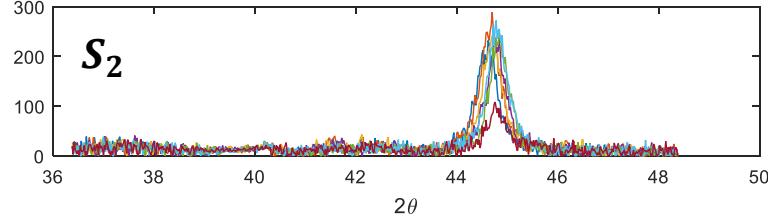
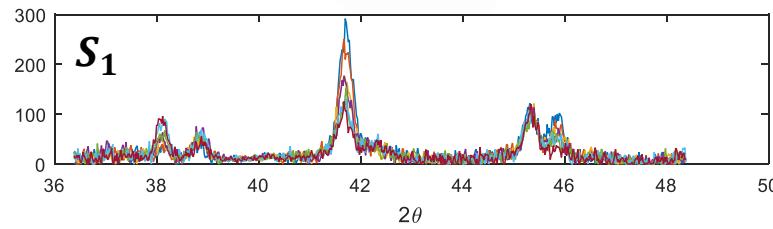
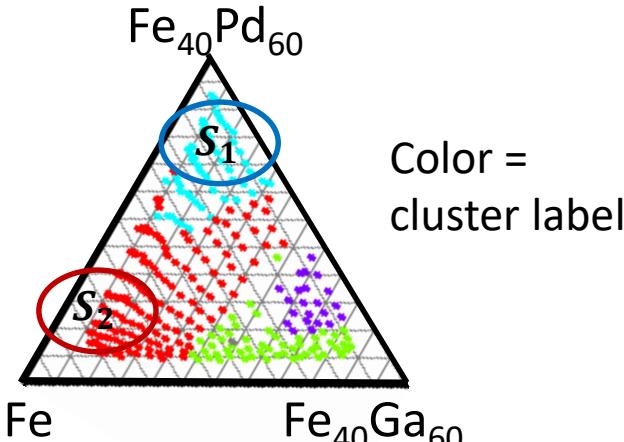
# Jupyter: K-Means

- Example



# Jupyter: K-Means for Selecting Representative Sample

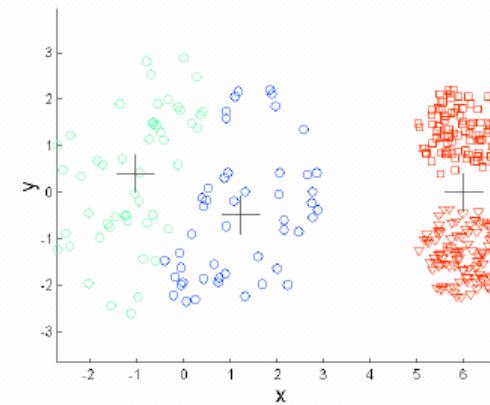
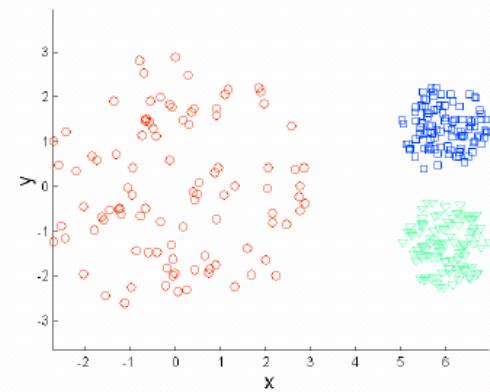
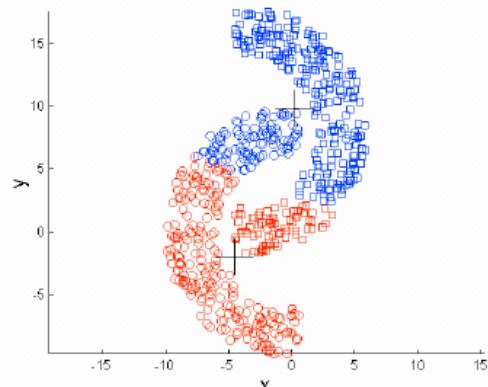
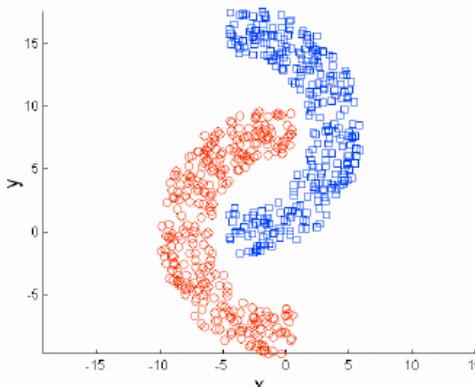
- Example



Samples closest to cluster center

# Combinatorial Method: K-Means Algorithm

- Issues:
  - Different cluster density
  - Different cluster size
  - Non-spherical cluster
  - Outliers
  - Empty cluster (e.g. due to sampling)



# K-Medoids

- Similar to k-means
- Useful when only the dissimilarity measure is available (doesn't have to be a metric).
- Initialize prototypes by picking data points.

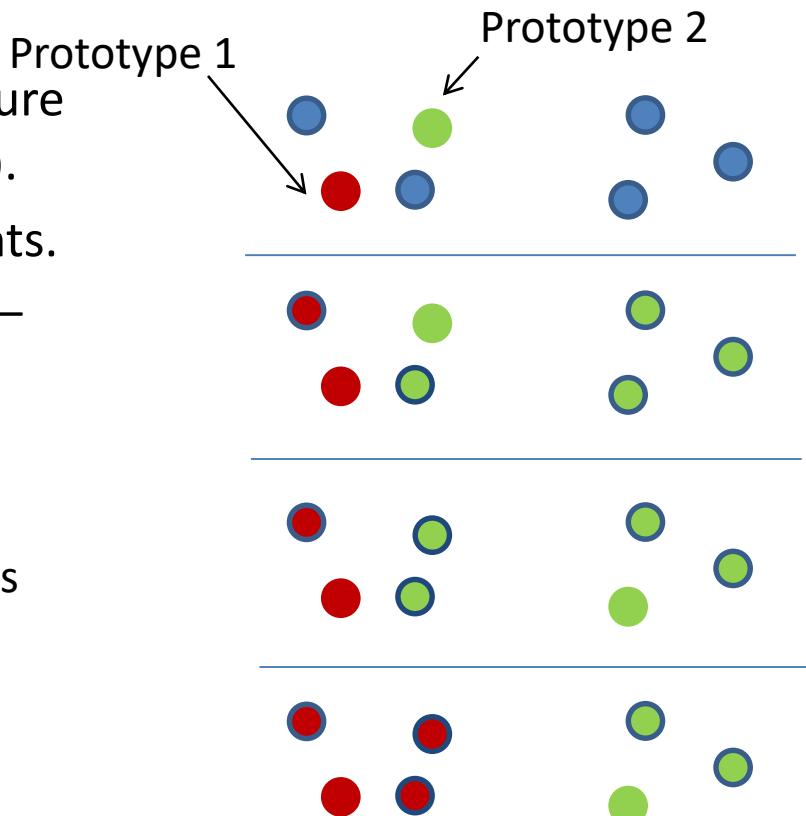
1. Assign each point to the nearest prototype – defining clusters.

$$C(i) = \operatorname{argmin}_k d(x_i, x_{i_k^*})$$

2. Find the point in each cluster that minimizes total distance to other points in cluster.

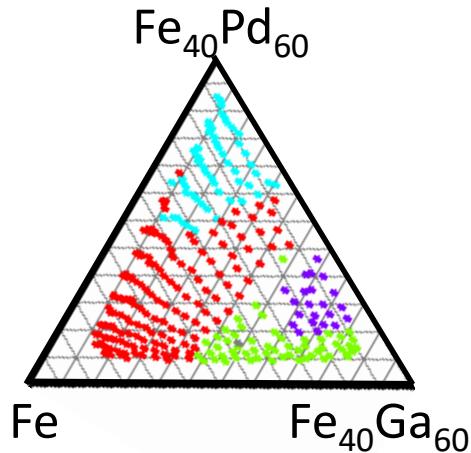
$$i_k^* = \operatorname{argmin}_{\substack{\{i: C(i)=k\} \\ C(i')=k}} \sum d(x_i, x_{i'})$$

- Iterate (1) and (2) until convergence.

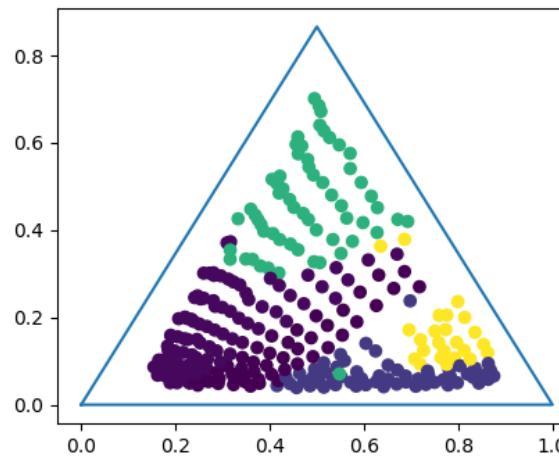
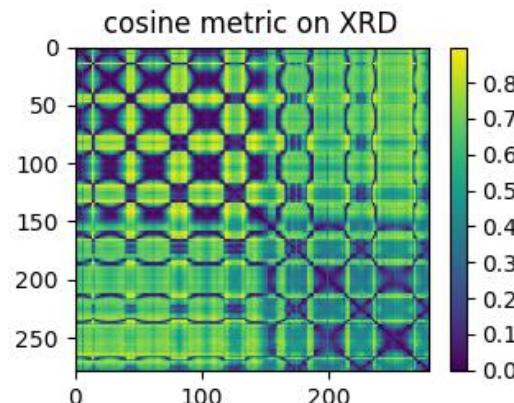
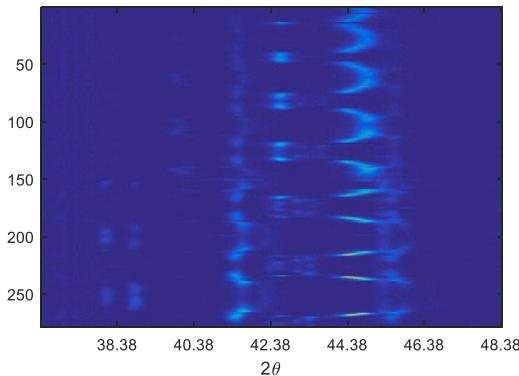


# K-Medoids

- Demo - Cluster original XRD data using Cosine Dissimilarity



X-ray Diffraction

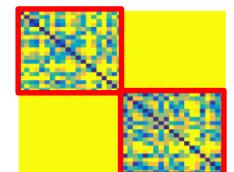


# Parameter Selection: How To Pick The Number of Clusters

- Classification
  - Can use training data to evaluate performance
  - Pick classification parameters that optimize performance
- Clustering
  - No training data
  - Performance often is subjective – requires expert's eye
- How to evaluate clustering results?
  - Combinatorial method: minimize within cluster scatter

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} d(x_i, x_j)$$

$$W(C)$$

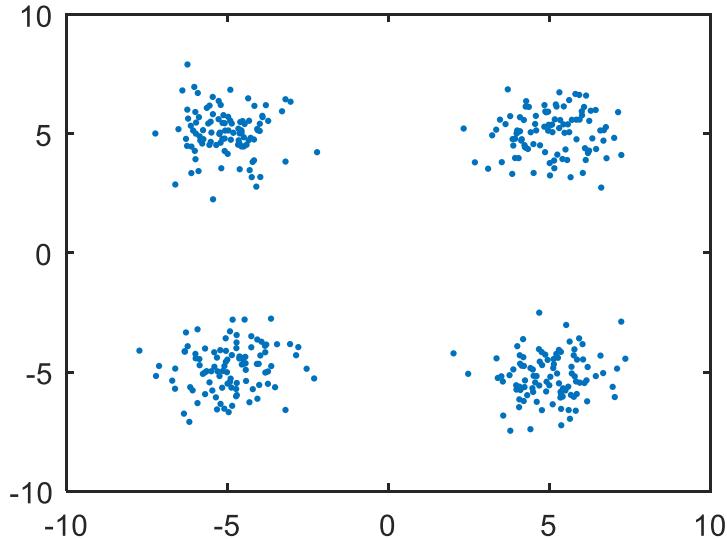


- Repeat clustering N times and return result with minimum W(C).

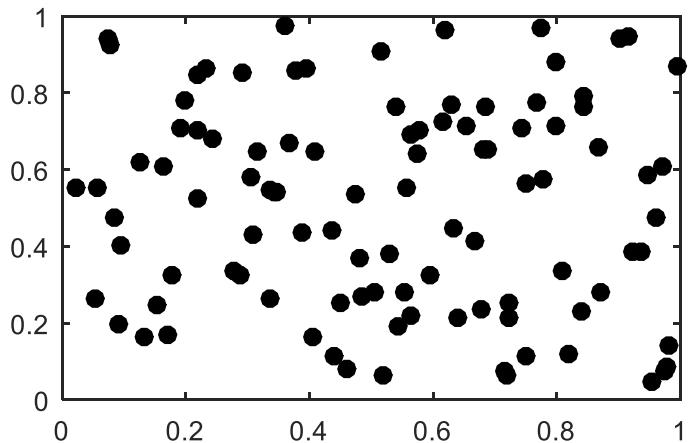
# Parameter Selection: How To Pick The Number of Clusters

- Gap Statistic
  - Heuristic based on within cluster scatter  $W_k$
  - K-Means!

