

Aplikácie audio rozhrania

Juraj Kačur

UMIKT, FEI STU, Bratislava

Obsah

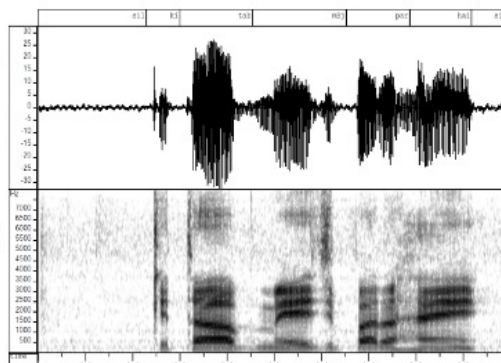
- Detekcia reči
- Potlačenie šumov
- Kompresia reči a audio signálov
- Identifikácia hovoriaceho
- Rozpoznávanie reči
- Syntéza reči
- Ďalšie použitia

2

Obsahom tejto state je vymenovanie základných aplikácií, ktoré môžu byť využité na strane stroja pri modalitách rozhrania využívajúcej zvukové signály, teda zmysel sluchu. Pre každú aplikáciu sú uvedené jej výhody, nevýhody, najčastejšie formy použitia, problémy s jej nasadením, rôzne druhy jej realizácie- typy postupov, atď..

Detekcia reči

- Úloha: zistiť či je v signály obsiahnutá reč a zdetegovať jej polohu
- Detekcia reči (VAD)
- Algoritmy využívajú známe vlastnosti reči
 - Nestacionárny signál
 - Frekv. Pásmo 300-4000 Hz
 - Periodické a neznělé časti
 - Dynamika reči 50dB
 - Rôzne spektrá
 - Plynú za sebou v čase
 - Dôležitá časová postupnosť



3

Najjednoduchšia a najzákladnejšia úloha je aby vôbec stroj vedel že na vstupe je rečový signál a odkiaľ po kľak sa tento signál v čase nachádza. Existujú rôzne algoritmy ktoré sa opierajú o základne vlastnosti rečových signálov.

Detekcia reči

- Využitie VAD
- Prenos rečových signálov
 - Systém prerušovaného vysielania (napr. GSM)
 - Systém variabilnej rýchlosti prenosu
- Systémy automatického rozpoznávania reči
 - Pre klasickú DTW je to nutnosť
 - Pre HMM VAD znamená: zrýchlenie, spresnenie a zjednodušenie celého procesu
- Zlepšenie kvality rečových signálov
 - Jednokanálový prístup: odhad šumu počas páuz a jeho následná filtrácia
- Analýza rečových signálov

4

Prehľad kde sa takéto algoritmy najčastejšie nachádzajú

Odstránenie šumov

- Šumy

- Aditívne

$$y=x+\text{šum}$$

- Odhad spektra šumu z nerečových časti
- Viacero prístupov
 - „Odčítanie“ priemerného šumu zo spektra zašumeného signálu
 - Periodický signál: priemerovanie periód
 - Existencia viacerých mikrofónov: zašumený a požadovaný(menej zašumený) signál, adaptívna filtrácia

5

Ďalšou úlohou je potlačenie rôznych druhov šumov ktoré sa v reálnom prostredí môžu do užitočného audio signálu dostať. Tu sú spomenuté 2 základné typy šumov- skreslení.

Odstránenie šumov

- Šumy
 - Konvolútorne
 - Vplyv prostredia - „ozveny“
- $Y = \text{conv}(x, \text{šum})$
- Odhad impulzovej odpovede (filtra) miestnosti
 - Viacero prístupov
 - Inverzný filter
 - Podelenie spektra odhadom priemerného spektra signálu
 - Modifikácia príznakov „zvuku“ aby štatistické vlastnosti zodpovedali nejakému prostrediu, modelu

Odstránenie šumov

- Metódy
 - Nezávislé od dát:
 - Zmeny parametrov v reči sú v rozsahu 1-20Hz
 - Filtrácie pre zvýraznenie modulačného spektra reči
 - Aplikuje sa na rečové príznaky: výkonové spektrá, kepstrálne koeficienty
 - Dátovo závislé (adaptívne)
 - Bez učiteľa (neexistujú triedy dát)
 - výkonová normalizácia, odčítanie stredných hodnôt
 - S učiteľom (existujú triedy dát : originál, nové)
 - mapovanie príznakov

7

Na potlačenie týchto šumov sa používa veľa metód, ktoré sa podľa dostupných dát delia do rôznych tried.

