



AGH

AKADEMIA GÓRNICZO-HUTNICZA IM. STANISŁAWA STASZICA W KRAKOWIE

Wydział Fizyki i Informatyki Stosowanej

Praca inżynierska

Paweł Pęksa

kierunek studiów: informatyka stosowana

Wykorzystanie algorytmów uczenia maszynowego w celu rozpoznawania nastroju muzyki

Opiekun: dr Marcin Wolter

Kraków, styczeń 2015

Oświadczam, świadomy(-a) odpowiedzialności karnej za poświadczenie nieprawdy, że niniejszą pracę dyplomową wykonałem(-am) osobiście i samodzielnie i nie korzystałem(-am) ze źródeł innych niż wymienione w pracy.

.....
(czytelny podpis)

Merytoryczna o cena pracy przez opiekuna:

Ocena:

Data:

Podpis:

Merytoryczna o cena pracy przez recenzenta:

Ocena:

Data:

Podpis:

Spis treści

1	Wprowadzenie	6
1.1	Przedmiot pracy	6
1.2	Problematyka pracy	6
2	Sztuczne sieci neuronowe	7
2.1	Budowa sieci neuronowej	7
2.1.1	Budowa neuronu	7
2.1.2	Topologia sieci	7
2.2	Uczenie sieci neuronowej	9
2.2.1	Reguła delta	9
2.2.2	Algorytm wstecznej propagacji błędów	10
3	Ekstrakcja cech dźwiękowych	10
3.1	Cyfrowa reprezentacja sygnału audio	11
3.2	Spektralna reprezentacja sygnału audio	11
3.2.1	Transformacja Fouriera	11
3.3	Wstępna obróbka sygnału	12
3.3.1	Okno czasowe	13
3.3.2	Algorytm wyrównywania poziomu głośności dźwięku	13
3.4	Cechy dźwięku bazujące na czasowej reprezentacji dźwięku	14
3.4.1	Wskaźnik zmiany znaku (<i>Zero Crossing Rate</i>)	14
3.4.2	Wskaźnik zmian (<i>Onset rate</i>)	14
3.5	Cechy dźwięku bazujące na spektralnej reprezentacji dźwięku	14
3.5.1	Złożoność spektralna (<i>Spectral complexity</i>)	14
3.5.2	Kształt spektralny (<i>Spectral shape</i>)	14
3.5.3	Płaskość spektralna <i>Spectral flatness</i>	16
3.5.4	Dysonans (<i>Dissonance</i>)	16
3.5.5	Skala	16
4	Matematyczny model emocji	17
5	Opis stworzonego systemu	18
5.1	Schemat systemu	18
5.2	Użyte narzędzia	18
6	Wyniki	18
7	Wnioski	18

1 Wprowadzenie

1.1 Przedmiot pracy

Muzyka towarzyszyła człowiekowi od czasów prehistorycznych. Z czasem stała się jedną z form sztuki. Niewątpliwie, gdy słyszymy jakąś melodię, nie sprawia nam większego kłopotu, aby określić emocje z nią związane. Należy jednak mieć na uwadze źródło emocji. Rozróżnić możemy emocje wyrażane przez muzykę oraz przez nią indukowane. Tematem niniejszej pracy jest rozpoznawanie nastroju muzyki przy użyciu uczenia maszynowego. Przez określenie "nastrój muzyki" mamy tutaj na myśli emocje, które ta muzyka reprezentuje. Konkretnym narzędziem wybranym dla tego przedsięwzięcia jest sztuczna sieć neuronowa. Idea klasyfikacji utworów muzycznych pod tym kątem jest względnie nowa, lecz można zauważyć wzrastające zainteresowanie tym tematem[?]. Do tej pory powstał szereg różnych prac podejmujących to zadanie[?][?][?]. Niektóre korzystających nie tylko z samego sygnału audio, ale także np. z tekstu utworu[?]. W tej pracy jednak pod uwagę brany jest jedynie sygnał audio z którego podjęto próbę wyekstrahowania odpowiednich cech audio, które mogłyby pozwolić sieci neuronowej rozpoznawać emocje reprezentowane przez dany utwór muzyczny według wybranego modelu.

W kolejnych rozdziałach zostały opisane niezbędne podstawy teoretyczne, których przyswojenie pozwala zrozumienie w jaki sposób postawione zadanie jest realizowane. W rozdziale 2 można zapoznać się z krótkim opisem sieci neuronowych, w rozdziale 3 czytelnik dowie się o tym jakie konkretnie cechy dźwiękowe zostały rozważane, natomiast w rozdziale 4 opisany jest sposób w jaki matematyka pomaga nam w opisanu emocji. Opis stworzonego systemu można znaleźć w rozdziale 5. Wyniki oraz wnioski wynikające z podjętej próby podejścia do opisywanego zagadnienia zostały przedstawione w rozdziale kolejno 6 oraz 7.