Data from 4 Gaussians



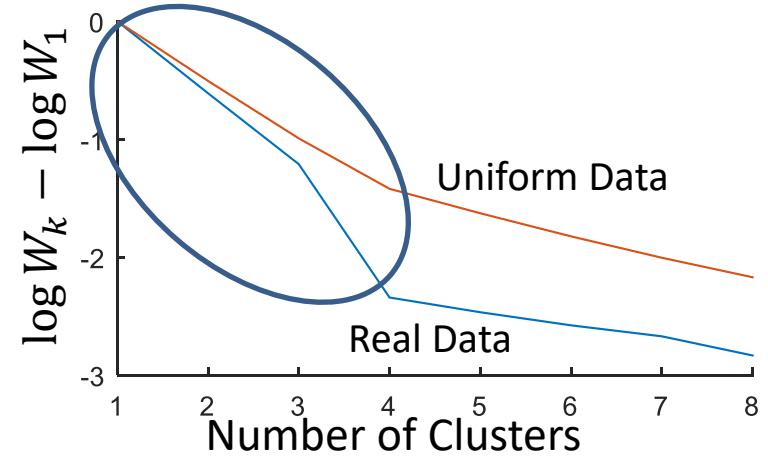
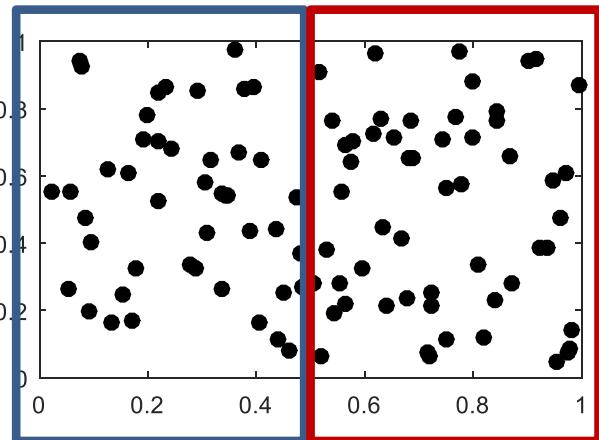
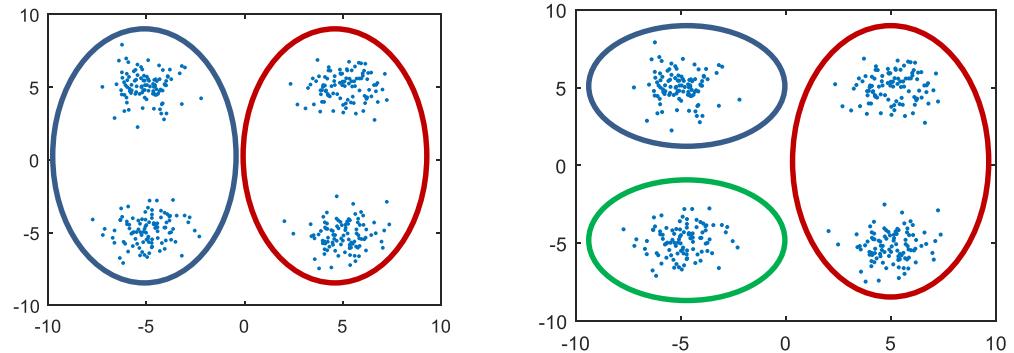
Reference data: Uniform distribution



# Parameter Selection: How To Pick The Number of Clusters

- Gap Statistic
  - Heuristic based on within cluster scatter  $W_k$

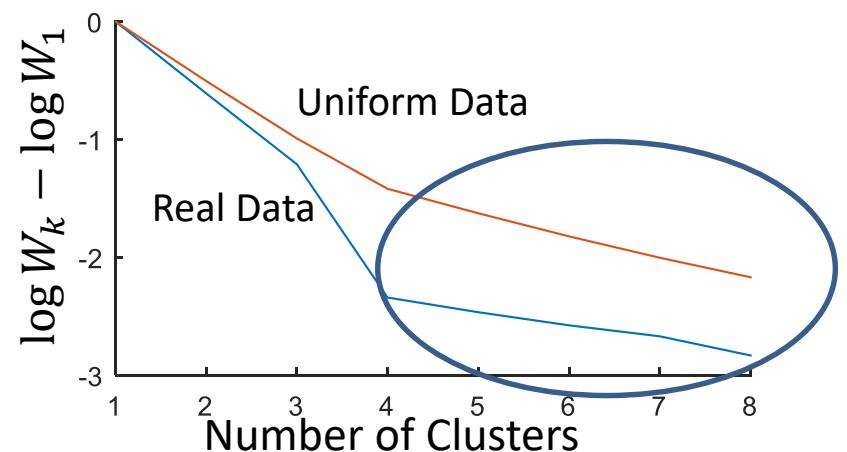
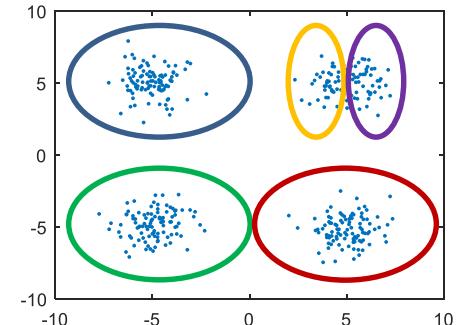
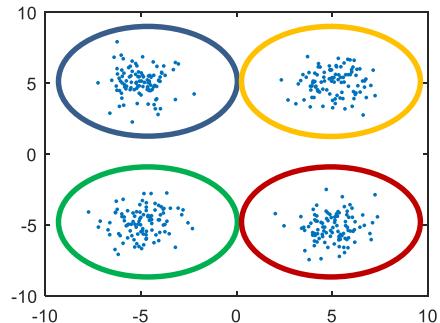
- $\hat{k} < k$ :
- as  $\hat{k}$  increases, more cluster centers  $\rightarrow W_k$  decreases
- Clustering becomes more accurate  $\rightarrow W_k$  decreases rapidly



# Parameter Selection: How To Pick The Number of Clusters

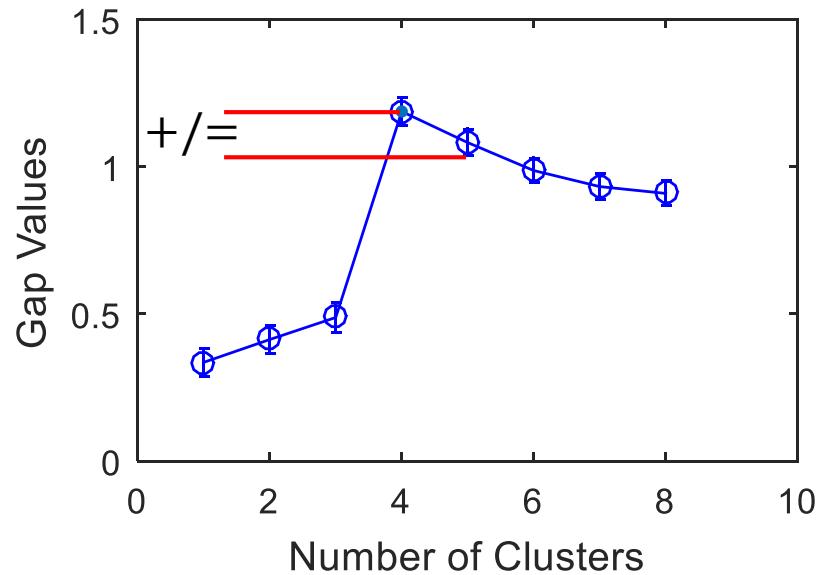
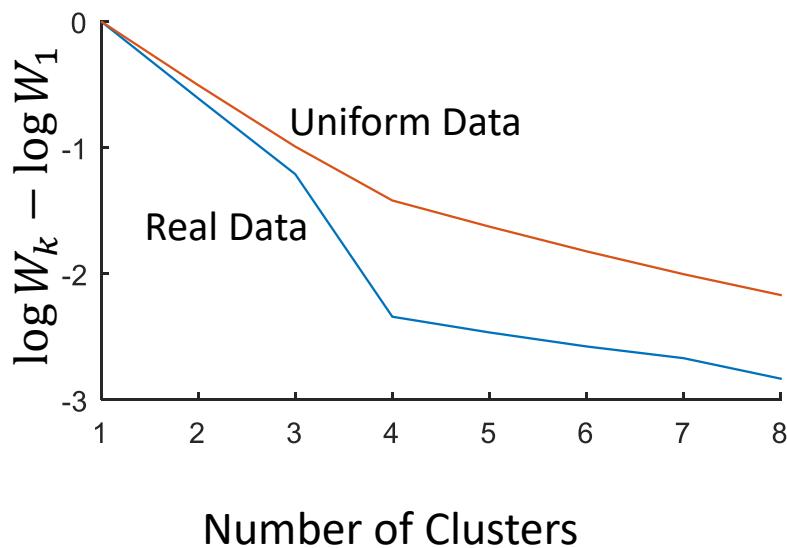
- Gap Statistic
  - Heuristic based on within cluster scatter  $W_k$

- $\hat{k} > k$ :
- as  $\hat{k}$  increases, more cluster centers  $\rightarrow W_k$  decreases
- Clustering is less accurate  $\rightarrow W_k$  decreases slower



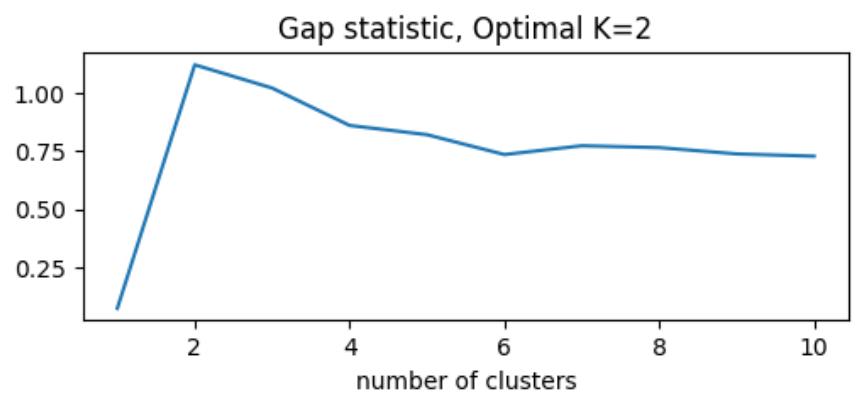
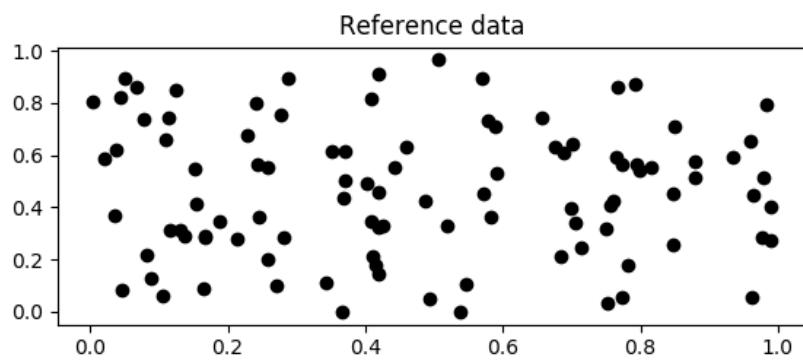
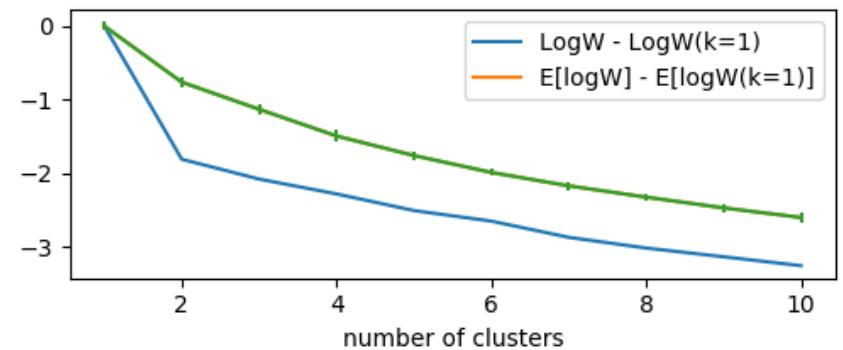
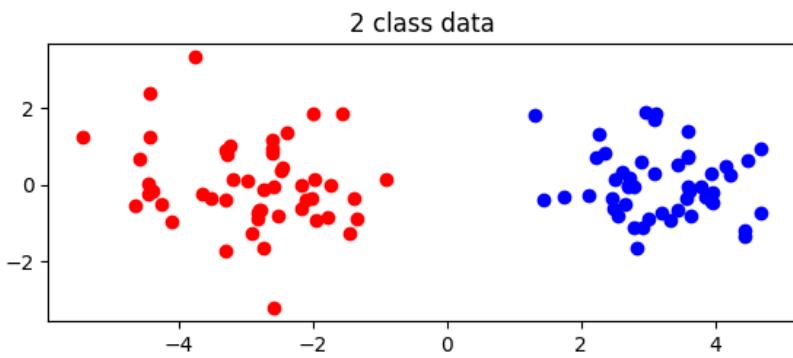
# Parameter Selection: How To Pick The Number of Clusters

- Compare Within cluster scatter  $W_k$  between clustering on real data and reference data.
- Compute “Gap”
- $K^* = \operatorname{argmin}_K [G(K) \geq G(K + 1) - \text{std}(\log W_{K+1})]$



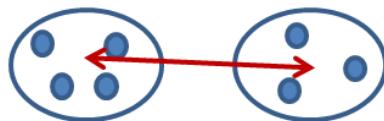
# Jupyter: Gap Statistic

- Demo – 2 Gaussian clusters

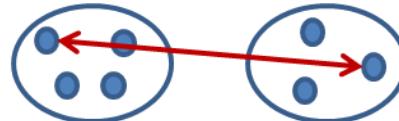


# Hierarchical cluster analysis

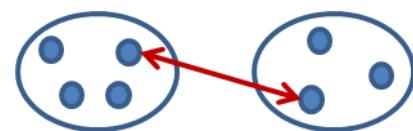
- Agglomerative hierarchical clustering
  - Don't need to define number of clusters  $k$
  - Need dissimilarity matrix
  - Define dissimilarity between clusters



