

Sebastian Dziadzio

Sprawozdanie z laboratoriów HTK

1. Opis projektu i zastosowanej gramatyki

Projekt polegał na stworzeniu systemu rozpoznawania mowy, który w uproszczony sposób symulował sterowanie samochodem. Zaimplementowane funkcje to: zwiększanie i redukcja biegu, zmiana prędkości o zadaną wartość, skręt o zadany kąt, włączanie lub wyłączenie świateł, silnika i tempomatu. Podstawowym kryterium podczas projektowania gramatyki była jej naturalność pod względem językowym. W tym celu konieczne było na przykład ograniczenie dokładności przy zmianie prędkości czy skręcie do pięciu jednostek. Udało się jednak otrzymać gramatykę nie wymagającą od użytkownika specjalnych przygotowań do pracy z systemem. Komendy są intuicyjne i poprawne językowo, przykładowo: „przyspiesz (trzydzieści kilometrów/mil)”, „zwiększ/zredukuj bieg”, „(wyłącz/włącz) światła drogowe”, „skręć (dwadzieścia stopni) w lewo”.

2. Opis nagrań

Nagrania sporządzono w pokoju o powierzchni 15m² za pomocą mikrofonu wbudowanego w kamerkę internetową Logitech C310, wyposażonego w funkcję redukcji szumów. Hałas tła był znikomy, najgłośniejszym słyszalnym zakłóceniem był pracujący komputer stacjonarny. Rejestrowano mowę ciągłą z przerwami pomiędzy kolejnymi wypowiedziami. Następnie przekonwertowano pliki na format mono, podpróbkowano do 16 000 Hz i zapisano jako waveform audio file. Operacje te wykonano z użyciem programu Audacity. Początkowy łączny czas nagrań treningowych to 3 minuty i 26 sekund. Ze względu na stosunkowo małą skuteczność działania programu, inkorporowano nagrania pochodzące z innego systemu. Wówczas czas nagrań treningowych wzrósł do 6 minut i 41 sekund. Nagrania testowe sporządzono w tych samych warunkach. Nagrania w folderach *test*, *test_snr20*, *test_snr30* i *test_snr40* pochodzą od tego samego mówcy co nagrania treningowe, przy czym w trzech ostatnich folderach zostały one sztucznie zniekształcone za pomocą pakietu Matlab poprzez dodanie szumu białego o stosunku sygnału do szumu wartości odpowiednio 20, 30 i 40 dB. Folder *test_inny_mowca* zawiera nagrania sporządzone przez innego mówcę.

3. Wyniki testów

Początkowo system przetestowano za pomocą dwudziestu czterech nagrań testowych pochodzących od tego samego mówcy co nagrania treningowe. Wykonano 20 reestymacji i nie uwzględniono modelu ciszy. System nie działał poprawnie, rozpoznanie słów było znikome (poniżej 10%). Po zamodelowaniu ciszy (jako osobnego fonemu sp), RR% dla słów wzrósł do około 20%. Ze względu na tak niski wynik, zdecydowano się na wcielenie nagrań treningowych pochodzących z innego projektu (dzięki uprzejmości Krzysztofa Drewnianego), czego śladem jest plik *phones0_dr.mlf* w folderze *mlf*. Wbrew oczekiwaniom, nie wpłynęło to znacznie na skuteczność systemu. Podjęto próbę zmiany parametrów *word insertion penalty* oraz *grammar scale factor*, jednak nawet znaczne zmiany ich wartości pozostawały bez wpływu na działanie programu, więc zaprzestano na kilkunastu próbach. System nie działał nawet gdy testowano go bezpośrednio z wykorzystaniem nagrań treningowych,