1.2 Problematyka pracy

Analiza muzyki pod kątem emocji jest zadaniem, które nie wiąże się tylko z przetwarzaniem sygnałów oraz uczeniem maszynowym, ale także z psychologią muzyki oraz jej teorią. Biorąc to pod uwagę jest to bardzo wymagający problem. Nastrój utworu muzycznego może być wysoce subiektywnym odczuciem. Wpływ na ocenę mogą mieć także wspomnienia danej osoby, nastrój w danej chwili, indywidualne preferencje czy poziom wykształcenia muzycznego. Wszystkie wspomniane jednak kwestie nie są na tyle znaczące, aby podjęcie pracy nad tym tematem było niemożliwe czy całkowicie nieskuteczne. Oczywiście nie da stworzyć się systemu, który będzie działał niezawodnie, bo tak jak zostało to wspomniane, zbyt wiele indywidualnych czynników ma wpływ na percepcje człowieka. Dotychczasowe badania udowadniają, że jest możliwe sprostanie temu zadaniu w zadowalającym stopniu[?] i taka też próba zostaje podjęta w niniejszej pracy.

2 Sztuczne sieci neuronowe

Sztuczną siecią neuronową, dalej zwaną po prostu siecią neuronową, określamy model matematyczny służący jako system przetwarzania informacji. Źródłem inspiracji dla tegoż modelu była biologia, a mianowicie sieci neuronowe istniejące w ludzkim mózgu. Dużą zaletą sieci neuronowych jest fakt, iż nie działają jak tradycyjne algorytmy, lecz mają one zdolność do uczenia się. Jest to kolejna analogia związana z pracą ludzkiego mózgu, co powodowało wzrastające zainteresowanie tym tematem. Pomimo tego, że nie udało się z ich pomocą odtworzyć pracy tego niewątpliwie niesamowitego narządu, znajdują one zastosowanie w wielu dziedzinach nauki. Najbardziej podstawowym rodzajem sieci neuronowej jest sieć jednokierunkowa tj. czyli taka, w której nie występują sprzężenia zwrotne i takie też sieci są wykorzystane w tej pracy oraz opisane w tym rozdziale.

2.1 Budowa sieci neuronowej

2.1.1 Budowa neuronu

Sieć neuronowa zbudowana jest z neuronów, które są odpowiednikami komórek nerwowych. Synapsy łączące poszczególne komórki modelowane są przez wagi liczbowe, których wielkości z kolei można interpretować jako wpływ jednej komórki na drugą. Matematyczny model neuronu użyty w sieciach neuronowych jest przedstawiony na rysunku 1. Składa się on z $n + 1$ wejść oraz jednego wyjścia. Dodatkowym wejściem neuronu jest tzw. bias, który przyjmuje stałą wartość, ale jest także modyfikowany w procesie uczenia (podobnie jak pozostałe wagi). Najczęściej przyjmuje się, że przed rozpoczęciem uczenia jego wartość równa jest 1. Zależność pomiędzy sygnałami wejściowymi x_i , wagami w_i , a tzw. sygnałem sumarycznego pobudzenia φ w najprostszym przypadku może być określana przez wzór:

$$\varphi = w_0 + \sum_{i=1}^n w_i x_i. \quad (1)$$

Sposób obliczenia sygnału wyjściowego neuronu określany jest przez funkcję aktywacji:

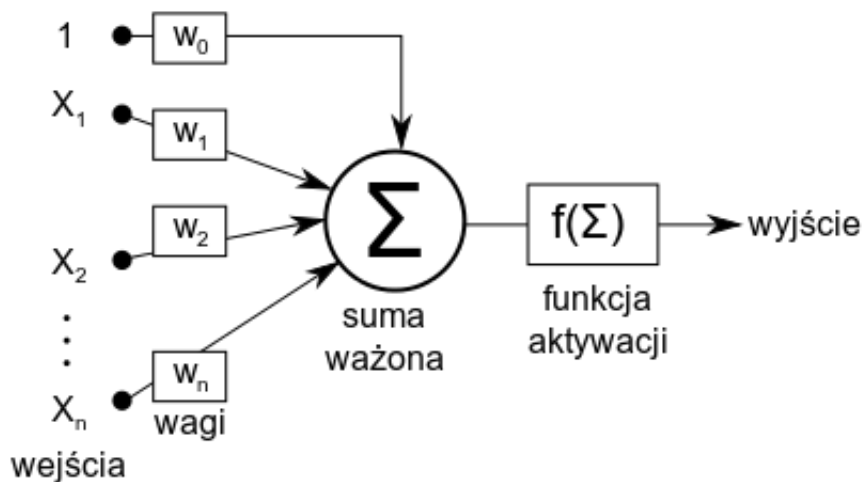
$$y = f(\varphi). \quad (2)$$

Rozważane były różne funkcje aktywacji, ale ostatecznie powszechnie używa się czterech z nich: funkcji liniowej, funkcji sigmoidalnej, funkcja tangens hiperboliczny oraz funkcji Gaussa[?]. Poszczególne funkcje zostały przedstawione na rysunku 2