Kompresia audio signálov

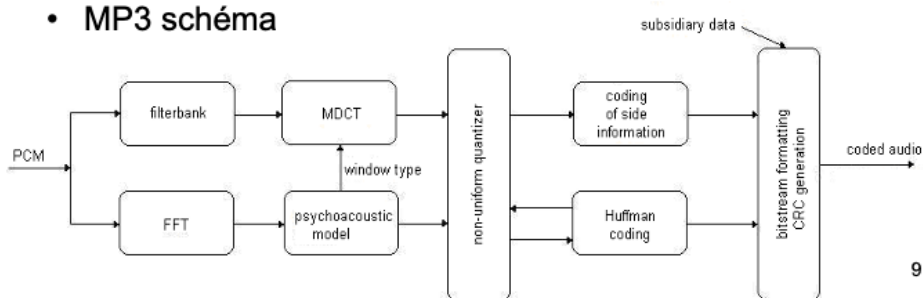
- Úloha: vyjadriť signál pomocou menšieho počtu bitov tak, aby nebol počuť rozdiel v kvalite, resp. nebol horší ako nastavená kvalita
- Kompakcia – menší počet technických bitov, ich lepšie využitie, nedochádza k strate informácie
- Kompresia – dochádza k strate informácie, ale aby bola dodržaná istá kvalita
- Meranie kvality
 - Subjektívne metódy
 - Objektívne metódy
- Audio- hudba: aby nebol rozdiel bežne počuteľný- rušivý
- Reč (telefónia): aby bola zachovaná zrozumiteľnosť

8

Veľmi významnou skupinou aplikácií pre rečové ako aj audio rozhranie sú algoritmy na kompresiu reči a kompresiu audio signálov. Tieto postupy značne (napr. Pre audio-hudbu až cca 12 krát) redukujú množstvo dátového priestoru na uchovanie a prenos týchto signálov. V nasledovnom sú spomenuté niektoré ich delenia a ich vlastnosti.

Kompresia audio signálov

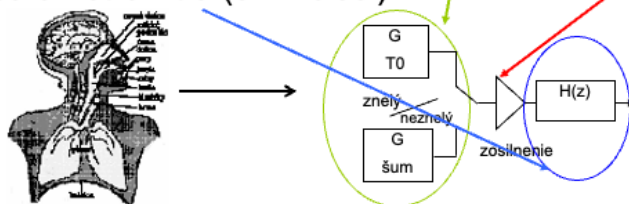
- Audio
- Je založená na využívaní psychoakustických princípov:
 - Absolútny prah počuteľnosti, frekvenčné maskovanie, časové maskovanie, nelineárne- psychoakustické frekvenčné mierky, nelineárne (log) vnímanie hlasitosti
- Najvýznamnejšie postupy: MP3, OGG, WMA; dosahujú cca 12 násobnú kompresiu
- MP3 schéma



Základom pre kompresiu audio signálov (hudba) je využitie nedokonalosti ucha teda odstránenie tých častí signálu ktoré ucho nevie vnímať. Vyžívajú sa tu poznatky z psychoakustiky najmä fenomén frekvenčného maskovania.

Kompresia audio signálov

- Reč
- Je založená na odhade parametrov lineárneho modelu produkcie reči
- Parametre: budenie (periodické/ "šum"), zosilnenie, parametre filtra (8-12 čísel)



- Najvýznamnejšie postupy: CELP, ACELP, ... Používajú sa vo VoIP a mobilných telefónoch
- Prenosová rýchlosť cca 9kbit/s až 800bit/s (hranica zrozumiteľnosti)

10

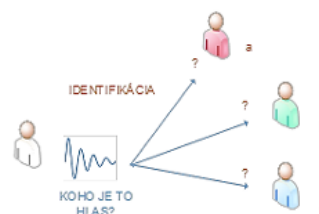
Pri kompresii rečových signálov, napr. Telefónia, sa naopak využíva známy model ako je reč produkovaná hlasovými organmi človeka. Potom parametre takéhoto modelu sú posielane, resp. uchovávané a nie samotný signál. Ten sa pomocou modelu a týchto parametrov vie jednoducho znova vygenerovať. Tým sa dosahuje ešte výraznejšia kompresia.

