



MSc in AI  
NCSR Demokritos - University of Piraeus

Course: **Machine Learning for Multimodal Data**

# Lesson 2

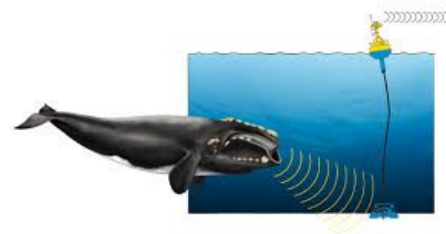
## Audio Representations and Feature Extraction

Theodoros Giannakopoulos



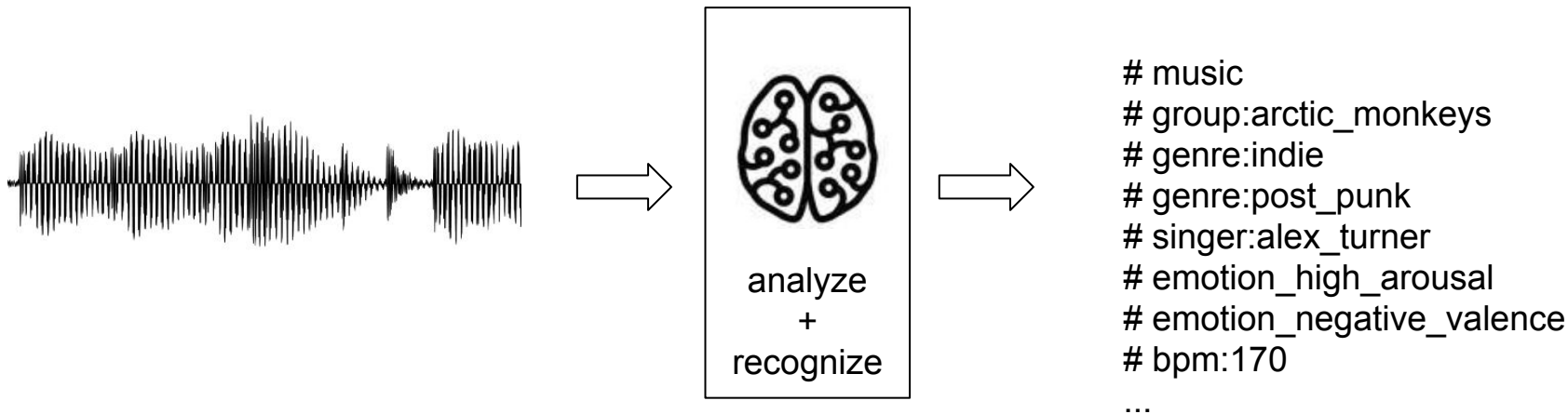
# Audio Analysis Applications

Applications	Audio	Speech	Music
Automatic Speech Recognition (ASR)			
Virtual Assistants			
Music Search			
Surveillance			
Environmental Monitoring			
Recommendation			



# Audio Analysis Goal & Applications

- Goal: extract high-level descriptions from raw audio signals (sounds)
- Using:
  - signal analysis: to extract features and representations
  - machine learning (supervised or unsupervised) to train models and to discover patterns
- Speech / Music / Audio
- Also referred to as “machine listening”



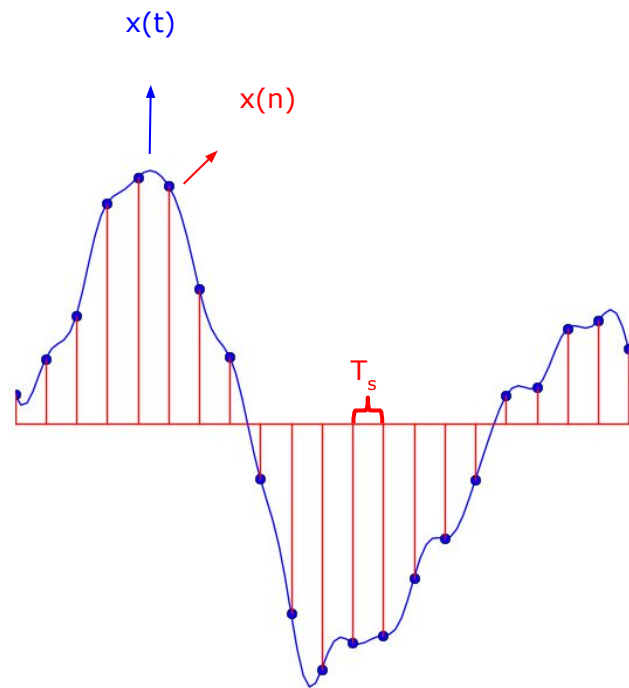
# Demo code

**Course's code samples will be available at this github repo:**

<https://github.com/tyiannak/multimodalAnalysis>

# What is sound?

- sound (physics):
  - a travelling vibration (wave)
  - through a medium (e.g. air)
  - transfers energy (particle to particle)
  - until “perceived” by our ears
- amplitude - loudness
- frequency - vibrations per sec
- analog sound → digital sound
  - sampling (sampling freq),  $f_s$
  - quantization (bits per sample)
  - example:
    - 44100 Hz
    - 16 bits per sample (sample resolution)
    - ~8 million integers for an average song! (single channel....)



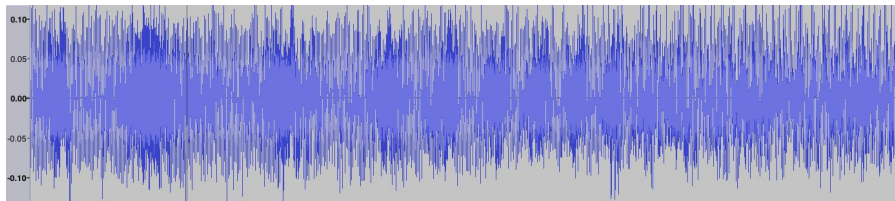
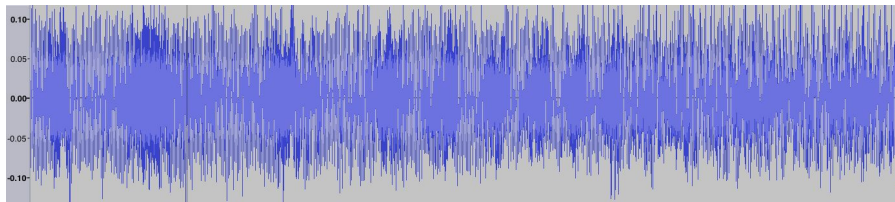
# DFT

- Discrete Fourier Transform (DFT)
- A frequency domain representation of the signal (time domain)
- FFT: efficient implementation of the DFT
- Given  $x(n)$  (signal), DFT is 
$$X(k) = \sum_{n=0}^{N-1} x(n) \exp(-j \frac{2\pi}{N} kn), k = 0, \dots, N-1$$
- Inverse  $\rightarrow x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k) \exp(j \frac{2\pi}{N} kn), n = 0, \dots, N-1,$
- Can be re-written in the form  $\rightarrow x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k) y_k(n), n = 0, \dots, N-1, \quad y_k(n) = \exp(j \frac{2\pi}{N} kn), n = 0, \dots, N-1$
- I.e. original signal can be written as a weighted average of complex exponentials (weights are DFT coefficients)

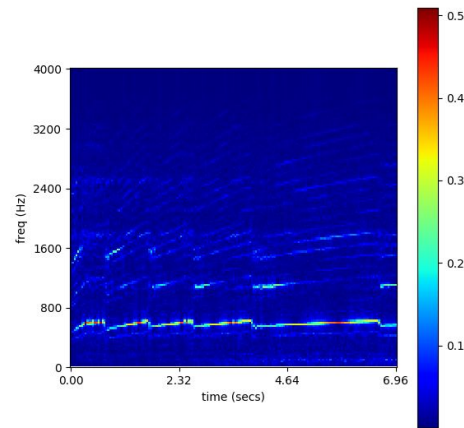
*DFT is:*

- *Defined in the range  $0..fs$*
- *Symmetric (center  $0..fs/2$ )*

# Representation: Time Vs Frequency

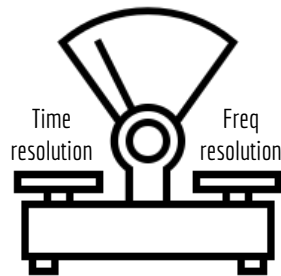


?