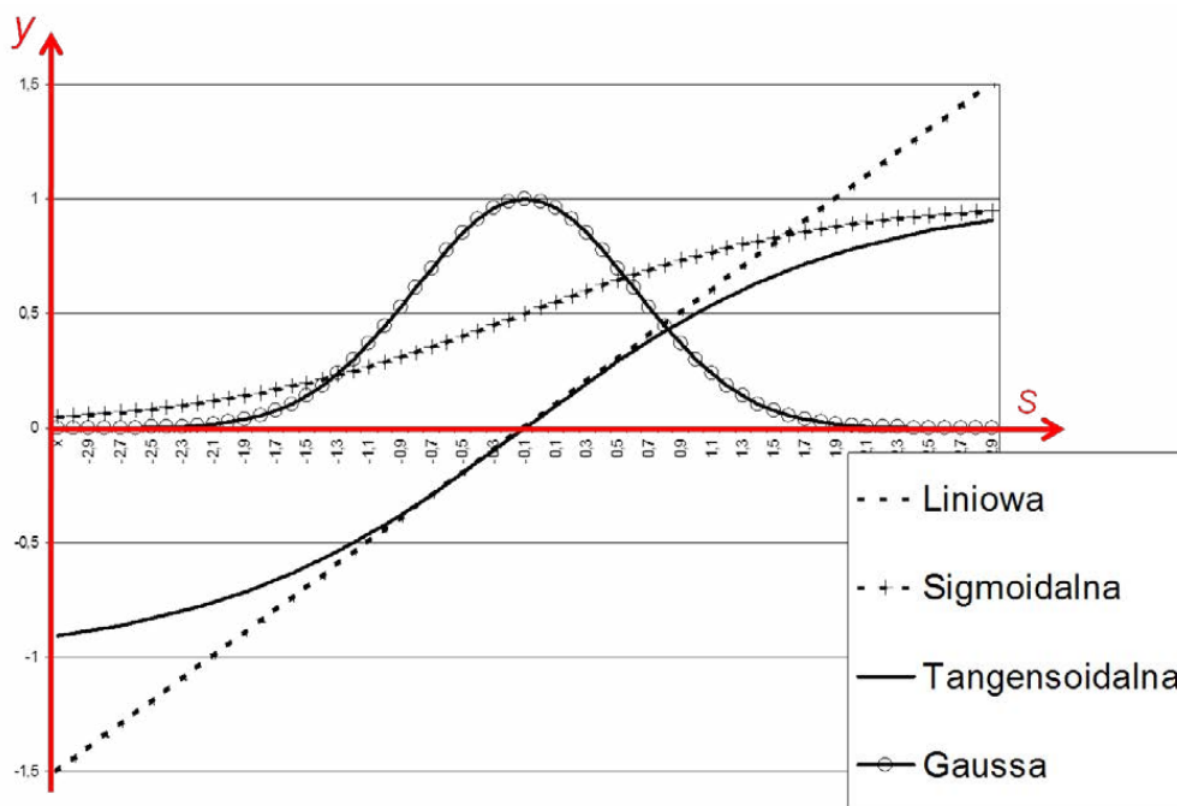
2.1.2 Topologia sieci

Wszystkie neurony zgrupowane są w warstwy z których możemy wyróżnić:

- warstwę wejściową



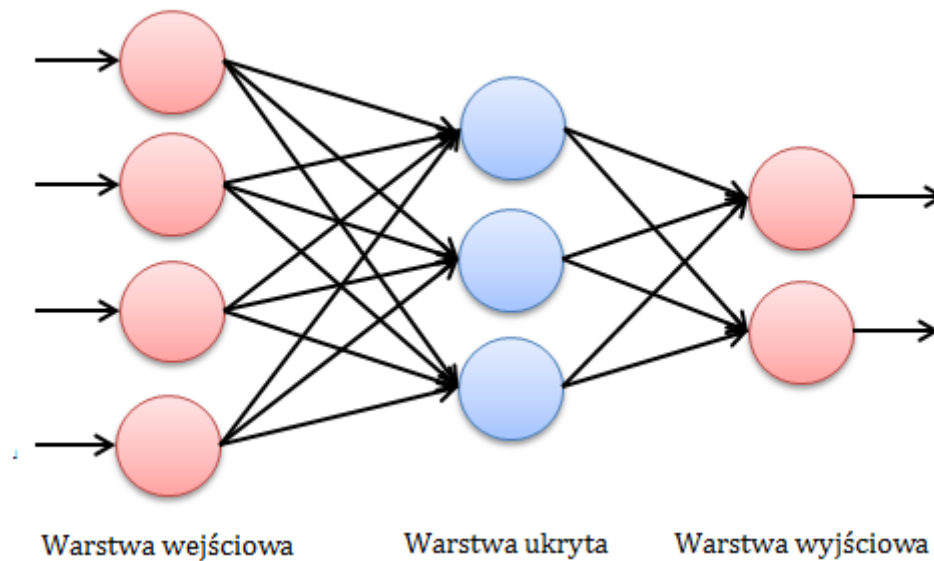
Rysunek 1: Matematyczny model neuronu



Rysunek 2: Najpowszechniejsze funkcje aktywacji

- jedną lub więcej warstw ukrytą
- warstwę wyjściową

W każdej z warstw znajduje się dowolna liczba neuronów, który posiada połączenie do wszystkich neuronów znajdujących się w warstwie kolejnej. Przykład sieci neuronowej składającej się z dokładnie trzech warstw zawierającej różną ilość neuronów w każdej z nich został przedstawiony na rysunku 3.



