# Music Information Retrieval in Polyphonic Mixtures
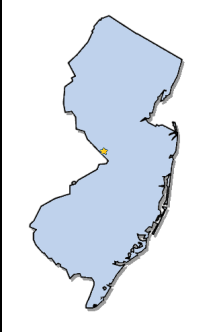
Steve Tjoa

MIR Workshop
CCRMA, Stanford University
iZotope, Inc.
San Francisco, CA, USA
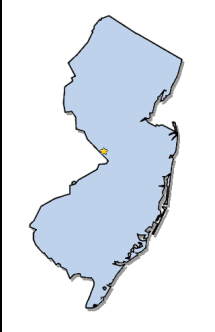
June 27, 2012

**iZotope**

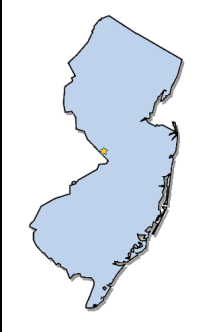A bit about myself...

## Quick Review

- What are the three main components of any classification system?

## Quick Review

- What are the three main components of any classification system?
- What are some useful features for MIR?
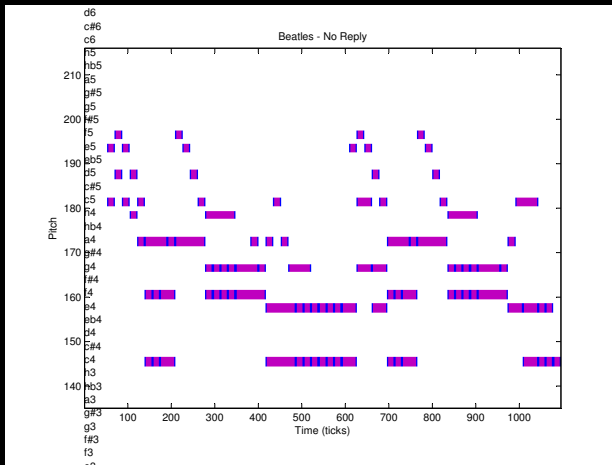
## Quick Review

- What are the three main components of any classification system?
- What are some useful features for MIR?
- What are some problems and applications addressed by MIR?

## Music Transcription

From this song...

## Music Transcription

From this song... get the "piano roll":

## Music Source Separation

Isolate, amplify, or suppress a musical voice/instrument.

Example: From these beats...

## Music Source Separation

Isolate, amplify, or suppress a musical voice/instrument.

Example: From these beats... isolate the kick drum and snare drum.

## A Really Special Tool

**Nonnegative Matrix Factorization** (NMF):

- Given $\mathbf{X}$ nonnegative, find $\mathbf{W}$ and $\mathbf{H}$, both nonnegative, that minimize some distance $d(\mathbf{X}, \mathbf{WH})$.

## A Really Special Tool

**Nonnegative Matrix Factorization** (NMF):

- Given $\mathbf{X}$ nonnegative, find $\mathbf{W}$ and $\mathbf{H}$, both nonnegative, that minimize some distance $d(\mathbf{X}, \mathbf{WH})$.
- Easy! And it works.
- *Meaningful* to humans.
- Widely used.

## Why NMF?

Energy of musical events are **nonnegative**.

## Brief Refresher: Matrix Multiplication

$$\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = a + 2b$$

## Brief Refresher: Matrix Multiplication

$$\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = a + 2b$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} a & b & c \end{bmatrix} = \begin{bmatrix} 3a & 3b & 3c \\ 4a & 4b & 4c \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = a + 2b$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} a & b & c \end{bmatrix} = \begin{bmatrix} 3a & 3b & 3c \\ 4a & 4b & 4c \end{bmatrix}$$

$$\mathbf{w} \begin{bmatrix} a & b & c \end{bmatrix} = \begin{bmatrix} a\mathbf{w} & b\mathbf{w} & c\mathbf{w} \end{bmatrix}$$