Rozpoznávanie rečníka

- Rozpoznávanie rečníka

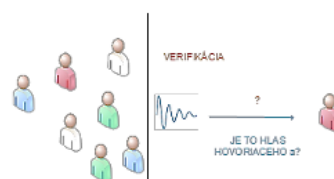
- Identifikácia rečníka

- Kto zo skupiny hovorí?
 - Uzavretá skupina
 - Nekooperujúci používatelia



- Verifikácia rečníka

- Hovorí ten kto to tvrdí?
 - Otvorená množina
 - Kooperujúci používatelia



11

Často krát je potrebné na rečovom rozhraní riešiť problém identity hovoriacich. To riešia aplikácie z oblasti rozpoznávanie rečníka. Tie sa ďalej delia na problém identifikácie rečníka a autentifikácie. Pri identifikácii je potrebné určiť kto zo skupiny rozpráva a pri autentifikácii je potrebné potvrdiť autentitu daného jedinca na základe jeho hlasu. To je možné preto lebo rečový signál nesie okrem lexikálnej informácie aj biometrickú informáciu (zdedenú aj naučenú)

Rozpoznávanie rečníka

- Problémy spojené s identifikáciou
 - Prítomnosť aditívnych šumov
 - Prítomnosť konvolútorých šumov
 - Nahrávacie zariadenie
 - Prostredie záznamu
 - Rozdielne podmienky tréovania a nasadenia
 - „session variability“
 - Zmena vlastností hlasu
 - Chcená: zmena spôsobu rozprávania
 - Nechcená: aktuálny stav (zdravotný, emocionálny)

12

V ďalšom sú rozobrané problémy, výhody a nevýhody využívania aplikácií rozpoznávania rečníka, ako aj rôzne spôsoby ich realizácie.

Rozpoznávanie rečníka

- Použitie systémov rozpoznania rečníka
 - Identifikácia rečníka
 - Kriminalistika
 - Automatické indexovanie rečových databáz (tzv. meta dáta)
 - Personifikácia služieb, nastavení, atď.
 - Verifikácia rečníka
 - Prístup k zabezpečeným údajom pomocou jednoduchej a neinvazivnej biometrie

Rozpoznávanie rečníka

- Podľa požadovaného vstupu
 - Textovo závislé
 - Vyžadujú presný text - frázu
 - Najpresnejšie (až 99% úspešnosť)
 - Menšie použitie
 - Využívajú systémy rozpoznania reči (HMM, DTW)
 - Textovo nezávisle
 - Neočakávajú konkrétny text
 - Úspešnosti 70-95%
 - Širšie a jednoduchšie použitie

Rozpoznávanie rečníka

- Výhody
 - Ľahko získaný biometrický signál
 - Neinvazívna metóda
 - Obsahuje viacero špecifických informácií o užívateľovi (fyzické, naučené)
- Nevýhody
 - Veľká variabilita rečníka
 - Nestálosť parametrov
 - Prítomnosť šumov (rozdiel medzi tréňovaním a nasadením)
 - Ľahká reprodukcia (nahrávacie zariadenia)

15

Rozpoznávanie rečníka

- Existuje veľa príznakov rôznych vlastností
- Delenie:
 - Akustické
 - Prozodické
 - Príznačky vyššej úrovne

Rozpoznávanie rečníka

- Akustické príznaky (1. úroveň)
 - Odrážajú fyzické vlastnosti hlasových orgánov
 - Ľahko sa počítajú na úsekoch cca 10-30ms
 - Sú náchylné na zmeny prostredia a hovoriaceho (kompenzácia prostredia)
 - Najčastejšie
 - Obálka spektra