Rysunek 3: Przykładowa jednokierunkowa sieć neuronowa

2.2 Uczenie sieci neuronowej

Wyróżniamy dwa podstawowe sposoby uczenia sieci:

- uczenie z nauczycielem (uczenie nadzorowane)
- uczenie bez nauczyciela (uczenie nienadzorowane)

Dzięki możliwości uczenia sieci neuronowej nie jest konieczne projektowanie algorytmu, który przetwarza dla nas informacje w oczekiwany sposób. Sieć neuronowa korzystając z odpowiedniego algorytmu uczenia sama modeluje ten algorytm poprzez modyfikację wag. Należy też wspomnieć, że początkowe wagi sieci neuronowej zwykle inicjalizowane są wartościami losowymi.

2.2.1 Reguła delta

Proces uczenia, tak jak zostało to wspomniane, polega na modyfikowaniu jej współczynników wagowych. Opisane w tym podrozdziale zostanie uczenie z nauczycielem, ponieważ takie właśnie zostało wykorzystane w tej pracy. W przypadku uczenia z nauczycielem potrzebujemy zbioru uczącego składającego się z wejścia, które podajemy sieci i oczekiwanego rezultatu dla tego wyjścia, co możemy oznaczyć jako pary (x_i, z_i) , gdzie z_i jest oczekiwaną odpowiedzią dla sygnału wejściowego x_i . Zadaniem sieci jest wymodelowanie funkcji:

$$h(x) = z. \quad (3)$$

Uczenie jest procesem iteracyjnym, gdzie w każdej iteracji modyfikujemy wagi sieci. Liczba iteracji N równa jest liczbie par (x_i, z_i) . W każdym kroku j procesu uczenia możemy zdefiniować

wielkość błędu neuronu wyjściowego jako:

$$\delta_i^j = |z_i^j - y_i^j|. \quad (4)$$

Cel procesu uczenia można określić jako minimalizowanie funkcji:

$$Q = \frac{1}{2} \sum_{j=1}^N (\delta_i^j)^2. \quad (5)$$

Korzystając z metody gradientowej możemy zdefiniować poprawkę δw dla wagi w neuronu i jako:

$$\Delta w_i = -\eta \frac{\partial Q}{\partial w_i}, \quad (6)$$

oraz, idąc dalej, zdefiniować wzór dla korygowania wagi w w kolejnych krokach:

$$w_i^{j+1} = w_i^j + \Delta w_i \quad (7)$$

przy czym η jest dodatkowym współczynnikiem liczbowym, który decyduje o szybkości uczenia. W ten sposób poprawiamy wagi sieci j razy.

Problemem z którym borykano się do połowy lat 80-tych, a który spostrzegawczy czytelnik mógł już zauważyć, jest fakt, iż tym sposobem niemożliwe jest uczenie sieci, która składa się z więcej niż jednej warstwy, ponieważ nieznana jest oczekiwana odpowiedź neuronów warstwy innej niż wyjściowej. Wzór 6 jest jednak punktem wyjściowym większości algorytmów automatycznego uczenia[?] i jest on określany w literaturze regułą delta[?].

2.2.2 Algorytm wstecznej propagacji błędów

W celu przeprowadzenia uczenia dla sieci wielowarstwowej spotykamy się z potrzebą określenia błędu δ dla neuronów, które należą do warstw ukrytych sieci neuronowej. Umożliwia nam to algorytm wstecznej propagacji błędów. Błąd dla takiego neuronu obliczamy korzystając ze wszystkich błędów neuronów do których wysłał on swój sygnał. Uwzględniane są także wagi połączeń. Mowa jest o wstecznej propagacji, ponieważ odbywa się ona przeciwnie do przepływów sygnałów w sieci. Błąd dla neuronów znajdujących się w warstwie innej niż wejściowa możemy określić jako:

$$\delta_i = f'(\varphi) \sum_k w_k \delta_k, \quad (8)$$

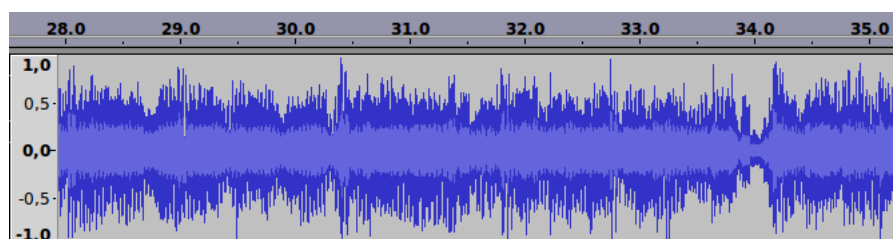
przy czym w_k oraz δ_k są kolejno wagami oraz błędami neuronów do których analizowany neuron wysyłał swój sygnał.

3 Ekstrakcja cech dźwiękowych

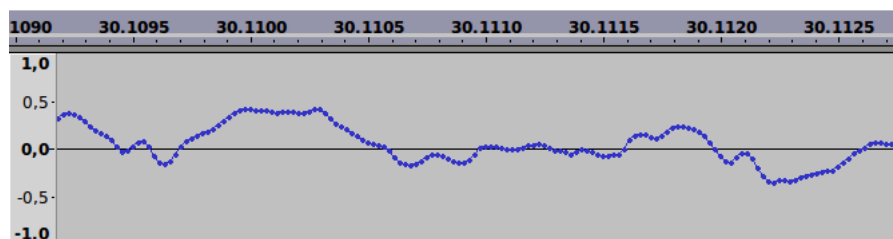
Cechami dźwiękowymi nazywamy zmienne wyekstrahowane z sygnału audio, które opisują ten sygnał i pozwalają uzyskać dodatkowe informacje na jego temat[?]. Opisane tym w rozdziale zostały cechy sygnału audio, które są wejściami dla użytego klasyfikatora tj. sieci neuronowej. Poszczególne cechy bazowały zarówno na czasowej jak i spektralnej reprezentacji sygnału.