# Frequency Representation: Spectrogram

- Spectrogram: Time - Frequency 2D representation
- 1st step: Windowing
  - Signal broken into short-term windows (or frames)
  - Typically 20 to 100 mseconds
  - (Non)overlapping
  - E.g.: 50 msec frame size, 10 msec step: 80% overlap
- 2nd step: FFT
  - Fast Implementation of the Discrete Fourier Transform (DFT)
  - Not to be confused with Discrete Time Fourier Transform which:
    - Is continuous in the frequency domain
    - Is periodic in the range  $0..f_s$
- FFT Resolution & window length:
  - longer windows  $\rightarrow$  more dense representations  $\rightarrow$  better frequency resolutions
  - but also losing time resolution (because of signals' non - stationarity)
  - trade off between time and frequency resolution



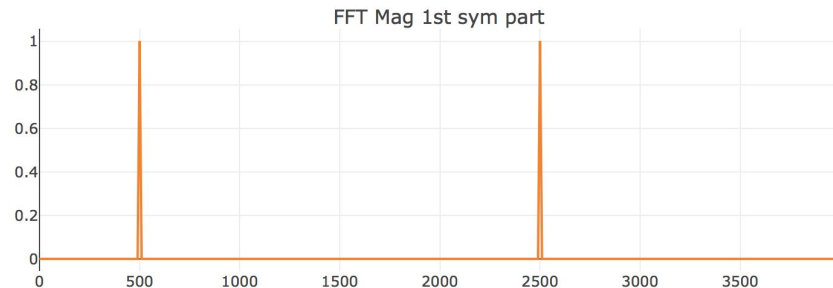
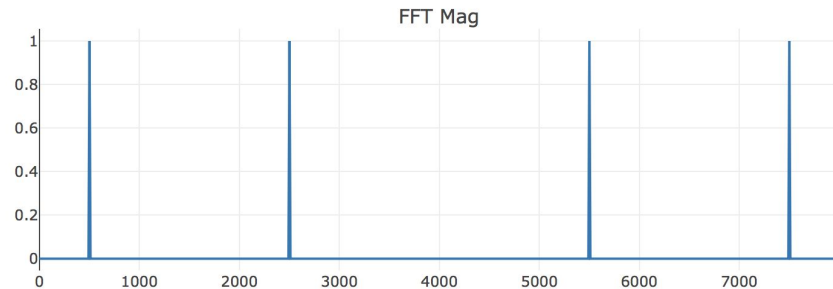


# Frequency Representation: FFT Example

```

1  """
2  @brief Example 01
3  @details FFT computation example
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import scipy.fftpack as scp
7  import numpy as np
8  import plotly
9  import plotly.graph_objs as go
10 if __name__ == '__main__':
11     f1, f2, fs, duration = 500, 2500, 8000, 0.1
12     # define time range
13     t = np.arange(0, duration, 1.0/fs)
14     # define signal as sum of cosines
15     x = np.cos(2 * np.pi * t * f1) + np.cos(2 * np.pi * t * f2)
16     # get mag of fft
17     X = np.abs(scp.fft(x))
18     # normalize FFT mag
19     X = X / X.max()
20     freqs = np.arange(0, 1, 1.0/len(X)) * fs
21     # get 1st symmetric part
22     freqs_1 = freqs[0:int(len(freqs)/2)]
23     X_1 = X[0:int(len(X)/2)]
24     figs = plotly.tools.make_subplots(rows=2, cols=1,
25                                       subplot_titles=["FFT Mag",
26                                                         "FFT Mag 1st sym part"])
27     figs.append_trace(go.Scatter(x=freqs, y=X, showlegend=False), 1, 1)
28     figs.append_trace(go.Scatter(x=freqs_1, y=X_1, showlegend=False), 2, 1)
29     plotly.offline.plot(figs, filename="temp.html", auto_open=True)

```



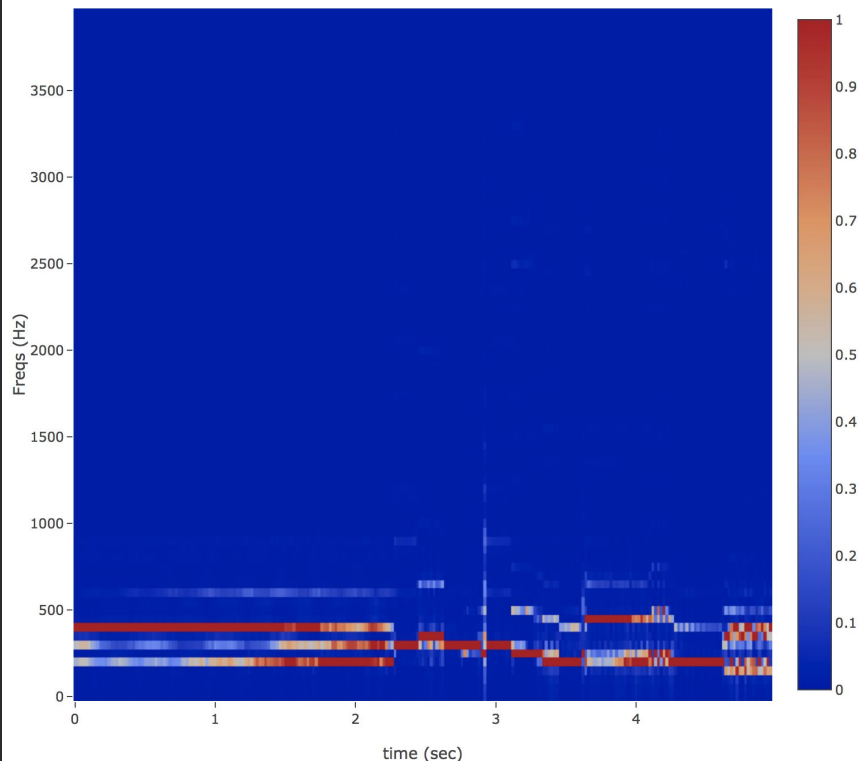
# Frequency Representation: Spectrogram Example

```

1  """
2  @brief Example 02
3  @details Example of spectrogram computation for a wav file, using only scipy
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import scipy.fftpack as scp
7  import numpy as np
8  import scipy.io.wavfile as wavfile
9  import plotly
10 import plotly.graph_objs as go
11 layout = go.Layout(title='Spectrogram Calculation Example',
12                    xaxis=dict(title='time (sec)',),
13                    yaxis=dict(title='Freqs (Hz)',))
14
15
16 def get_fft_spec(signal, fs, win):
17     frame_size, signal_len, spec, times = int(win * fs), len(signal), [], []
18     # break signal into non-overlapping short-term windows (frames)
19     frames = np.array([signal[x:x + frame_size] for x in
20                        np.arange(0, signal_len - frame_size, frame_size)])
21     for i_f, f in enumerate(frames): # for each frame
22         times.append(i_f * win)
23         # append mag of fft
24         X = np.abs(scp.fft(f)) ** 2
25         freqs = np.arange(0, 1, 1.0/len(X)) * (fs/2)
26         spec.append(X[0:int(len(X)/2)] / X.max())
27     return np.array(spec).T, freqs, times
28
29 if __name__ == '__main__':
30     [Fs, s] = wavfile.read("../data/sample_music.wav")
31     S, f, t = get_fft_spec(s, Fs, 0.02)
32     heatmap = go.Heatmap(z=S, y=f, x=t)
33     plotly.offline.plot(go.Figure(data=[heatmap], layout=layout),
34                          filename="temp.html", auto_open=True)

