

Invariances and Data Augmentation for Supervised Music Transcription

Published in: 2018 IEEE International Conference on
Acoustics, Speech and Signal Processing (ICASSP)

7107018026

統研一 劉俊廷

Contents

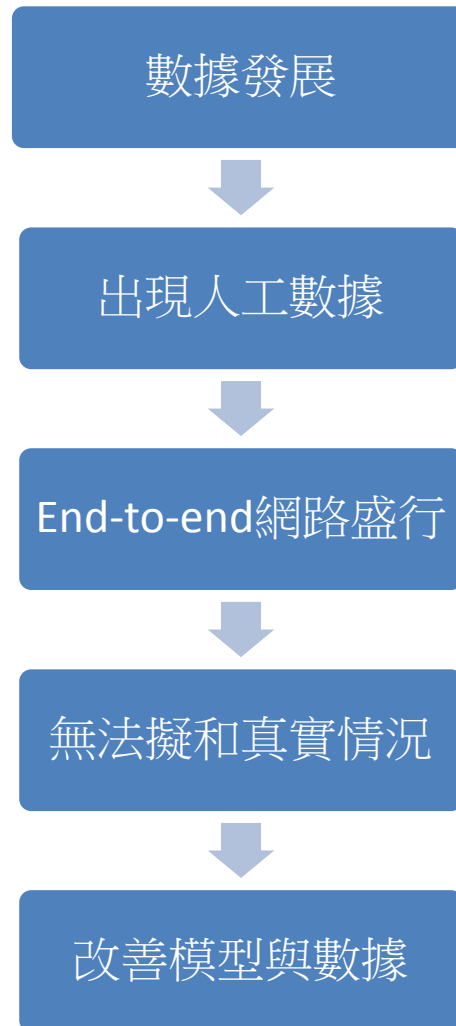
Introduction

Methods

Training

Results

Introduction



Introduction

Model	Prec.	Rec.	Acc.	Etot
MIREX 2009 Dataset				
THK1	82.2	78.9	72.0	.316
KD1	72.4	81.1	66.9	.419
MHMTM1	72.7	78.2	65.5	.441
WCS1	64.0	80.6	59.3	.569
ZCY2	62.7	56.2	50.6	.601
Su Dataset				
THK1	70.1	54.6	51.0	.529
KD1	45.9	45.0	38.1	.745
WCS1	63.6	39.7	35.7	.700
MHMTM1	61.2	36.8	35.2	.676
ZCY2	40.9	28.2	26.2	.799

Methods

本文考慮的模型可以分成三大類:

1. Filter bank:

概念:藉由filter去篩選特徵值，以減少網路參數

1. **Short-time Fourier transform**
2. **Log-spaced filterbank**
3. **Windowed filterbank(cosine window:1-cos(t))**

$$filter_k = (w_{k,\sin}^T x_t)^2 + (w_{k,\cos}^T x_t)^2$$

Methods

2.Translation-invariant:

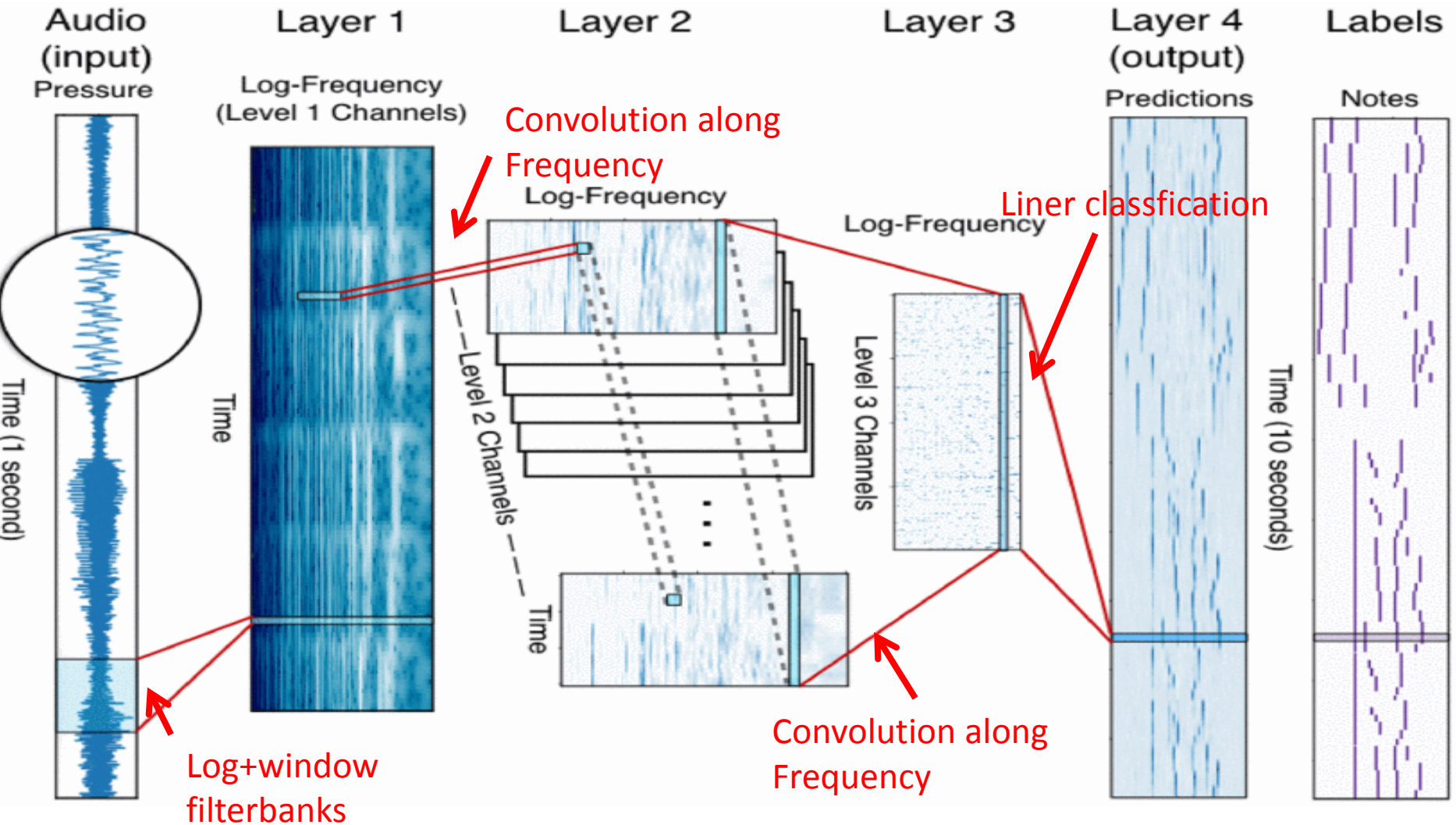
概念:依照頻率的平移不變性概念改進網路，使網路保留頻率的空間關係

1.baseline

2.pitch-shift（使用音高變換優化）

3.wide layer 3(增加第三層節點量)

Methods



Methods

3.end-to-end:

1.three layer network

2.channel convolution(類似Translation-invariant
的結構)

Training

- dataset:
MusicNet
- 學習模式:
momentum($p=.95$)
- 資料增強:
對頻域中的音調做偏移(± 5 半音)，藉此增加資料，可有效強化Translation-invariant network，也對每個數據點移動 $[-1, 1]$ 的位置

Results

Model	Avg. Prec.	Acc.	Err.
filterbanks			
STFT (no compress)	40.4	15.9	.860
STFT	60.4	36.2	.681
log frequencies	62.7	39.8	.646
cosine windows	66.1	38.7	.637
log + windows	66.7	38.9	.633
three layer network	73.8	51.4	.541
end-to-end			
learned filterbank [7]	67.8	48.9	.634
three layer network	70.8	48.8	.558
deep complex [12]	72.9	-	-
channel convolution	73.3	50.4	.531
translation-invariant			
baseline	76.5	53.2	.496
pitch-shift	77.1	54.5	.482
wide layer 3	77.3	55.3	.474
commercial software			
Melodyne [22]	58.8	41.0	.760

Thanks for listening