3.1 Cyfrowa reprezentacja sygnału audio

Dźwięk jest sygnałem analogowym. W celu przechowywania go na cyfrowych nośnikach spotkano się z potrzebą jego digitalizacji tj. reprezentacji w postaci cyfrowej. Najczęściej stosowaną w tym celu jest metoda PCM¹ w której to sygnał analogowy jest próbkowany w równych odstępach czasu i zapisywany cyfrowo. Powszechna częstotliwość próbkowania wynosi 44 100 Hz ze względu na zakres słyszalnych częstotliwości człowieka, który wynosi około 20 000 Hz, a zgodnie z prawem Nyquista sygnał powinien być próbkowany z dwa razy wyższą częstotliwością niż maksymalna częstotliwość sygnału w celu uzyskania dokładnej reprezentacji bez zniekształceń[?]. Przykładowy, reprezentowany cyfrowo, sygnał audio przedstawia rysunek 4 oraz 5. Na rysunku 5 przedstawiony został ten sam sygnał, który widzimy na 4, lecz w bardzo dużym przybliżeniu. Oba rysunki zostały wygenerowane za pomocą programu Audacity² z pliku audio w formacie mp3³.



Rysunek 4: Przykładowy sygnał audio



Rysunek 5: Przykładowy sygnał audio - przybliżenie

3.2 Spektralna reprezentacja sygnału audio

3.2.1 Transformacja Fouriera

Każdy sygnał, który jest reprezentowany jako zmieniająca się w czasie amplituda posiada też odpowiadające spektrum częstotliwościowe, które wymiennie nazywa się też widmem. Dotyczy to także sygnału audio. Dzięki przedstawieniu sygnału dźwiękowego w ten sposób możliwe jest uzyskanie dodatkowych informacji na jego temat[?]. Spektrum przedstawia skład

¹PCM - Pulse Code Modulation

²Program służący do edycji plików dźwiękowych

³Format cyfrowego zapisu i kompresji plików dźwiękowych

częstotliwościowy dźwięku. Możemy wyróżnić spektrum amplitudowe oraz spektrum fazowe. Przy obliczaniu spektrum z pomocą przychodzi transformacja Fouriera. Oznaczając $x(t)$ jako sygnał, a $X(f)$ jako wynik transformacji, możemy zapisać:

$$X(f) = \int_{-\infty}^{\infty} x(t)e^{-j2\pi ft} dt, \quad (9)$$

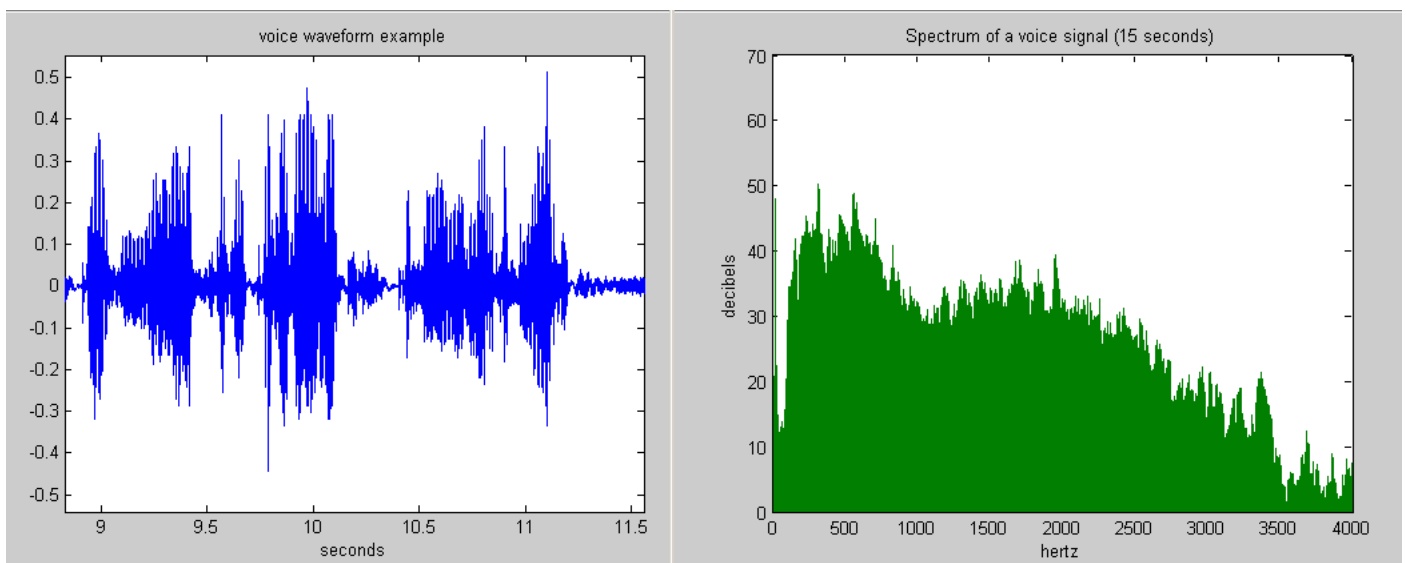
przy czym f to częstotliwość, natomiast t oznacza czas. W celu otrzymania spektrum amplitudowego, z którego szeroko korzystano w niniejszej pracy, należy z otrzymanego wyniku obliczyć wartość bezwzględną tj:

