

ENGR-E 511; ENGR-E 399

# “Machine Learning for Signal Processing”

## Module 03: Lecture 01: Clustering

**Minje Kim**

Department of Intelligent Systems Engineering

Email: [minje@indiana.edu](mailto:minje@indiana.edu)

Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>



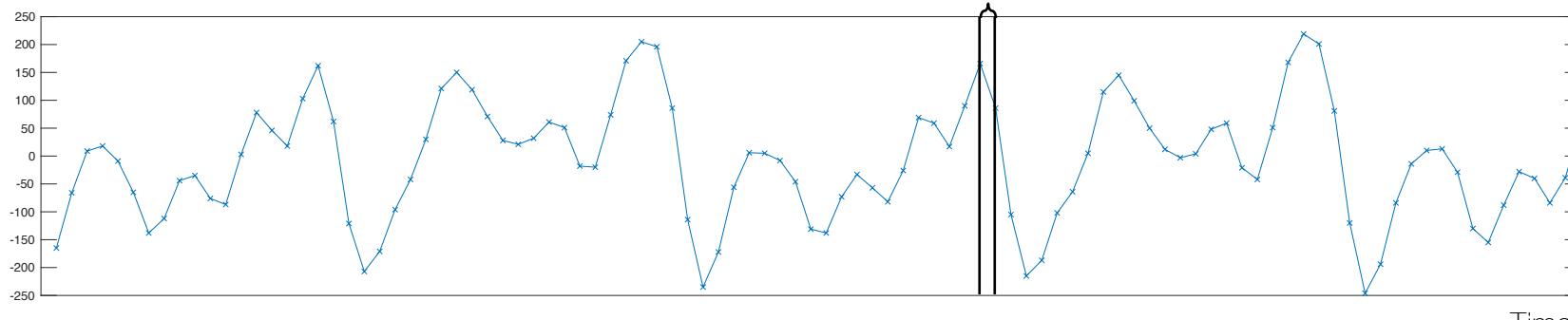
INDIANA UNIVERSITY  
**SCHOOL OF INFORMATICS,  
COMPUTING, AND ENGINEERING**



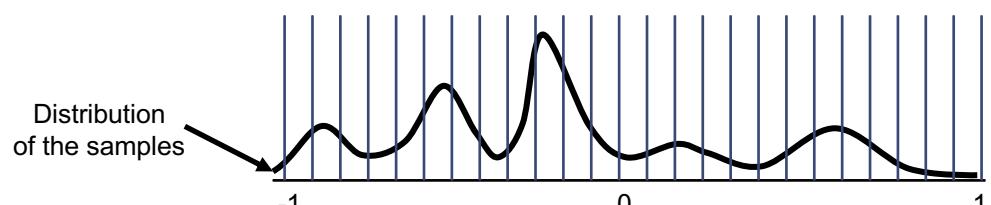
# Motivating Problems

## - CD

- How do we represent music in a CD?
  - What? What is CD?
  - 44.1 kHz, 16 bit, LPCM
  - What does it mean?



- We sample from the (continuous) waveform at every  $1/44100$  second
- Each sample is represented with one of  $2^{16}=65536$  values
  - e.g. 0000 0000 0000 1000 : -0.9998
  - e.g. 1111 1111 1101 1100 : +0.9989
  - e.g. 1000 0000 0000 0000 : 0
- Can we do better?
  - (Real-valued) samples in the same block is represented by the same value (and I don't like it)



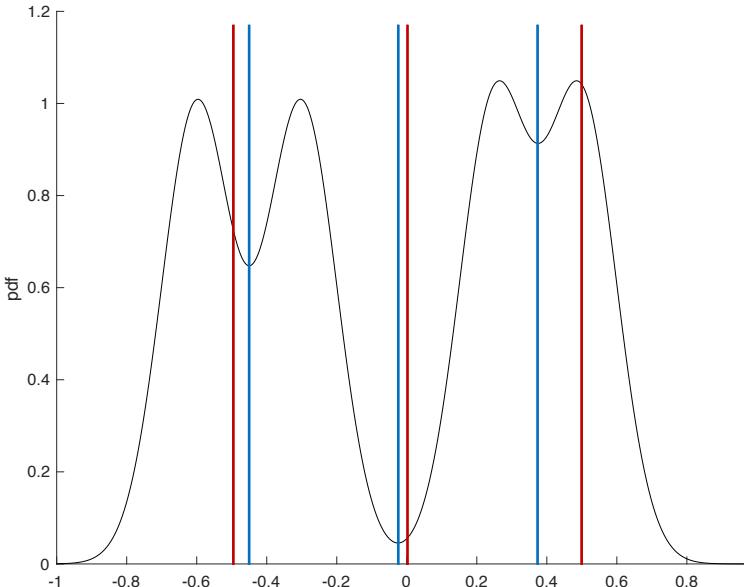
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Motivating Problems

## - The Lloyds-Max algorithm

- First of all, we need to save the bits
  - What if there's a lack of bits?
    - Some values are misrepresented
- Assume that the raw audio samples are following a distribution like:



- Which do you prefer between the red boundaries and blue boundaries (they are all for 2 bit encoding)?



8bits

4bits

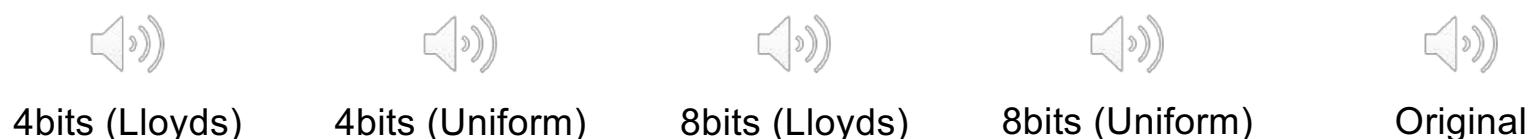
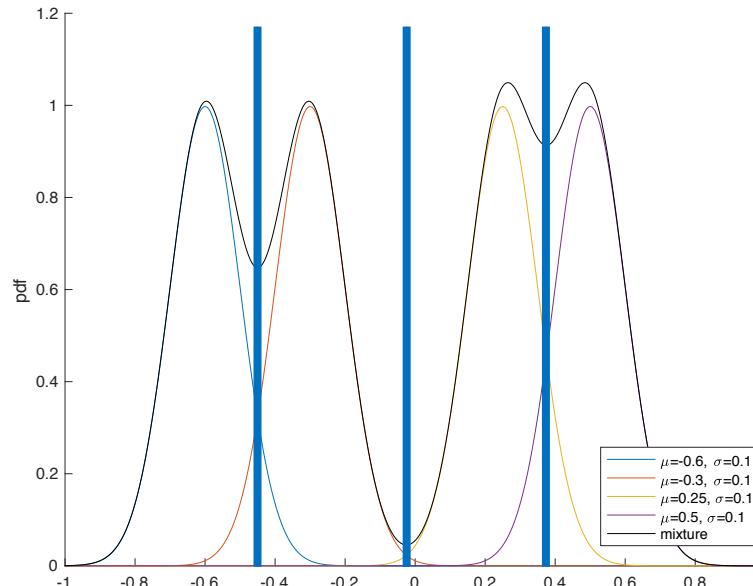


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Motivating Problems

- The Lloyds-Max algorithm
  - You prefer the blue boundaries because they go well with the underlying structure of the sample distribution
    - Underlying structure?
  - And they sound better!



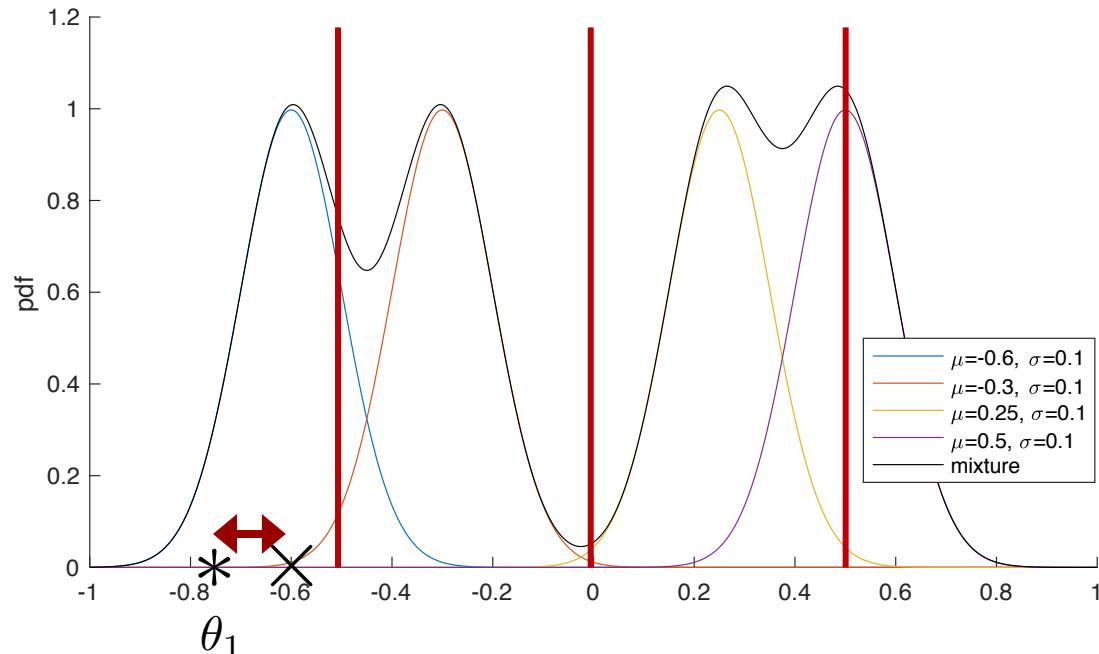
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# k-Means Clustering

## - A scalar case

- Where's the distortion from?
- Let's start from the red boundaries we don't like
  - The discrepancy between the representative and the actual samples
  - We want to find the representative that creates the least discrepancy
    - For each quantization level
- How do we measure the amount of discrepancy?