Rozpoznávanie rečníka

- Prozodické príznaky (2. úroveň)
 - Charakterové vlastnosti (naučené)
 - Dynamika, prízvuk, intonácia, rýchlosť, pauzy v reči, atď.
 - Ťažšie sa extrahujú a kvantifikujú
 - Počítajú sa z dlhších časových intervalov
 - Najúspešnejšie sú založené na vyhodnocovaní priebehu hlasivkovej periódy

Rozpoznávanie rečníka

- Príznamy vyššej úrovne (3. úroveň)
 - Jazyk
 - Dialekt
 - Slovná zásoba
 - Využívajú metódy rozpoznávania reči
 - Nie sú tak náchylné na vplyv prostredia
 - Málo diskriminujúce

Rozpoznávanie reči

- Automatický prepis rečového signálu do textu
- Vstup rečový signál (PCM vzorky)
- Výstup: prepis čo bolo povedané, text



20

Asi najzložitejšou a najpraktickejšiu úlohou pre spracovanie reči je automatické rozpoznanie reči, teda úloha aby počítač vedel prepísať (do formy textu) čo bolo povedané. Pozor to ešte nie je úloha aby významu aj porozumel a na to aj vhodne reagoval. Na to sú iné aplikácie ktoré spadajú do oblasti umelej inteligencie.

Rozpoznávanie reči

- Sú rozpoznané slová len zo slovníka
- Postupnosť slov je obmedzená zvolenou gramatikou
- Úloha nerieši problém porozumenia významu neseného v jazykovej informácii
 - umelá inteligencia (vstup text po rozpoznaní)

21

Zvyčajne takéto systémy majú rôzne očakávania čo sa týka vstupnej reči a obmedzenia na spôsob činnosti. Môžu preto existovať veľmi jednoduché systémy až po veľmi komplexne riešenia využívajúce veľké množstvo výpočtových a pamäťových prostriedkov. Preto sú j tie najzakladanejšie aspekty ďalej spomenuté ako aj spôsoby delenia systémov automatického rozpoznávania reči.

Rozpoznávanie reči

- Existuje veľa aplikácií, systémov a postupov riešenia
- Každé so svojimi predpokladmi, obmedzeniami na činnosť, formou vstupu a výstupu
- Majú rôzne zložitosti: systémové, časové, finančné,...

Rozpoznávanie reči

- Na ich korektné porovnanie je potrebná ich kategorizácia, zavedenie tried-kategórii
- Rôzne kritéria na ich delenie:
 - veľkosť slovníka
 - typ vstupnej reči
 - čas odpovedi systému
 - atď.

Rozpoznávanie reči

- Veľkosť slovníka
 - Veľmi dôležité vzhľadom na:
 - úspešnosť
 - čas výpočtu
 - zdroje systému
 - Rozlišuje sa:
 - Malý ~ (desiatky až stovky) slov
 - Stredný ~ (stovky až tisíce) slov
 - Veľký ~ sto tisíce slov

Pozn. počty sa menia vzhľadom na pokroky v technológiách

Rozpoznávanie reči

• Typ vstupnej reči	Úrovně zložitosti
• Rozpoznanie izolovaných slov – Na vstupe je len jedno slovo	– ľahká
• Diktované prehovorenie – Slova idú za sebou, ale je medzi nimi dostatočne dlhá oddeľujúca pauza	
• Kontinuálna reč – Slová idú plynulo za sebou, predpokladá sa však obmedzenie na prostredie a typ reči	– Stredná
• Prirodzená reč – Bežná konverzácia: šumy, chyby výslovnosti, artefakty rečníka, opakovanie slov, atď.	– Najťažšia