Average: distance  
between cluster  
centers



Complete:  
distance between  
furthest points



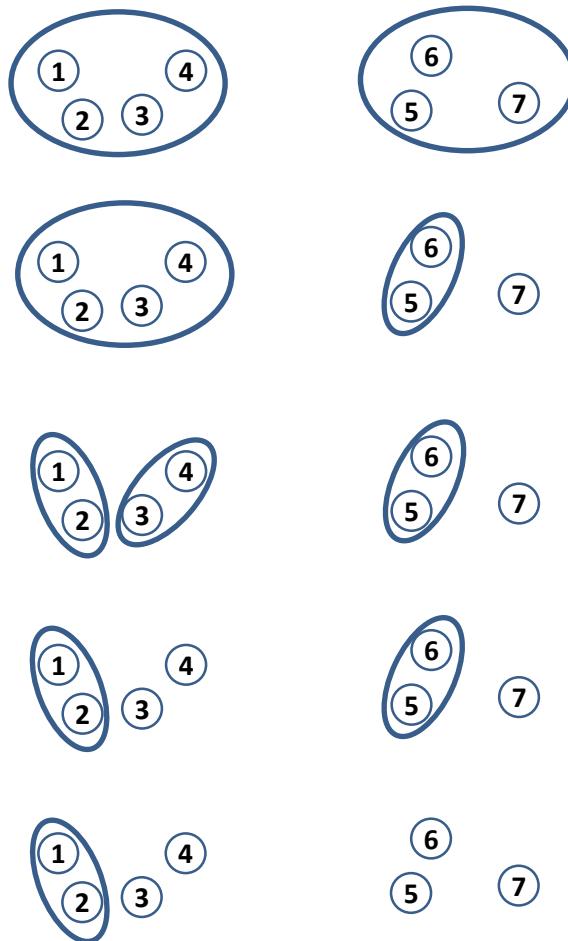
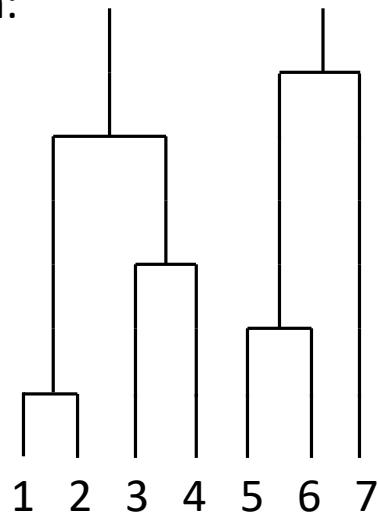
Single: distance  
between closest  
points

$$\text{Ward: } d(A, B) = W_{A \cup B} - W_A - W_B$$

# Hierarchical cluster analysis

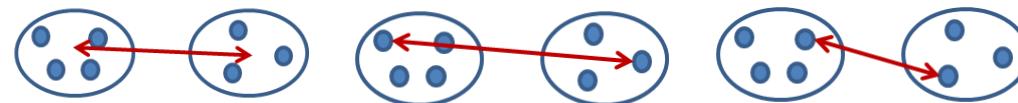
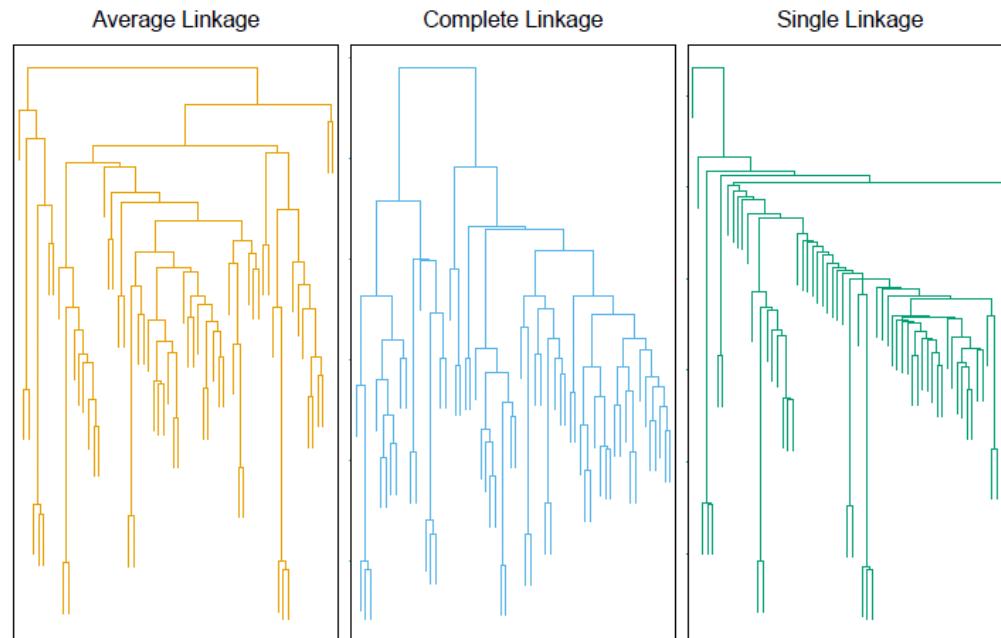
- Agglomerative hierarchical clustering

Dendrogram:



# Agglomerative Hierarchical Cluster Analysis

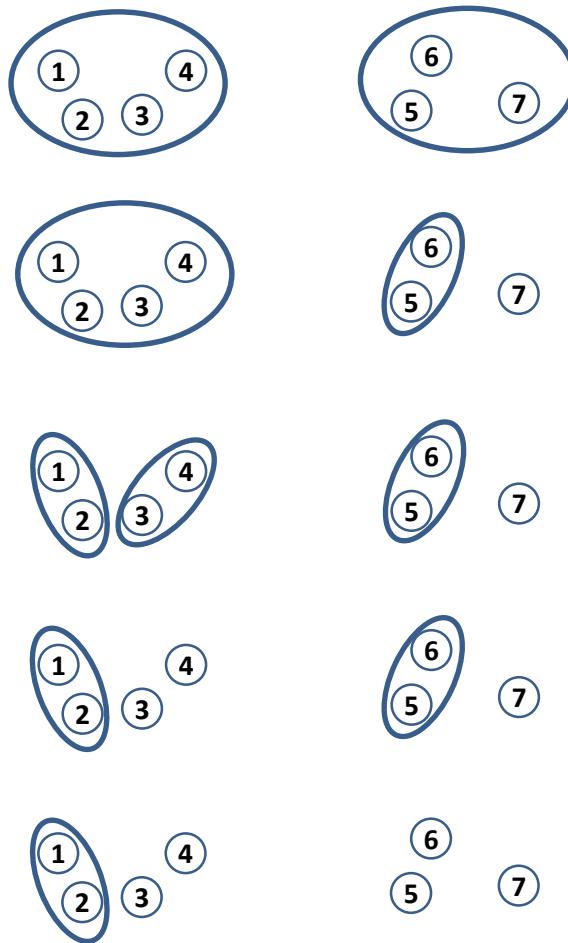
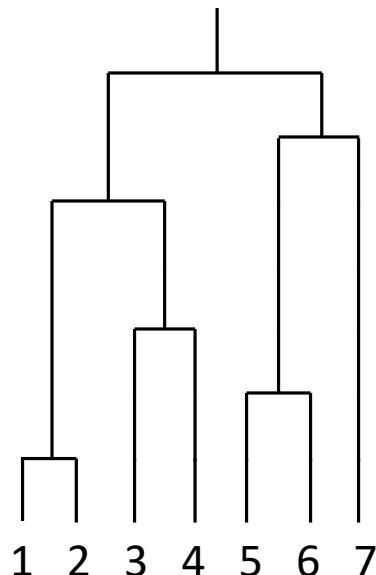
- Different types of Agglomerative HCA



Comparison: Hastie, et al. Elements of Statistical Learning

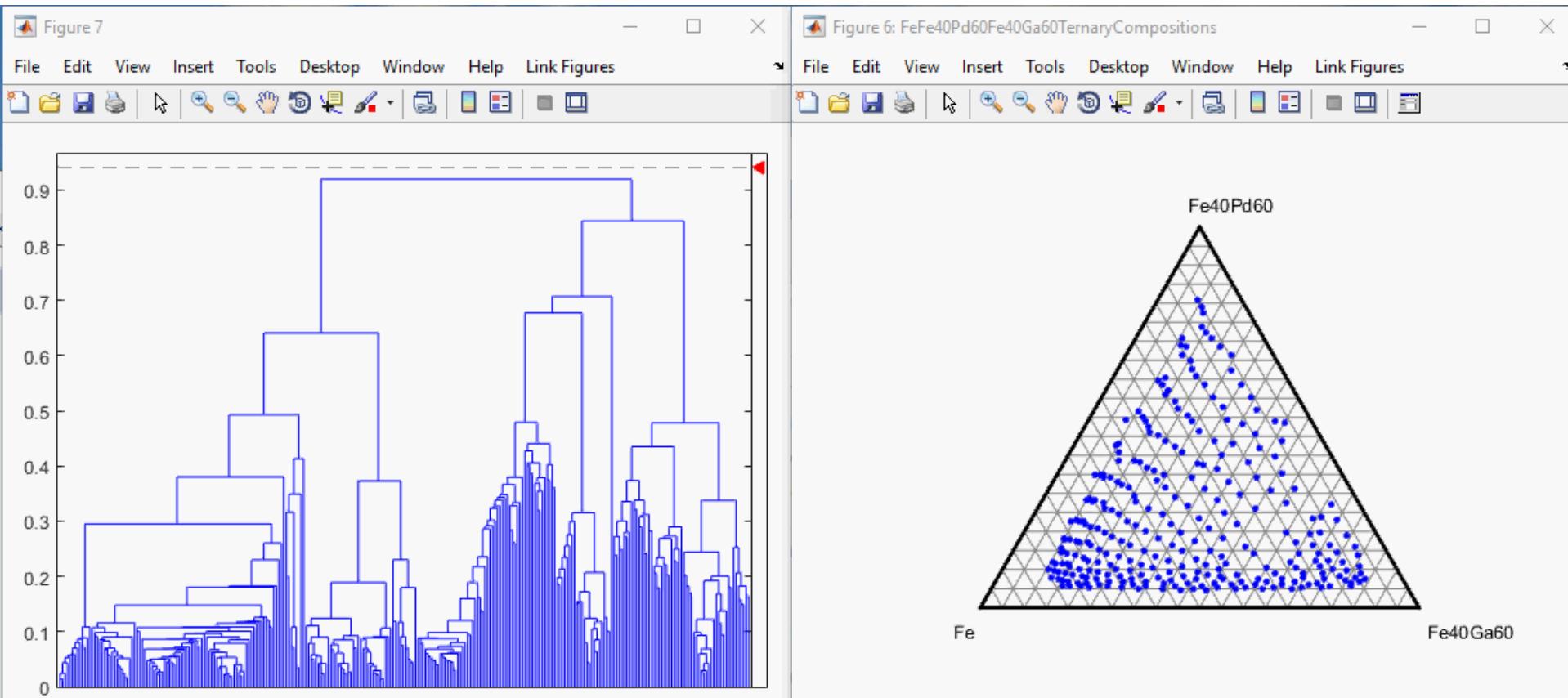
# Hierarchical cluster analysis

- Cluster based on dissimilarity
  - Great way to visualize choice of dissimilarity measure.
- Don't need to choose number of clusters



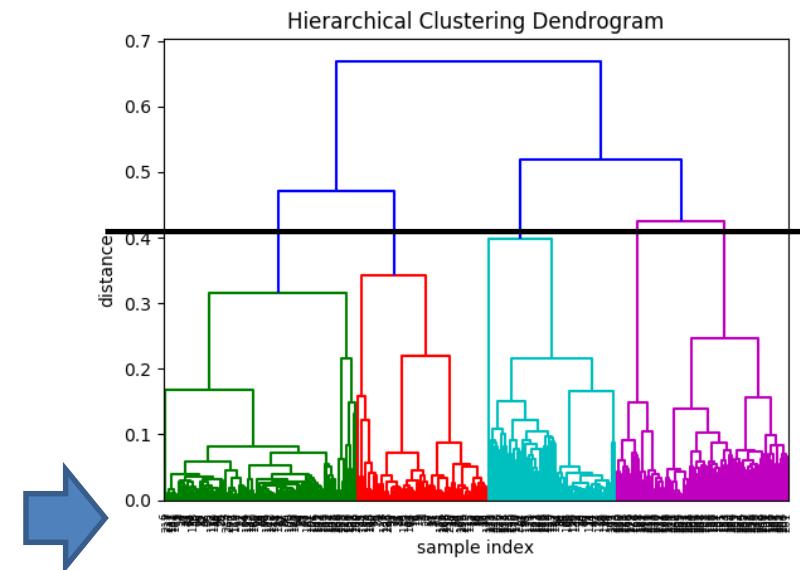
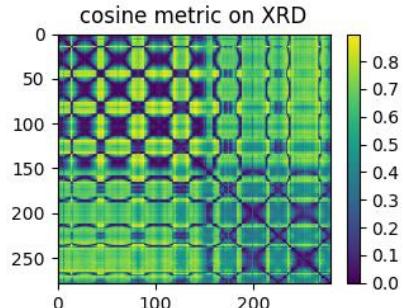
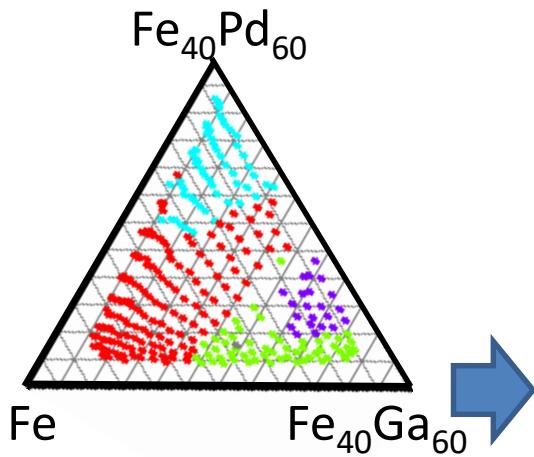
# Hierarchical Cluster Analysis