# k-Means Clustering

## - A scalar case

- So, the objective for j-th range is to find the representative value that minimizes the error

$$\arg \min_{\theta_j} \sum_{i \in \mathcal{C}_j} \|x_i - \theta_j\|^2$$

- For all examples that belong to the j-th range

- $\mathcal{C}_j$  holds the indices for j-th range

- So, what's the solution? (why?)  $\theta_j = \frac{1}{|\mathcal{C}_j|} \sum_{i \in \mathcal{C}_j} x_i$

- Eventually, we need to do this for all J ranges

$$\arg \min_{\theta_1, \theta_2, \dots, \theta_J} \sum_{j=1}^J \sum_{i \in \mathcal{C}_j} \|x_i - \theta_j\|^2$$

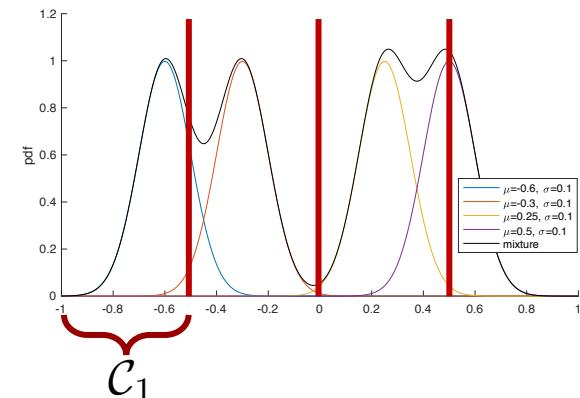
- Let me put it in another way

$$\arg \min_{\theta_1, \theta_2, \dots, \theta_J} \sum_{j=1}^J \sum_{i=1}^N u_{ij} \|x_i - \theta_j\|^2$$

- Where  $u_{ij}$  is a matrix that indicates the membership

- e.g.  $u_{ij} = 1$  then  $i \in \mathcal{C}_j$

- The membership matrix  $u_{ij}$  should also meet another constraint:  $\sum_j u_{ij} = 1$



	j=1		j=J	
i=1	1	0	0	0
	0	0	0	1
	1	0	0	0
	0	1	0	0
	0	0	1	0
	...	...	...	...
i=N	0	1	0	0

4<sup>th</sup> sample belongs to the 2<sup>nd</sup> cluster



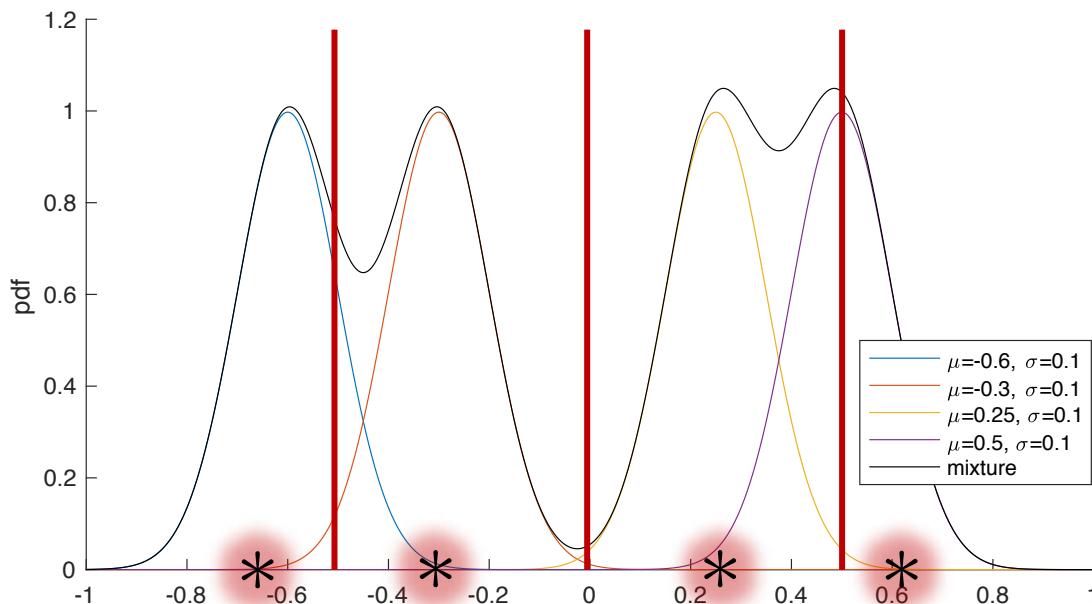
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# k-Means Clustering

## - A scalar case

- Are we done?
  - No, we assumed that the boundaries are correct, but they aren't



- In other words, we don't know if the  $u_{ij}$  matrix is correct



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# k-Means Clustering

## - A scalar case

- We need to optimize w.r.t. the membership matrix as well

$$\arg \min_{\theta, U} \sum_{j=1}^J \sum_{i=1}^N u_{ij} \|x_i - \theta_j\|^2$$

- The k-means clustering algorithm on scalar samples
  - Initialize the means  $\theta = \{\theta_1, \theta_2, \dots, \theta_J\}^\top$  with random numbers
  - Update the membership matrix
    - This is actually a complicated optimization problem (why?), but the solution is simple
    - For a fixed set of means

$$u_{ij} = \begin{cases} 1 & \text{if } j = \arg \min_{j'} \|x_i - \theta_{j'}\|^2 \\ 0 & \text{otherwise} \end{cases}$$

- Update the means
$$\begin{aligned}\theta_j &= \frac{1}{|\mathcal{C}_j|} \sum_{i \in \mathcal{C}_j} x_i \\ &= \frac{1}{|\mathcal{C}_j|} \sum_{i=1}^N u_{ij} x_i\end{aligned}$$



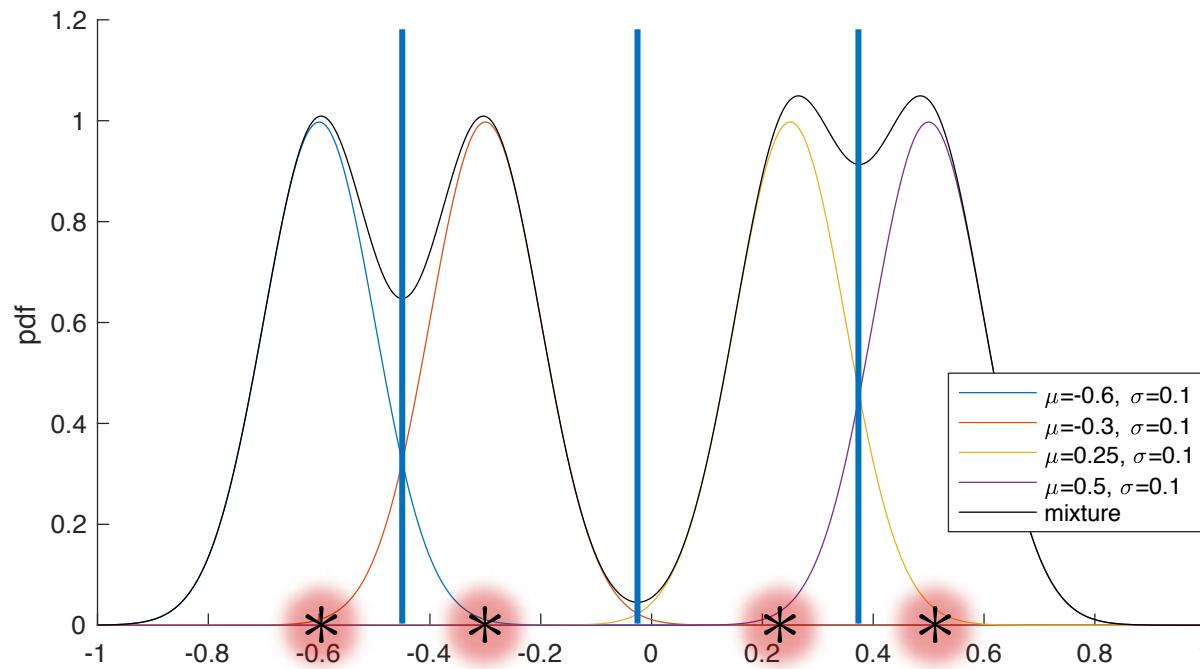
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# k-Means Clustering

## - A scalar case

- Let's get back to the CD encoding problem (Lloyd-Max algorithm)
- Now instead of all the possible real values between -1 and +1
- We replace the values within each range with their corresponding representatives, i.e. the means.