25

Rozpoznávanie reči

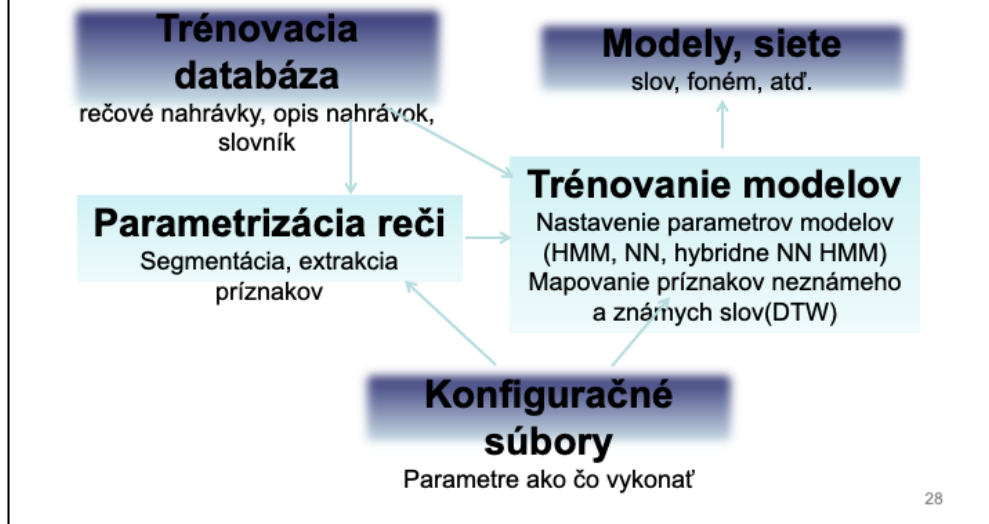
- Závislosť na hovoriacom
 - Závislé na hovoriacom
 - Vyššia úspešnosť
 - Menšia robustnosť pri zmene hlasu
 - Vhodný kompromis: adaptácia všeobecných modelov (natrénované na veľkom množstve hovoriacich) na konkrétneho užívateľa (existuje malé množstvo tréovacích dát)
 - Nezávislé na hovoriacom
 - Potenciálne menej presné
 - Majú väčšiu robustnosť

Rozpoznávanie reči

- Odozva systému
 - V reálnom čase
 - Tolerované oneskorenie cca sekunda/y
 - Je len jeden prechod cez všetky hypotézy
 - Zobrazujú sa postupne len čiastočné hypotézy; ešte pred skončením prehovorenia
 - Neoptimálne riešenie, menšie úspešnosti
 - „Off line“
 - Nemajú stanovenú časovú odozvu
 - Môžu čakať na skončenie prehovorenia
 - Môžu brať do úvahy všetky (viacero) hypotézy
 - potenciálne presnejšie
 - Môžu vykonať viacero prechodov cez sieť hypotéz⁷

Rozpoznanie reči

– Fáza tréovania (klasické systémy)