$$|X(f)|. \quad (10)$$

Trzeba jednak pamiętać, że w przypadku sygnału audio zapisanego w pamięci komputera mamy do czynienia z sygnałem dyskretnym, więc należy użyć DFT⁴ - transformaty Fouriera dla sygnałów dyskretnych wyrażającej się wzorem:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-i2\pi \frac{k}{N} n}, \quad (11)$$

przy czym x_n to kolejne wartości próbkowanego sygnału, X_k to wartości transformaty. Na rysunku 6 przedstawiony jest przykładowy sygnał wraz z odpowiadającym mu spektrum amplitudowym.



Rysunek 6: Sygnał audio wraz z odpowiadającym mu spektrum amplitudowym

3.3 Wstępna obróbka sygnału

Przed analizą sygnału dźwiękowego, jakim są utwory muzyczne, w celu poprawy jakości danych przeprowadza się obróbkę wstępną, co pozwala na bardziej efektywną ich analizę. Przed

⁴Discrete Fourier Transform

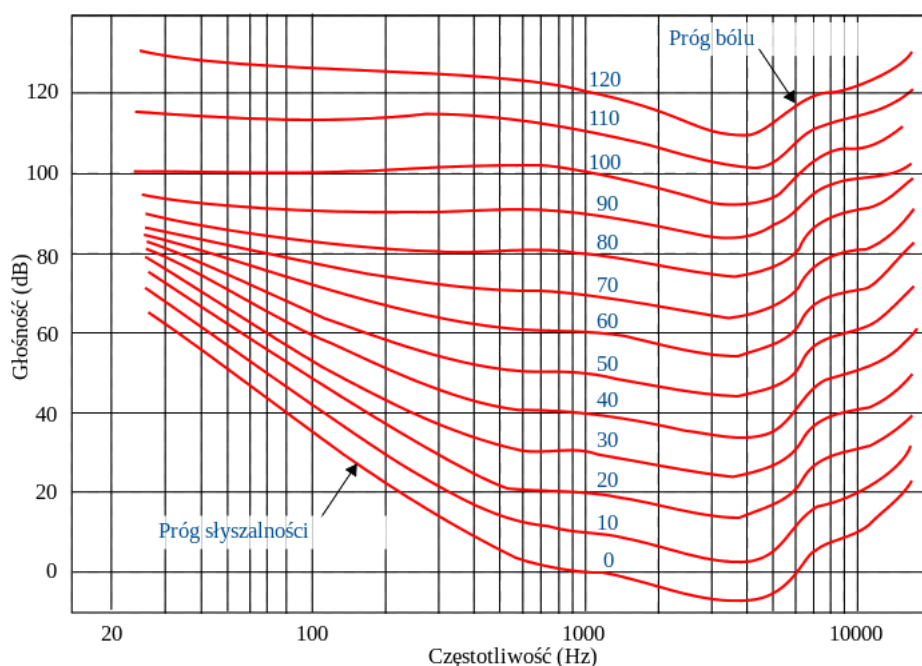
ekstrakcją cech dźwiękowych wykorzystana została funkcja okna czasowego oraz algorytm wyrównywania poziomu głośności dźwięku opisane w kolejnych dwóch podrozdziałach.

3.3.1 Okno czasowe

Ważną rolę przy korzystaniu z transformaty Fouriera odgrywa okresowość. W przypadku, gdy mamy do czynienia z danymi, które występują niecałkowitą ilość razy, na końcach analizowanych danych występują nieciągłości, które powodują, że otrzymane spektrum jest zniekształcone. Rozwiązaniem tego problemu jest okno czasowe. Jest to funkcja, którą mnożymy przy sygnale w celu zmniejszenia zniekształceń[?].

3.3.2 Algorytm wyrównywania poziomu głośności dźwięku

Ludzkie ucho nie słyszy dźwięków o wszystkich częstotliwościach jako dźwięków o tym samym poziomie głośności. Problem ten można rozwiązać częściowo przy użyciu algorytmu wyrównywania poziomu głośności, który filtruje dźwięk korzystając z krzywych izofonicznych. Pojęcie izofony lub też krzywej izofonicznej określa w fonach⁵ słyszalne natężenie dźwięku w zależności od częstotliwości[?]. Ze względu na subiektywność postrzegania głośności nie istnieją ściśle określone krzywe izofoniczne. Istnieje jednak model zaproponowany przez Fletchera oraz Munsona[?], który jest przedstawiony na rysunku 7.