```

Spectrogram Calculation Example



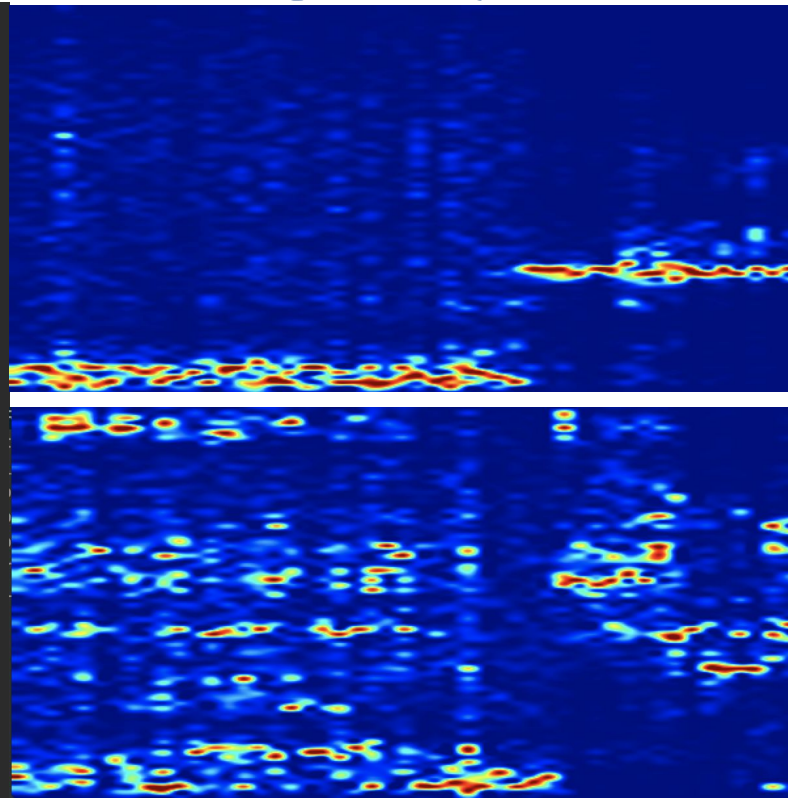
- scipy.fftpack used for fft

# Frequency Representation: Spectrogram Recording Example

```

1  """
2  @brief Example 03
3  @details Example of audio recording and spectrogram computation
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import scipy.io.wavfile as wavfile
8  import example02, pyaudio, struct, cv2, signal, sys
9
10 fs = 8000
11 bufSize = int(fs * 1.0)
12 def signal_handler(signal, x):
13     wavfile.write("output.wav", fs, np.int16(aggregated_buf))
14     sys.exit(0)
15 signal.signal(signal.SIGINT, signal_handler)
16
17 if __name__ == '__main__':
18     global aggregated_buf; aggregated_buf = np.array([])
19     pa = pyaudio.PyAudio() # initialize recording
20     stream = pa.open(format=pyaudio.paInt16, channels=1, rate=fs,
21                      input=True, frames_per_buffer=bufSize)
22     while 1:
23         # read recorded data, convert bytes to samples and then to numpy array
24         block = stream.read(bufSize, exception_on_overflow=False)
25         s = np.array(list(struct.unpack("%dh"%(len(block)/2), block))).astype(float)
26         aggregated_buf = np.concatenate((aggregated_buf, s))
27         s /= (2 ** 15)
28         # get spectrogram and visualize it using opencv
29         specgram, f, t = example02.get_fft_spec(s, fs, 0.02)
30         iSpec = np.array(specgram[:-1] * 255, dtype=np.uint8)
31         iSpec2 = cv2.resize(iSpec, (600, 350), interpolation=cv2.INTER_CUBIC)
32         iSpec2 = cv2.applyColorMap(iSpec2, cv2.COLORMAP_JET)
33         cv2.imshow('Spectrogram', iSpec2)
34         cv2.moveWindow('Spectrogram', 100, 100)
35         ch = cv2.waitKey(10)

```



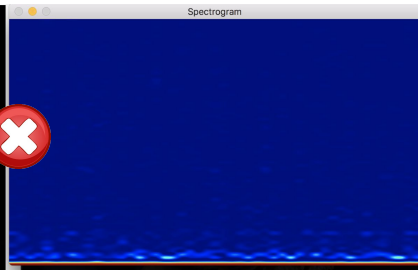
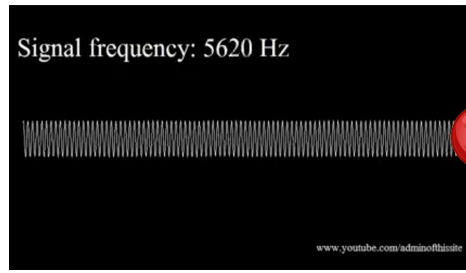
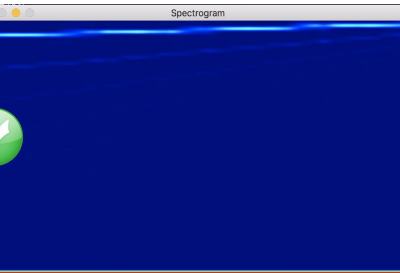
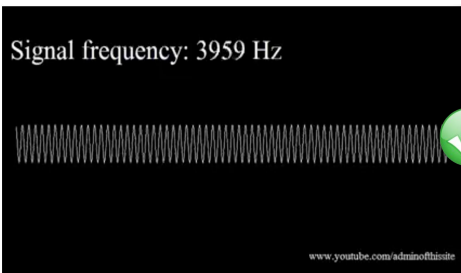
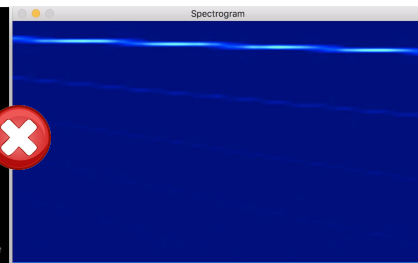
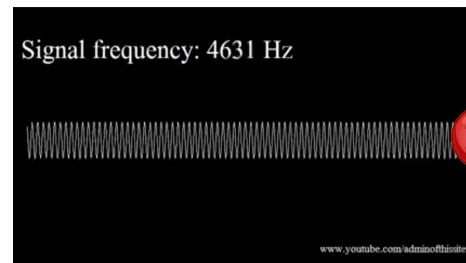
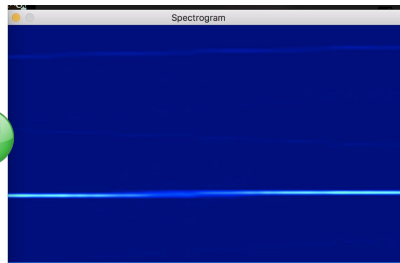
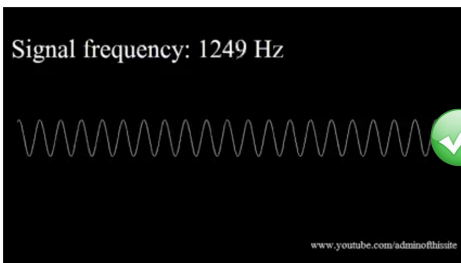
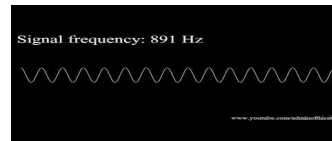
**scipy.fftpack** used for fft (imported from example02)

**pyaudio** used for data acquisition (recording)

**opencv** used for fast online, non-blocking visualization (matplotlib slower)

# Frequency Representation: Spectrogram Recording Example

Sampling freq / Aliasing ?  
Check the example03.py script on this sound:



- $F_s = 8000 \text{ Hz} \rightarrow F_{\text{nyquist}} = 4000 \text{ Hz}$
- Freqs after the Nyquist Freq are not captured
  - Run experiment on paura

# Open-source libraries for audio signal analysis

- librosa (Python)
  - <https://librosa.github.io>
  - Implements various audio features (mfccs, chroma, beat, etc)
  - Audio fx (e.g. pitch shift)
  - Some ML components
- pyAudioAnalysis (Python)
  - <https://github.com/tyiannak/pyAudioAnalysis>
  - Implements various audio features
  - Built-in training/testing of classifiers (using scikit-learn)
  - Clustering (speaker diarization)
  - Visualization
- Opensmile (C++)
  - <https://audeering.com/technology/opensmile/>
  - Richest set of audio features

# Spectrogram calculation using pyAudioAnalysis

```

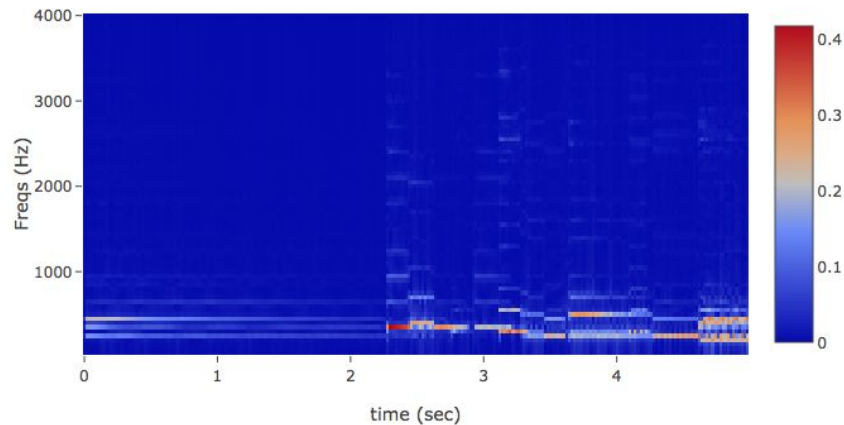
1  """
2  @brief Example 04
3  @details pyAudioAnalysis spectrogram calculation and visualization example
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import scipy.io.wavfile as wavfile
8  import plotly
9  import plotly.graph_objs as go
10 from pyAudioAnalysis import ShortTermFeatures
11 as aF
12 layout = go.Layout(title='Spectrogram Extraction Example using pyAudioAnalysis',
13                    xaxis=dict(title='time (sec)',),
14                    yaxis=dict(title='Freqs (Hz)',))
15
16 def normalize_signal(signal):
17     signal = np.double(signal)
18     signal = signal / (2.0 ** 15)
19     return (signal - signal.mean()) / ((np.abs(signal)).max() + 0.000000001)
20
21 if __name__ == '__main__':
22     [Fs, s] = wavfile.read("../data/sample_music.wav")
23     s = normalize_signal(s)
24     [S, t, f] = aF.spectrogram(s, Fs, int(Fs * 0.020), int(Fs * 0.020))
25     heatmap = go.Heatmap(z=S.T, y=f, x=t)
26     plotly.offline.plot(go.Figure(data=[heatmap], layout=layout),
27                          filename="temp.html", auto_open=True)

```

Only for 16-bit sample resolution

Also returns time and frequency  
scales (in Hz and secs  
respectively)

Spectrogram Extraction Example using pyAudioAnalysis



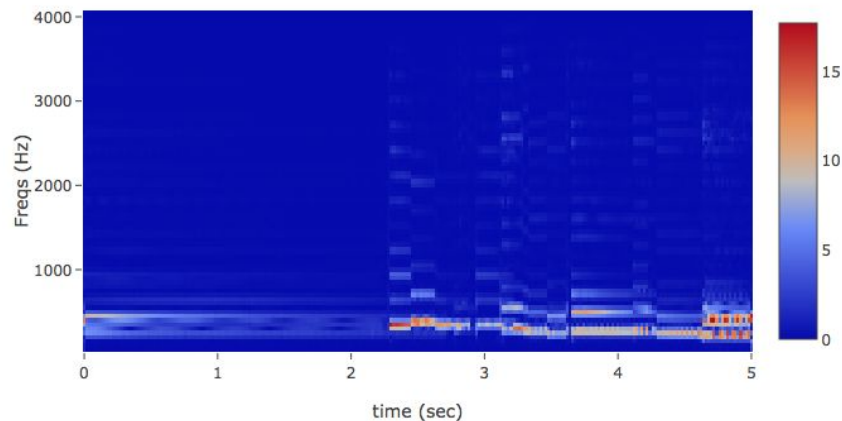
# Spectrogram calculation using librosa

```

1  """
2  @brief Example 05
3  @details Librosa spectrogram calculation and visualization example
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import scipy.io.wavfile as wavfile
8  import plotly
9  import librosa
10 import plotly.graph_objs as go
11 layout = go.Layout(title='Spectrogram Extraction Example using librosa',
12                    xaxis=dict(title='time (sec)',),
13                    yaxis=dict(title='Freqs (Hz)',))
14
15 def normalize_signal(signal):
16     signal = np.double(signal)
17     signal = signal / (2.0 ** 15)
18     return (signal - signal.mean()) / ((np.abs(signal)).max() + 0.000000001)
19
20 if __name__ == '__main__':
21     [Fs, s] = wavfile.read("../data/sample_music.wav")
22     s = normalize_signal(s)
23     S = np.abs(librosa.stft(s, int(Fs * 0.020), int(Fs * 0.020)))
24     # create frequency and time axes
25     f = [float((f + 1) * Fs) / (int(Fs * 0.020)) for f in range(S.shape[0])]
26     t = [float(t * int(Fs * 0.020)) / Fs for t in range(S.shape[1])]
27     heatmap = go.Heatmap(z=S, y=f, x=t)
28     plotly.offline.plot(go.Figure(data=[heatmap], layout=layout),
29                          filename="temp.html", auto_open=True)

```

Spectrogram Extraction Example using librosa



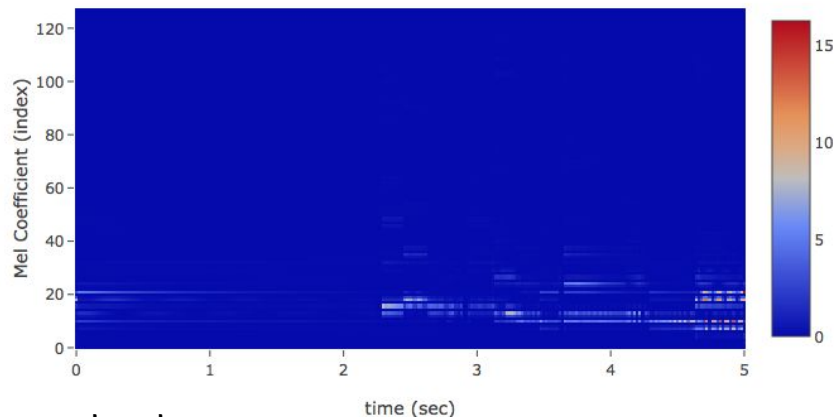
# Melgram calculation using librosa

```

1  """
2  @brief Example 06
3  @details Librosa spectrogram calculation and visualization example
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import scipy.io.wavfile as wavfile
8  import plotly
9  import librosa
10 import plotly.graph_objs as go
11 layout = go.Layout(title='Melgram Extraction Example using librosa',
12                    xaxis=dict(title='time (sec)',),
13                    yaxis=dict(title='Mel Coefficient (index)',))
14
15 def normalize_signal(signal):
16     signal = np.double(signal)
17     signal = signal / (2.0 ** 15)
18     return (signal - signal.mean()) / ((np.abs(signal)).max() + 0.000000001)
19
20 if __name__ == '__main__':
21     [Fs, s] = wavfile.read("../data/sample_music.wav")
22     s = normalize_signal(s)
23     S = librosa.feature.melspectrogram(s, Fs, None, int(Fs * 0.020),
24                                       int(Fs * 0.020), power=2)
25     # create frequency and time axes
26     f = range(S.shape[0])
27     t = [float(t * int(Fs * 0.020)) / Fs for t in range(S.shape[1])]
28     heatmap = go.Heatmap(z=S, y=f, x=t)
29     plotly.offline.plot(go.Figure(data=[heatmap], layout=layout),
30                        filename="temp.html", auto_open=True)

```

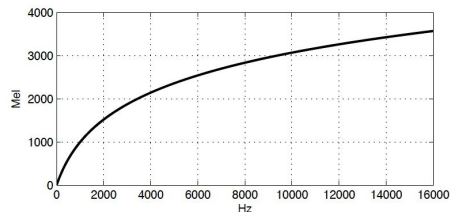
Melgram Extraction Example using librosa



## mel- scale

- Conform with psychoacoustic observations
- The human auditory system can distinguish neighboring frequencies more easily in the low frequency regions

$$f_w = 1127.01048 * \log(f/700 + 1)$$





# Why Mel?

```

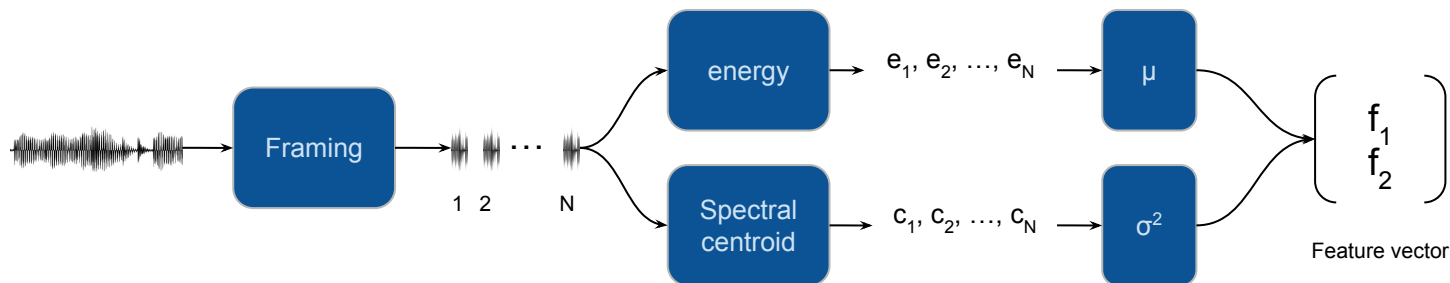
1  """
2  @brief Example 07
3  @details Frequency perceived discrimination experiment
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  from __future__ import print_function
7  import os, time, scipy.io.wavfile as wavfile, numpy as np
8  from random import randint
9
10 def play_sound(freq, duration, fs):
11     t = np.arange(0, duration, 1.0/fs); x = 0.5*np.cos(2 * np.pi * t * freq)
12     wavfile.write("temp.wav", fs, x); os.system("play temp.wav -q")
13
14 if __name__ == '__main__':
15     freqs, thres, n_exp, fs = [250, 500, 1000, 2000, 3000], [2, 5, 10, 20], 10, 16000
16     answers = [[] for i in range(len(freqs))]
17     for i_f, f in enumerate(freqs):
18         for t in thres:
19             answers[i_f].append(0)
20             for i in range(n_exp):
21                 sequel = randint(1, 2)
22                 if sequel == 2:
23                     play_sound(f, 0.5, fs); time.sleep(0.5); play_sound(f+t, 0.5, fs)
24                 else:
25                     play_sound(f+t, 0.5, fs); time.sleep(0.5); play_sound(f, 0.5, fs)
26                 ans = int(raw_input('Which was higher (1/2):'))
27                 if ans == sequel: answers[i_f][-1] += 1
28     print("Freq\t", end='')
29     for t in thres: print("{0:.1f}\t".format(t), end='')
30     print("")
31     for i_f, f in enumerate(freqs):
32         print("{} Hz\t".format(f), end='')
33         for i_t, t in enumerate(thres):
34             print("{0:.1f}\t".format(answers[i_f][i_t] / float(n_exp)), end='')
35     print("")

```

	Thresholds			
Freq	2 Hz	5 Hz	10 Hz	20 Hz
250 Hz	0.7	1	1	1
500 Hz	0.4	0.8	0.9	1
1000 Hz	0.6	0.8	1	0.9
2000 Hz	0.5	0.4	0.9	1
3000 Hz	0.5	0.5	0.6	1

# Audio segment feature extraction

- Short-term windowing:
  - “frames”
  - extract features per frame (such as energy, or spectral centroid)
  - result: sequence of vectors (one vector for each frame)
- Segment windowing:
  - segments are either predefined or applied on long recordings (e.g. fix-sized)
  - each segment corresponds to a sequence of short-term feature vectors
  - common practice
    - extract segment (mid-term) statistics ( $\mu, \sigma^2$ )
    - applied per sequence of short-term feature sequence (in the segment)



# Audio features: Segment Statistics

- Each feature is extracted in a short-term basis
- Segment feature statistics capture temporal changes in short-term feature sequences
- Statistics:
  - mean value
  - std/var
  - percentiles
  - max / min
  - Skewness
- Examples:
  - average zero crossing rate
  - deviation of the spectral centroid

# Time-domain features

- Energy

- usually normalized by window length
- high variation over successive speech frames (std statistic)

$$E(i) = \frac{1}{W_L} \sum_{n=1}^{W_L} |x_i(n)|^2$$

- Zero Crossing Rate

- rate of sign changes during the frame
- measure of noisiness
- high values for noisy signals

$$\text{sgn}[x_i(n)] = \begin{cases} 1, & x_i(n) \geq 0 \\ -1, & x_i(n) < 0 \end{cases} \quad Z(i) = \frac{1}{2W_L} \sum_{n=1}^{W_L} | \text{sgn}[x_i(n)] - \text{sgn}[x_i(n-1)] |$$

- Energy Entropy

- measure of abrupt changes in the signal's energy
- divide frames to K sub-frames and compute (normalized) sub-energies ( $e_{\text{subframe}_k}$ )
- compute entropy of  $e_{\text{subframe}_k}$  sequence

$$e_j = \frac{E_{\text{subFrame}_j}}{E_{\text{shortFrame}_i}} \quad E_{\text{shortFrame}_i} = \sum_{k=1}^K E_{\text{subFrame}_j}$$

$$H(i) = - \sum_{j=1}^K e_j^2 \cdot \log_2(e_j^2)$$

# Time-domain features - Example

```

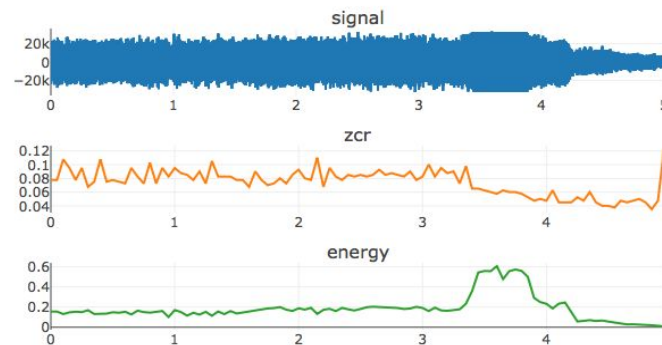
1  """
2  @brief Example 08
3  @details pyAudioAnalysis feature extraction example 1
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import plotly
8  import plotly.graph_objs as go
9  from pyAudioAnalysis import ShortTermFeatures as aF
10 from pyAudioAnalysis import audioBasicIO as aIO
11
12
13 if __name__ == '__main__':
14     # read machine sound
15     fs, s = aIO.read_audio_file("../data/general/objects/1-46744-A.ogg.wav")
16     duration = len(s) / float(fs)
17     # extract short term features and plot ZCR and Energy
18     [f, fn] = aF.feature_extraction(s, fs, int(fs * 0.050), int(fs * 0.050))
19     figs = plotly.tools.make_subplots(rows=3, cols=1,
20                                     subplot_titles=["signal", fn[0], fn[1]])
21     time = np.arange(0, duration - 0.050, 0.050)
22     time_s = np.arange(0, duration, 1/float(fs))
23     figs.append_trace(go.Scatter(x=time_s, y=s, showlegend=False), 1, 1)
24     figs.append_trace(go.Scatter(x=time, y=f[0, :], showlegend=False), 2, 1)
25     figs.append_trace(go.Scatter(x=time, y=f[1, :], showlegend=False), 3, 1)
26     plotly.offline.plot(figs, filename="temp.html", auto_open=True)

```

Also returns list of  
short-term feature  
names

This is a  $\#n\_frames \times \#n\_wins$  matrix

- using pyAudioAnalysis
- short-term feature sequences for ZCR / Energy
- sound: **vacuum cleaner**
- zcr is higher for “noisy” sounds



# Audio features: Frequency Domain (Spectral) Features

- Let  $X$  be the abs(FFT)
- Spectral Centroid
  - Center of gravity of the spectrum
- Spectral spread
  - 2nd central moment of the spectrum
- Spectral entropy
  - Divide spectrum into  $L$  sub-bands
  - Compute normalized sub-band energies ( $E_f$ )
  - Compute entropy

$$C_i = \frac{\sum_{k=1}^{Wf_L} (k+1) X_i(k)}{\sum_{k=1}^{Wf_L} X_i(k)}$$

$$S_i = \sqrt{\frac{\sum_{k=1}^{Wf_L} ((k+1) - C_i)^2 X_i(k)}{\sum_{k=1}^{Wf_L} X_i(k)}}$$

$$n_f = \frac{E_f}{\sum_{f=0}^{L-1} E_f}, f = 0, \dots, L-1. \quad H = - \sum_{f=0}^{L-1} n_f \cdot \log_2(n_f)$$

- Spectral Flux
  - Spectral change between two successive frames

$$EN_i(k) = \frac{X_i(k)}{\sum_{l=1}^{Wf_L} X_i(l)} \quad Fl_{(i,i-1)} = \sum_{k=1}^{Wf_L} (EN_i(k) - EN_{i-1}(k))^2$$

- Spectral Rolloff
  - Freq below which a percentage of the mag distribution of the spectrum is concentrated
  - If the  $m$ -th DFT coefficient is the spectral rolloff  $\rightarrow$

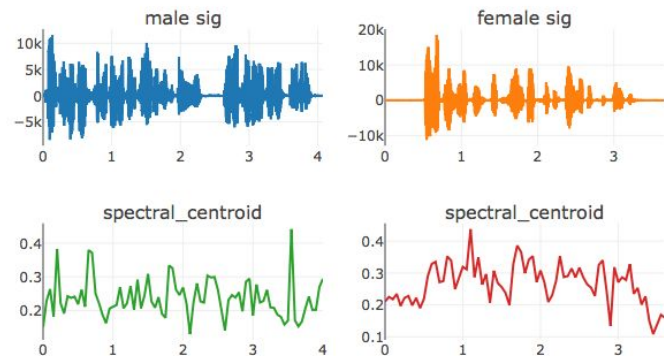
$$\sum_{k=1}^m X_i(k) = C \sum_{k=1}^{Wf_L} X_i(k)$$

# Audio features: Frequency Domain (Spectral) Features - Example

```

1  """
2  @brief Example 09
3  @details pyAudioAnalysis feature extraction example for male / female speeches
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import plotly
8  import plotly.graph_objs as go
9  from pyAudioAnalysis import ShortTermFeatures as aF
10 from pyAudioAnalysis import audioBasicIO as aIO
11
12
13 if __name__ == '__main__':
14     win = 0.05
15     fp1 = "../data/general/speech/m1_neu-m1-l1.wav.wav" # male
16     fp2 = "../data/general/speech/f1_neu-f1-l2.wav.wav" # female
17     # read machine sound
18     fs1, s1 = aIO.read_audio_file(fp1)
19     fs2, s2 = aIO.read_audio_file(fp2)
20     dur1, dur2 = len(s1) / float(fs1), len(s2) / float(fs2)
21     # extract short term features
22     [f1, fn] = aF.feature_extraction(s1, fs1, int(fs1 * win), int(fs1 * win))
23     [f2, fn] = aF.feature_extraction(s2, fs2, int(fs2 * win), int(fs2 * win))
24     figs = plotly.tools.make_subplots(rows=2, cols=2,
25                                       subplot_titles=["male sig", "female sig",
26                                                         fn[3], fn[3]])
27     t1 = np.arange(0, dur1 - win, win)
28     ts_1 = np.arange(0, dur1, 1/float(fs1))
29     t2 = np.arange(0, dur2 - win, win)
30     ts_2 = np.arange(0, dur2, 1/float(fs2))
31     figs.append_trace(go.Scatter(x=ts_1, y=s1, showlegend=False), 1, 1)
32     figs.append_trace(go.Scatter(x=ts_2, y=s2, showlegend=False), 1, 2)
33     figs.append_trace(go.Scatter(x=t1, y=f1[3, :], showlegend=False), 2, 1)
34     figs.append_trace(go.Scatter(x=t2, y=f2[3, :], showlegend=False), 2, 2)
35     plotly.offline.plot(figs, filename="temp.html", auto_open=True)

```



# Audio features: Cepstral Domain

- Mel-Frequency Cepstral Coefficients
  - Compute DFT
  - Mel-scale filter bank application
  - Compute  $O_k$  as the power of the output of each filter
  - Compute MFCCs as the discrete cosine transform coefficients of the mel-scaled log-power spectrum
- Usually select the first 13 MFCCs (considered to carry enough discriminative information especially for speech classification tasks)
- Cepstrum in general (not mel):
  - Inverse fft of the log fft

$$f_w = 1127.01048 * \log(f/700 + 1)$$

$$c_m = \sum_{k=1}^L (\log \widetilde{O}_k) \cos[m(k - \frac{1}{2})\frac{\pi}{L}], \quad m = 1, \dots, L$$

$$: \left| \mathcal{F}^{-1} \left\{ \log \left( |\mathcal{F}\{f(t)\}|^2 \right) \right\} \right|^2$$



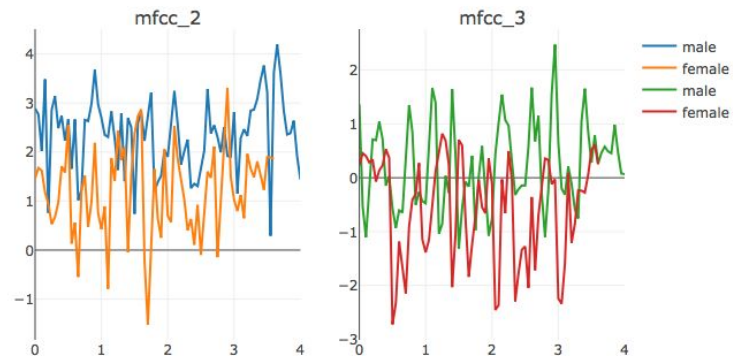


# Audio features: Cepstral Domain - Example

```

1  """
2  @brief Example 10
3  @details pyAudioAnalysis feature extraction example - mfccs for male/female
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import plotly
8  import plotly.graph_objs as go
9  from pyAudioAnalysis import audioFeShortTermFeaturesatureExtraction as aF
10 from pyAudioAnalysis import audioBasicIO as aIO
11
12
13 if __name__ == '__main__':
14     win = 0.05
15     fp1 = "../data/general/speech/m1_neu-m1-l1.wav.wav" # male
16     fp2 = "../data/general/speech/f1_neu-f1-l2.wav.wav" # female
17     # read machine sound
18     fs1, s1 = aIO.read_audio_file(fp1)
19     fs2, s2 = aIO.read_audio_file(fp2)
20     dur1, dur2 = len(s1) / float(fs1), len(s2) / float(fs2)
21     # extract short term features
22     [f1, fn] = aF.feature_extraction(s1, fs1, int(fs1 * win), int(fs1 * win))
23     [f2, fn] = aF.feature_extraction(s2, fs2, int(fs2 * win), int(fs2 * win))
24     figs = plotly.tools.make_subplots(rows=1, cols=2,
25                                     subplot_titles=[fn[9], fn[10]])
26     t1 = np.arange(0, dur1 - 0.050, 0.050)
27     t2 = np.arange(0, dur2 - 0.050, 0.050)
28     figs.append_trace(go.Scatter(x=t1, y=f1[9, :], name="male"), 1, 1)
29     figs.append_trace(go.Scatter(x=t2, y=f2[9, :], name="female"), 1, 1)
30     figs.append_trace(go.Scatter(x=t1, y=f1[10, :], name="male"), 1, 2)
31     figs.append_trace(go.Scatter(x=t2, y=f2[10, :], name="female"), 1, 2)
32
33     plotly.offline.plot(figs, filename="temp.html", auto_open=True)

```



# Audio features: Chroma Vector

- 12-element frequency representation
- In music applications
- Group the DFT coefficients of a window into 12 bins
- Each bin represents the 12 equal-tempered classes of western-type music
- Bins in semitone spacing
- $S_k$  is the set of frequencies for the  $k$ -th bin (representing DFT coefficients)

$$v_k = \sum_{n \in S_k} \frac{X_i(n)}{N_k}, k \in 0..11$$

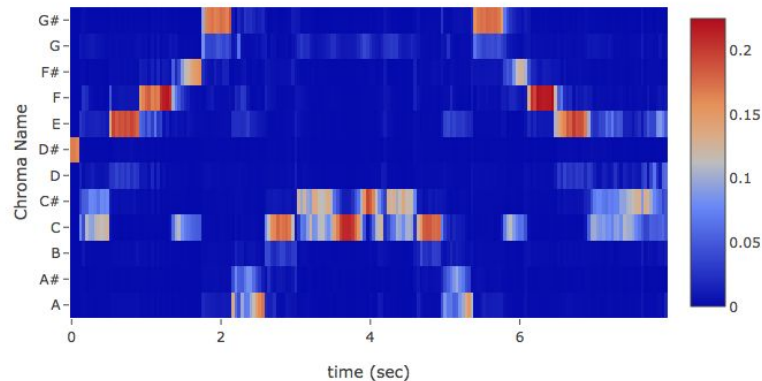
# Audio features: Chroma Vector - Example

```

1  """
2  @brief Example 11
3  @details pyAudioAnalysis chromagram example
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import plotly
7  import plotly.graph_objs as go
8  from pyAudioAnalysis import ShortTermFeatures as aF
9  from pyAudioAnalysis import audioBasicIO as aIO
10 layout = go.Layout(title='Chromagram example for doremi.wav signal',
11                    xaxis=dict(title='time (sec)',),
12                    yaxis=dict(title='Chroma Name',))
13
14
15 if __name__ == '__main__':
16     win = 0.04
17     fp = "../data/doremi.wav" # music sample
18     # read machine sound
19     fs, s = aIO.read_audio_file(fp)
20     fs = float(fs)
21     dur1 = len(s) / float(fs)
22     spec, time, freq = aF.chromagram(s, fs, int(fs * win),
23                                     int(fs * win), False)
24     heatmap = go.Heatmap(z=spec.T, y=freq, x=time)
25     plotly.offline.plot(go.Figure(data=[heatmap], layout=layout),
26                         filename="temp.html", auto_open=True)

```

Chromagram example for doremi.wav signal



# Plotly histogram representation (function in utilities.py)

```

6  import plotly
7  import plotly.graph_objs as go
8  import numpy as np
9  import matplotlib.pyplot as plt
10
11 def plot_feature_histograms(list_of_feature_mtr, feature_names,
12                             class_names, n_columns=5):
13     '''
14     Plots the histograms of all classes and features for a given
15     classification task.
16     :param list_of_feature_mtr: List of feature matrices
17                               (n_samples x n_features) for each class
18     :param feature_names:     List of feature names
19     :param class_names:      List of class names, for each feature matr
20     '''
21     clr_map = plt.cm.get_cmap('jet')
22     n_features = len(feature_names)
23     n_bins = 12
24     n_classes = len(class_names)
25     n_rows = int(n_features / n_columns) + 1
26     figs = plotly.tools.make_subplots(rows=n_rows, cols=n_columns,
27                                       subplot_titles=feature_names)
28     figs['layout'].update(height=(n_rows * 250))
29     range_cl = range(int(int(255/n_classes)/2), 255, int(255/n_classes))
30     clr = []
31     for i in range(n_classes):
32         clr.append('rgba({}, {}, {}, {})'.format(clr_map(range_cl[i])[0],
33                                                  clr_map(range_cl[i])[1],
34                                                  clr_map(range_cl[i])[2],
35                                                  clr_map(range_cl[i])[3]))

```

```

36     for i in range(n_features):
37         # for each feature get its bin range (min:(max-min)/n_bins:max)
38         f = np.vstack([x[:, i:i + 1] for x in list_of_feature_mtr])
39         bins = np.arange(f.min(), f.max(), (f.max() - f.min()) / n_bins)
40         for fi, f in enumerate(list_of_feature_mtr):
41             # load the color for the current class (fi)
42             mark_prop = dict(color=clr[fi], line=dict(color=clr[fi], width=3))
43             # compute the histogram of the current feature (i) and normalize:
44             h, _ = np.histogram(f[:, i], bins=bins)
45             h = h.astype(float) / h.sum()
46             cbins = (bins[0:-1] + bins[1:]) / 2
47             scatter_1 = go.Scatter(x=cbins, y=h, name=class_names[fi],
48                                   marker=mark_prop, showlegend=(i == 0))
49             # (show the legend only on the first line)
50             figs.append_trace(scatter_1, int(i/n_columns)+1, i % n_columns+1)
51     for i in figs['layout']['annotations']:
52         i['font'] = dict(size=10, color='#224488')
53     plotly.offline.plot(figs, filename="report.html", auto_open=True)

```

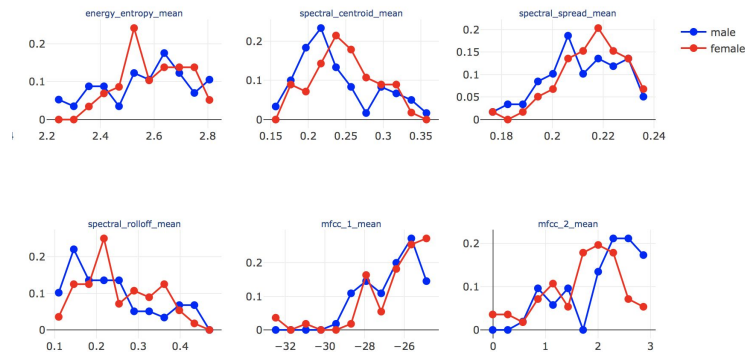
<https://plot.ly/python/>

# Feature discrimination example: male vs female segments

```

1  """
2  @brief Example 12
3  @details pyAudioAnalysis feature extraction for classes organized in folders
4  and feature histogram representation (per feature and class).
5  Binary classification task: male vs female speech segments
6  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
7  """
8  from pyAudioAnalysis import MidTermFeatures as aF
9  import os.path
10 import utilities as ut
11
12 if __name__ == '__main__':
13     dirs = ["../data/gender/male",
14            "../data/gender/female"]
15     class_names = [os.path.basename(d) for d in dirs]
16     m_win, m_step, s_win, s_step = 1, 1, 0.1, 0.05
17     features = []
18     for d in dirs:
19         # get feature matrix for each directory (class)
20         f, files, fn = aF.directory_feature_extraction(d, m_win, m_step, s_win,
21                                                       s_step)
22         features.append(f)
23     ut.plot_feature_histograms(features, fn, class_names)

```

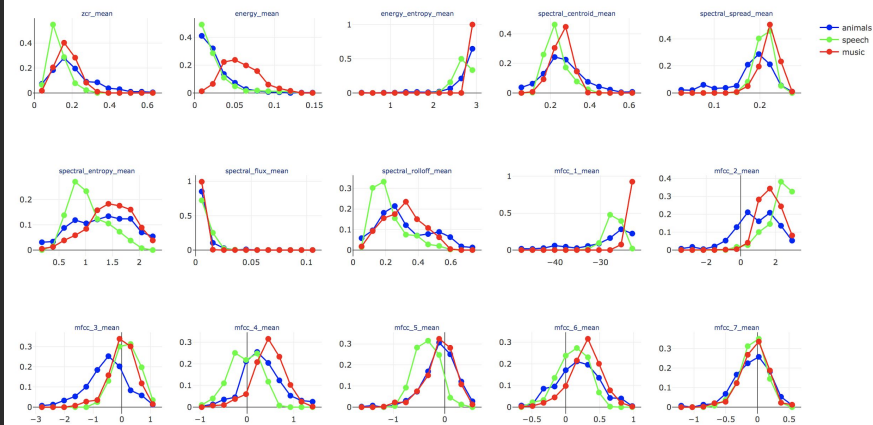


# Feature discrimination example: 3-class task

```

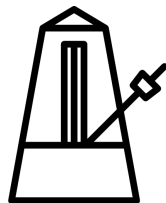
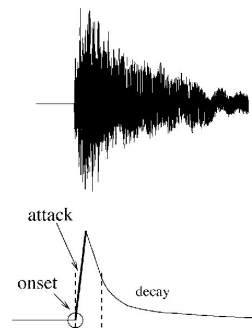
1  """
2  @brief Example 13
3  @details pyAudioAnalysis feature extraction for classes organized in folders
4  and feature histogram representation (per feature and class).
5  3-class classification task: animals vs speech vs music segments
6  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
7  """
8  from pyAudioAnalysis import MidTermFeatures as aF
9  import os.path
10 import utilities as ut
11
12 if __name__ == '__main__':
13     dirs = ["../data/general/animals",
14             "../data/general/speech",
15             "../data/general/music"]
16     class_names = [os.path.basename(d) for d in dirs]
17     m_win, m_step, s_win, s_step = 1, 1, 0.1, 0.05
18     features = []
19     for d in dirs:
20         # get feature matrix for each directory (class)
21         f, files, fn = aF.directory_feature_extraction(d, m_win, m_step, s_win,
22                                                       s_step)
23         features.append(f)
24     ut.plot_feature_histograms(features, fn, class_names)

```

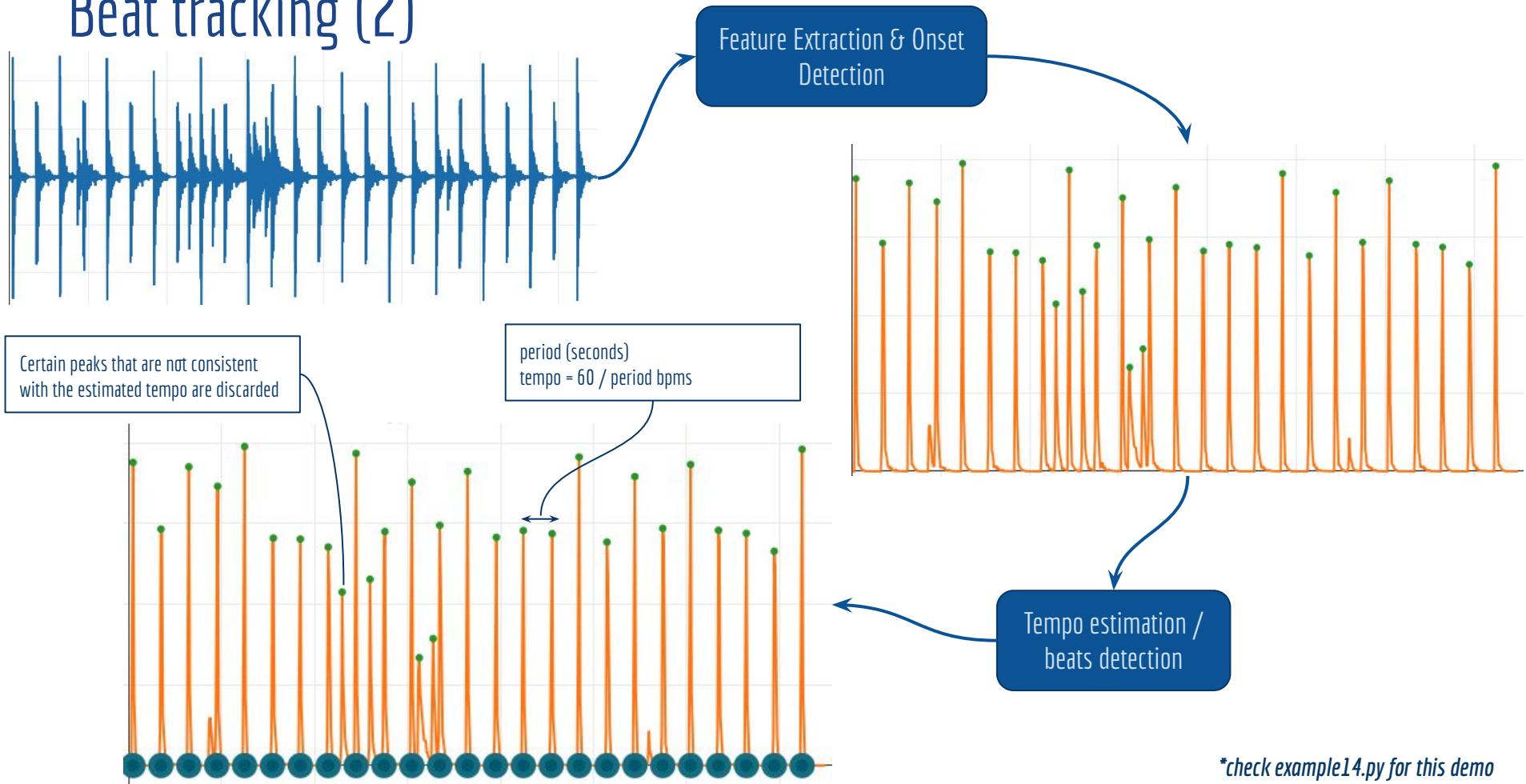


# Beat tracking (1)

- Tempo / beat: fundamental properties of music
- Beat: a steady “pulse” that provides the temporal framework of a song
- Beat tracking: “tapping the foot when music plays”
- **onset** detection:
  - onset: time position of a significant signal change (e.g note)
  - change in signal's energy or frequency distribution
  - onset → attack → decay
- **tempo** estimation & beat **peaks** selection:
  - detect peaks that are (almost) equally spaced in time
  - detected peaks are (almost) consistent with est. tempo



## Beat tracking (2)



*\*check example14.py for this demo*



# Beat tracking (3)

- Perception of beat is
  - hierarchical
  - ambiguous
- Estimated tempo can be multiple of the “real” tempo
- E.g. rap song tracked at 170 bpms: true tempo is  $170/2 = 85$  bpms!
- External expert knowledge may be needed to put constraints in the last step of the tempo extraction method (see prev slide)

J. P. Bello, L. Daudet, S. Abdallah, C. Duxbury, M. Davies, M. B. Sandler, “A Tutorial on Onset Detection in Music Signals,” *IEEE Tr. Speech and Audio Proc.*, vol. 13, no. 5, pp. 1035-1047, September 2005

P. Desain & H. Honing, “Computational models of beat induction: The rule-based approach,” *J. New Music Research*, vol. 28 no. 1, pp. 29-42, 1999.

Eric. D. Scheirer, “Tempo and beat analysis of acoustic musical signals,” *J. Acoust. Soc. Am.*, vol. 103, pp. 588-601, 1998.

Davies, M. E., & Plumbley, M. D. (2007). Context-dependent beat tracking of musical audio. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(3), 1009-1020.

McKinney, M. F., & Moelants, D. (2006). Ambiguity in tempo perception: What draws listeners to different metrical levels?. *Music Perception: An Interdisciplinary Journal*, 24(2), 155-166.

# Beat tracking - Example

```

1  """
2  @brief Example 15
3  @details Librosa beattracking example
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import numpy as np
7  import scipy.io.wavfile as wavfile
8  import sys
9  import librosa
10
11  if __name__ == '__main__':
12      # needs filepath as main argument:
13      if len(sys.argv) != 2:
14          sys.exit()
15      # load file and extract tempo and beats:
16      [Fs, s] = wavfile.read(sys.argv[1])
17      tempo, beats = librosa.beat.beat_track(y=s, sr=Fs, units="time")
18      beats -= 0.05
19      # add small 220Hz sounds on the 2nd channel of the song on each beat
20      s = s.reshape(-1, 1)
21      s = np.array(np.concatenate((s, np.zeros(s.shape)), axis=1))
22      for ib, b in enumerate(beats):
23          t = np.arange(0, 0.2, 1.0 / Fs)
24          amp_mod = 0.2 / (np.sqrt(t)+0.2) - 0.2
25          amp_mod[amp_mod < 0] = 0
26          x = s.max() * np.cos(2 * np.pi * t * 220) * amp_mod
27          s[int(Fs * b):
28             int(Fs * b) + int(x.shape[0]), 1] = x.astype('int16')
29      wavfile.write("output.wav", Fs, np.int16(s))

```

## Usage example:

```

audio|master ⚡ ⇒ python3 example15.py
./data/musical_genres_small/hiphop/nwa_straght_outta_campton.wav

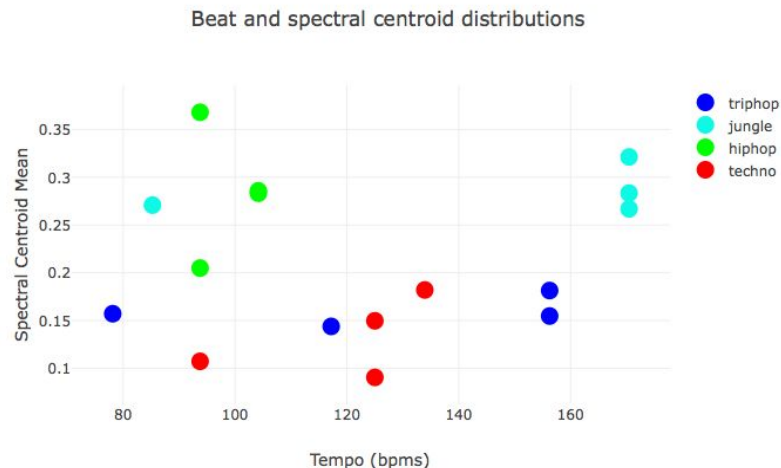
```

# Beat tracking - Discrimination Example

```

1  """
2  @brief Example 16
3  @details librosa beattracking example: extract tempo and spectral centroid for
4  songs from different musical genres
5  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
6  """
7  import scipy.io.wavfile as wavfile, utilities as ut
8  import glob, os, librosa, plotly, numpy as np, plotly.graph_objs as go
9  from pyAudioAnalysis.MidTermFeatures import mid_feature_extraction as mt
10
11 layout = go.Layout(title='Beat and spectral centroid distributions',
12                    xaxis=dict(title='Tempo (bpms)'),
13                    yaxis=dict(title='Spectral Centroid Mean'),)
14
15 def get_dir_features(dir_name):
16     feats = []
17     for f in glob.glob(os.path.join(dir_name, "*.wav")):
18         [Fs, s] = wavfile.read(f)
19         tempo, _ = librosa.beat.beat_track(y=s, sr=Fs)
20         f, _, fn = mt(s, Fs, int(1.0*Fs), int(1.0*Fs), int(0.1*Fs), int(0.1*Fs))
21         feats.append([tempo, np.mean(f[fn.index("spectral_centroid_mean")],
22                                     axis=0)])
23     return np.array(feats)
24
25 if __name__ == '__main__':
26     g_paths = glob.glob("../data/musical_genres_small/*/*")
27     g_names = [p.split('/')[-2] for p in g_paths]
28     clr = ut.get_color_combinations(len(g_paths))
29     features = [get_dir_features(g) for g in g_paths]
30     plots = [go.Scatter(x=features[i][:, 0], y=features[i][:, 1],
31                        mode='markers', name=g_names[i], marker=dict(
32                            color=clr[i], size=15))
33              for i in range(len(g_paths))]
34     plotly.offline.plot(go.Figure(data = plots, layout=layout),
35                        filename="temp.html", auto_open=True)

```



"True" tempo values:

- Triphop: 90-110
- Jungle: 160-170
- Hiphop: 80-100
- Techno: 120-130

- Tempo is not always discriminative
- Tempo estimation has errors (e.g. triphop estimated values are double)

*Need for multiple features to achieve accurate discrimination*

# Pitch tracking

- **f0:**
  - fundamental frequency
  - a **physical** property of sound:
    - speech: glottal pulses freq
    - music: most dominant freq of a note (eg freq of vibration of a string)
- **pitch**
  - a **subjective** phenomenon (f0 open to measurement)
  - perceptual
  - follows f0
  - speech:
    - not always clear
    - vad required
  - music:
    - note transcription
    - polyphony

## Pitch tracking:

- Time / spectral domain
- Spectral:
  - simple argmax?
  - no! f0 not always the freq with the max freq in spectrogram

# Pitch Tracking - Example using librosa

```

1  """
2  @brief Example 17
3  @details Librosa pitch tracking example
4  @author Theodoros Giannakopoulos {tyiannak@gmail.com}
5  """
6  import scipy.io.wavfile as wavfile
7  import librosa
8  import plotly
9  import numpy as np
10 import plotly.graph_objs as go
11 from scipy.signal import medfilt as mf
12
13 layout = go.Layout(title='Librosa pitch estimation',
14                    xaxis=dict(title='time frame',),
15                    yaxis=dict(title='freq (Hz)',))
16
17 def get_librosa_pitch(signal, fs, window):
18     pitches, magnitudes = librosa.piptrack(y=signal, sr=fs, n_fft=int(window),
19                                           hop_length=int(window/10))
20     pitch_pos = np.argmax(magnitudes, axis=0)
21     pitches_final = []
22     for i in range(len(pitch_pos)):
23         pitches_final.append(pitches[pitch_pos[i], i])
24     pitches_final = np.array(pitches_final)
25     pitches_final[pitches_final > 2000] = 0 # cut high pitches
26     return mf(pitches_final, 3) # use medfilt for smoothing
27
28 if __name__ == '__main__':
29     [fs, s] = wavfile.read("../data/acapella.wav")
30     p = get_librosa_pitch(s, fs, fs/20)
31     plt = go.Scatter(x=np.arange(len(p)), y=p, mode='markers', showlegend=False)
32     plotly.offline.plot(go.Figure(data=[plt], layout=layout),
33                         filename="temp.html", auto_open=True)

```