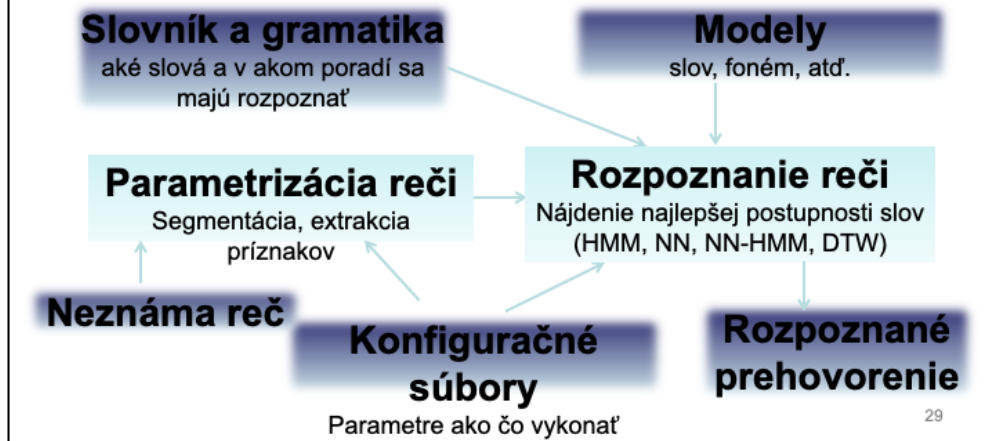


Pred samotným nasadením systému rozpoznávania reči je ho najprv potrebné natréovať na danú reč, prípadne prostredie a rečníka. Preto existujú dve fázy a to fáza tréovania a fáza nasadenia (rozpoznávania) systému. V závislosti od komplexnosti systému fáza tréovania predstavuje najzložitejšiu úlohu tak z pohľadu teoretického, t.j. nájdenie vhodného modelu reči a jeho následne nastavenie na reálnych dátach, tak z pohľadu praktickej realizácie, množstvo zložitých výpočtov, dát, parametrov, nastavení. Existuje viacero fundamentálnych metód rozpoznávania reči, napr. jednoduché porovnávanie vzorov (DTW), štatistické modelovanie reči (HMM) s modelovaním jazykogramatiky, metódy založené na rôznych typoch neurónových sietí až po tzv. end 2 end systémy. Každá z nich vyžaduje iné vstupné dáta, nastavenia a spôsoby tréovania. Všeobecná schéma tréovania je zobrazená na nasledovnom grafe.

Rozpoznanie reči

– Rozpoznávanie (recognition)

- Nájdenie najväčšej zhody, resp. pravdepodobnosti medzi neznámym prehovorením a vhodnou postupnosťou modelov slov



Pri fáze rozpoznávania sa už použije natrénovaný model a na základe neho sa zo vstupného signálu vyberie najpravdepodobnejšia hypotéza, ktorá zodpovedá danému prehovoreniu. To ako sa to konkrétne robí závisí od metódy. Všeobecná schéma fázy rozpoznávania je zobrazená na nasledovnom obrázku.

Rozpoznávanie reči

- DTW (dynamic time warping)
 - Rozpoznávanie izolovaných slov
 - Porovnáva slovo zo slovami zo slovníka
- HMM – Skryté Markovové modely
 - Spojitá reč
 - Štatistická technika modelovania dynamických systémov
 - Každá reči je modelovaná “HMM” modelmi

30

V ďalšom sú v stručnosti spomenuté najznámejšie používané postupy; niektoré sú už skôr historické.

Rozpoznávanie reči

- **Hybridné neurónové siete a HMM**
 - NN akustické modelovanie, HMM časove
 - najbežnejšie a najlepšie výsledky
 - Komplikované- znalosti z viacerých oblastí
- **Neurónové siete**
 - Najnovšie, veľa dát, dlhé tréovanie
 - Jednoduchšie použitie bez hlbších znalostí (jazyk, fonológia, modely- jazyka, gramatiky)
 - End 2 end systémy: vstup signál alebo príznaky a výstup text bez ďalších medzistupňov, napr. Transformery, conformery, attention-encoder decoder, rekurentné siete

-31

Syntéza reči

- Vstup: text, výstup: prehovorenie (PCM signál)
- Požiadavky
 - Zrozumiteľnosť (nutná)
 - Prirodzenosť (veľmi ťažké)
 - Možnosť voľby dodatočných parametrov:
 - Typ hlasu, rýchlosť prehovorenia, emocie, atď.

32

Ďalšou veľkou oblasťou na audio rozhraní je syntéza reči. Teda umele generovanie prehovorenia (akustická forma) počítačom. Využíva sa vtedy keď počítač prezentuje užívateľovi informáciu formou reči. V ďalšom sú spomenuté základné požiadavky na syntézu, jej vlastnosti, rôzne spôsoby realizácie so svojimi výhodami a nevýhodami.

Syntéza reči

- Existuje viacero prístupov
 - Syntéza spájaním jednotiek
 - Nájdenie a spojenie vhodných časti reči z nahovorenej databázy
 - najviac prirodzene znejúca syntetizovaná reč
 - Typy
 - Výber jednotiek: fonémy, slabiky, slová, vety. Spájanie jednotiek z databázy tak, aby plynulo nadväzovali, veľká databáza
 - Difónová: len dvojice foném, prechod medzi nimi je v strede segmentu, malá databáza
 - Syntéza na špecifickej oblasti: z nahratých slov a fráz vytvára reč. Používa sa v aplikáciách s limitovaný rečovým výstupom na určitú oblasť

