SUPERVISED LEARNING FOR GENRE CLASSIFICATION OF AUDIO TRACKS

Author: Angel Bergantiños Yeste Director: Dr. Sergio Escalera Guerrero Deparment of Applied Mathematics and Analysis

Context and motivation



History



Telephones were the first to transmit audio through electricity



Thomas Stockham with the first digital recorder (1976)



First iPod released by Apple in 2001

State of the art



MFCC to extract information of a song's frequencies



Deep learning techniques are being used in bigger studies

Method whole overview





Method: MFCC

Mel Scale

$$m = 2595 \cdot \log_{10} \left(1 + \frac{F}{7000} \right)$$





Method: MFCC

Discrete Cosine Transform:

$$C_n = \sum_{k=1}^{k} (log D_k) \cos\left[m(k-\frac{1}{2})\frac{\pi}{k}\right]$$



-5-449.0	-2.5929	9.5790	3.73.90	1+0974	8-198.019	0.4002	-0.4141	-0.5532	0.2255	-0.1020	-0.0399	11	
-3.6378	-3.1142	1.0405	1.11403	1.8470	1.2766	0.0804	0.2627	-0.8690	0.5633	-0.1017	-0.6296		
-5.0660	-4.6436	6.0634	-2.2439	1.4021	0.5963	-1.4794	1.5902	-0.2153	0.0075	-0.1500	-0.3393		
-1,0055	-9.9207	3.7205	-0.9701	1.0237	-1.2302	-1.4634	0.9555	-1.2980	0.2904	-0.6725	0.1645		
-0.3097	-9.7017	3.6233	0.0970	1.7642	-0.5255	-1.0207	8.4669	-0.4703	0,2964	-0.1535	0,2327		and the second se
-1.3033	-5.2693	2.0101	-2,3445	-0.9608	0.9933	-0.5446	0.971	0.1454	0.0297	~0,1504	-0.3407		Each worr
-1.7296	-6,4530	4.1630	-2.0912	-0.5213	0,9857	-0.7430	0.5864	0.2250	1.1920	0.1467	-0.0459		Lach row
-1,9640	-3.0545	3.5566	-0.9636	0.9436	0.3203	-1.1640	1.0310	-0.4067	0.1223	0.2637	-0.4239		
-0.4012	-3.3034	1.3002	-0.4901	-0.2171	0.1642	-0.3015	1.2705	0.0540	0.6925	0.3036	-0.1465		a state and an and a state of the
-0.5713	-2.4569	8.5435	-0.1623	0.0595	0.1440	-0.2646	2.2.441	0.2524	0.0075	-0.4515	-0.4493		renresent
0.4967	-6.4202	0.9275	-2.34848	-1.1777	1.1550	0.4456	2.2468	0.0612	-0.0079	0.6328	-0.5521		represent
0.1287	-3.3791	0.7737	-0.5596	-0.3027	0.9926	0.2468	1.0093	0.2848	0.6419	0.5000	-0.2926		-