4bits (Lloyds)



4bits (Uniform)



8bits (Lloyds)



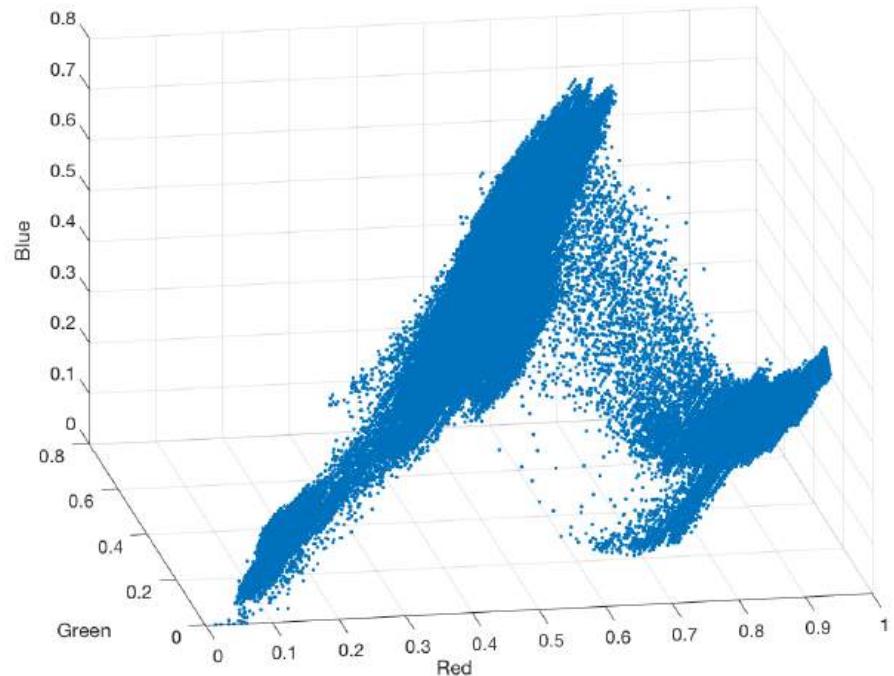
8bits (Uniform)



Original

# Motivating Problems

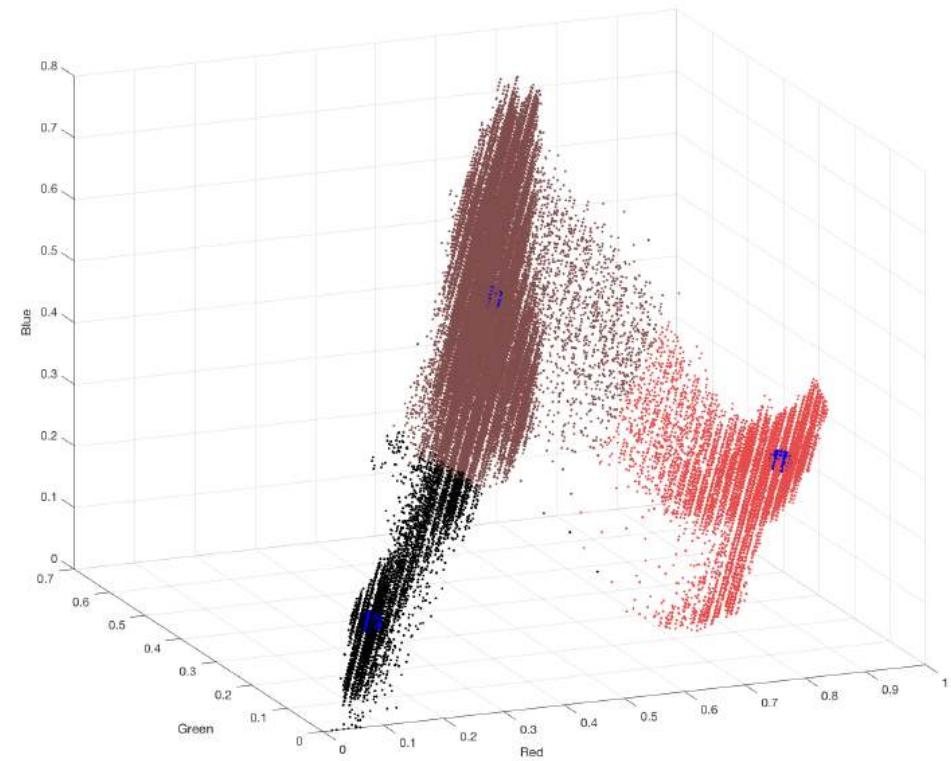
- Black cat, red wall, gray ground
  - Now let's move on to the multi-dimensional case
  - In general, how do we quantize a vector?
  - First off, can you (verbally) describe this picture?



# Motivating Problems

- Black cat, red wall, gray ground

- k-means with three clusters

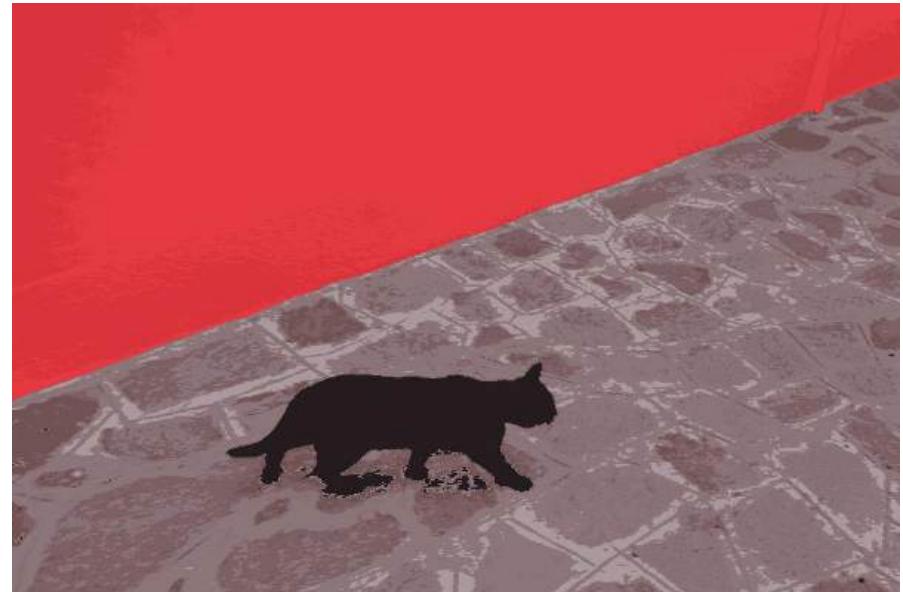


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Motivating Problems

- Black cat, red wall, gray ground
- k-means with 8 clusters and 16 clusters



- Algorithm-wise everything is the same except for the fact that the input samples are 3D (RGB) vectors

$$\arg \min_{\boldsymbol{\theta}, \mathbf{U}} \sum_{j=1}^J \sum_{i=1}^N u_{ij} \|x_i - \theta_j\|^2 \longrightarrow \arg \min_{\boldsymbol{\Theta}, \mathbf{U}} \sum_{j=1}^J \sum_{i=1}^N u_{ij} \|\mathbf{x}_i - \boldsymbol{\theta}_j\|^2$$

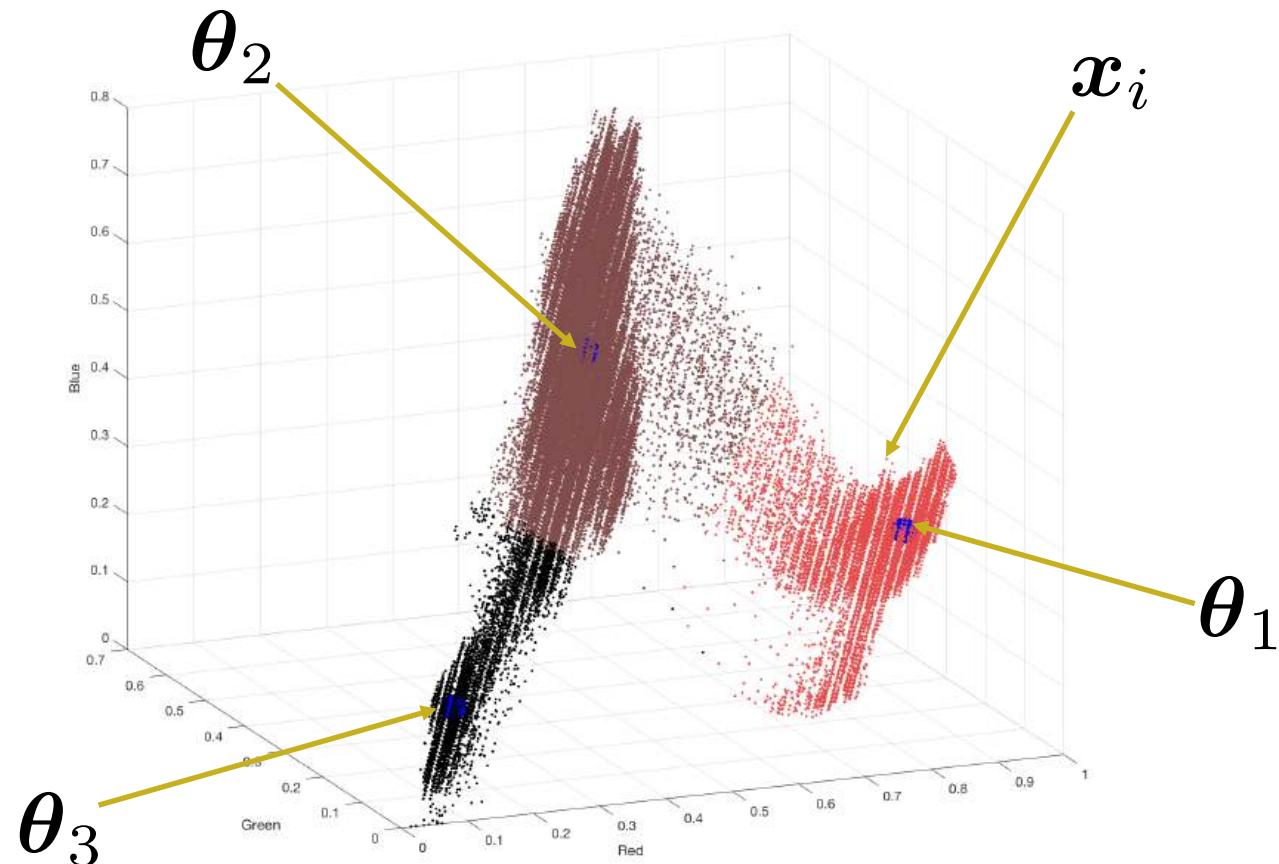


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Motivating Problems

- Black cat, red wall, gray ground



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Vector Quantization

- Clustering on multi-dimensional samples
- What we did is something called **Vector Quantization (VQ)**
  - Do clustering
  - Replace vector samples with the mean of the cluster they belong to
- In image, we had 3D vectors (RGB)
  - If we assume  $J$  clusters, we need  $\log_2 J$  bits to encode a pixel
    - Rather than 8 bits per pixel
- What we need:
  - A good clustering
    - A small number of means that are representative enough
  - Dictionary (codebook)
    - A codeword corresponds to one of the means
  - Index to the codebook
    - Index of the pixel-wise membership