33

Syntéza reči

– Formantová syntéza

- Nepoužíva vzorky ľudskej reči
- Reč vytvorená zo zvukových modelov
- Parametre: základná hlasivková frekvencia, počet tvar a poloha rezonančných frekvencií (vrcholov spektra), časy trvania
- Dobrá kontrola nad vlastnosťami reči
- Problémy aby to znelo prirodzene

Syntéza reči

– Artikulačná syntéza

- Je založená na modeli ľudského hlasového traktu a artikulácie

– HMM syntéza

- Parametre ako frekvenčné spektrum, základná hlasivková frekvencia, prozódia reči sú modelované pomocou HMM
- Priebeh reči je generovaný pomocou HMM a maximálnej pravdepodobnosti generovania parametrov za predpokladu textu

Ostatné aplikácie

- Rozpoznávanie emócií z reči
 - Call centra, odozvy, reakcie zákazníkov
 - Medicína: detekcia porúch správania
 - Starostlivosť napr. o pacientov, hendikepovaných, starých
 - Využívajú najmä prozodické vlastnosti:
 - dynamiku reči, tempo reči, pauzy v reči, moduláciu, atď..
 - Pracujú na dlhších intervaloch 1-5 s

36

Okrem týchto základných aplikácií existuje viacero menej častých resp. novo skúmaných a nasadzovaných aplikácií. Tieto budú jednotne spomenuté v tomto odseku. Zaujímavá a aj potrebná úloha je detegovať emócie z reči. Jej výhody, nevýhody, a problémy s ňou spojené sú uvedené v ďalšom texte.

Ostatné aplikácie

- Problémy:

- viacero emócií naraz
- Rôzny prejav emócií- kultúrny, individuálny
- Predstieranie, resp. hranie emócií: trénovacie databázy
- Rozlišujeme cca 5 až 7 základných emócií: neutrálna, hnev, prekvapenie, radosť, odpor, nuda
- Niekedy sa presne nemeria konkrétna emócia ale skôr trieda emócií

- Rozpoznávanie udalosti v audio signáloch

- Rozbite sklo, výstrel, krik, plač
- Automatické bezpečnostne dohľadové systémy (kamery všetko nevidia, zvuky sa šíria ďalej aj do nepokrytých miest)

37

Ďalšou dôležitou úlohou ktorá sa bude postupnejšie masívnejšie nasadzovať je detekcia resp. rozpoznanie rôznych udalosti zo zvukových signáloch. Jej využitie bude najmä v monitorovacích a bezpečnostných systémoch.

Ostatné aplikácie

- Monitorovacie systémy:
 - Živé, neživé a umelo vytvorené procesy
 - Či výrobné procesy prebiehajú tak ako majú
 - výskyt prírodných dejov, zvierat, atď..
 - Využívajú sa aj dlhšie časové intervaly: trvanie daného deju
 - Širšie frekvenčné pásma ako len reč, napr. 20- 20kHz, niekedy aj viac
 - Vzhľadom na rôznosť dejov sa môžu sledovať rôzne parametre signálu
 - Momentálne najlepšie systémy využívajú NN

Ostatné aplikácie

- Rozpoznávanie prostredia z audio signálu
 - Z nahrávky zistiť prostredie kde to bolo nahraté
 - Istým spôsobom je to pod úloha detekcie udalosti
- Detekcia príbuzenstva
 - Určiť zo zvukových – rečových nahrávok jedincov ktorí z nich sú v príbuzenskom vzťahu
 - Využívajú sa nie len biologické- fyzické podobnosti ale aj naučené správanie, ako sa ľudia v spoločnom vzťahu podobne vyjadrujú, atď..
 - Kvôli komplikovanosti úlohy, kde nie je úplne zrejmé čo skúmať a vyhodnocovať sa s výhodou používajú NN prístupy

39

Ďalšie aplikácie môžu byť rozpoznávanie prostredia, určenie príbuzenstva užívateľov na základe reči, vkladanie skrytej informácie do rečových resp. audio signálov (watermarking), atď.

Ďakujem za pozornosť

juraj.kacur@stuba.sk