- CombiView video of changing # of clusters

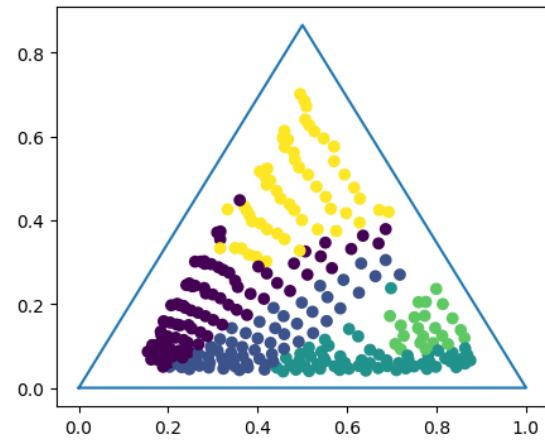
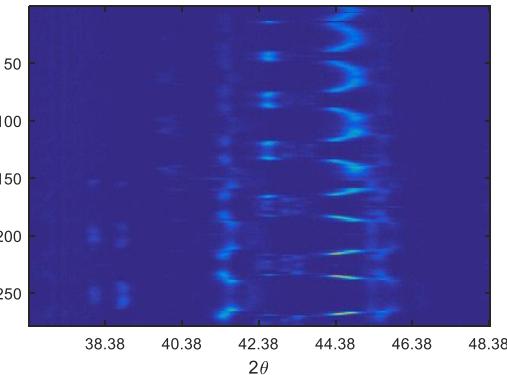


# Jupyter: Hierarchical Cluster Analysis

- Demo

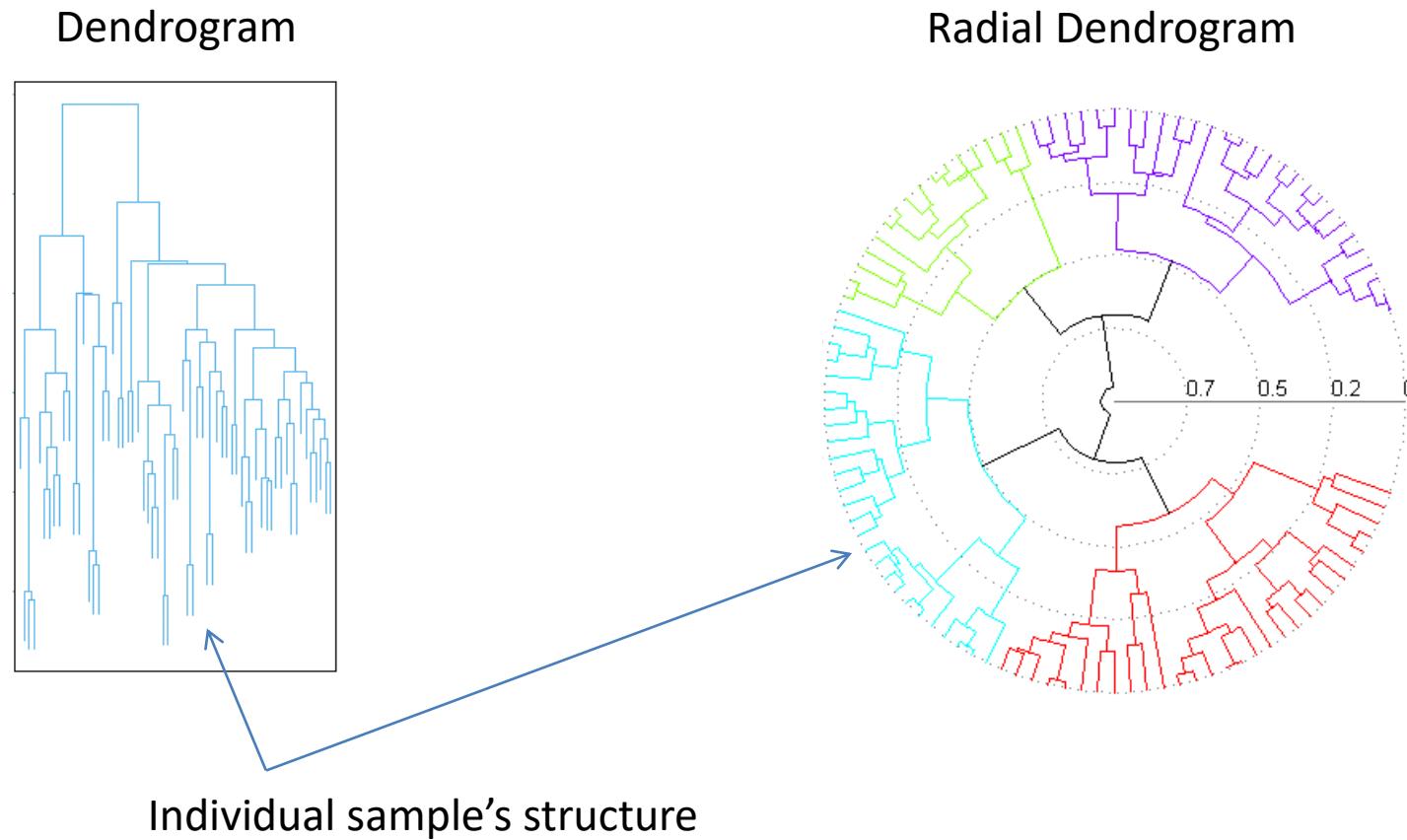


## X-ray Diffraction



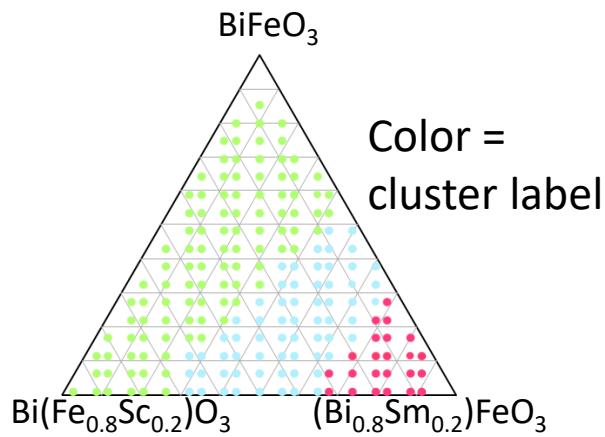
# Hierarchical Cluster Analysis

- Visualization Tool

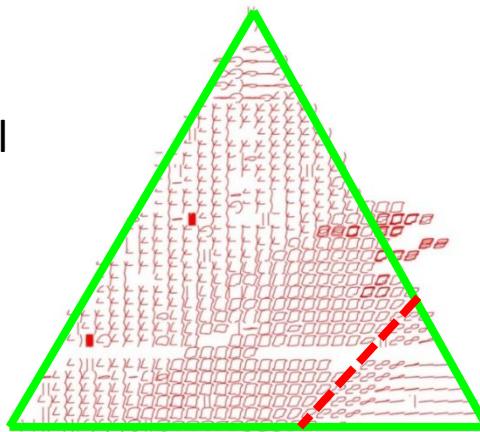


# Radial Dendrogram: Ferroelectrics Relating Structure to Properties

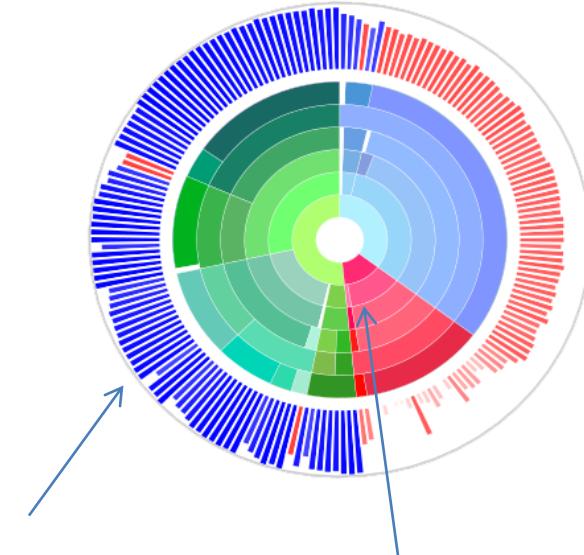
Cluster Raman Data



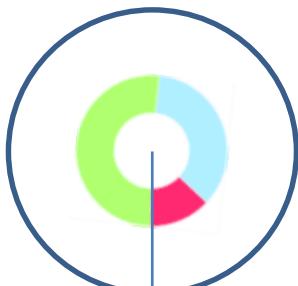
Functional Property:  
Ferroelectric Coercivity



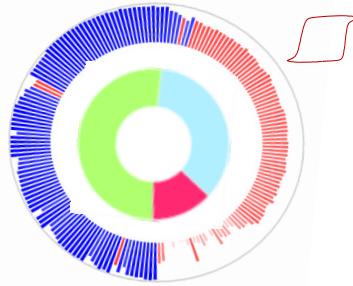
Sunburst



Visualize Samples on a Ring



Sample 1 @ 0°



Functional  
property

Samples  
clustered by  
structure

Structure is correlated to FE Coercivity!

# Mixture Model: Gaussian

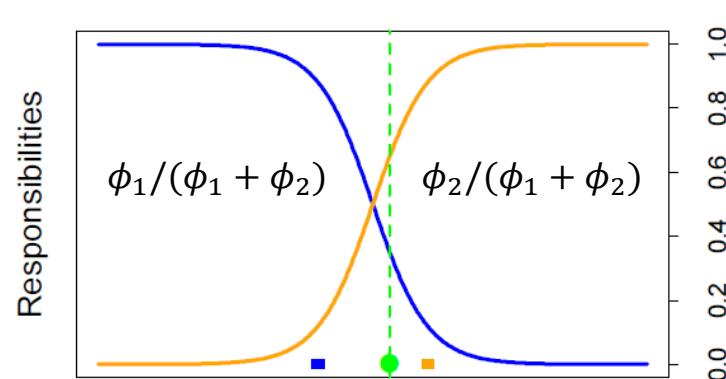
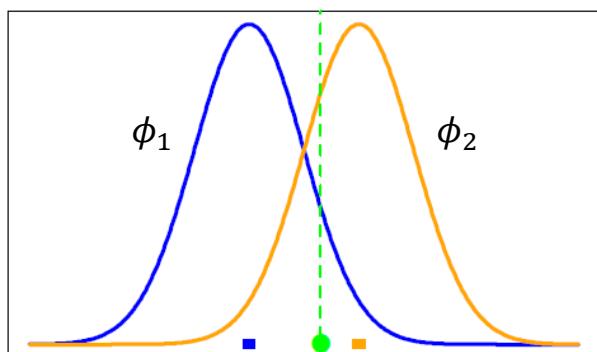
- Assume that the data is drawn from a set of  $k$  PDFs. In this case, Gaussian.
- Each PDF is associated with a cluster. ( 1 Cluster  $\leftrightarrow$  1 PDF )

$$p(x) = \sum_{j=1}^k \phi_j = \sum_{j=1}^k \pi_j N(\mu_j, \sigma_j^2)$$

- Match each point to its generating PDF  $\rightarrow$  Cluster labels
- Example – 2 Gaussians:

$$p(x) = \pi_1 N(\mu_1, \sigma_1^2) + \pi_2 N(\mu_2, \sigma_2^2) \quad \pi_1 + \pi_2 = 1$$

- Responsibility  $w_{i,1} = \frac{\phi_1}{\sum_{l=1}^k \phi_l} = \phi_1(x)/(\phi_1(x) + \phi_2(x))$



# Mixture Model

- Solve with Iterative greedy decent.

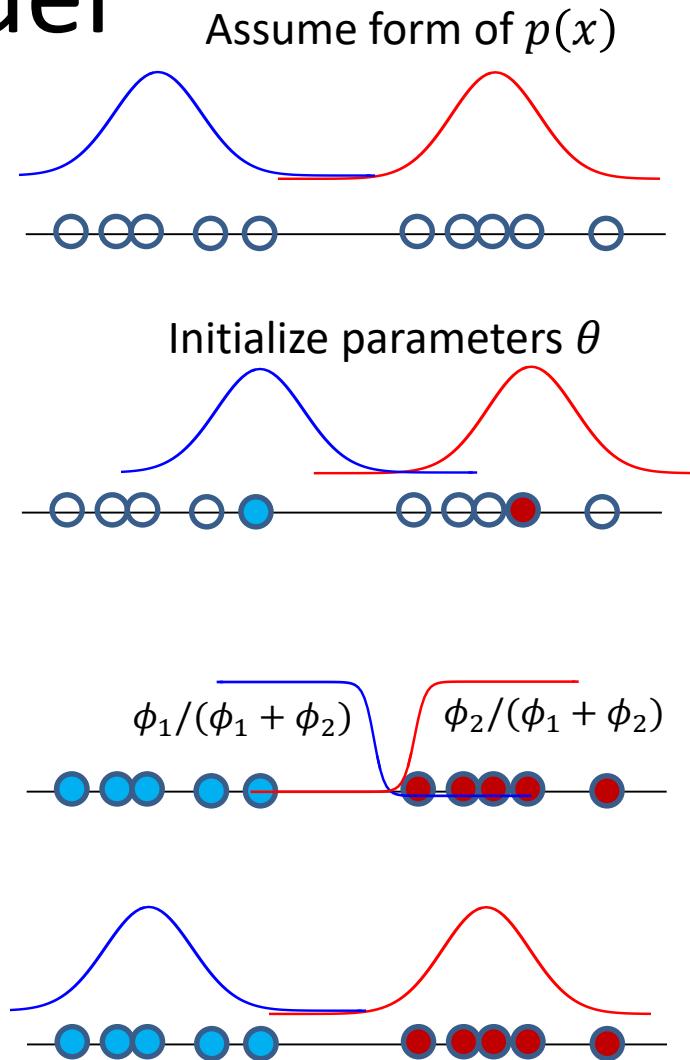
$$p(x) = \pi_1\phi_1(x) + \pi_2\phi_2(x)$$

- Unknowns:  $\theta = \{\pi_1, \mu_1, \sigma_1, \pi_2, \mu_2, \sigma_2\}$
- Initialize:  $\pi_j = .5, \sigma_j = \sigma_T, \mu_j \rightarrow$  random pts
- 1. Compute PDF responsibility for each point & Assign each point to PDF with greatest responsibility.
- 2. Recompute  $\theta$ s using responsibility  $w_{i,j}$  as weight:

$$\pi_j = \frac{1}{N} \sum_{i=1}^N w_{i,j} \quad \mu_j = \frac{\sum_{i=1}^N w_{i,j} x_i}{\sum_{i=1}^N w_{i,j}}$$

etc.