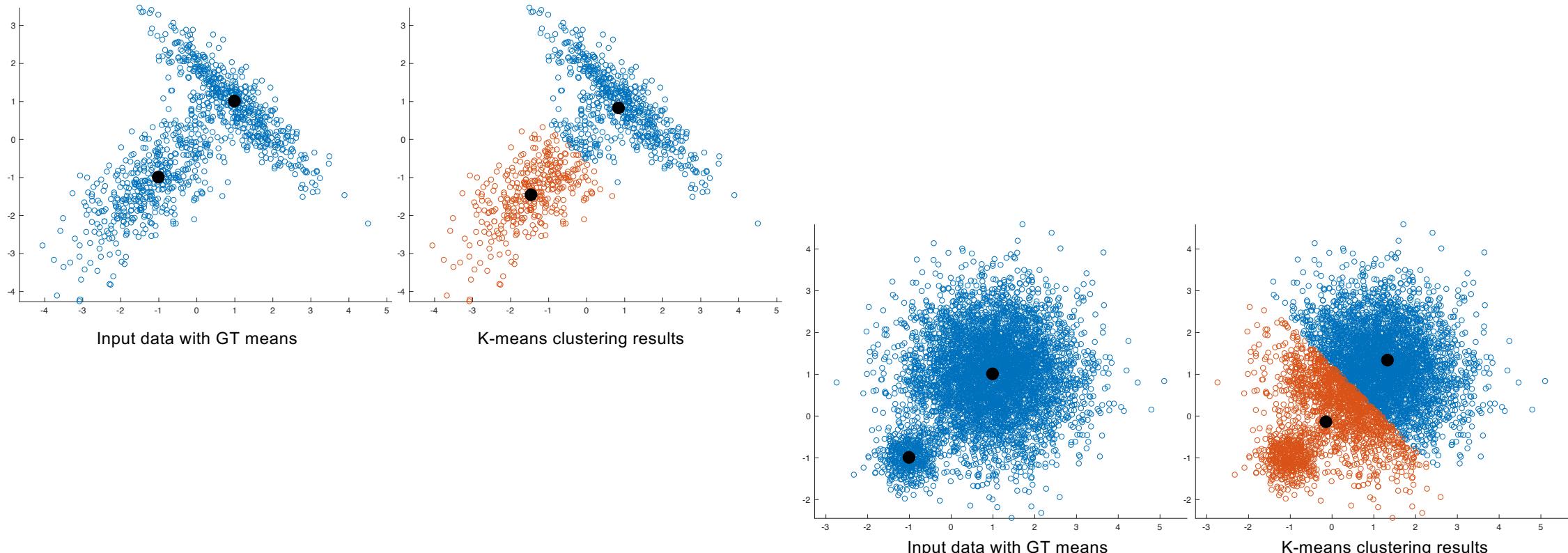


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

- What's wrong with k-means?
  - What I don't like about k-means
    - Euclidean distance, hard decision, equiprobable clusters, diagonal cov...



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

- An alternative: Mahalanobis distance
  - Let's tweak k-means
  - First, let's take variance into account for the distance metric
    - Mahalanobis distance:

$$\mathcal{D}_M(x_i || \mu_j) = \sqrt{\frac{(x_i - \mu_j)^2}{\sigma_j^2}}$$

$\mu_j$  Mean of j-th cluster  
 $\sigma_j$  Standard dev. of j-th cluster

- Multi-dimensional cases with covariance:

$$\mathcal{D}_M(x_i || \mu_j) = \sqrt{(x_i - \mu_j)^\top \Sigma^{-1} (x_i - \mu_j)}$$

- For a 2D Gaussian with

$$\Sigma = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}$$

- (1,1)

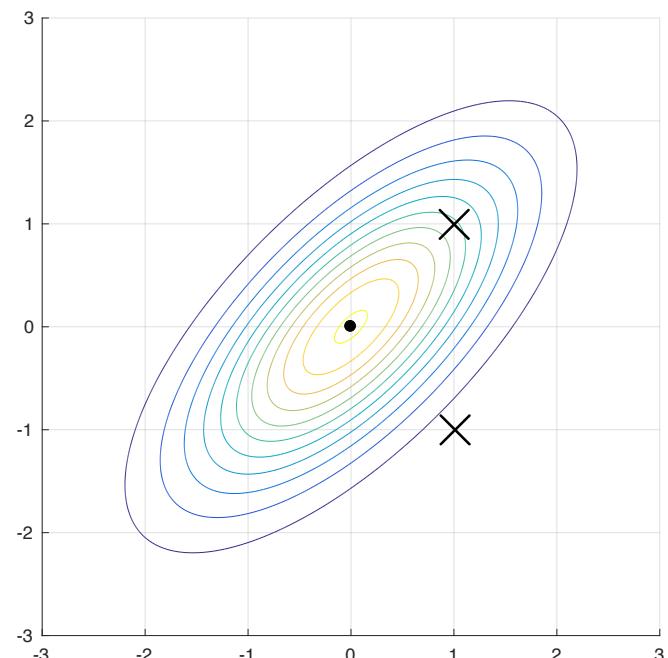
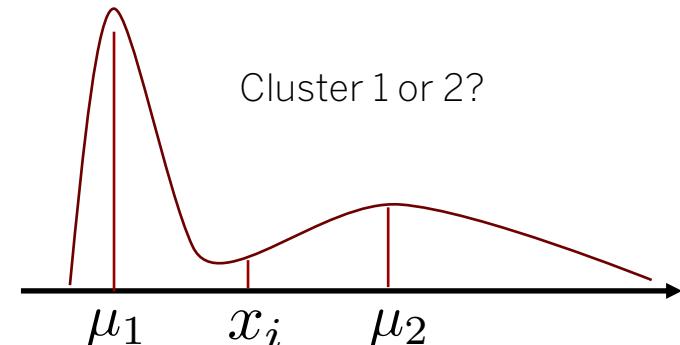
- Euclidean:  $\sqrt{2}$

- Mahalanobis: 1.0847

- (1,-1)

- Euclidean:  $\sqrt{2}$

- Mahalanobis: 2.5820



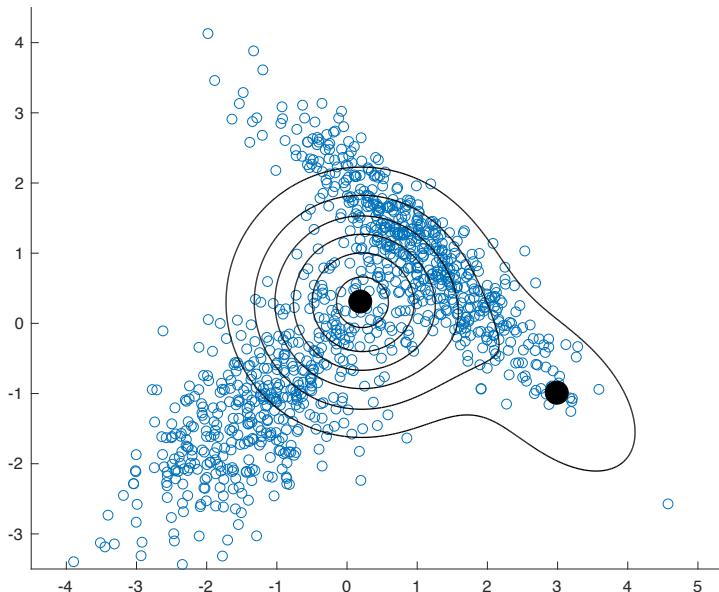
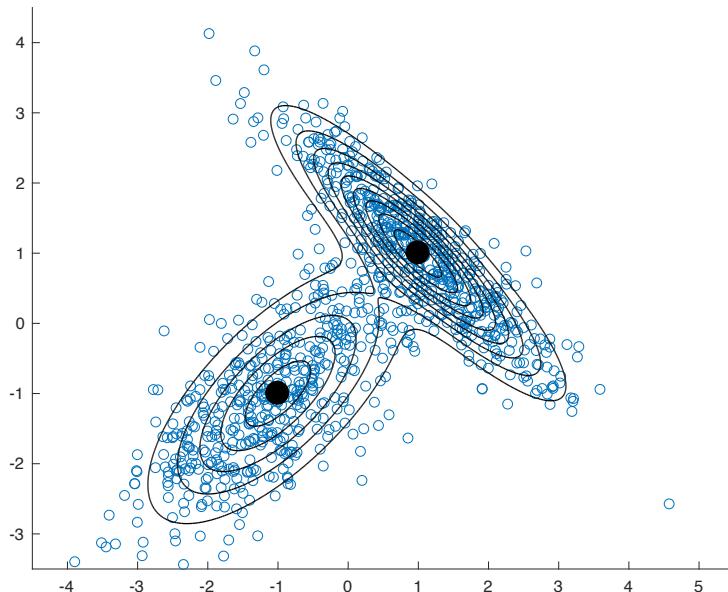
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

## - Maximum Likelihood

- Mixture of Gaussians (MoG) or Gaussian Mixture Model (GMM)
  - A maximum likelihood problem
    - Given the data, find the best fit among the family of prob. distributions with a certain parametric form
- Which one is a better fit?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

## - Maximum Likelihood

- We know how to solve a maximum likelihood problem:

$$\arg \max_{\Theta} \prod_{i=1}^N p(x_i; \Theta)$$

- For the GMM case,  
we can break down the likelihood as follows:

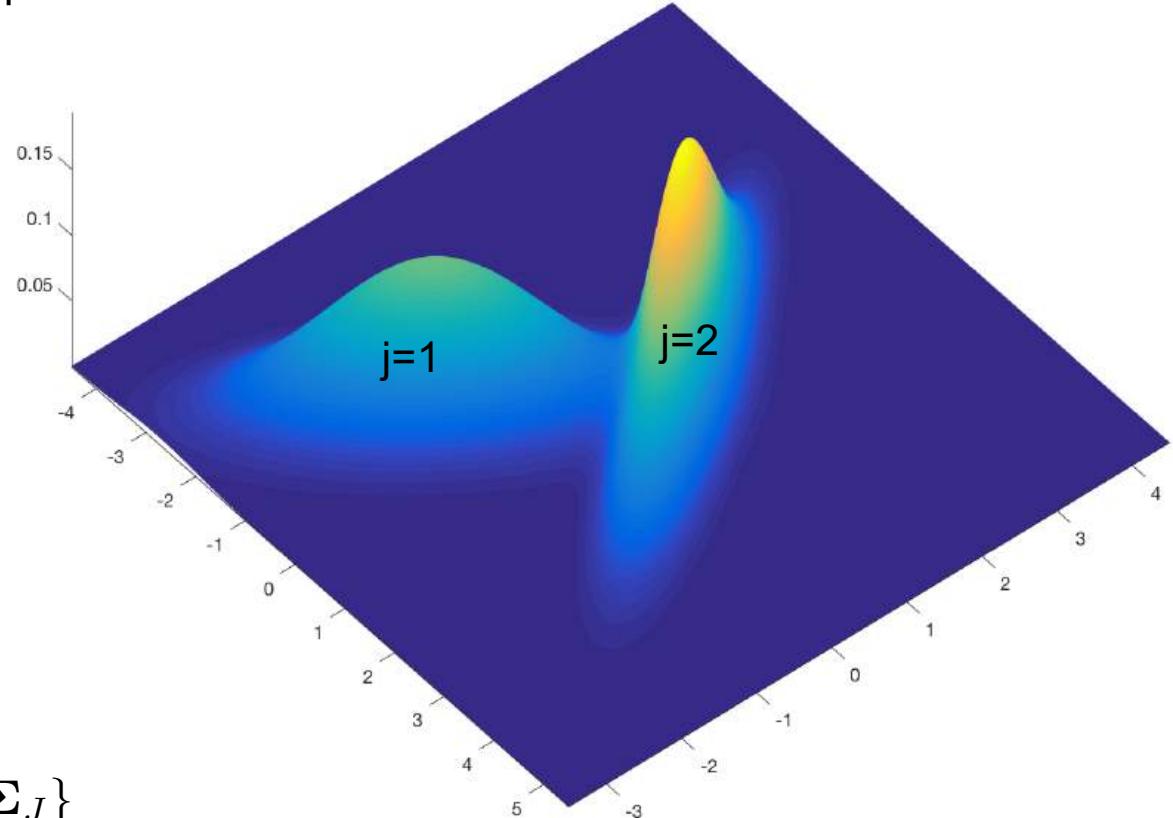
$$\mathcal{L} = \prod_{i=1}^N \sum_{j=1}^J P_j \mathcal{N}(x_i; \mu_j, \Sigma_j)$$

- Because,

$$p(x_i; \Theta) = \sum_{j=1}^J P_j \mathcal{N}(x_i; \mu_j, \Sigma_j)$$

- Note that:

$$\Theta = \{P_1, \mu_1, \Sigma_1, P_2, \mu_2, \Sigma_2, \dots, P_J, \mu_J, \Sigma_J\}$$



# Gaussian Mixture Model

## - Maximum Likelihood

- We also know the p.d.f. of a Gaussian:

$$\mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}_j|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)\right)$$

- Therefore, the likelihood is

$$\begin{aligned}\mathcal{L} &= \prod_{i=1}^N \sum_{j=1}^J P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \\ &= \prod_{i=1}^N \sum_{j=1}^J P_j \left( \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}_j|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)\right) \right)\end{aligned}$$

- Then the log-likelihood is:

$$\begin{aligned}\mathcal{LL} &= \sum_{i=1}^N \log \left( \sum_{j=1}^J P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \right) \\ &= \sum_{i=1}^N \log \left( \sum_{j=1}^J P_j \left( \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}_j|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)\right) \right) \right)\end{aligned}$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

## - Maximum Likelihood

- The objective function:

$$\arg \max_{\Theta} \mathcal{LL} + \lambda \left( \sum_{j=1}^J P_j - 1 \right)$$

$$\arg \max_{\Theta} \sum_{i=1}^N \log \left( \sum_{j=1}^J P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \right) + \lambda \left( \sum_{j=1}^J P_j - 1 \right)$$

- Differentiation is difficult

- Why?
    - Because of the summation inside the logarithm

- What should we do?

- Jensen's inequality

$$\log \left( \sum_{j=1}^J \frac{P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \mathbf{U}_{ij}}{\mathbf{U}_{ij}} \right) \geq \sum_{j=1}^J \mathbf{U}_{ij} \log \left( \frac{P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\mathbf{U}_{ij}} \right)$$

- What? Why?



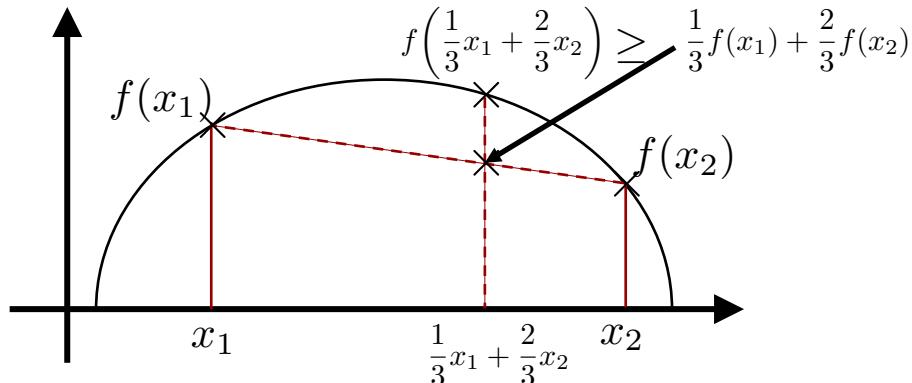
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

- Jensen's Inequality

- For example



- For a concave function  $f: f\left(\frac{\sum a_i x_i}{\sum a_i}\right) \geq \frac{\sum a_i f(x_i)}{\sum a_i}$

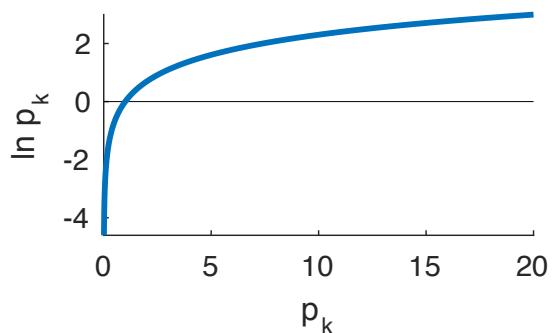
- Or:  $f\left(\sum a_i x_i\right) \geq \sum a_i f(x_i)$  if  $\sum a_i = 1$  and  $a_i \geq 0$

- Logarithmic functions are concave

- Why?

$$\log'(x) = \frac{1}{x}$$

$$\log''(x) = -\frac{1}{x^2}$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

## - Expectation Maximization (EM)

- Let's get back to the ML problem for GMM

$$\begin{aligned}\mathcal{LL} &= \sum_{i=1}^N \log \left( \sum_{j=1}^J \frac{P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \mathbf{U}_{ij}}{\mathbf{U}_{ij}} \right) & \sum_j \mathbf{U}_{ij} = 1 \text{ and } \mathbf{U}_{ij} \geq 0 \\ &\geq \sum_{i=1}^N \sum_{j=1}^J \mathbf{U}_{ij} \log \left( \frac{P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\mathbf{U}_{ij}} \right) & \text{Jensen's inequality} \\ &= \sum_{i=1}^N \sum_{j=1}^J \mathbf{U}_{ij} \log \left( \frac{p(j|\mathbf{x}_i)p(\mathbf{x}_i)}{\mathbf{U}_{ij}} \right) & \because p(j|\mathbf{x}_i) = \frac{p(\mathbf{x}_i|j)p(j)}{\sum_j p(\mathbf{x}_i|j)p(j)} = \frac{\mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) P_j}{p(\mathbf{x}_i)} \\ &= \sum_{i=1}^N \sum_{j=1}^J \mathbf{U}_{ij} \log \left( \frac{p(j|\mathbf{x}_i)}{\mathbf{U}_{ij}} \right) + \sum_{j=1}^J \mathbf{U}_{ij} \log p(\mathbf{x}_i) \\ &= \sum_{i=1}^N \sum_{j=1}^J \mathbf{U}_{ij} \log \left( \frac{p(j|\mathbf{x}_i)}{\mathbf{U}_{ij}} \right) + \log p(\mathbf{x}_i) & = \sum_{i=1}^N -\mathcal{D}_{KL}(\mathbf{U}_{ij} || p(j|\mathbf{x}_i)) + \log p(\mathbf{x}_i)\end{aligned}$$

- What if we fix all the other parameters and maximize  $\mathcal{LL}$  w.r.t.  $\mathbf{U}_{ij}$ ?
  - When the KL divergence is minimal:  $\mathbf{U}_{ij} = p(j|\mathbf{x}_i)$
- This procedure is the **E-step** of the EM algorithm for GMM