Rysunek 7: Izofony „normalnego” ucha według Fletchera i Munsona, wartości fonów dla poszczególnych krzywych określone są liczbami w kolorze niebieskim

⁵Fon - jednostka poziomu głośności wg słownika języka polskiego PWN

3.4 Cechy dźwięku bazujące na czasowej reprezentacji dźwięku

3.4.1 Wskaźnik zmiany znaku (*Zero Crossing Rate*)

Wskaźnik zmiany znaku (*Zero Crossing Rate*) jest jedną z najprostszych cech dźwiękowych obliczanych korzystając z wykorzystaniem dźwięku reprezentowanym jako zmiana amplitudy w czasie. Wyraża on liczbę zmiany znaków w fali dźwiękowej. Możemy określić go wzorem:

$$ZCR = \frac{1}{2} \sum_{n=1}^N |sgn(x[n]) - sgn(x[n-1])|. \quad (12)$$

Jest to deskryptor często stosowany w pozyskiwaniu informacji z muzyki, ale także stosowany w rozpoznawaniu mowy. Zawdzięcza to łatwości jego obliczania, a także faktowi, że przechowuje informację o szumach występujących w dźwięku[?].

3.4.2 Wskaźnik zmian (*Onset rate*)

Wskaźnik zmian (*Onset rate*) jest podstawowym wskaźnikiem rytmu utworu muzycznego mającym duży wpływ na postrzeganie emocji reprezentowanych przez muzykę. Określa on zmienność dźwięku. Jest on określany liczbą ekstremów obwiedni dźwięku⁶. Zakłada się, że różnica czasowa pomiędzy dwoma zmianami brany pod uwagę w zliczaniu musi wynosić przynajmniej 60 ms[?].

3.5 Cechy dźwięku bazujące na spektralnej reprezentacji dźwięku

3.5.1 Złożoność spektralna (*Spectral complexity*)

Złożoność spektralna (*spectral complexity*) jest liczbą ekstremów w widmie amplitudowym sygnału dźwiękowego. Opisuje ona złożoność tego widma. Utwory z większą średnią złożonością spektralną charakteryzują się większą energicznością[?].

3.5.2 Kształt spektralny (*Spectral shape*)

Kształt spektralny jest kształtem widma amplitudowego danego sygnału dźwiękowego. W jego skład wchodzi m.in. środek masy widma sygnału (*spectral centroid*), współczynnik skośności widma sygnału (*spectral skewness*), kurtoza widma sygnału (*spectral kurtosis*) oraz tzw. *spectral roll-off*. Wszystkie te cechy mają wpływ na odbiór utworu muzycznego przez słuchacza pod kątem reprezentowanych przez utwór emocji[?].

Moment centralny

W celu określenia kolejnych wzorów określających kształt spektralny użyteczne jest zdefiniowa-

⁶Krzywa opisująca zmianę amplitudy sygnału[?]

nie momentu centralnego. Dla zmiennej dyskretnej możemy moment centralny określić wzorem:

$$\mu = \sum_{i=1}^N \frac{(f_i - \bar{f})^r}{N}, \quad (13)$$

gdzie f_i są kolejnymi częstotliwościami występującymi w widmie, \bar{f} średnią częstotliwością, natomiast N to liczba wszystkich częstotliwości. Moment centralny rzędu drugiego, korzystając ze wzoru 13 obliczamy w następujący sposób:

$$\sigma^2 = \sum_{i=1}^N \frac{(f_i - \bar{f})^2}{N} \quad (14)$$

i nazywany wariancją. Pierwiastek kwadratowy z wariancji σ określany jest mianem odchylenia standardowego.

Środek masy widma

Środek masy widma wyrażamy wzorem:

$$SC = \frac{\sum f_i a_i}{a_i}, \quad (15)$$

gdzie f_i jest częstotliwością, natomiast a_i amplitudą dla poszczególnych częstotliwości[?].

Współczynnik skośności widma

Współczynnik skośności mówi o asymetryczności widma. W przypadku, gdy jest on mniejszy od zera, więcej danych znajduje się po lewej stronie widma, w przypadku, gdy jest on większy od zera, więcej danych znajduje się po prawej stronie widma. Możemy określić go wzorem:

$$\gamma = \frac{\mu_3}{\sigma^3} \quad (16)$$

Kurtoza widma

Kurtoza widma jest miarą jego spłaszczenia, wyraża się wzorem:

$$K = \frac{\mu_4}{\sigma^4} \quad (17)$$

Roll-off widma

Cechą dźwięku określana mianem *Roll-off*'u widma jest częstotliwość, która dzieli widmo sygnału na dwie części według ustalonego progu T , który zazwyczaj wynosi 0.95[?]. Omawianą cechę możemy określić wzorem:

$$\sum_{i=1}^{R_t} f_i = T \sum_{i=1}^N f_i[?]. \quad (18)$$

3.5.3 Płaskość spektralna *Spectral flatness*

Płaskość spektralna jest stosunkiem średniej arytmetycznej do średniej geometrycznej widma amplitudowego wyrażonym w decybelach:

$$spectralFlatness = 10\log_{10} \frac{G}{A} \quad (19)$$

przy czym G jest średnią geometryczną:

$$G = \sqrt[n]{\prod_{i=1}^N a_i} \quad (20)$$

oraz A jest średnią arytmetyczną:

$$A = \frac{\sum_{i=1}^N a_i}{N}. \quad (21)$$

Cecha ta określa jak bardzo spłaszczony jest wykres widma amplitudowego. Wraz ze wzrostem tego wskaźnika dźwięk bardziej przypomina szum. Wartość bliska 1 jest przyjmowana dla tej cechy dźwiękowej w przypadku białego szumu⁷[?].