## Brief Refresher: Matrix Multiplication

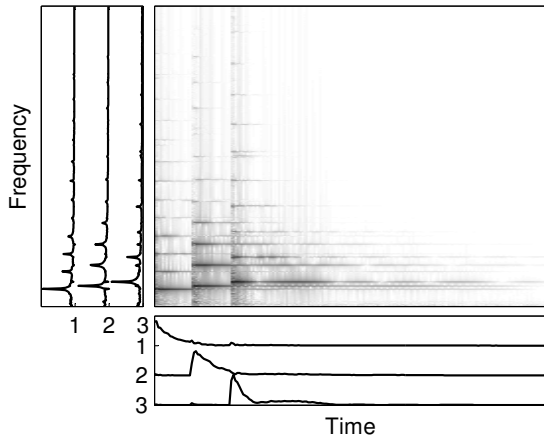$$\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = a + 2b$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} a & b & c \end{bmatrix} = \begin{bmatrix} 3a & 3b & 3c \\ 4a & 4b & 4c \end{bmatrix}$$

$$\mathbf{w} \begin{bmatrix} a & b & c \end{bmatrix} = \begin{bmatrix} a\mathbf{w} & b\mathbf{w} & c\mathbf{w} \end{bmatrix}$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix} \mathbf{h} = \begin{bmatrix} 3\mathbf{h} \\ 4\mathbf{h} \end{bmatrix}$$
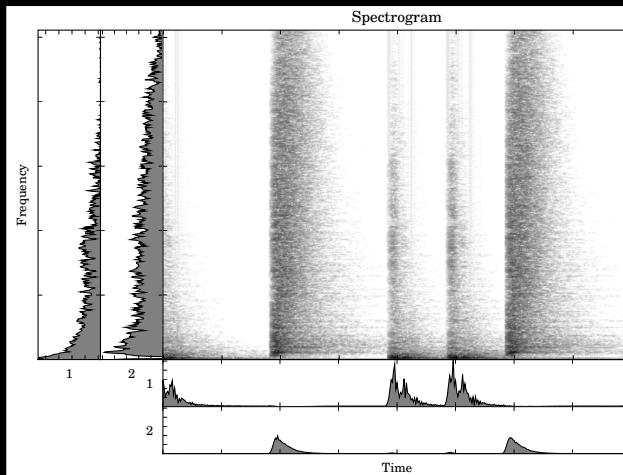
**Nonnegative Matrix Factorizaton**

Top right: **X**. Left: **W**. Bottom: **H**. Three piano notes:

## Nonnegative Matrix Factorizaton

Top right: **X**. Left: **W**. Bottom: **H**. Kick and snare:

## NMF Algorithms

Multiplicative update rules:

$$\mathbf{W} \leftarrow \mathbf{W} \cdot \frac{\mathbf{X}\mathbf{H}^T}{\mathbf{W}\mathbf{H}\mathbf{H}^T} \qquad \mathbf{H} \leftarrow \mathbf{H} \cdot \frac{\mathbf{W}^T\mathbf{X}}{\mathbf{W}^T\mathbf{W}\mathbf{H}}$$

See [Lee and Seung, NIPS 2001].

## NMF Algorithms

Easy to implement!
Python:

```python
1  for iter in range(maxiter):
2      W = multiply(W, (X*H.T)/(W*H*H.T))
3      H = multiply(H, (W.T*X)/(W.T*W*H))
```

## NMF Algorithms

Easy to implement!
Python:

```
1  for iter in range(maxiter):
2      W = multiply(W, (X*H.T)/(W*H*H.T))
3      H = multiply(H, (W.T*X)/(W.T*W*H))
```

Matlab:

```
1  for iter=1:maxiter
2      W = W.*(X*H')./(W*H*H');
3      H = H.*(W'*X)./(W'*W*H);
4  end
```

## Example: Source Separation

kick and snare:

- [kick drum] and [snare drum]

## Example: Source Separation

kick and snare:

- [kick drum] and [snare drum]

oboe and horn:

- Duan et. al: [oboe] and [horn]
- Wang et. al: [oboe] and [horn]
- Tjoa and Liu: [oboe] and [horn]

## Example: Source Separation

kick and snare:

- [kick drum] and [snare drum]
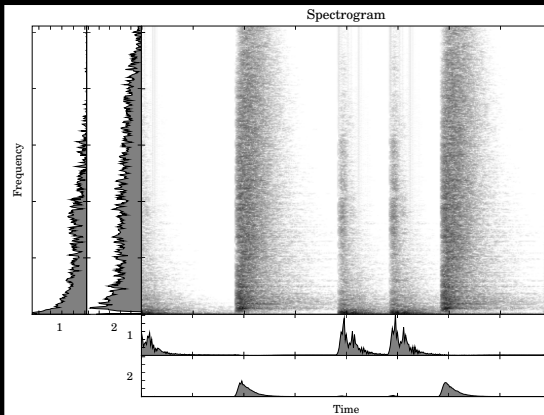
oboe and horn:

- Duan et. al: [oboe] and [horn]
- Wang et. al: [oboe] and [horn]
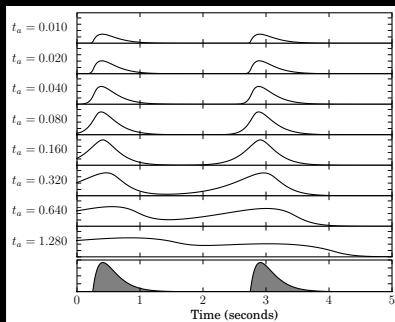- Tjoa and Liu: [oboe] and [horn]

Vivaldi, *Winter, Four Seasons*:

- [solo] and [accompaniment]

## Example: Instrument Recognition

Use NMF to identify the instruments in a musical signal.
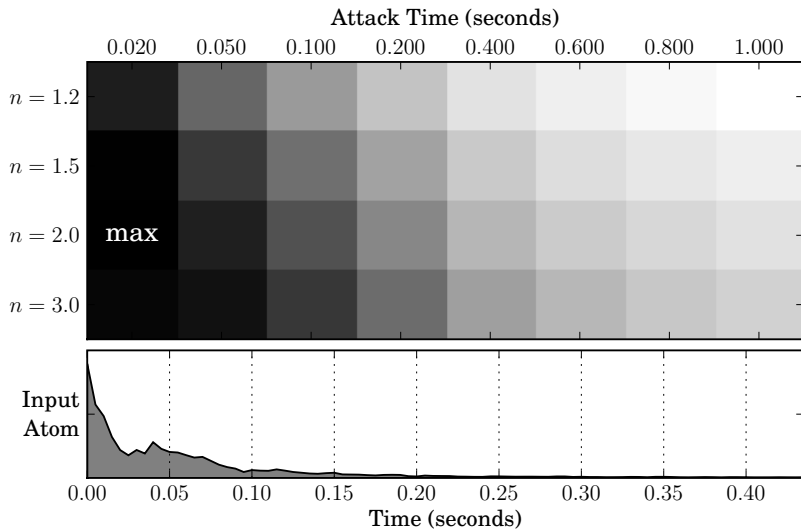Observe these atoms:

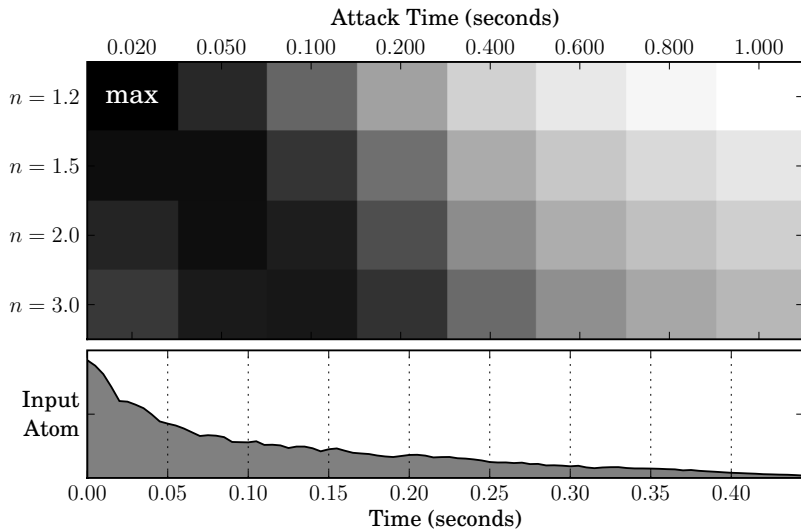Filter the temporal atoms from NMF [Tjoa and Liu, 2010]:



- Use support vector machine (SVM) to classify the processed spectral and temporal atoms.
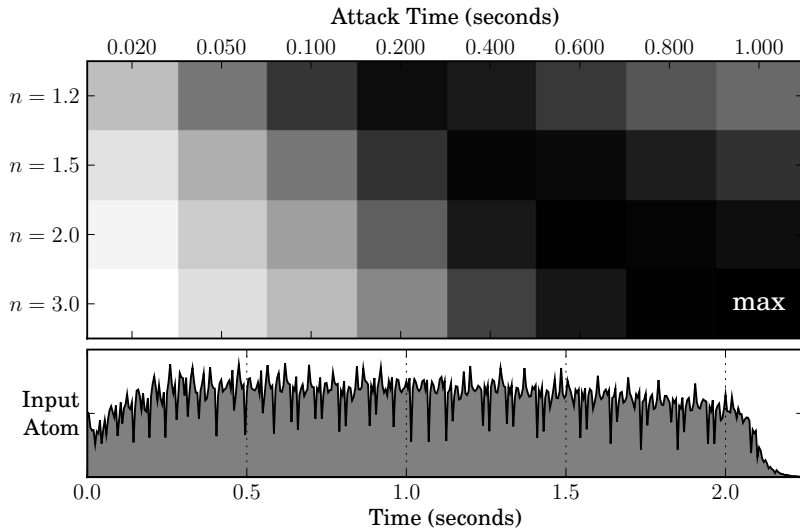
Feature Vector of Kick Drum

**Feature Vector of Snare Drum**

Attack Time (seconds)

|  | 0.020 | 0.050 | 0.100 | 0.200 | 0.400 | 0.600 | 0.800 | 1.000 |

$n = 1.2$ — max

$n = 1.5$

$n = 2.0$

$n = 3.0$

Input Atom

Time (seconds)

0.00   0.05   0.10   0.15   0.20   0.25   0.30   0.35   0.40

## Feature Vector of Trumpet

**Feature Vector of Violin**

Attack Time (seconds)

0.020   0.050   0.100   0.200   0.400   0.600   0.800   1.000

$n = 1.2$

$n = 1.5$        max

$n = 2.0$

$n = 3.0$

Input Atom

0   1   2   3   4   5   6

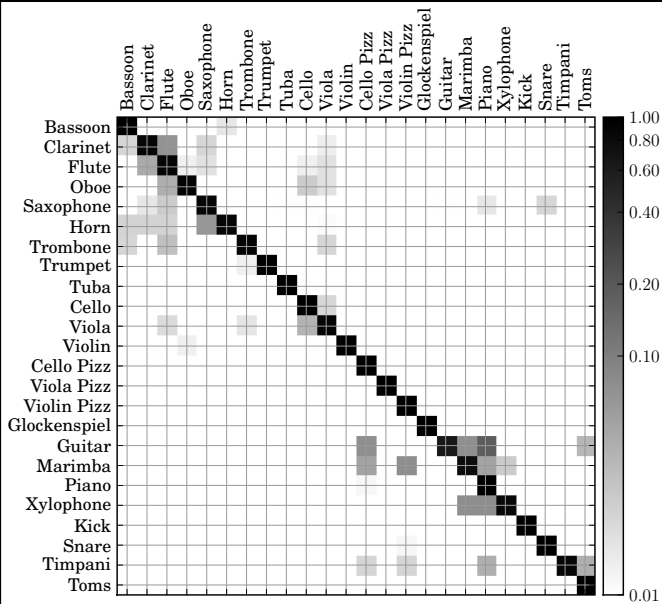Time (seconds)

## Results: Isolated Instrument Recognition

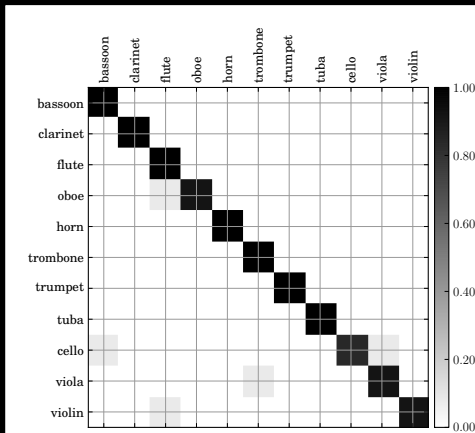Experiments on isolated instrument sounds:

- Accuracy: **92.3%**
- Reflect state-of-the-art performance for isolated instrument recognition among as many as 24 classes.

## Results: Solo Melodic Phrases
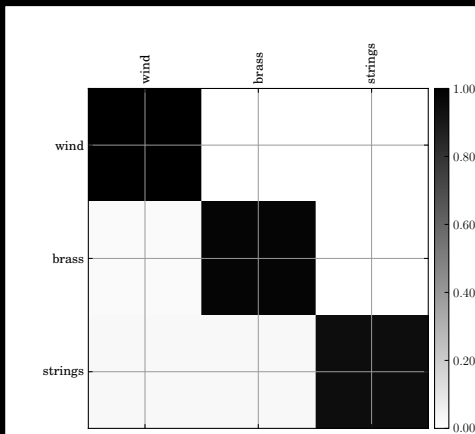
Instrument classifications. One decision per signal.
Accuracy: **96.2%**.

## Results: Solo Melodic Phrases

Family classifications. One decision per signal.
Accuracy: **97.4%**.

Existing algorithms cannot handle "**complicated**" music.

Related work:

- smoothness
- harmonicity
- statistical priors

## Sparse Coding

**What if you already have a large dictionary?**

$$\min_{\mathbf{s}} d(\mathbf{x}, \mathbf{As})$$

- Solution: Impose sparsity on $\mathbf{s}$.
- Benefits: guaranteed spectral structure; labels already known.

## Sparse Coding

Related work:

- matching pursuit (MP)
- orthogonal matching pursuit (OMP)
- basis pursuit (BP)

Disadvantages:

- Complexity that is **linear** in the dictionary size.
- **Neither fast nor scalable.**

# Example: Orthogonal Matching Pursuit

OMP [Pati et al., 1993]:

- Input: $\mathbf{x} \in \mathbb{R}^M$; $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, ..., \mathbf{a}_K] \in \mathbb{R}^{M \times K}$ s.t. $||\mathbf{a}_k||_2 = 1$ for all $k$.
- Output: $\hat{\mathbf{s}} \in \mathbb{R}^K$
- Initialize: $\mathcal{S} \leftarrow \emptyset$; $\mathbf{s} \leftarrow \mathbf{0}$; $\mathbf{r} \leftarrow \mathbf{x}$; $\epsilon > 0$.
- While $||\mathbf{r}|| > \epsilon$:
  1. $k \leftarrow \text{argmax}_j \, \mathbf{a}_j^T \mathbf{r}$
  2. $\mathcal{S} \leftarrow \mathcal{S} \cup k$
  3. Solve for $\{s_j | j \in \mathcal{S}\}$: $\min_{s_j | j \in \mathcal{S}} ||\mathbf{x} - \sum_{j \in \mathcal{S}} \mathbf{a}_j s_j||$
  4. $\mathbf{r} \leftarrow \mathbf{x} - \mathbf{A}\mathbf{s}$
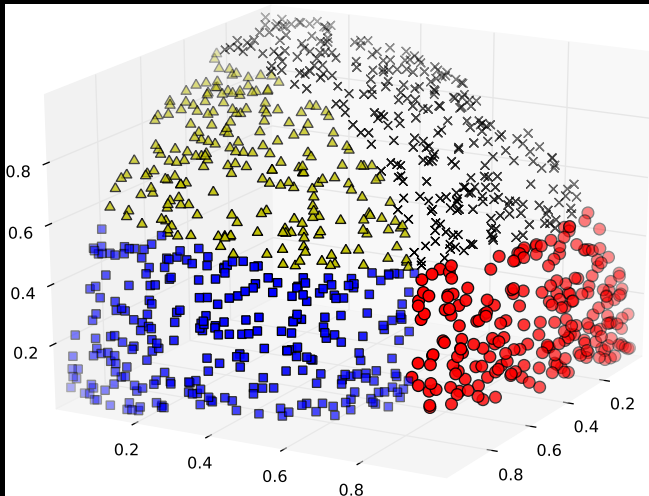- $\hat{\mathbf{s}} \leftarrow \mathbf{s}$

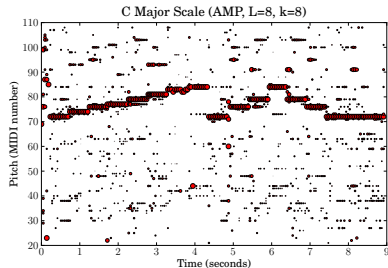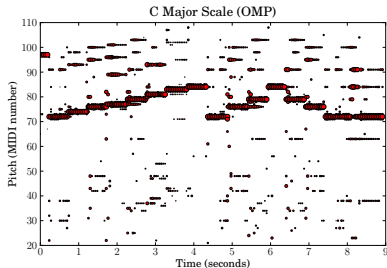## Proposed Algorithm: Approximate Matching Pursuit

AMP [Tjoa and Liu]:

- Input: $\mathbf{x} \in \mathbb{R}^M$; $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, ..., \mathbf{a}_K] \in \mathbb{R}^{M \times K}$ s.t. $||\mathbf{a}_k||_2 = 1$ for all $k$.
- Output: $\hat{\mathbf{s}} \in \mathbb{R}^K$
- Initialize: $\mathcal{S} \leftarrow \emptyset$; $\mathbf{s} \leftarrow \mathbf{0}$; $\mathbf{r} \leftarrow \mathbf{x}$; $\epsilon > 0$.
- While $||\mathbf{r}|| > \epsilon$:
  1. Find any $k$ such that $\mathbf{a}_k$ and $\mathbf{r}$ are near neighbors.
  2. $\mathcal{S} \leftarrow \mathcal{S} \cup k$
  3. Solve for $\{s_j | j \in \mathcal{S}\}$: $\min_{s_j | j \in \mathcal{S}} ||\mathbf{x} - \sum_{j \in \mathcal{S}} \mathbf{a}_j s_j||$
  4. $\mathbf{r} \leftarrow \mathbf{x} - \mathbf{A}\mathbf{s}$
- $\hat{\mathbf{s}} \leftarrow \mathbf{s}$

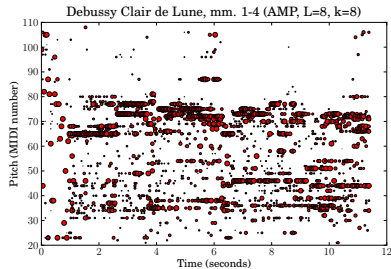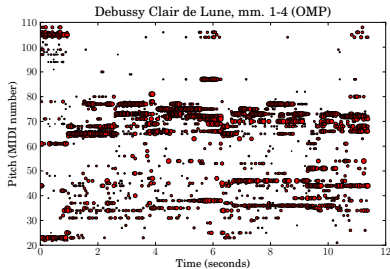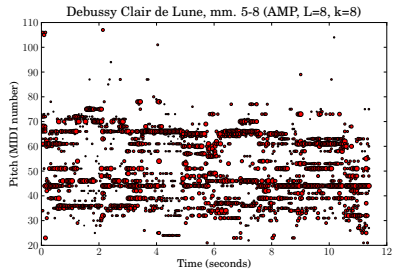## Locality Sensitive Hashing

Idea: Hash nearby points into the same bin.

# Experiments: Music Transcription

# Experiments: Music Transcription

## Experiments: Music Transcription

Execution times in seconds.

| Song | OMP | $AMP_{8,8}$ | $AMP_{10,10}$ |
|---|---|---|---|
| C-major scale | 81.05 | 43.63 | 21.03 |
| Debussy mm. 1-4 | 118.57 | 88.45 | 29.01 |
| Debussy mm. 5-8 | 123.05 | 121.73 | 121.84 |

## Where to Learn More

Conferences:

- Int. Society of Music Information Retrieval (ISMIR)
- MIR Evaluation Exchange (MIREX)
- Int. Computer Music Conference (ICMC)
- IEEE Int. Conf. Audio, Speech, Signal Processing (ICASSP)
- ACM Multimedia

Journals:

- IEEE Trans. Audio, Speech, Language, Processing
- Journal of New Music Research
- Computer Music Journal

# Lab 3

## Lab 3: Summary

Summary:

## Lab 3: Matlab Programming Tips

- Pressing the up and down arrows let you scroll through command history.
- A semicolon at the end of a line simply means "suppress output".
- Type `help <command>` for instant documentation. For example, `help wavread`, `help plot`, `help sound`. Use `help` liberally!

## Lab 3.1: Source Separation

1. In Matlab: Select File $\rightarrow$ Set Path.
   Select "Add with Subfolders".
   Select /usr/ccrma/courses/mir2011/lab3skt.

2. As in Lab 1, load the file, listen to it, and plot it.

```
1  [x, fs] = wavread('simpleLoop.wav');
2  sound(x, fs)
3  t = (0:length(x)-1)/fs;
4  plot(t, x)
5  xlabel('Time (seconds)')
```

3. Compute and plot a short-time Fourier transform, i.e., the Fourier transform over consecutive frames of the signal.

```
1  frame_size = 0.100;
2  hop = 0.050;
3  X = parsesig(x, fs, frame_size, hop);
4  imagesc(abs(X(200:-1:1,:)))
```

Type `help parsesig`, `help imagesc`, and `help abs` for more information.

This step gives you some visual intuition about how sounds (might) overlap.

4. Let's separate sources!

```
1  K = 2;
2  [y, W, H] = sourcesep(x, fs, K);
```

Type `help sourcesep` for more information.

5. Plot and listen to the separated signals.

```
1  plot(t, y)
2  xlabel('Time (seconds)')
3  legend('Signal 1', 'Signal 2')
4  sound(y(:,1), fs)
5  sound(y(:,2), fs)
```

Feel free to replace Signal 1 and Signal 2 with Kick and Snare (depending upon which is which).

6. Plot the outputs from NMF.

```
1 figure
2 plot(W(1:200,:))
3 legend('Signal 1', 'Signal 2')
4 figure
5 plot(H')
6 legend('Signal 1', 'Signal 2')
```

What do you observe from `W` and `H`?
Does it agree with the sounds you heard?

## Lab 3.1: Source Separation

7. Repeat the earlier steps for different audio files.
    - `125BOUNC-mono.WAV`
    - `58BPM.WAV`
    - `CongaGroove-mono.wav`
    - `Cstrum chord_mono.wav`

    ... and more.

    Experiment with different values for the number of sources, `K`.

    Where does this separation method succeed?

    Where does it fail?

Begin with simpleLoop.wav. Then try others.

1. Add noise to the input signal, plot, and listen.

```
1  xn = x + 0.01*randn(length(x),1);
2  plot(t, xn)
3  sound(xn, fs)
```

2. Separate, plot, and listen.

```
1 [yn, Wn, Hn] = sourcesep(xn, fs, K);
2 plot(t, yn)
3 sound(yn(:,1), fs)
4 sound(yn(:,2), fs)
```

How robust to noise is this separation method?
Compared to the noisy input signal, how much noise is left in the output signals?
Which output contains more noise? Why?

## Lab 3.3: Classification

Follow the K-NN example in Lab 1, but classify the *separated* signals.

1. As in Lab 1, extract features from each training sample in the kick and snare drum directories.

2. Train a K-NN model using the kick and snare drum samples.

```
1 labels=[[ones(10,1) zeros(10,1)];
2        [zeros(10,1) ones(10,1)]];
3 model_snare =
4   knn(5, 2, 1, trainingFeatures, labels);
5 [voting, model_output] =
6   knnfwd(model_snare, featuresScaled)
```

3. Extract features from the drum signals that you separated in Lab 3.1.
   Classify them using the K-NN model that you built.
   Does K-NN accurately classify the separated signals?
   Repeat for different numbers of separated signals (i.e., the parameter $K$ in NMF).

4. Overseparate the signal using $K = 20$ or more. For those separated components that are classified as snare, add them together using sum. The listen to the sum signal. Is it coherent, i.e., does it sound like a single separated drum?

## ...and more!

- If you have another idea that you would like to try out, please ask me!
- Please collaborate with a partner.
  Together, brainstorm your own problems, if you want!

**Good luck!**

# Music Information Retrieval in Polyphonic Mixtures

Steve Tjoa

MIR Workshop
CCRMA, Stanford University
iZotope, Inc.
San Francisco, CA, USA