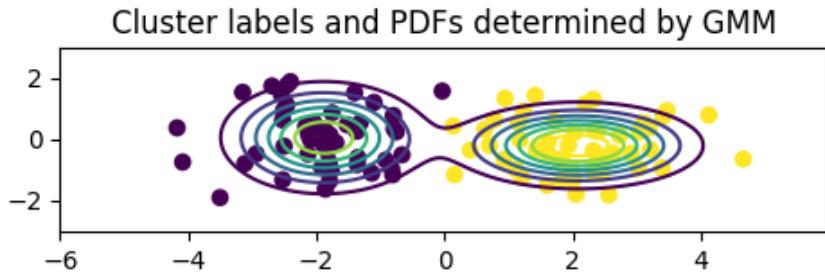
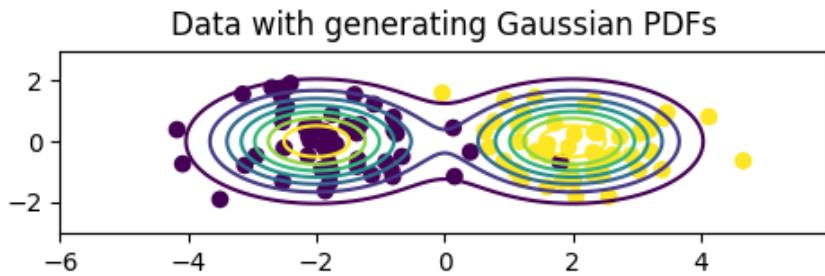
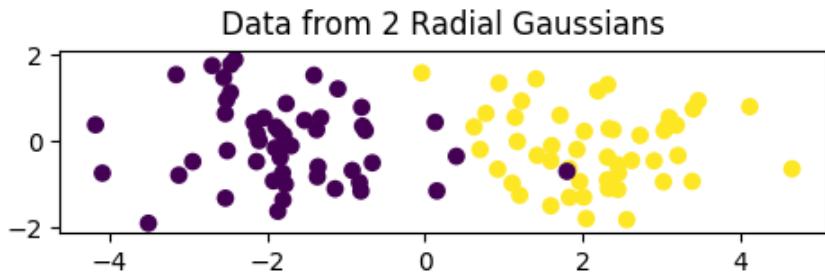
- Repeat 1 & 2 until convergence.



# Jupyter: Gaussian Mixture Model

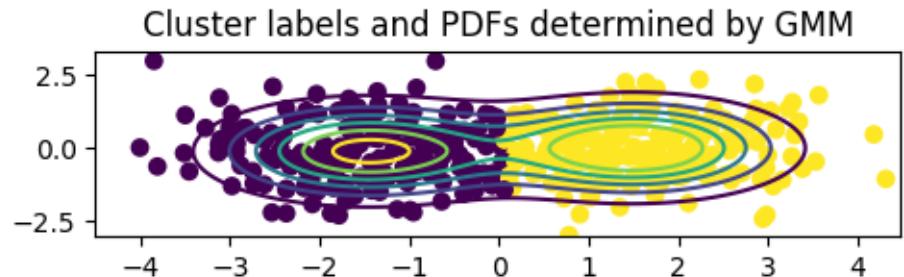
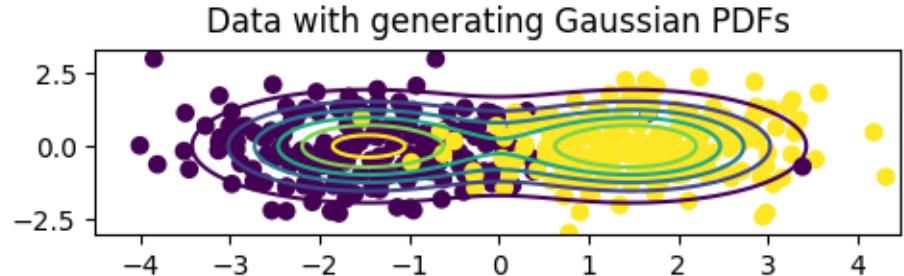
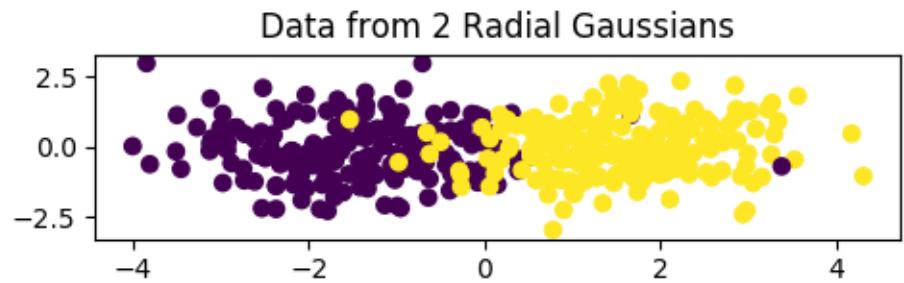
100 samples:

$$p(x) = .5N(-2,1) + .5N(2,1)$$



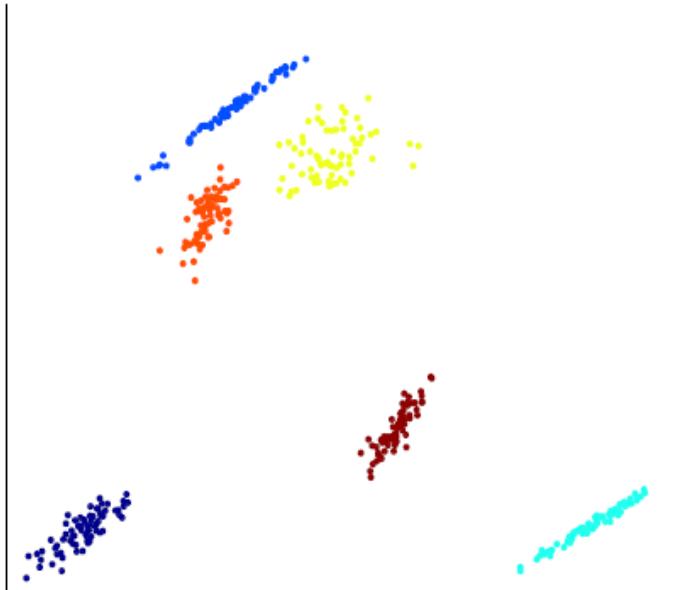
400 samples:

$$p(x) = .5N(-1.5,1) + .5N(1.5,1)$$

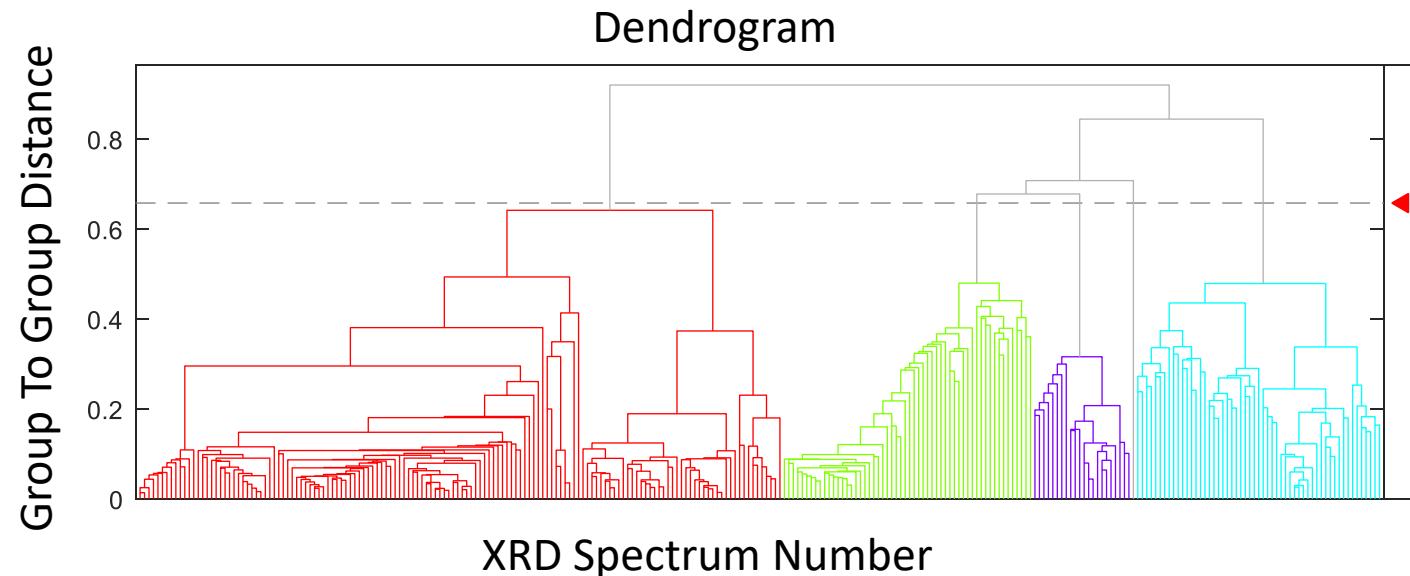


# Combinatorial Method: Spectral Clustering

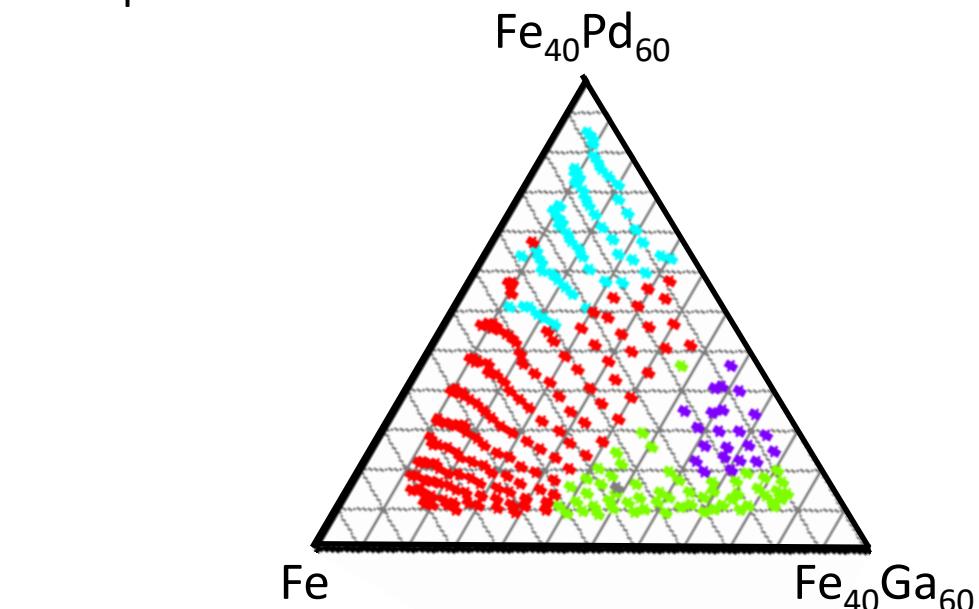
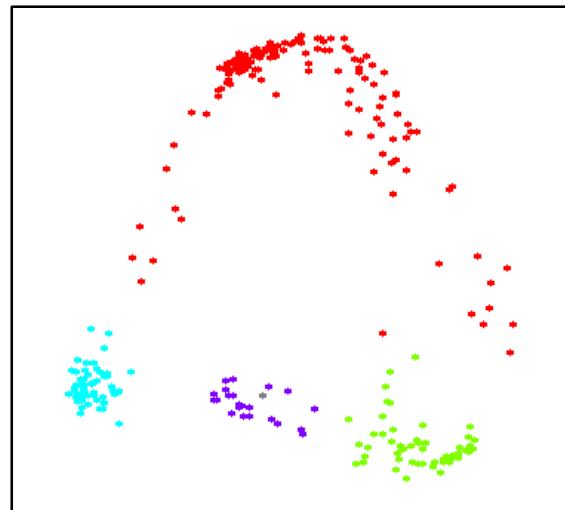
- Commonly more complex clustering methods are required.



# Clustering: Hierarchical Cluster Analysis

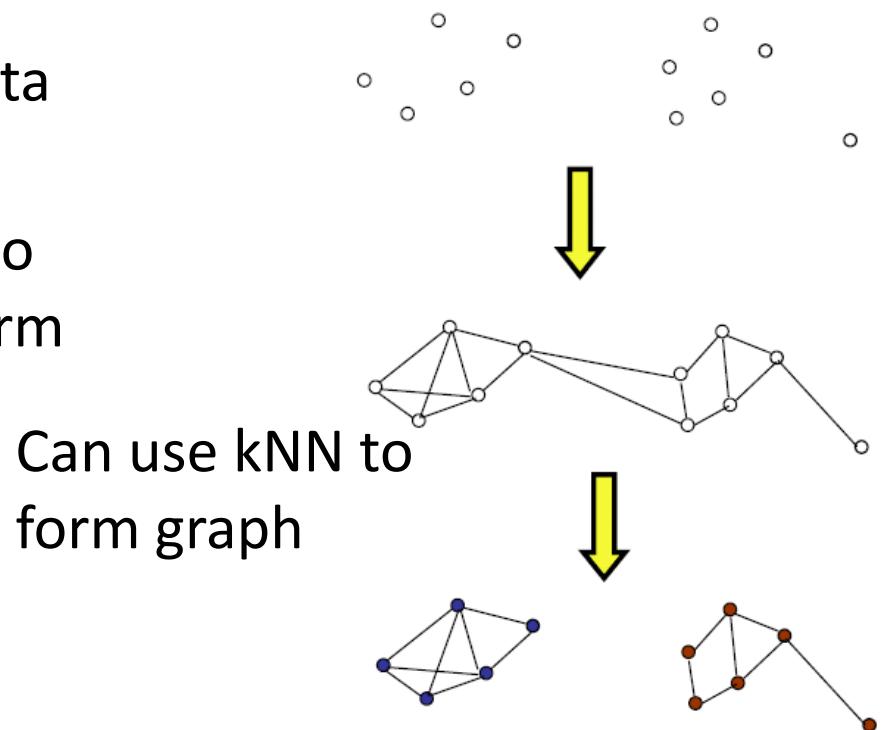


Multi-Dimensional Data Scaling



# Combinatorial Method: Spectral Clustering

- Graph cut clustering
  - Identify a graph for the data points
  - Identify the optimal ‘cut’ to segment the graph and form clusters