więc skonstruowano go od nowa. Zrezygnowano z wykorzystania dodatkowych nagrań treningowych. Wykryto przy tym nieznaczne błędy w gramatyce oraz nieprawidłowe parametry w nagraniach testowych (ich częstotliwość próbkowania wynosiła 44100 Hz zamiast 16000 Hz). Po podpróbkowaniu oraz zmianie rozdzielczości bitowej, RR% dla słów przekroczył 30%. Kluczowa dla poprawy skuteczności działania programu była zmiana liczby reestymacji. Najprawdopodobniej doszło do przetrenowania modeli, gdyż okazało się, że przyjęcie danych otrzymanych po dwóch reestymacjach skutkuje znacznie lepszym działaniem programu niż wykorzystanie modeli otrzymanych po piętnastu czy dwudziestu reestymacjach. W ten sposób udało się otrzymać wynik 77%:

```
C:\Users\Sebastian\Desktop\Sebastian\Studio\TH\HTK\Projekt\Projekt>HResults -I nlf/testref.nlf txt/monophones0.txt nlf/recout.nlf
===== HTK Results Analysis =====
Date: Mon Feb 18 01:50:36 2013
Ref : nlf/testref.nlf
Rec : nlf/recout.nlf
----- Overall Results -----
SENT: %Correct=54.17 [H=13, S=11, N=24]
WORD: %Corr=77.78, Acc=76.19 [H=49, D=6, S=8, I=1, N=63]
=====
```

Współczynnik rozpoznania dla zdań jest stosunkowo wysoki (54%), co wynika z faktu, że zdania w systemie są bardzo krótkie, składają się maksymalnie z pięciu-sześciu słów. Łącznie wykonano testy z użyciem 120 nagrań testowych, o różnych parametrach akustycznych. Zostało to bardziej szczegółowo omówione w podpunkcie piątym.

4. Analiza błędów rozpoznania

Najczęstsze zaobserwowane błędy w działaniu systemu to zamiana wyrazów o podobnej strukturze fonetycznej lub identycznej funkcji gramatycznej, np. włącz-wyłącz, trzydzieści-czterdzieści, pięćdziesiąt-sześćdziesiąt, zwolnij-przyspiesz, lewy-prawy. Można poradzić sobie z tym problemem na etapie projektowania gramatyki, poprzez unikanie wyrazów o podobnym brzmieniu, jeśli mają pełnić w zdaniu tę samą funkcję (na przykład *włącz* można zastąpić słowem *uruchom*). Niestety zmniejsza to naturalność gramatyki oraz intuicyjność interfejsu głosowego, która jest przecież jedną z jego podstawowych zalet.

5. Analiza różnych rozwiązań

Jak napisano w punkcie trzecim, łączna liczba nagrań testowych to 120: 24 nagrania mówcy A, 24 nagrania mówcy A z szumem o SNR 20 dB, 24 nagrania mówcy A z szumem o SNR 30 dB, 24 nagrania mówcy A z szumem o SNR 40 dB oraz 24 nagrania mówcy B. Poniżej zaprezentowano wyniki oraz ich zgrubną analizę w zależności od liczby reestymacji.

a) SNR 20dB

```
C:\Users\Sebastian\Desktop\Sebastian\Studio\TH\HTK\Projekt\Projekt>HResults -I nlf/testref_snr20.nlf txt/monophones0.txt nlf/recout_snr20.nlf
===== HTK Results Analysis =====
Date: Mon Feb 18 03:04:07 2013
Ref : nlf/testref_snr20.nlf
Rec : nlf/recout_snr20.nlf
----- Overall Results -----
SENT: %Correct=8.33 [H=2, S=22, N=24]
WORD: %Corr=14.29, Acc=14.29 [H=9, D=15, S=39, I=0, N=63]
=====
```

b) SNR 30 dB

```
C:\Users\Sebastian\Desktop\Sebastian\Studia\TH\HTK\Projekt\Projekt>HResults -I nlf/testref_snr30.nlf txt/nonophones0.txt nlf/recout_snr30.nlf
===== HTK Results Analysis =====
Date: Mon Feb 18 03:12:13 2013
Ref : nlf/testref_snr30.nlf
Rec : nlf/recout_snr30.nlf
----- Overall Results -----
SENT: %Correct=8.33 [H=2, S=22, N=24]
WORD: %Corr=17.46, Acc=17.46 [H=11, D=15, S=37, I=0, N=63]
=====
```