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

## - Expectation Maximization (EM)

### ○ M-step

- We find  $\Theta$  that maximizes  $\mathcal{L} + \lambda(\sum_{j=1}^J P_j - 1)$

$$\begin{aligned}\mathcal{L} &\geq \sum_{i=1}^N \sum_{j=1}^J U_{ij} \log \left( \frac{P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{U_{ij}} \right) \\ &= \sum_{i=1}^N \sum_{j=1}^J U_{ij} \log (P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)) - \sum_{i=1}^N \sum_{j=1}^J U_{ij} \log U_{ij} \quad \text{constant}\end{aligned}$$

- Therefore the final objective function for the M-step is

$$\arg \max_{\Theta} \sum_{i=1}^N \sum_{j=1}^J U_{ij} \log (P_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)) + \lambda \left( \sum_j P_j - 1 \right)$$

- More specifically (since we're going to do the partial differentiation)

$$\arg \max_{\boldsymbol{\mu}_j} \mathcal{J}_{\boldsymbol{\mu}_j}, \quad \mathcal{J}_{\boldsymbol{\mu}_j} = \sum_{i=1}^N -\frac{1}{2} \mathbf{U}_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) + const.$$

$$\arg \max_{P_j} \mathcal{J}_{P_j}, \quad \mathcal{J}_{P_j} = \sum_{i=1}^N \mathbf{U}_{ij} \log P_j + \lambda (\sum_j P_j - 1) + const.$$

$$\arg \max_{\boldsymbol{\Sigma}_j} \mathcal{J}_{\boldsymbol{\Sigma}_j}, \quad \mathcal{J}_{\boldsymbol{\Sigma}_j} = \sum_{i=1}^N \mathbf{U}_{ij} \left( -\frac{1}{2} \log |\boldsymbol{\Sigma}_j| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) \right) + const.$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

## - Expectation Maximization (EM)

- M-step

- Partial differentiation w.r.t. the parameters and find the local maxima

- For the means:

$$\frac{\partial \mathcal{J}_{\boldsymbol{\mu}_j}}{\partial \boldsymbol{\mu}_j} = \sum_{i=1}^N -\mathbf{U}_{ij}(\mathbf{x}_i - \boldsymbol{\mu}_j) = 0, \quad \boldsymbol{\mu}_j = \frac{\sum_{i=1}^N \mathbf{U}_{ij} \mathbf{x}_i}{\sum_{i=1}^N \mathbf{U}_{ij}}$$

$$\arg \max_{\boldsymbol{\mu}_j} \mathcal{J}_{\boldsymbol{\mu}_j}, \quad \mathcal{J}_{\boldsymbol{\mu}_j} = \sum_{i=1}^N -\frac{1}{2} \mathbf{U}_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) + const.$$

- For the priors:

$$\frac{\partial \mathcal{J}_{P_j}}{\partial \lambda} = (\sum_j P_j - 1) = 0 \quad \frac{\partial \mathcal{J}_{P_j}}{\partial P_j} = \frac{\sum_{i=1}^N \mathbf{U}_{ij}}{P_j} + \lambda = 0$$

$$\arg \max_{P_j} \mathcal{J}_{P_j}, \quad \mathcal{J}_{P_j} = \sum_{i=1}^N \mathbf{U}_{ij} \log P_j + \lambda (\sum_j P_j - 1) + const.$$

$$\Leftrightarrow \sum_{j=1}^J \sum_{i=1}^N \mathbf{U}_{ij} = -\lambda \sum_{j=1}^J P_j$$

$$\Leftrightarrow \lambda = -\sum_{j=1}^J \sum_{i=1}^N \mathbf{U}_{ij} = -N$$

$$\Leftrightarrow P_j = \frac{\sum_{i=1}^N \mathbf{U}_{ij}}{N}$$

- For the covariance (see matrixcookbook 2.1.2 and 2.2):

$$\boldsymbol{\Sigma}_j = \frac{\sum_i \mathbf{U}_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j) (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top}{\sum_i \mathbf{U}_{ij}}$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

## - Expectation Maximization (EM)

- E-step: calculate posterior probabilities

$$U_{ij} = p(j|x_i) = \frac{P_j \mathcal{N}(x_i; \mu_j, \Sigma_j)}{\sum_j P_j \mathcal{N}(x_i; \mu_j, \Sigma_j)}$$

- M-step: update parameters

$$\mu_j = \frac{\sum_{i=1}^N U_{ij} x_i}{\sum_{i=1}^N U_{ij}}$$

$$P_j = \frac{\sum_{i=1}^N U_{ij}}{N}$$

$$\Sigma_j = \frac{\sum_i U_{ij} (x_i - \mu_j)(x_i - \mu_j)^\top}{\sum_i U_{ij}}$$

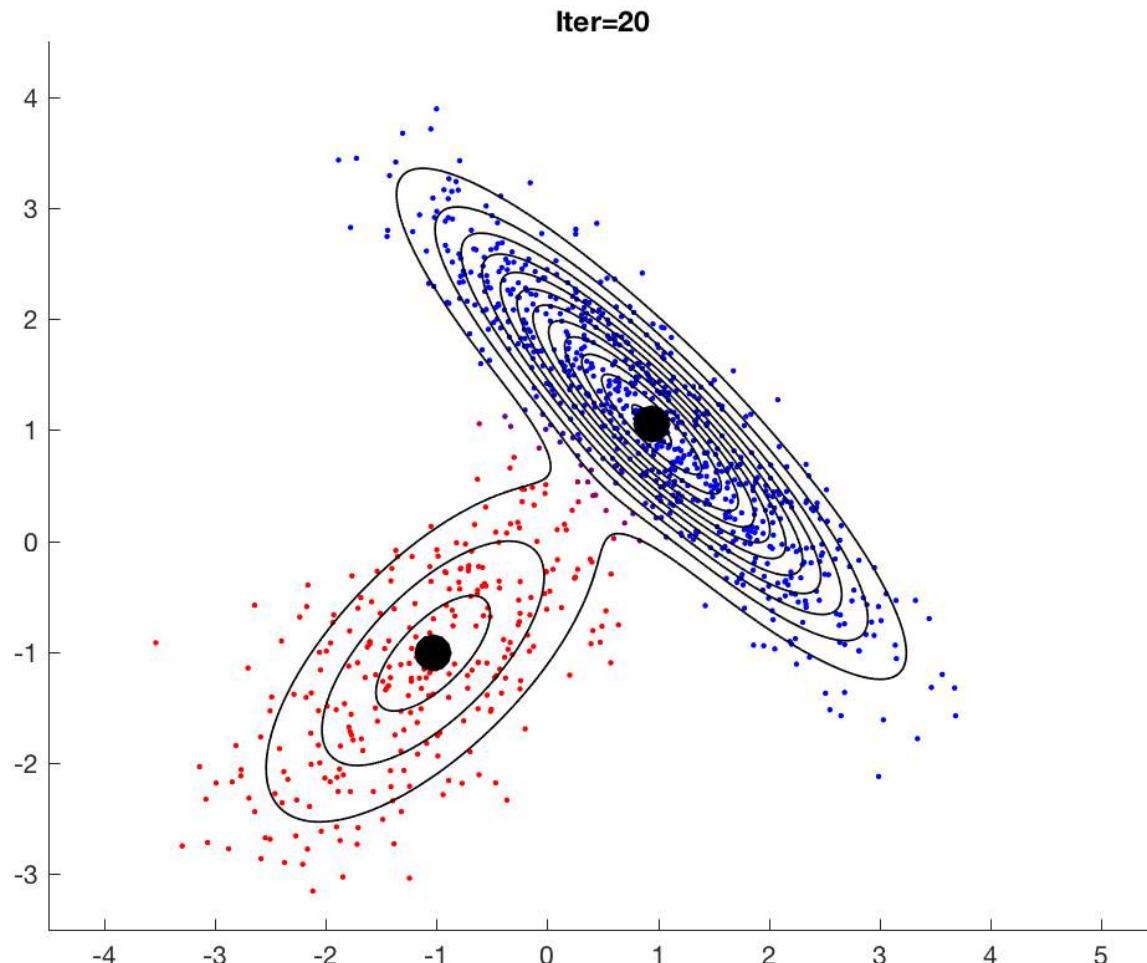


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Gaussian Mixture Model

- Too much math?

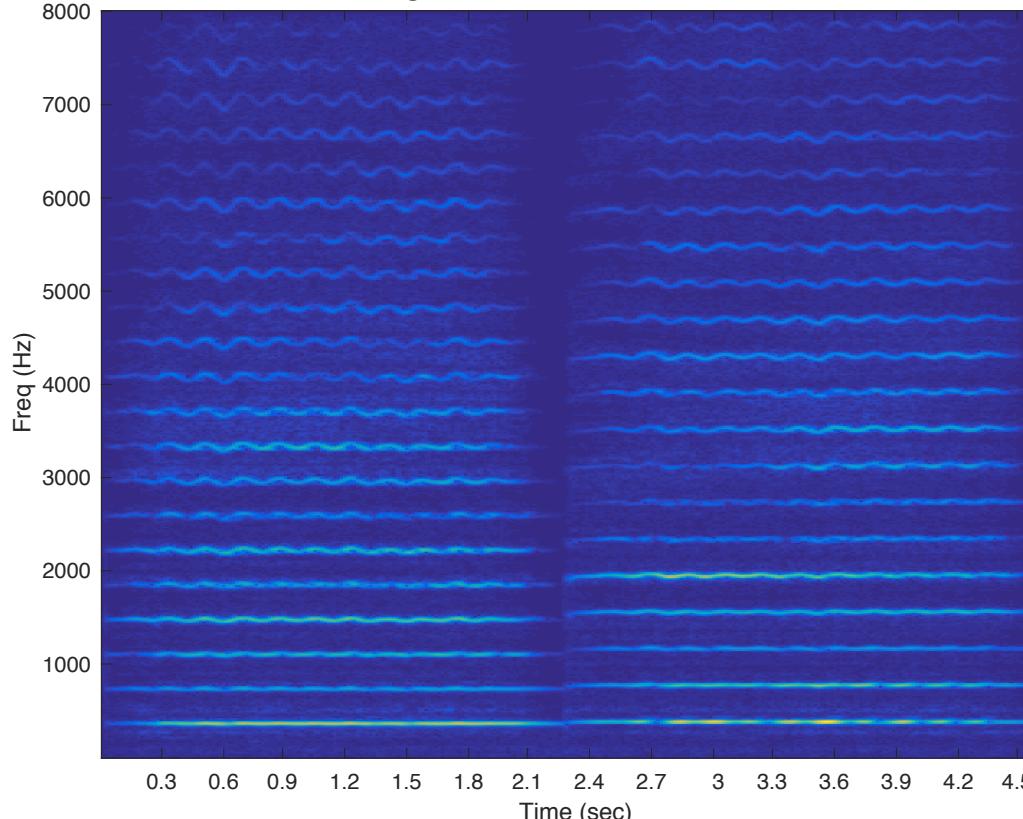


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Mixture of Multinomial Distributions