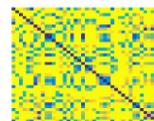
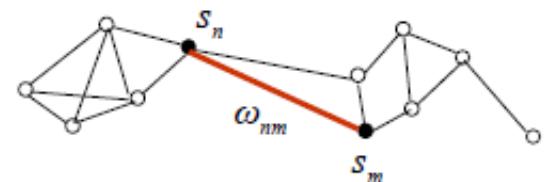


# Combinatorial Method: Spectral Clustering

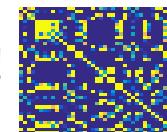
- Once we have a graph, we need to define some terms.

pairwise affinity (similarity)

$$\omega_{nm} = e^{-\frac{d(x_i, x_j)^2}{\sigma^2}}$$



Dissimilarity -> Similarity!



node volume (degree)

$$D_n = \sum_{m=1}^N \omega_{nm}$$



volume of a set (cluster)

$$Vol(C) = \sum_{n \in C} D_n$$

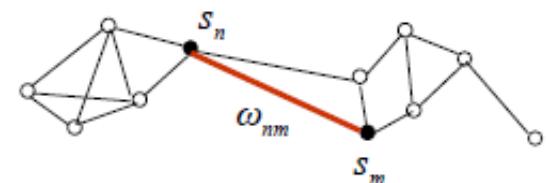


# Spectral Clustering

- Optimal cut will identify block diagonal similarity matrix

pairwise affinity (similarity)

$$\omega_{nm} = e^{-\frac{d(x_i, x_j)^2}{\sigma^2}}$$

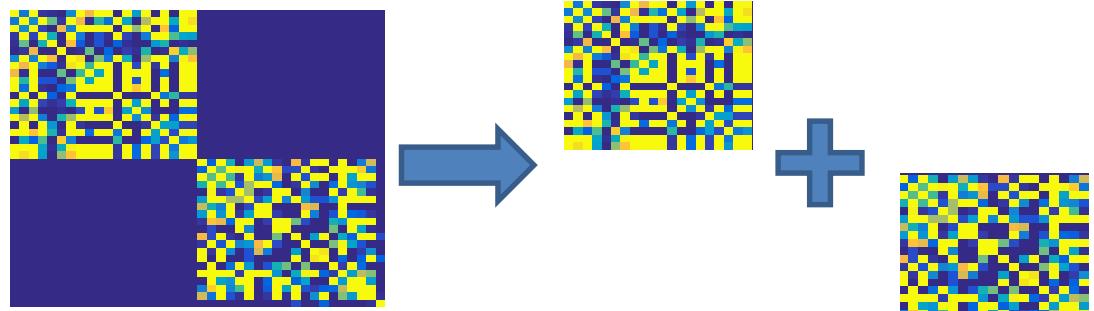


cut between 2 sets

$$Cut(C_1, C_2) = \sum_{n \in C_1} \sum_{m \in C_2} \omega_{nm}$$



Optimal cut separates the similarity matrix into two.

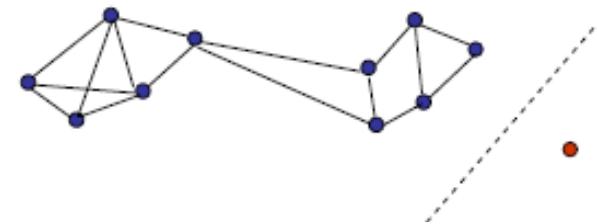


# Spectral Clustering

- Different ways to cut a graph:

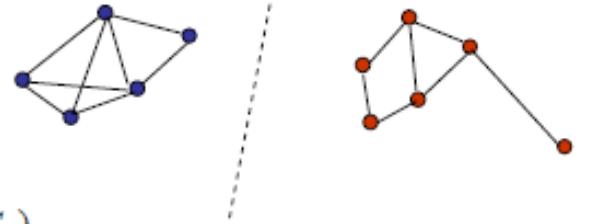
minimal bipartition cut

$$\min \text{Cut}(C_1, C_2)$$



minimal bipartition normalized cut

$$\min \frac{\text{Cut}(C_1, C_2)}{\text{Vol}(C_1)} + \frac{\text{Cut}(C_2, C_1)}{\text{Vol}(C_2)} = \min \left( \frac{1}{\text{Vol}(C_1)} + \frac{1}{\text{Vol}(C_2)} \right) \text{Cut}(C_1, C_2)$$



- Finding the optimal normalized cut is NP-hard.
- Approximation using Spectral Clustering

# Spectral Clustering

$$D = \begin{bmatrix} D_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & D_N \end{bmatrix}$$

- First, form the Graph Laplacian.

$$L = D - W$$

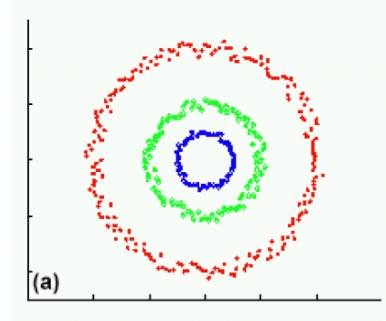
Node volume matrix      Similarity matrix

$$L_{rw} = D^{-1}L = I - D^{-1}W$$

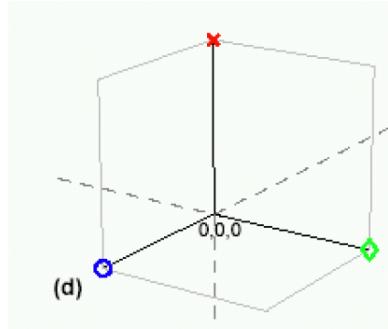
Shi and Malik (2000)

- Eigen analysis:  $L_{rw}\nu = \lambda\nu$        $V = [\nu_1, \nu_2, \dots \nu_k]$       k eigenvectors with smallest eigenvalues
- The eigen vectors with smallest eigen values are used to map data to new highly clustered space.

Original data



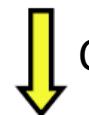
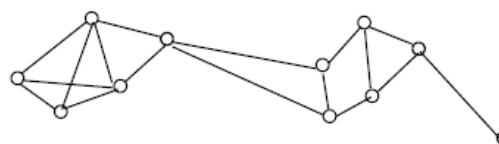
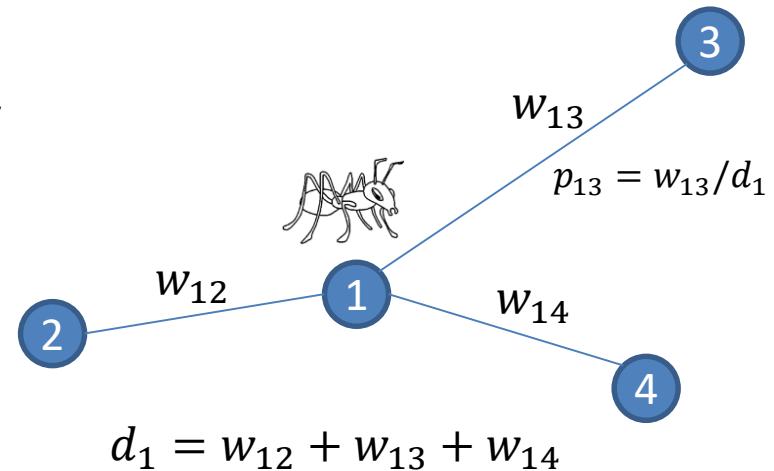
Projected data



# Spectral Clustering

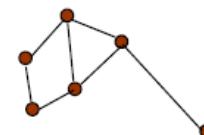
- Random Walk: Identify partition so that a random walk usually remains within a cluster and rarely jumps between clusters.

- Ant starts at a vertex  $i$ .
- Ant jumps to neighbor  $j$  with probability  $p_{ij} = w_{ij}/d_i$
- $P = D^{-1}W$
- $\sum_j p_{ij} = 1$



Cut least likely walking paths

Ant will tend  
to roam within  
clusters

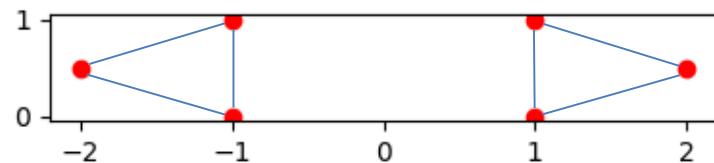


# Spectral Clustering

- Understanding the normalized Graph Laplacian.

$$L_{rw} = I - D^{-1}W = I - P$$

$$\sum_j p_{ij} = 1$$



$$I \quad -p_{ij}$$

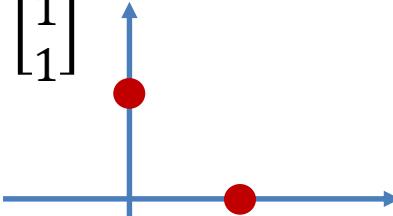
$$L_{rw} = \begin{array}{|c|c|c|c|c|c|} \hline & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \hline \text{---} & 1.00 & -0.50 & -0.50 & 0.00 & 0.00 & 0.00 \\ \hline -0.50 & 1.00 & -0.50 & -0.50 & 0.00 & 0.00 & 0.00 \\ \hline -0.50 & -0.50 & 1.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ \hline 0.00 & 0.00 & 0.00 & 1.00 & -0.50 & -0.50 & \\ \hline 0.00 & 0.00 & 0.00 & -0.50 & 1.00 & -0.50 & \\ \hline 0.00 & 0.00 & 0.00 & -0.50 & -0.50 & 1.00 & \\ \hline \end{array}$$

$$L_{rw}v = \lambda v, \quad \lambda = 0$$

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

If fully connected,  
first eigenvector:

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$



# Spectral Clustering

$$D = \begin{bmatrix} D_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & D_N \end{bmatrix}$$

- First, form the Graph Laplacian.

$$L = D - W$$

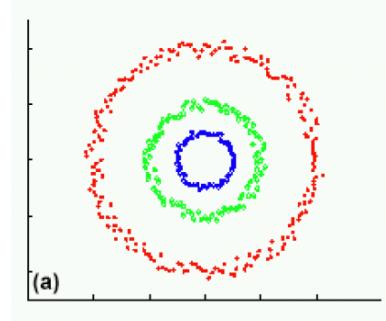
Node volume matrix      Similarity matrix

$$L_{rw} = D^{-1}L = I - D^{-1}W$$

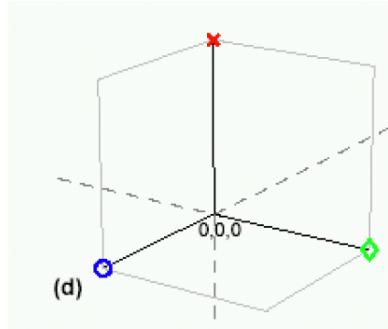
Shi and Malik (2000)

- Eigen analysis:  $L_{rw}\nu = \lambda\nu$        $V = [\nu_1, \nu_2, \dots \nu_k]$       k eigenvectors with smallest eigenvalues
- The eigen vectors with smallest eigen values are used to map data to new highly clustered space.

Original data



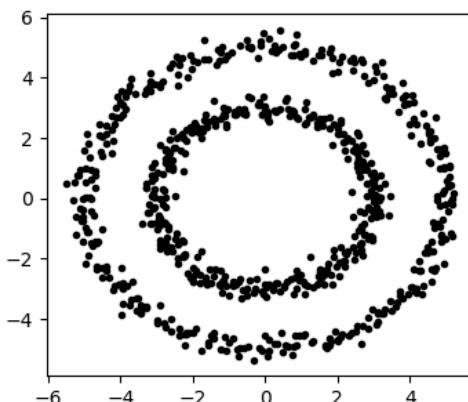
Projected data



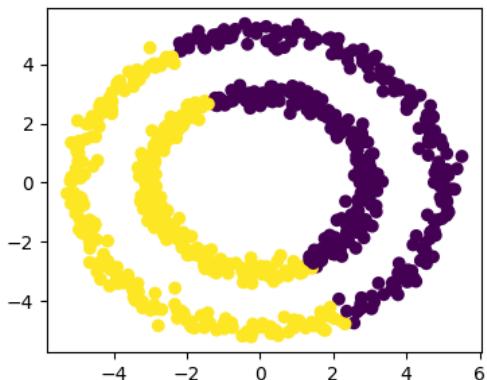
# Jupyter: Spectral Clustering

Demo:

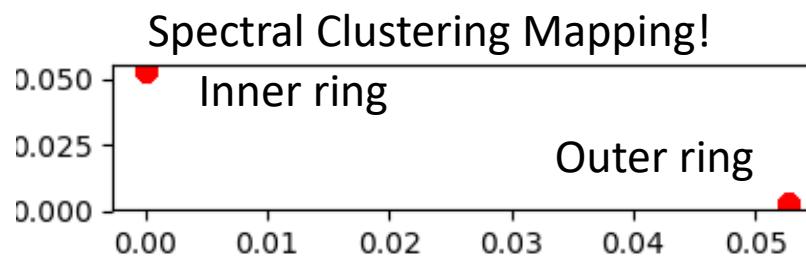
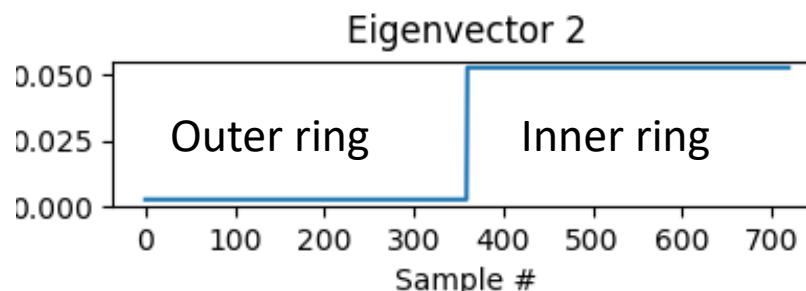
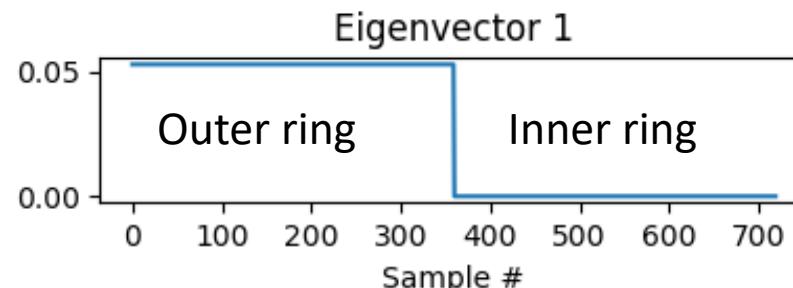
Demo Data



kmeans



Spectral Clustering Eigenvectors



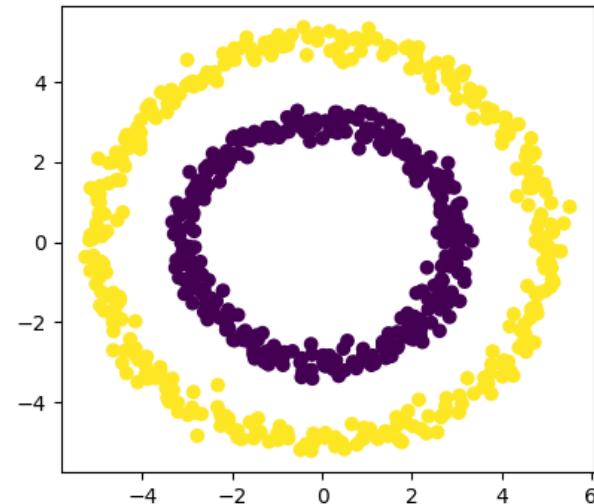
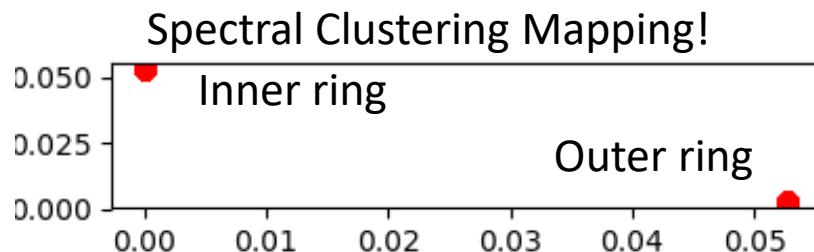
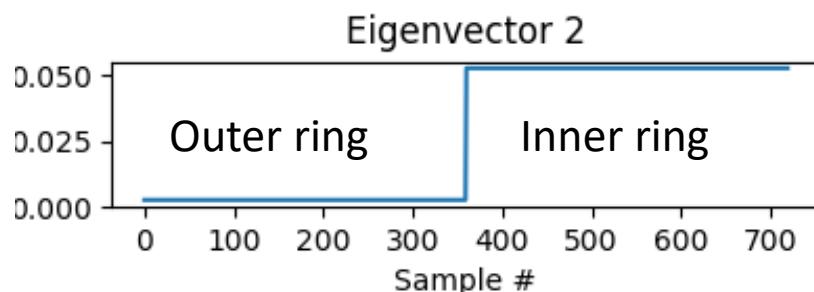
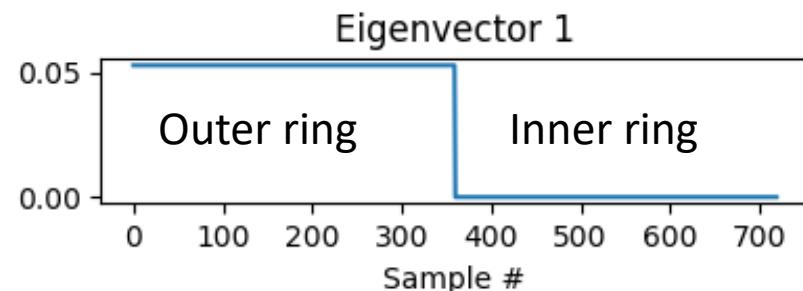
Sample 50  
mapped to  
[ 0.05, 0 ]

Sample 700  
mapped to  
[ 0, 0.05 ]

# Jupyter: Spectral Clustering

Demo:

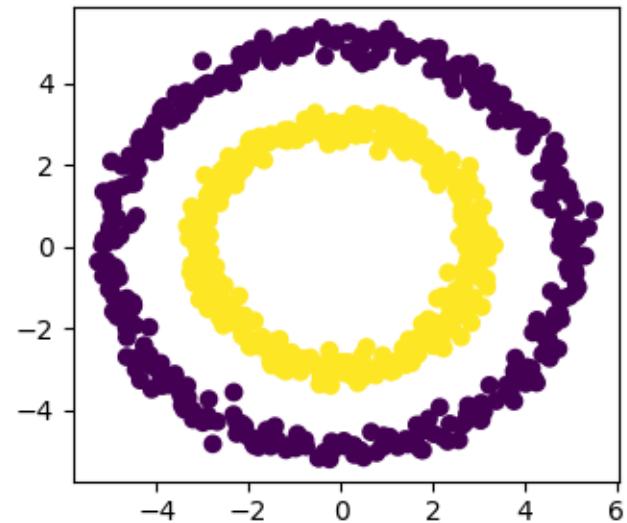
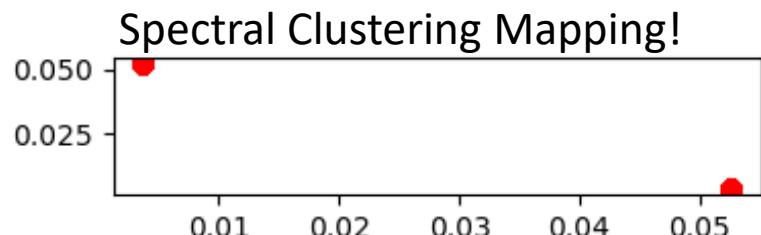
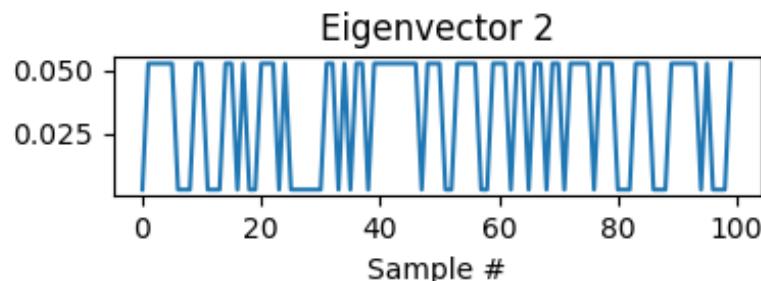
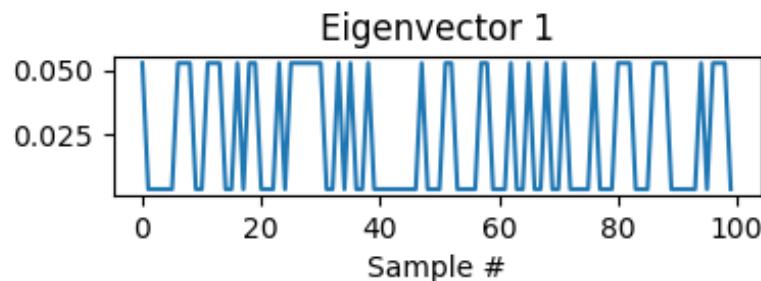
Spectral Clustering Eigenvectors

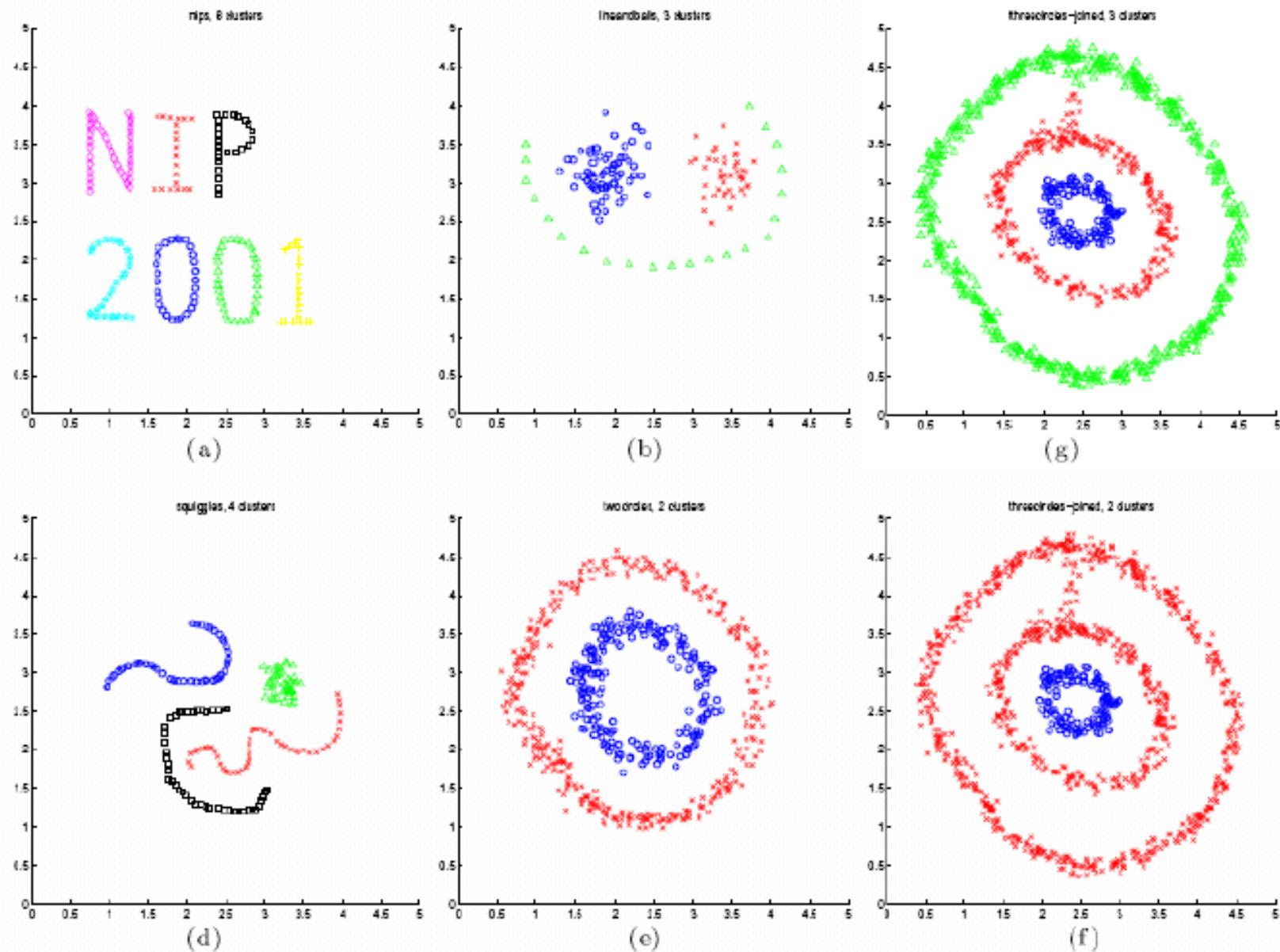


# Jupyter: Spectral Clustering

What if the samples had their order were mixed?

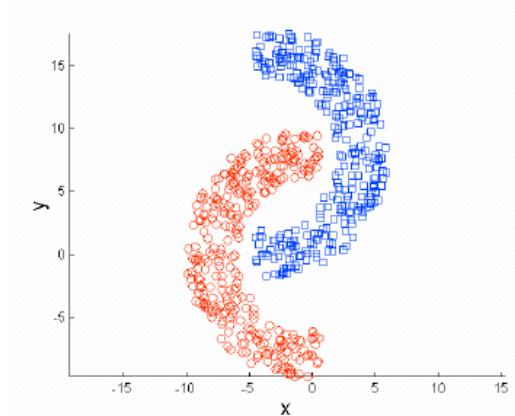
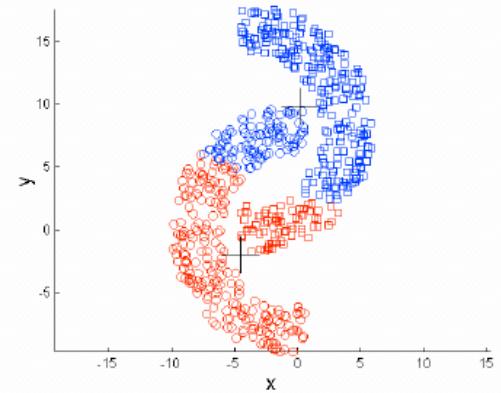
Spectral Clustering Eigenvectors



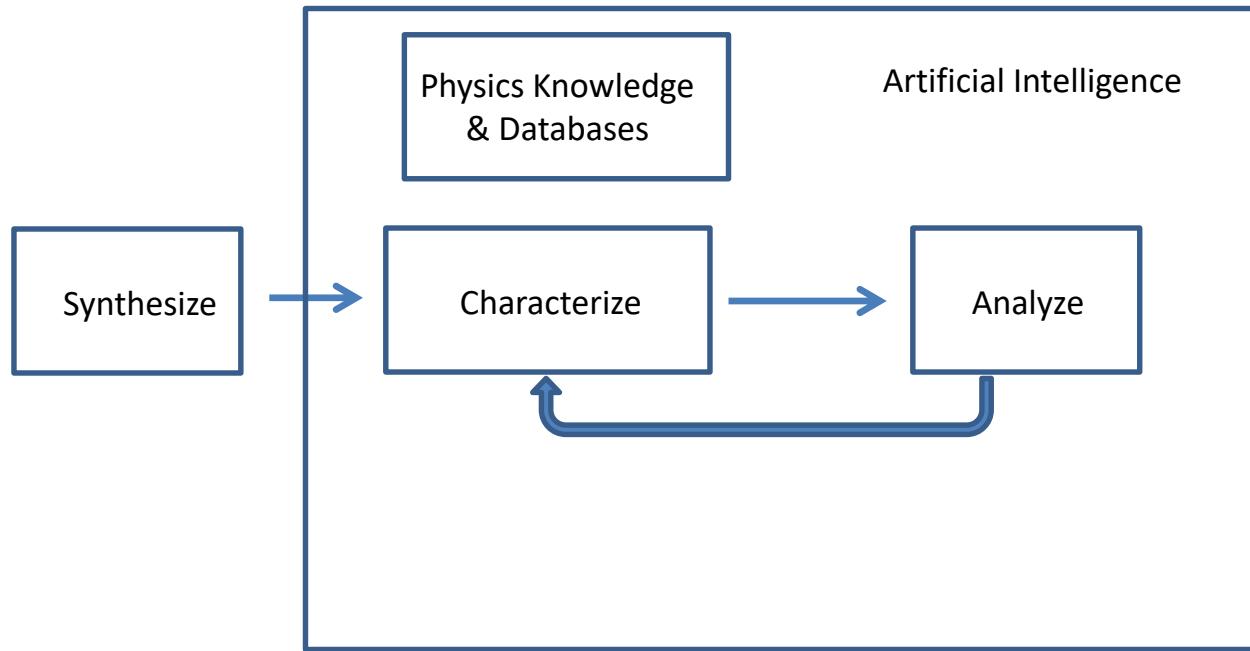


# Clustering: Challenges

- Data pre-processing
- Define appropriate dissimilarity
- Selecting appropriate method
  - What do you know about the data?
  - What methods have been successfully used on similar data?
  - Try multiple methods and visualize results.
- Number of clusters
  - Heuristic, e.g. Gap statistic
- Visualize Results!
  - The most important check is to have an expert look at the results.

