

LEARNING PITCH INVARIANTS FOR INSTRUMENT RECOGNITION

Names should be omitted for double-blind reviewing
Affiliations should be omitted for double-blind reviewing

ABSTRACT

Musical performance combines a wide range of pitches, nuances, and expressive techniques. Audio-based classification of musical instruments thus requires to build signal representations that are invariant to such transformations. Focusing on pitch invariance, this article investigates the construction of multi-stage architectures for instrument recognition. We show that Mel-frequency cepstral coefficients (MFCC) lack invariance with respect to realistic pitch shifts. In turn, a convolutional neural network (ConvNet) in the time-frequency domain is able to disentangle pitch variability from timbral information in a subtler way. We further improve the ConvNet architecture by limiting weight sharing to octave-wide frequency bands at the first layer, while allowing full weight sharing at deeper layers. We extend our method to the recognition of multiple instruments playing simultaneously.

1. INTRODUCTION

Among the cognitive attributes of musical tones, pitch is distinguished by a combination of three properties. First, it is relative: ordering pitches from low to high gives rise to intervals and melodic patterns. Secondly, it is intensive: multiple pitches heard simultaneously produce a chord, not a single unified tone – contrary to loudness, which adds up with the number of sources. Thirdly, it is invariant to instrumentation: this makes possible the transcription of polyphonic music under a single symbolic system.

Besides this invariance property, understanding the influence of pitch in audio streams is paramount to the design of an efficient system for automated classification, tagging, and similarity retrieval in music.

Section 2 demonstrates that pitch is the major factor of variability among musical notes of a given instrument, if described by their Mel-frequency cepstra. Section 3 describes a typical deep learning architecture for spectrogram-based classification, consisting of two convolutional layers and one densely connected layer. Section 4 improves the aforementioned architecture by splitting spectrograms into octave-wide frequency bands, training specific convolutional layers over each band in parallel, and

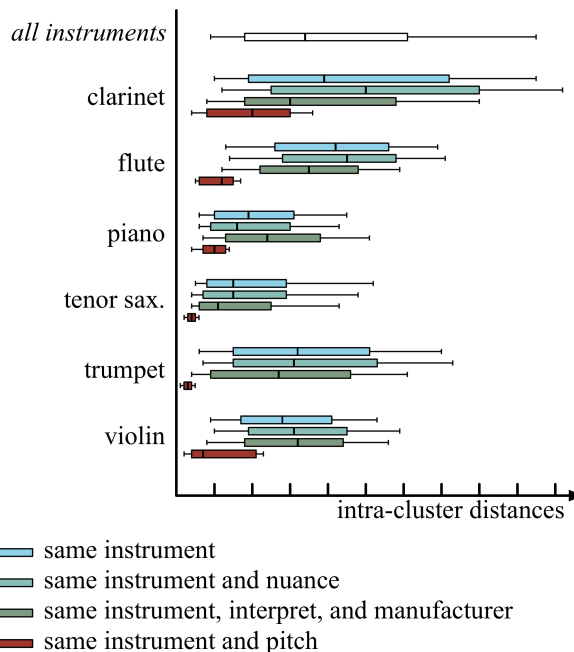


Figure 1: Distributions of squared Euclidean distances among various clusters in the RWC dataset. Whisker ends correspond to lower and upper deciles. See text for details.

gathering feature maps at a later stage. Section 5 discusses the effectiveness of the presented systems on a challenging dataset for music instrument recognition.

Time-frequency representations, such as the constant-Q wavelet scalogram, are a useful first step for the construction of pitch-adaptive features.

2. HOW INVARIANT IS THE MEL CEPSTRUM ?

The MFCCs were extracted from a filterbank of 40 Mel-frequency bands and 13 discrete cosine transform coefficients.

3. DEEP CONVOLUTIONAL NETWORKS

A deep learning system for classification is built by stacking multiple layers of weakly nonlinear transformations, whose parameters are jointly optimized such that the top-level layer fits a training set of labeled examples. This section introduces a typical deep learning architecture for audio classification and describes the functioning of each layer.



© Names should be omitted for double-blind reviewing.

Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** Names should be omitted for double-blind reviewing. “Learning pitch invariants for instrument recognition”, 16th International Society for Music Information Retrieval Conference, 2015.