- Clustering Musical Notes
  - How many clusters? (samples are magnitude spectra)
    - I'm curious what kind of notes are there in the signal



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Mixture of Multinomial Distributions

## - Clustering Musical Notes

- For i-th spectrum  $\mathbf{x}_i$  :

- Although the magnitudes are real numbers we can scale them to convert them into integers:

$$[1.2, 1.35, 0.1, 5.525] = 0.025 \times [48, 54, 4, 221]$$

- Then, we can think of this as an observation from a multinomial dist.

$$\mathcal{M}(\mathbf{x}_i; \boldsymbol{\theta}) = \frac{N!}{\prod_d x_{id}!} \prod_d \theta_d^{x_{id}}$$

- EM for mixture of multinomial distributions

- Initialize two mean spectra (random numbers)  $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}_+^D$
  - Initialize two prior prob (random numbers that sum to one)  $P_1 + P_2 = 1$
  - Calculate posterior prob (E-step)

$$\begin{aligned} p(j=1|\mathbf{x}_i) &= \frac{P_1 \mathcal{M}(\mathbf{x}_i; \boldsymbol{\theta}_1)}{P_1 \mathcal{M}(\mathbf{x}_i; \boldsymbol{\theta}_1) + P_2 \mathcal{M}(\mathbf{x}_i; \boldsymbol{\theta}_2)} \\ &= \frac{P_1 \frac{N!}{\prod_d x_{id}!} \prod_d \theta_{1d}^{x_{id}}}{P_1 \frac{N!}{\prod_d x_{id}!} \prod_d \theta_{1d}^{x_{id}} + P_2 \frac{N!}{\prod_d x_{id}!} \prod_d \theta_{2d}^{x_{id}}} = \frac{P_1 \prod_d \theta_{1d}^{x_{id}}}{P_1 \prod_d \theta_{1d}^{x_{id}} + P_2 \prod_d \theta_{2d}^{x_{id}}} \end{aligned}$$

It's usually fine not to convert the spectra into integers, although it's not strictly correct.  
In this example I just normalized each spectrum.

- Update means (M-step)

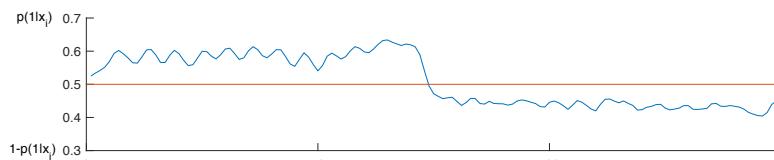
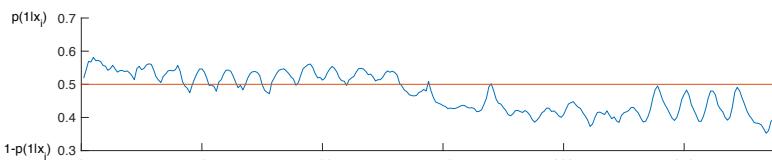
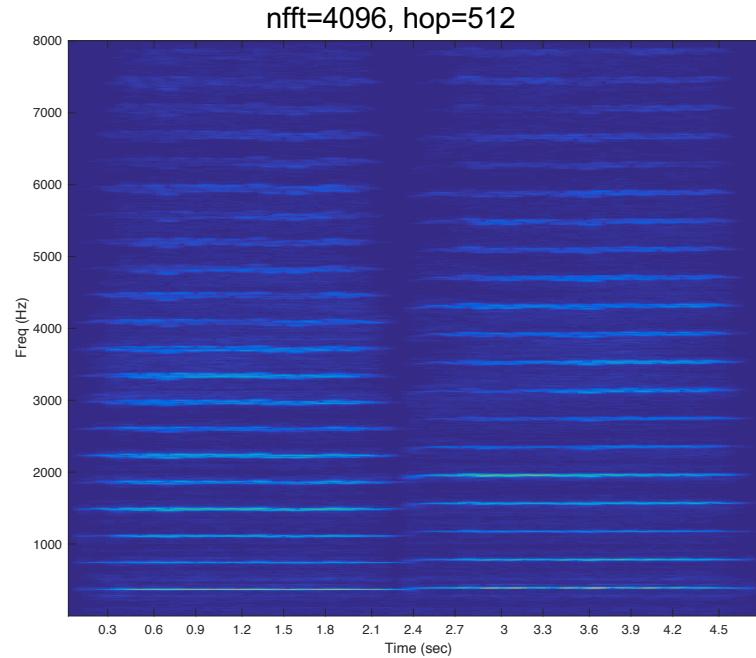
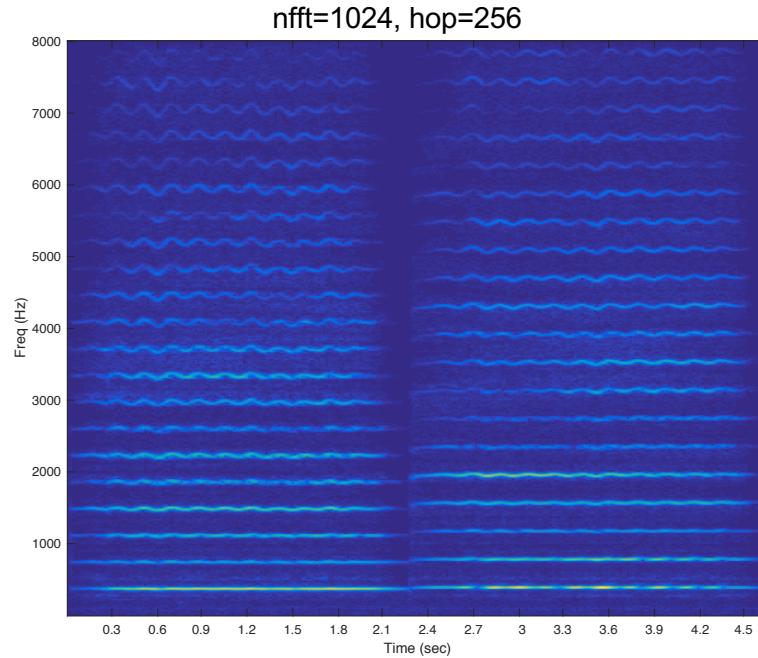


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Mixture of Multinomial Distributions

- Clustering Musical Notes
- It's actually a difficult clustering task with a lot of spurious local minima



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Locality Sensitive Hashing

## - From clustering to hashing

- Hashing is a popular concept in database
  - There is a query to the database
  - Instead of comparing the original representations, a hash function maps the query down to an integer (binary) address
  - The address is associated with a **bucket**. It can contain a few different database records
    - We say that those records collide
  - Then, we refine the search inside the bucket
  - This is cheaper than seeing the entire database
- Traditional challenges
  - Records are better off evenly distributed (for the speed)
  - Overflow



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Locality Sensitive Hashing

- From clustering to hashing
  - There's another hashing concept in machine learning
    - Locality sensitive hashing or semantic hashing
  - For the data points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  in a D-dimensional space
    - If they are close enough  $\mathcal{D}(\mathbf{x}_i||\mathbf{x}_j) < \tau$
    - Then, the Hamming distance between them after hashing is zero  $\mathcal{H}(\phi(\mathbf{x}_i)||\phi(\mathbf{x}_j)) = 0$  with probability  $p$
    - If they are far enough  $\mathcal{D}(\mathbf{x}_i||\mathbf{x}_j) \geq c\tau$
    - Then, the Hamming distance between them after hashing is zero  $\mathcal{H}(\phi(\mathbf{x}_i)||\phi(\mathbf{x}_j)) = 0$  with probability  $q$
    - $p > q$  !
  - The hash function  $\phi(\mathbf{x}_i)$  that meets above conditions are said to be in the locality sensitive hash function family
  - In other words...
    - Originally similar items collide in the same bucket
    - Share the same address
    - Quantized using the same binary string
    - Are in the same cluster



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Locality Sensitive Hashing

## - From clustering to hashing

- How do we find the hash function?
  - Well, it's not easy
- One way is to rely on a bunch of random projections

$$K \begin{bmatrix} +1 \\ -1 \\ \vdots \\ -1 \end{bmatrix} = \text{sign} \left( \begin{bmatrix} -1 & +1 & \cdots & -1 \\ +1 & -1 & \cdots & +1 \\ \vdots & \vdots & \ddots & \vdots \\ +1 & +1 & \cdots & -1 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \right)$$

- With  $K$  projections, we can represent  $D$  dimensional data with  $K$  bits
  - e.g. 513 magnitude Fourier coefficients (real) with  $K=128$  bits
- Why does it work?
  - Johnson-Lindenstrauss Theorem
- It's also related to sparse coding and compressive sensing