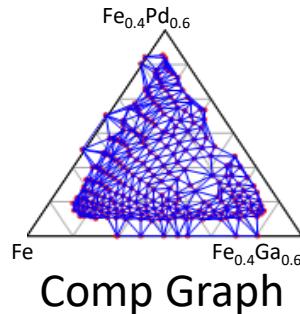


# Autonomous Metrology using the techniques you've learned

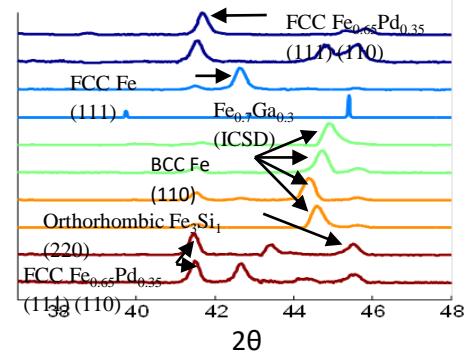
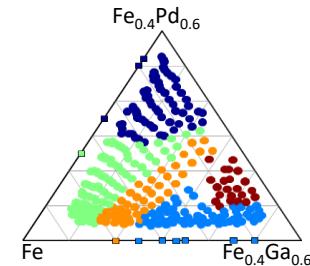


# Autonomous HiTp XRD

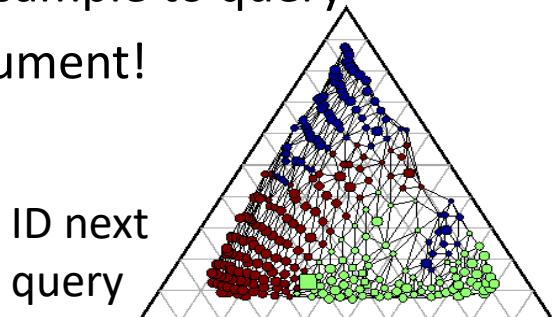
1. Pick a metric -> Cosine metric
2. Form a Graph to represent samples  
(here, using composition for graph)
3. Graph cut clustering -> phase regions.
4. NMF -> constituent spectra for each phase region.
  - Iterate (2) and (3) until convergence.
4. Graph based semi-supervised learning ->  
extrapolate from known samples to unknown.
5. Identify optimal sample to query
6. Talk to XRD instrument!
7. Repeat 2-5



Phase Diagram



Constituent Phases



# Phase Mapping: High-Throughput Approach

- Measurement is a time / resource sink
- For wafer of 500+ samples:
- In Lab: Takes weeks-months
- Synchrotron: Takes 5+ Hours (Every second counts)



Mn-Ni-Ge library  
535 samples



Bruker D8  
30 Minutes per sample  
2 weeks!



Stanford Synchrotron Radiation Lightsource  
30 seconds per sample  
4.5 hours

# Autonomous Metrology: Motivation

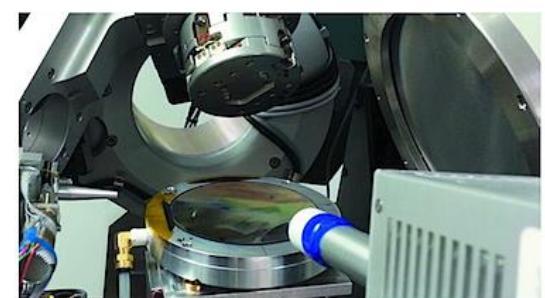
Why use AI to just analyze data? Put it in control of the equipment!

Instead of measuring all the samples, measure only the ones that count -> AI for optimal experiment design

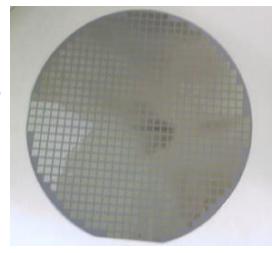
- Minimum measurements -> Maximum knowledge
- Save on worker hours and instrument time.
- Start it up and let it run.
- Minimize human bias: experiment design, execution, data analysis
- Replaced with traceable algorithmic bias
- Democratize Science
  - Simplify equipment use
  - Collaboratory



Bruker D8  
30 Minutes per sample  
2 weeks!

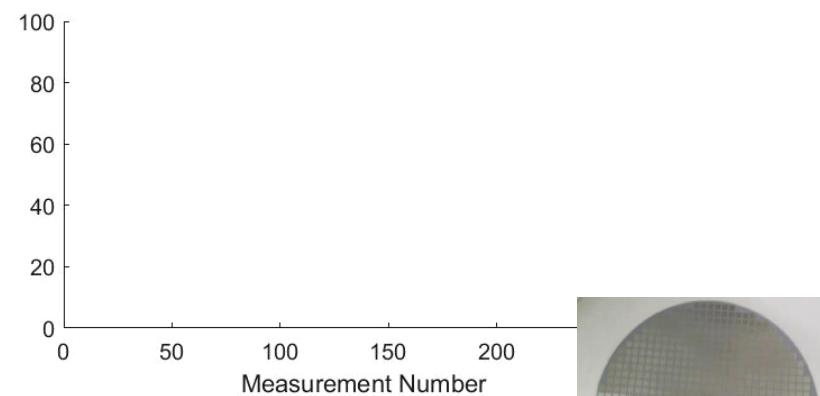
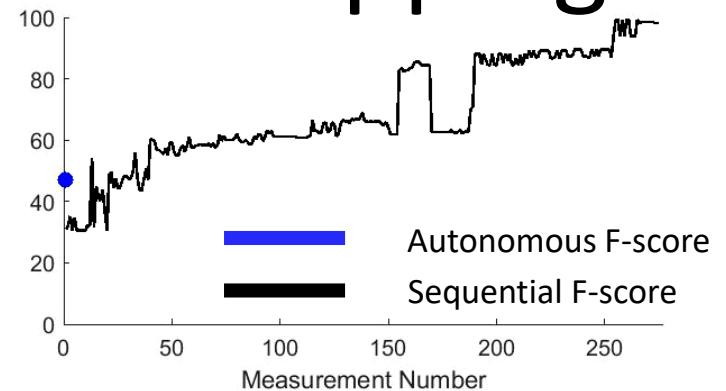
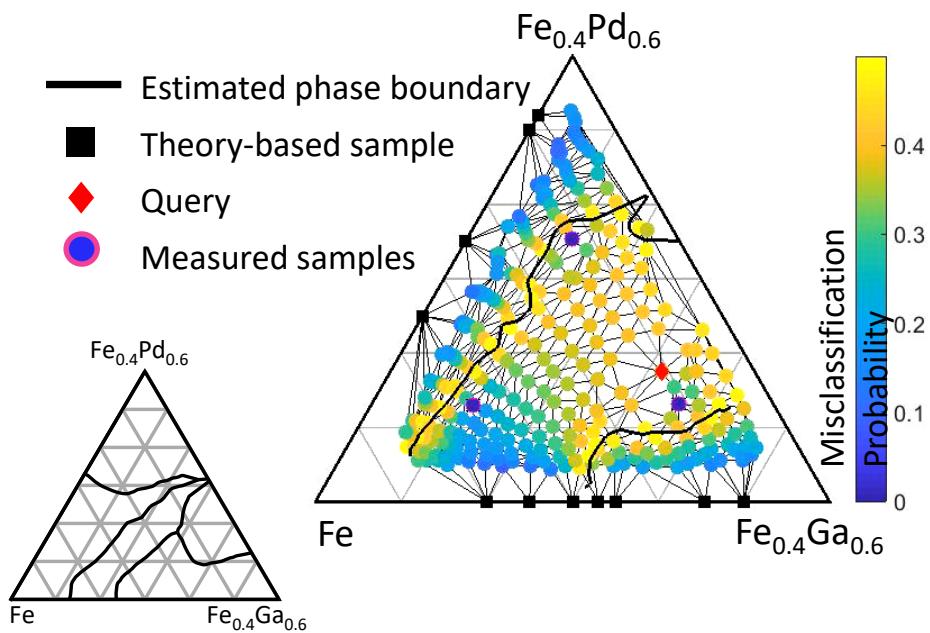


Stanford Synchrotron Radiation  
Lightsource  
30 seconds per sample  
4.5 hours



Mn-Ni-Ge library  
535 samples

# Autonomous Phase Mapping

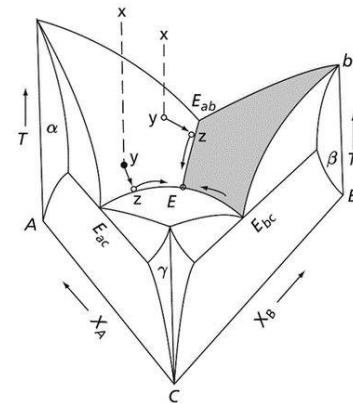
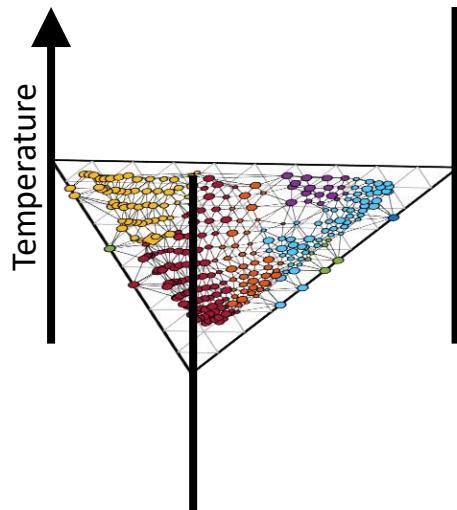
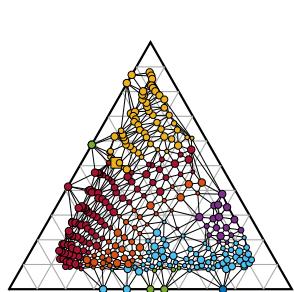


AI is controlling X-ray diffraction systems at SLAC & in the lab!



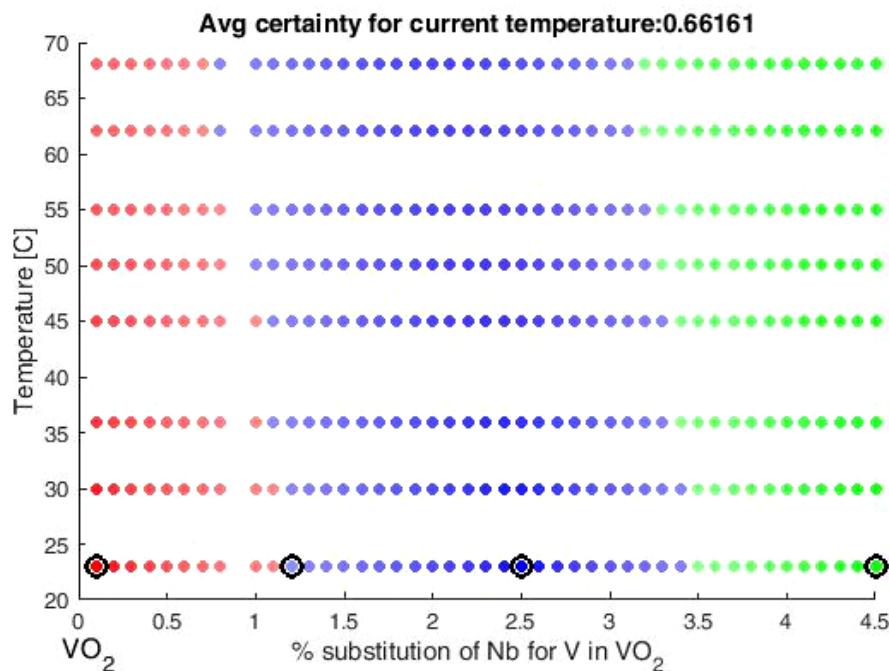
# Autonomous Temperature Phase Mapping

- Developing 2 systems for autonomous composition & temperature phase mapping.
- Minimum measurements for maximum knowledge.



Wikipedia: ternary\_cooling.jpg

# Composition-Temperature Phase Mapping of $V_{1-x}Nb_xO_2$



# The Challenge of Materials Discovery

- High dimensional space

- Exhaustive Search:

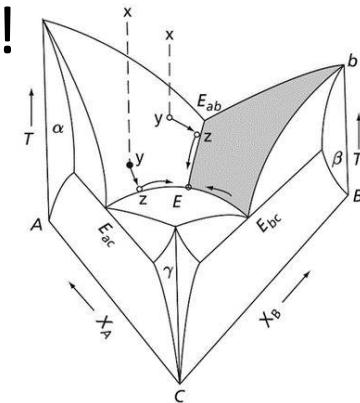
Assume: For each new parameter, 100 experiments over range.

For N parameters  $\rightarrow 10^{2(N-1)}$  experiments!

4 Parameters  $\rightarrow 10^6$  experiments

6 Parameters  $\rightarrow 10^{10}$  experiments

- Complex materials are out of reach!



4 Fabrication parameters:  
3 Elements  
1 Processing (temperature)  
 $\rightarrow 10$

# Any Questions?

- Feel free to contact me:

Gilad Kusne – [aaron.kusne@nist.gov](mailto:aaron.kusne@nist.gov)