3.5.4 Dysonans (*Dissonance*)

Dysonans jest deskryptorem dźwięku obliczanym na podstawie odstępów pomiędzy ekstremami widma amplitudowego. W przypadku utworów muzycznych cechujących się mniejszym rozdźwiękiem obserwuje się większą równomierność tychże odstępów. Matematycznie dysonans można określić wzorem:

$$dissonance = \frac{1}{H} \sum_{h=1}^H a(h) - SE(h), \quad (22)$$

gdzie H jest liczbą ekstremów, $a(h)$ amplitudą dla danego ekstremum oraz $SE(h)$ amplitudą obwiedni spektrum dla częstotliwości $f(h)$ [?].

3.5.5 Skala

Skala muzyczna jest składowa się z dźwięków o różnych częstotliwościach ułożonych według ustalonego schematu. Możemy wyróżnić dwa podstawowe skale: molową oraz durową. Powszechnie uznaje się, że utwory muzyczne bazujące na skali durowej mają radosne brzmienie, natomiast na skali molowej smutne brzmienie[?]. W celu wyekstrahowania skali muzycznej utworu należy najpierw obliczyć jego HPCP⁸, który obliczany jest na podstawie ekstremów widma amplitudowego według wzoru:

$$HPCP(n) = \sum_{i=1}^N w(n, f_i) a_i^2, \quad (23)$$

gdzie a_i oraz f_i są kolejno amplitudą oraz częstotliwością ekstremum, N jest liczbą wszystkich ekstremów, n kolejną wartością wektora HPCP, natomiast w funkcją wagową określającą w

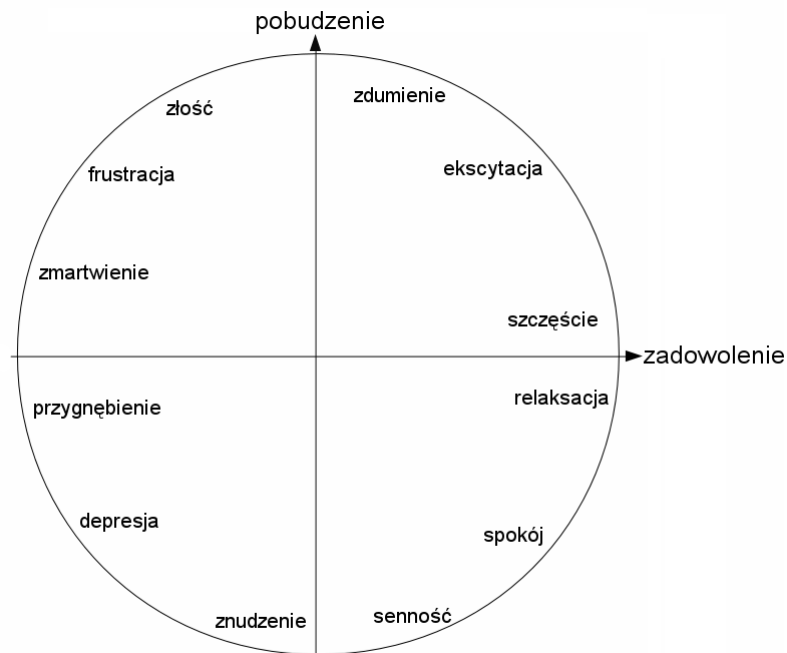
⁷Szum akustyczny o prawie płaskim widmie

⁸Harmonic Pitch Class Profile

jaki sposób poszczególne ekstremum wpływa na wartość n wartości wektora HPCP. W celu określenia skali muzycznej obliczana jest korelacja pomiędzy wektorem HPCP, a odpowiednimi profilami dla obu skali[?].

4 Matematyczny model emocji

W celu realizacji postawionego zadania należy zdefiniować emocje matematycznie. Jednym z podejść stosowanych w pracach o podobnej tematyce jest wykorzystanie etykiet mówiących o emocjach. Każdy z utworów etykietowano z wykorzystaniem przymiotników takich jak np. "smutny", "wesoły", "relaksujący". Innym rozwiązaniem jest określanie emocji z wykorzystaniem dwuwymiarowej przestrzeni opracowanej przez Russella i podobnie postąpiono w niniejszej pracy. W tym przypadku, emocje określone są przy pomocy dwóch parametrów: pobudzenia(ang. *arousal*) oraz zadowolenia(ang. *valence*), co przedstawia rysunek 8. Model ten zakłada, że wszystkie emocje wynikają z dwóch niezależnych systemów neurofizjologicznych⁹ w kontekście do poprzednich teorii zakładających istnienie niezależnych systemów dla każdej podstawowej emocji. Wyniki najnowszych badań są jednak bardziej konsistentne z nowszym podejściem[?].



Rysunek 8: Model emocji według Russela

⁹system neurofizjologiczny - system układu nerwowego

5 Opis stworzonego systemu

5.1 Schemat systemu

5.2 Użyte narzędzia

6 Wyniki

7 Wnioski

Literatura