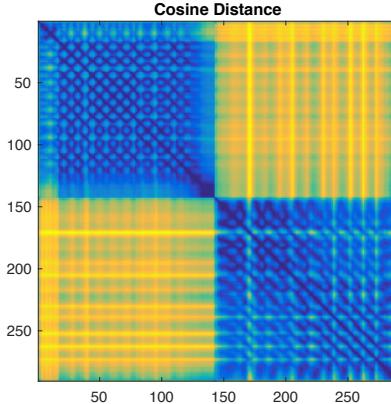
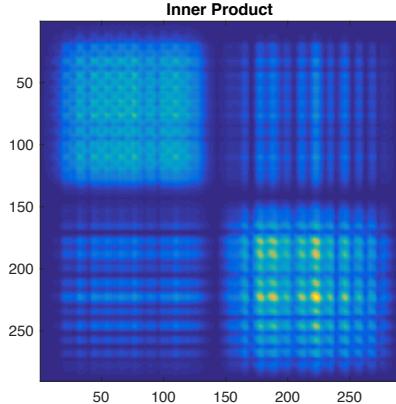
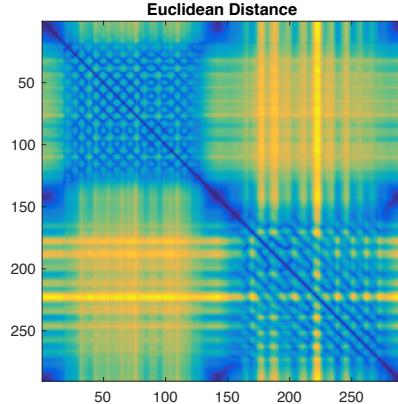
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

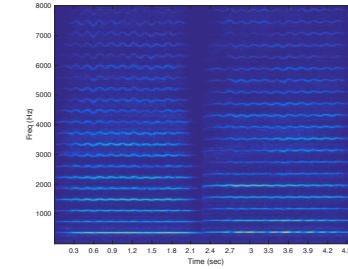
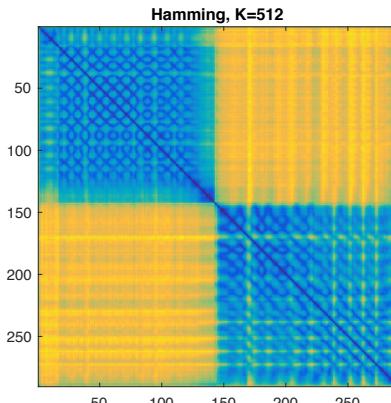
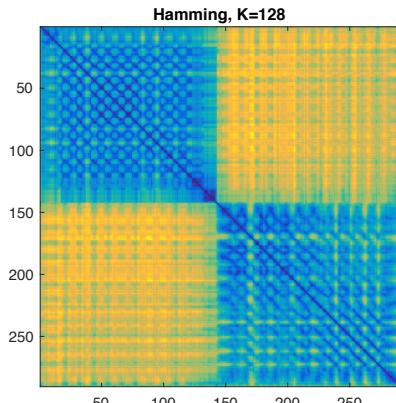
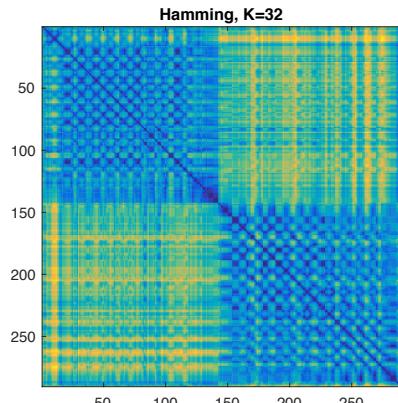
# Locality Sensitive Hashing

## - From clustering to hashing

- Let's see how the original spectra are similar to each other



- And in terms of Hamming distance after random projection



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Spectral Hashing

- More machine learning involved

- Just another hashing technique, but it tries to minimize the difference between

Original pairwise similarity

$$\min \sum_{ij} W_{ij} \mathcal{H}(y_i || y_j)$$

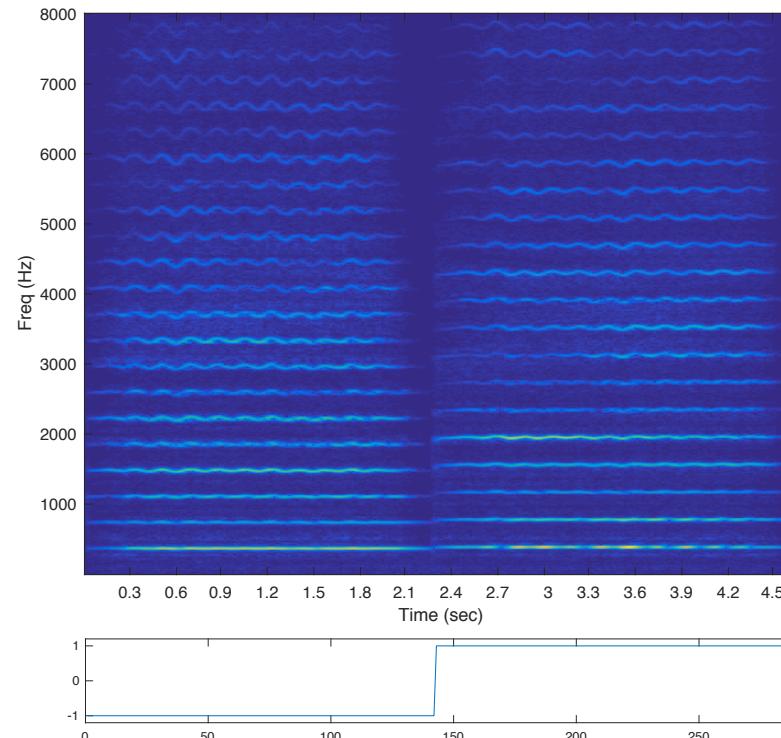
Pairwise Hamming distance

subject to:  $y_i \in \{-1, +1\}^K$

$$\sum_i y_i = \mathbf{0}$$

$$y_i^\top y_j = 0 \quad \text{if } i \neq j$$

- See the paper for the more optimization detail
  - But, the basic idea is to see the problem as an eigendecomposition problem
  - That's why it's called **spectral** hashing



INDIANA UNIVERSITY

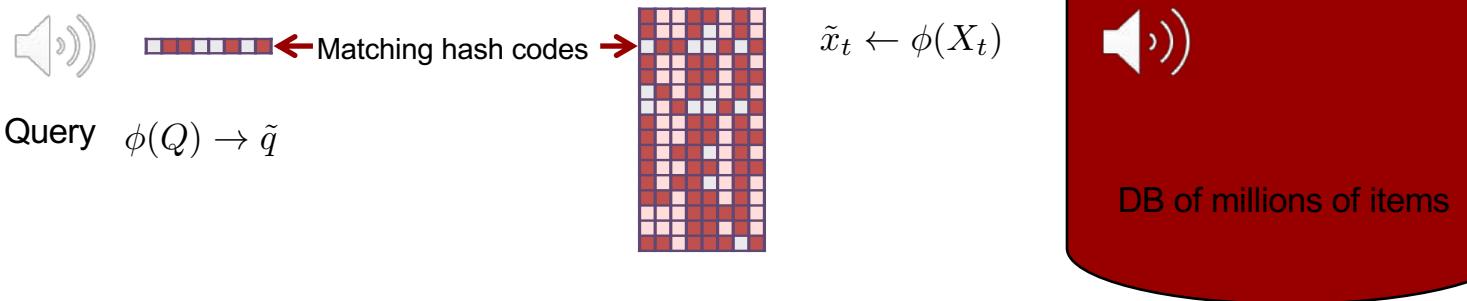
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Weiss, Yair, Antonio Torralba, and Rob Fergus. "Spectral hashing." *Advances in neural information processing systems*. 2009.

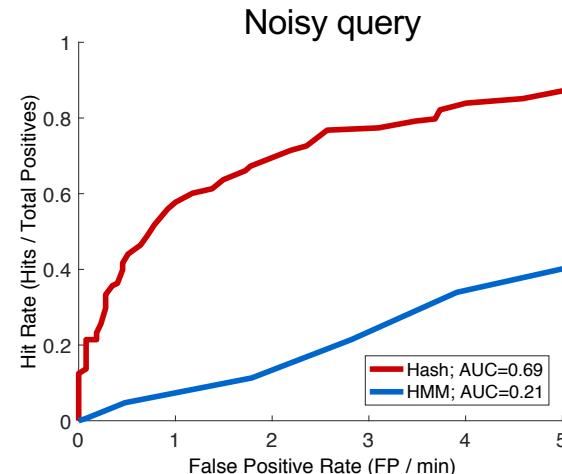
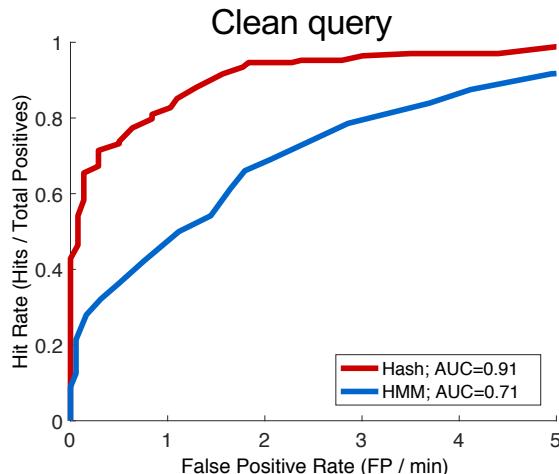
# Locality Sensitive Hashing

- Why is it useful?

- For faster detection



- Keyword spotting demo [HMM](#) versus [Spatial-Temporal WTA Hashing](#)



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

# Reading

- Textbook 6.8 – 6.16
- Textbook 2.5.5
- Bishop, “Pattern Recognition and Machine Learning” Chapter 9



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING



# Thank You!



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING