

Discriminating Sound Textures

Jayson Lynch*

Eric Mannes[†]

December 9, 2015

Abstract

We show sound textures are cool stuff.

Keywords: Machine Learning, Sound Textures, Audition

*MIT Computer Science and Artificial Intelligence Laboratory, 32 Vassar Street, Cambridge, MA 02139, USA, {jaysonl}@mit.edu.

[†]MIT Computer Science and Artificial Intelligence Laboratory, 32 Vassar Street, Cambridge, MA 02139, USA, {mannes}@mit.edu.

1 Introduction

Sound textures, such as rain or crackling fire, are the result of many similar acoustic events.

Recent work by McDermott and Simoncelli attempts to understand a special category of sounds called *sound textures* [MS11,MSS13]. They show that for many sound textures, the salient features that cause humans to classify these sounds are controlled by a number of time-averaged statistics of the sounds. This leads to an interesting question of why some sounds seem to be identified by time-averaged properties, when many other, such as speech or music, depend strongly on their temporal structure. McDermott et. al. conjecture that the ‘sparsity’ or ‘compressibility’ are related to the sounds that are perceived in this way. In addition, they mention that sound textures are ‘temporally homogeneous’ without giving formal definitions of these terms. This paper seeks to build upon these ideas both to investigate new features for audio classification as well as to lend evidence toward the previous conjectures about sound texture perception. We provide several formal definitions of audio compressibility, sparsity, and temporal homogeneity in Section 2. These are then extracted as a feature set and are used in a machine learning algorithm for sound textures, described in Section 3.

There has been other prior work on classifying sound textures. Saint-Arnaud’s Master’s Thesis attempts to extract the sound atoms that make up sound textures and use this as a basis for a classifier [SA95]. The thesis also discusses human perception of sound textures and the difficulties surrounding a good definition. A significant amount of the machine learning work in sound textures has been around synthesis often using wavelet hierarchies [?] or sound atoms [?]. Schwarz provides a general overview of methods as well as a classification scheme of the synthesis methods used [?].

[Cite McDermott’s papers [MS11,MSS13] and why we care about these features/- xxx
classification]

[Cite some other audio/sound texture classification work.] xxx

2 Definitions

We provide working definitions for examples of three main descriptors: sparsity, compressibility, and temporal homogeneity. Since these characteristics can be captured in different ways, which are not always equivalent, we choose to work with multiple formal definitions.

Compressibility was the simplest to work out. We chose a standard lossy and lossless compression algorithm and measured the compression ratio for the .wav files. In this case we used **[what xxx
audio compression did we use?]**

How to define sparsity was slightly less clear, since it usually refers to the quantity of zero entries with respect to some basis. To overcome this, we decided to use several natural representations. We looked at sparsity in the time domain with respect to the actual time series of the audio sample. We also looked at the sparsity in the frequency domain under both uniform and log-scale transforms. To be more precise, we took short-term Fourier transforms, constant q transforms, and Mel transforms of the time series (using the implementations in the Python library Librosa [MRL⁺15]) and counted the ratio of entries near zero. Due to noise and precision error, we assumed an entry was empty if the magnitude of the value was less than 10^{-5} . There are also algorithms for extracting low rank matrices assuming one has sampled from noisy data [?]. We would have liked to have implemented and run one of these algorithms, using the derived rank as a measure of sparsity, but we did not have time to extract this feature.

Finally temporal homogeneity refers to consistency in the signal over time. Obviously we

can't have everything be exactly the same, so one must pick specific features with which to check consistency. A paper on visual textures provides a precise definition: X is homogeneous with respect to the function $\Phi : \mathbb{R}^L \rightarrow \mathbb{R}$

[Give formal and informal definitions of the features we used] xxx

3 Methods

[Describe how we set up our feature extraction. How we set up our learning. What our dataset is. Any other important process things.] xxx

4 Results

[What correlated and what didn't? What feature (ensemble) lead to a good classifier?] xxx

5 Conclusion

[Speculate if these could be useful features in audio classification. Speculate about the implication to human audition. Give future directions.] xxx

Acknowledgments

We thank Prof Josh McDermott for his support in answering our questions about their research and providing us with their dataset. We also thank the course staff of 6.867 for their instruction and support this semester. In particular, we appreciate Marzyeh Ghassemi advice and guidance in shaping our project plan.

References

- [MRL⁺15] Brian McFee, Colin Raffel, Dawen Liang, Daniel PW Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto. librosa: Audio and music signal analysis in python. 2015.
- [MS11] Josh H. McDermott and Eero P. Simoncelli. Sound texture perception via statistics of the auditory periphery: Evidence from sound synthesis. *Neuron*, 71(5):926 – 940, 2011.
- [MSS13] Josh H McDermott, Michael Schemitsch, and Eero P Simoncelli. Summary statistics in auditory perception. *Nature neuroscience*, 16(4):493–498, 2013.
- [SA95] Nicolas Saint-Arnaud. *Classification of sound textures*. PhD thesis, Massachusetts Institute of Technology, 1995.