-0,0728	-3.7563	0.5699	-0.0489	0.2343	0.4756	0.4992	1.6032	0.1652	0.0122	-0.2101	-0.0362		the
-0.3395	-3.4437	0.7674	0.0454	-0.4404	0.2064	0.2407	0.0200	Q.7244	0,6887	0.3326	-0.0276		the
-0.0023	-2.6316	0.5554	0,2128	0.2285	8.7151	0.9026	0.9031	0.3237	0.6447	-0.2010	-0.7701		and the second second second
1.2008	-4,7609	1.4773	-1.9036	0.1806	0.7917	-1.9056	0.0885	-0.0204	~0.6530	-0.5596	-0.7898		C
4.2404	-7.0050	2.2855	-3.0301	0.4077	0.3009	-2.4106	1.0520	-0.4195	-0.3556	-1,4760	-1.0200		Coefficier
5.1101	-5-2439	5.9252	-2.0519	-0.0767	-0.5758	-4.5036	0.6527	-1.0000	-0.0240	-5.6930	-1.0995		Countrates
5.3670	-5.3009	1.7243	-3.2641	-0.4697	-0.5444	-2.4957	0.4939	-1,3370	-0.0723	-1.5724	-1.2014		
4.4299	-2.2795	2.0896	-1.7947	0.1299	-0.6778	-2.0999	0.7455	-0.5207	0.1553	-0.2849	-0.5669		te of one
9,3647	-2.2947	1.5319	-2.0383	1.2013	0.7549	-1.2512	0.7727	-0.7075	0.7140	0.2346	-0.4544		to or one
5.7591	1.1100	0.0734	-0.3302	1.2090	0,6303	-0.7911	0.0968	-0.674Z	-0.1100	0.3209	-0.6772		The second se
6,3980	-0.0449	3.2490	-0.448	0. <100	-0.0670	-2-9748	-0.2493	-1.0762	-0.5821	0.1700	-1.3295		fromo
7.4630	-1.4000	3.2646	-0.0418	0.2532	-1,0338	-3.1235	-0.3227	-2.4547	-0.5505	0.4764	-1.7012		ITame
7.5910	-0.0154	2.7045	-1,1042	0.4519	-1.0370	-2.9399	-0.4293	-2.2024	-0.4330	0.2524	-1.7090		100000000000000000000000000000000000000
7,4905	-0.1935	3.1.3.2.1	-9.3805	0.7255	-0.7637	-2.7753	-0.3251	-2.2230	-0.4442	0.4502	-1.3500		
4.6246	3.3056	3.723.52	0.2628	0.5791	-0.2289	-3.19043	-3.2537	-1.6025	-0.9995	0,1915	-1-3976		
6.1054	-3.0996	5.0034	-0.1610	-0.0508	0.2331	-0.0025	1.1105	-0.6608	0.0478	0.3511	-0.0502		
7.1377	-1.6330	4.9773	-0.2123	0.3024	-0.4199	0.2590	0.0320	-0.6091	0.3869	0.2132	0.2403		
5.2067	-1.4740	4.1267	-0.5511	-0.9703	0.3327	-0.1149	0.5340	-0.4765	0.8246	-0.0710	-0.5921	m	
											12		
1			F	ach co	humn	conros	ante on	10			12		
				acii co	iumm i	epres	ents on	le					
				Exte	antod (Cooffi	iont						
				EXI	acted	Coem	lent						

Method: PCA and Feature representation

Naive representation

Every frame is concatenated after another.

[Frame₁][Frame₂]...[Frame_N]

Where each Frame has as many values as MFCC. Frame1 = $[MFCC_1, ..., MFCC_M]$ Histogram representation Each histogram is formed with the MFCC of every frame.



Method: PCA and Feature representation

Naive: $[Frame_1][Frame_2]...[Frame_N]$ Histogram: $[Hist_1][Hist_2]...[Hist_M]$



Method: learning schemes

Support Vector Machine

Different kernels:

- Rbf
- Linear
- Poly
- Sigmoid



Results: Data

Marsyas dataset: http://marsyasweb.appspot.com/download/data_sets/

Genres		Properties
Blues	Jazz	• Sample rate: 22050Hz
Classical	Metal	• Channels: 1 (Mono)
Country	Рор	- Example rates 22050 frag
Disco	Reggae	• Frame rate: 22050 ips
Hip hop	Rock	

Results: Evaluation protocol and method parameters.



Confusion matrix



Gamma and C parameters

Results: Naive

Naïve representation without PCA

n_mfcc	rbf	Linear	Poly	Sigmoid
5	15	38	38	10
10	15	44	41	10
15	22	47	39	11
20	24	46	41	11



	rbf	Linear	Poly	Sigmoid		
n_mrcc	$\gamma = 10^{-7}$	C=1	$\gamma = 10^{-7}$	$\gamma = 10^{-9}$		
15	50	47	42	34		

Naïve representation with PCA

n_mfcc	rbf	Linear
5	12	34
10	12	42
15	12	44
20	12	41

6	rbf	Linear
n_mfcc	$\gamma = 10^{-7}$	$C = 10^{-2}$
15	46	45

Results: Histograms

Histogram representation



Results: Histograms

Histogram representation without PCA (linear)

n_mfcc	10	20	30	40	50	60	70	80	90
5	25	41	41	45	41	41	40	41	41
10	25	42	44	50	50	50	48	49	51
15	25	42	44	53	51	50	54	57	59
20	25	42	44	53	52	51	54	55	57

Histogram representation with PCA (linear)

n_mfcc	10	20	30	40	50	60	70	80	90
15	25	42	44	53	51	50	54	57	59

Histogram representation without PCA (rbf)

n_mfcc	10	20	30	40	50	60	70	80	90
5	24	37	33	35	34	35	31	32	32
10	24	37	37	33	35	35	33	31	31
15	24	37	37	35	37	33	31	31	29
20	24	37	37	35	37	32	30	31	30

Histogram representation with PCA (rbf)

n_mfcc	10	20	30	40	50	60	70	80	90
15	24	37	37	34	38	36	34	31	30

Results: Histograms

Best results

True label

	rbfLinear $\gamma = 1$ $C = 10$ 7065	
n_micc	$\gamma = 1$	C = 10
15	70	65

Confusion matrix, without normalization													
Blues -	5	0	2	0	0	0	2	0	1	0		9	
Classical -	0	9	0	0	0	0	0	0	0	1		- 8	1
Country -	0	0	8	0	0	1	0	1	0	0		- 7	
Disco -	0	0	2	6	0	0	0	1	1	0		- 6	
Hip hop -	0	0	0	1	5	0	1	1	2	0		- 5	
Jazz -	0	0	0	0	0	8	0	0	0	2		- 4	
Metal -	1	0	0	0	0	0	8	0	0	1		- 3	
Pop -	1	1	1	0	0	0	0	7	0	0		-2	
Reggae -	1	0	0	0	1	0	0	0	7	1			
Rock -	0	0	1	1	0	1	0	0	0	7			
	65		2.4	1.0	6	2.12	10	- 2		4.5		0	
	BINE	assic	ount	OFF	10 20	φ,	Meth	60. 62.	egge	600			
	0			Pre	dicte	ed la	bel						

set 1	set 2	set 3	set 4	set 5
61	69	74	64	75
set 6	set 7	set 8	set 9	set 10

Results: Feature relevance

Feature ranking: 1. feature 3 (0.011840) 2. feature 0 (0.011585) 3. feature 0 (0.010474) 4. feature 3 (0.010153) 5. feature 1 (0.010096) 6. feature 8 (0.010084) 7. feature 0 (0.009707) 8. feature 10 (0.009629) 9. feature 2 (0.009549) 10. feature 1 (0.009488) 11. feature 2 (0.009082) 12. feature 0 (0.008644) 13. feature 3 (0.008399) 14. feature 3 (0.008375) 15. feature 11 (0.008210 16. feature 9 (0.008207) 17. feature 0 (0.008100) 18. feature 0 (0.008054) 19. feature 0 (0.007974) 20. feature 1 (0.007854) 21. feature 7 (0.007851) 22. feature 2 (0.007696) 23. feature 2 (0.007666) 24. feature 3 (0.007388) 25. feature 1 (0.007354)



Results: Extra

AdaBoost: Worse performance (18%)

Naïve-Bayes: Lower accuracy, but similar to SVC default results (56%)

Conclussion and future work



Bigger is better: with more songs, we can learn better.

With a bigger dataset, deep learning can work better and allows the creation of a more accurate model.



The Elements of Music

Melody Rhythm Harmony Texture Form Tempo and Dynamics

Other features of music can be also used as input to work with the frequency for further audio analysis.



Thanks for your attention!

(20)