c) SNR 40 dB

```
C:\Users\Sebastian\Desktop\Sebastian\Studia\TH\HTK\Projekt\Projekt>HResults -I nlf/testref_snr40.nlf txt/nonophones0.txt nlf/recout_snr40.nlf
===== HTK Results Analysis =====
Date: Mon Feb 18 03:11:11 2013
Ref : nlf/testref_snr40.nlf
Rec : nlf/recout_snr40.nlf
----- Overall Results -----
SENT: %Correct=12.50 [H=3, S=21, N=24]
WORD: %Corr=20.63, Acc=14.29 [H=13, D=13, S=37, I=4, N=63]
=====
```

d) Mówca B

```
C:\Users\Sebastian\Desktop\Sebastian\Studia\TH\HTK\Projekt\Sebastian_Dziedzic\HTK>HResults -I nlf/testref_inny_mowca.nlf txt/nonophones0.txt nlf/recout_inny_mowca.nlf
===== HTK Results Analysis =====
Date: Mon Feb 18 15:30:23 2013
Ref : nlf/testref_inny_mowca.nlf
Rec : nlf/recout_inny_mowca.nlf
----- Overall Results -----
SENT: %Correct=33.33 [H=8, S=16, N=24]
WORD: %Corr=63.49, Acc=55.56 [H=40, D=0, S=23, I=5, N=63]
=====
```

Dla zaszumionych nagrań najlepszy współczynnik rozpoznania osiągnięto dla trzech reestymacji. Przy dziesięciu, piętnastu i dwudziestu skuteczność systemu gwałtownie spadała i wynosiła poniżej dziesięciu procent. Dla obcego mówcy optymalną liczbą reestymacji było sześć. Dla dziesięciu, piętnastu i dwudziestu skuteczność nieznacznie spadała (do około 50%), natomiast dla dwóch, trzech i czterech wynosiła poniżej 20%.

6. Wnioski

Działanie systemu w bardzo dużym stopniu zależy od warunków akustycznych. Ze względu na ewentualne zastosowanie tego typu oprogramowania, skuteczność rozpoznania zdań powinna wynosić co najmniej 95% dla różnych mówców, nawet w złych warunkach akustycznych (duża ilość zakłóceń). Odstęp sygnału od szumu rzędu 40 dB jest zapewne trudny do uzyskania we wnętrzu samochodu (pracujący silnik, hałas aerodynamiczny, muzyka, inni mówcy) a nawet dla tak "komfortowych" warunków akustycznych skuteczność systemu wynosi zaledwie 20%. W idealnych warunkach udało się uzyskać 77 % rozpoznawalności słów. Przy tak mało zróżnicowanym zbiorze treningowym (zaledwie jeden mówca) dość zaskakujący był stosunkowo niewielki spadek efektywności przy zmianie użytkownika. Należy jednak zaznaczyć, że mówca B to osoba o podobnej skali głosu, akcencie i intonacji oraz tej samej płci (brat mówcy A). Do stworzenia systemu niezależnego od mówcy potrzebny byłby dużo bardziej reprezentatywny korpus audio. Główną

przyczyną niewielkiej rozpoznawalności (zaledwie 77%) jest zapewne bardzo zła jakość nagrań treningowych. Jak okazało się podczas drugiego podejścia do stworzenia systemu, prawie wszystkie nagrania treningowe są zniekształcone (doszło do przesterowania mikrofonu). Również sprzęt rejestrujący nie był zbyt dobrej jakości.

Wyrażam zgodę na dołączenie moich nagrań do korpusu mowy AGH. Nagrania mogą być odtwarzane ale wyłącznie bez podawania tożsamości mówców (na przykład w celu prezentacji jakości, rodzaju nagrań itd.)