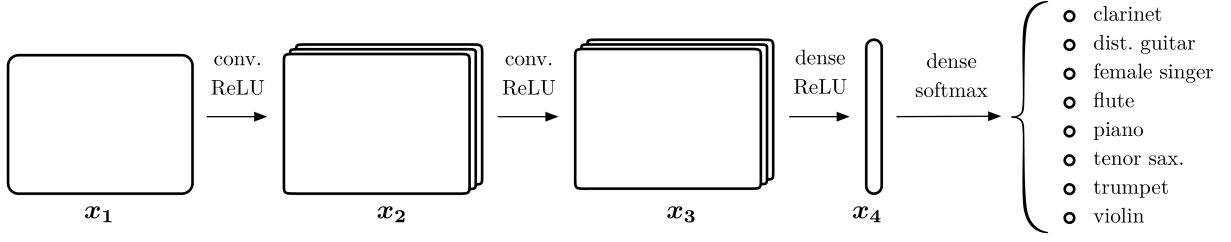


Figure 2: Architecture of a convolutional network with full weight sharing. See text for details.

3.1 Time-frequency representation

We used the implementation from the librosa package [5] with $Q = 12$ filters per octave, center frequencies ranging from 55 Hz to 14 kHz (8 octaves from A1 to A9), and a hop size of 23 ms. Furthermore, we applied nonlinear perceptual weighting of loudness in order to reduce the dynamic range between the fundamental partial and its upper harmonics. A 3-second sound excerpt $\mathbf{x}[t]$ is represented by a time-frequency matrix $\mathbf{x}_1[t, k_1]$ of width $T = 128$ samples and height $K_1 = 96$ MIDI indices, i.e. 8 octaves.

3.2 Architecture

Each layer in a convolutional network typically consists in the composition of three operations: two-dimensional convolutions, application of a pointwise nonlinearity, and subsampling.

First of all, we apply a family $\mathbf{W}_2[\tau, \kappa_1, k_2]$ of $K_2 = 32$ learned time-frequency convolutional operators, whose supports are constrained to have width Δt and height Δk_1 . Element-wise biases $\mathbf{b}_2[k_2]$ are added to the convolutions, resulting in the three-way tensor

$$\begin{aligned} \mathbf{y}_2[t, k_1, k_2] &= \mathbf{b}_2[k_2] + \mathbf{W}_2[t, k_1, k_2] \stackrel{t, k_1}{*} \mathbf{x}_1[t, k_1] \\ &= \mathbf{b}_2[k_2] + \sum_{\substack{0 \leq \tau < \Delta t \\ 0 \leq \kappa_1 < \Delta k_1}} \mathbf{W}_2[\tau, \kappa_1, k_2] \mathbf{x}_1[t - \tau, k_1 - \kappa_1]. \end{aligned} \quad (1)$$

The second step is the application of a pointwise nonlinearity. We have chosen the *rectified linear unit* [ReLU] because of its popularity in computer vision and its computational efficiency.

$$\mathbf{y}_2^+[t, k_1, k_2] = \max(\mathbf{y}_2[t, k_1, k_2], 0) \quad (2)$$

To achieve invariance to translation as well as frequency transposition, we pool neighboring units in the time-frequency domain (t, k_1) over non-overlapping rectangles of width Δt and height Δk_1 .

$$\mathbf{x}_2[t, k_1, k_2] = \max_{\substack{0 \leq \tau < \Delta t \\ 0 \leq \kappa_1 < \Delta k_1}} \left\{ \mathbf{y}_2^+[t - \tau, k_1 - \kappa_1, k_2] \right\} \quad (3)$$

We apply a family $\mathbf{W}_3[\tau, \kappa_1, k_2, k_3]$ of K_3 convolutional operators that perform a linear combination of time-frequency feature maps in \mathbf{x}_2 along the channel variable k_2 .

$$\begin{aligned} \mathbf{y}_3[t, k_1, k_3] &= \sum_{k_2} \mathbf{b}_3[k_2, k_3] + \mathbf{W}_3[t, k_1, k_2, k_3] \stackrel{t, k_1}{*} \mathbf{x}_2[t, k_1, k_2]. \end{aligned} \quad (4)$$

After nonlinear rectification and max-pooling, the layer \mathbf{y}_3 turns into a non-negative tensor $\mathbf{x}_3[t, k_1, k_3]$.

$$\mathbf{y}_4[k_4] = \sum_{t, k_1, k_3} \mathbf{W}_4[t, k_1, k_3, k_4] \mathbf{x}_3[t, k_1, k_3] \quad (5)$$

We apply a ReLU to \mathbf{y}_4 , yielding $\mathbf{x}_4[k_4] = \mathbf{y}_4^+[k_4]$. $\mathbf{y}_5[k_5] = \sum_{k_4} \mathbf{W}_5[k_4, k_5] \mathbf{x}_4[k_4]$.

$$\mathbf{x}_5[k_5] = \frac{\exp \mathbf{y}_5[k_5]}{\sum_{\kappa_5} \exp \mathbf{y}_5[\kappa_5]} \quad (6)$$

The above ensures that the coefficients of \mathbf{x}_5 are non-negative and sum to one, hence can be fit to a probability distribution.

$$\mathcal{L}(\mathbf{x}_5, \mathcal{I}) = - \sum_{k_5 \in \mathcal{I}} \log \mathbf{x}_5[k_5] + \sum_{m=1}^4 \lambda_m \|\mathbf{W}_m\|_2. \quad (7)$$

The goal is to minimize the average loss $\mathcal{L}(\mathbf{x}_5, \mathcal{I})$ for across all pairs $(\mathbf{x}, \mathcal{I})$ in the training set.

3.3 Training

The network is trained on categorical cross-entropy over shuffled mini-batches of size 512 with uniform class distribution. The learning rate policy for each scalar weight in the network is *Adam* [4], a state-of-the-art online optimizer for gradient-based learning.

3.4 Joint supervision

4. LIMITED WEIGHT SHARING

An Euclidean division of k_1 by Q yields $k_1 = j_1 \times Q + \chi_1$.

Representation	Error rate (%)
MFCC & random forest	—
ConvNet, full weight sharing	—
ConvNet, limited weight sharing	—

Table 1

Representation	Error rate (%)
MFCC & random forest	—
ConvNet, full weight sharing	—
ConvNet, limited weight sharing	—

Table 2

$$\mathbf{y}_2[t, k_1, k_2] = \mathbf{b}_2[j_1, k_2] + \mathbf{W}_2[t, \chi_1, j_1, k_2] \overset{t, \chi_1}{*} \mathbf{x}_1[t, \chi_1, j_1]. \quad (8)$$

Limited weight sharing has been introduced by Abdel-Hamid et al. [1].

5. SINGLE-INSTRUMENT CLASSIFICATION

5.1 Experimental design

In order to train the proposed algorithms, we used MedleyDB v1.1. [2], a dataset of 122 multitracks annotated with instrument activations as well as melodic f_0 curves when present. We extracted the monophonic stems corresponding to a selection of eight pitched instruments [see Figure 2]. Stems with leaking instruments in the background were discarded. The evaluation set consists of 120 recordings of solo music collected by Joder et al. [3]. We discarded recordings with extended instrumental techniques, since they are under-represented in MedleyDB.

5.2 Results

Results are charted in Table 1.

6. POLYPHONIC CLASSIFICATION

6.1 Experimental design

6.2 Results

Results are charted in Table 2.

7. CONCLUSIONS

8. REFERENCES

- [1] Ossama Abdel-Hamid, Abdel-rahman Mohamed, Hui Jiang, Li Deng, Gerald Penn, and Dong Yu. Convolutional neural networks for speech recognition. *Audio, Speech, and Language Processing, IEEE/ACM Transactions on*, 22(10):1533–1545, 2014.
- [2] Rachel Bittner, Justin Salamon, Mike Tierney, Matthias Mauch, Chris Cannam, and Juan Bello. Medleydb: a multitrack dataset for annotation-intensive mir research. *International Society for Music Information Retrieval Conference*, 2014.
- [3] Cyril Joder, Slim Essid, and Gaël Richard. Temporal integration for audio classification with application to musical instrument classification. *IEEE Transactions on Audio, Speech and Language Processing*, 17(1):174–186, 2009.
- [4] Diederik P. Kingma and Jimmy Lei Ba. Adam: a Method for Stochastic Optimization. *International Conference on Learning Representations*, pages 1–13, 2015.
- [5] Brian McFee, Matt McVicar, Colin Raffel, Dawen Liang, Oriol Nieto, Eric Battenberg, Josh Moore, Dan Ellis, Ryuichi Yamamoto, Rachel Bittner, Douglas Repetto, Petr Viktorin, Joo Felipe Santos, and Adrian Holovaty. librosa: 0.4.1. zenodo. 10.5281/zenodo.18369, October 2015.