

To Dr. G. van der Merwe, my mentor in this project.

Recurrent neural network-enhanced HMM speech recognition systems

by

Jan Willem Frederik Thirion

Submitted in partial fulfillment of the requirements for the degree

Master of Engineering (Electronics)

in the

Faculty of Engineering

UNIVERSITY OF PRETORIA

October 2001



Speech recognition using recurrent neural networks

Thesis submitted in partial fulfillment of the requirements for the degree of

Master of Science in Computer Science

Abstract

To Carina, who makes every moment in life special.

Keywords: speech recognition, speech recognition, recurrent neural networks, encoder-decoder, recurrent neural networks, hidden and convolutional networks, hidden state models.

Abstract

Speech recognition is fast becoming an attractive form of communication between humans and machines, due to the recent advances made in the field. The technology is however, far from being acceptable when used in its untrained form. Many attempts have been made to overcome some of the major difficulties involved in the automatic recognition of human speech by machine.

The dissertation attempts to provide better performance. The first contribution is an enhanced segmentation of speech signals, where speaker information is used on higher level layers to improve performance. The second is to introduce the hidden state segmentation layer and its variants. This method is to improve the performance and the accuracy of a state-of-the-art hidden Markov model (HMM) based recognition system using the segmentation information as new features to incorporate speaker information into the Viterbi decoding process or posterior probability. The third contribution technique outperforms other attempts in which the segmentations problems are solved.

The suggested solution results in a better overall performance of around 90% and is called an enhanced segmentation approach (ESA). The enhancement is implemented in current speech recognition that only uses spectrograms from the short-term analysis, namely a fixed window of the target syllable, and the long-term global features.

speech vectors. The problem of the choice of the fixed size window is thus eliminated, at the cost of difficult real-time implementation. The BRNN can be trained using a modified form of backpropagation through time (BPTT). The input to the BRNN is the entire speech sentence and the two BRNN outputs are the probability that a boundary between phonemes occurs in that frame of speech, as well as the probabilities of each phoneme occurring. The sequence of phonemes can then be recovered by scanning the speech signal into phonemes using only this output signal, and no higher level lexical knowledge such as the sequence of utterances.

Keywords: speech segmentation, speech recognition, recurrent neural networks, extended recurrent neural networks, bi-directional recurrent neural networks, hidden Markov models

Speech recognition is fast becoming an attractive form of communication between humans and machines, due to the recent advances made in the field. The technology is, however, far from being acceptable when used in an unlimited form. Many attempts have been made to overcome some of the many difficulties faced in the automatic recognition of human speech by machine.

This dissertation attempts to accomplish two ambitious goals. The first is the reliable, automatic segmentation of speech signals, in a speaker independent manner, where no higher level lexical knowledge is used in the process. The system is limited to segmentation in an off-line manner. The second is to improve the phoneme recognition accuracy of a state-of-the-art hidden Markov model (HMM) based recognition system, using the segmentation information. A new technique of incorporating segmentation information into the Viterbi decoding process is presented, and it is shown that this technique outperforms other attempts to include the segmentation probabilities.

The segmentation system consists of a bi-directional recurrent neural network (BRNN), also called an extended recurrent neural network (ERNN). In contrast to conventional recurrent neural networks, that only use speech vectors from the past, present and possibly a fixed window of the future, BRNNs use all of the past, present and future



speech vectors. The problem of the choice of the fixed size window is thus eliminated, at the cost of difficult real-time implementation. The BRNN can be trained using a modified form of backpropagation through time (BPTT). The input to the BRNN is the entire speech sentence and the two BRNN outputs are the probability that a boundary between phonemes occurs in that frame of speech, as well as the probability that no boundary occurs. The segmentation system segments the speech signal into phonemes, using only the speech signal, and no higher level lexical knowledge such as the sequence of phonemes.

The recognition system can incorporate the segmentation probabilities in one of two ways. The first is to modify the HMM transition probabilities by combining the HMM transition probabilities and the BRNN outputs. The second method, developed in this dissertation, involves the use of an adaptive word (phoneme) transition penalty. Previously, only a fixed transition penalty was used between words (phonemes). By making the transition penalty adaptive (based on segmentation information), the phoneme recognition performance can be significantly improved. It is also shown that the adaptive word transition penalty outperforms the HMM transition probability modification technique, used by others.

All of the experiments used in this dissertation are conducted on the TIMIT database, in order to provide a convenient way to compare the results to that of others in the field. The hidden Markov toolkit (HTK) from Cambridge University is used for all phoneme recognition experiments.

soos die volledige spraak gespreek is. Die vorm van 'n segmenteerstelsel wat die grootste tussen-silens en tussen-silens tydsonde maak, sou ook die waarskynlikheid dat geen groter verlies sou skep. Die segmentering moet respondeer op die spraak. In sommige dinge word die spraaksein gebruik te maak. Geen hoor-wiel behoeft daarom meer as die enigste van funksie, word gebruik nie.

Uittreksel

Die uitgangspunt van hierdie verhandeling was om 'n betroubare spraaksegmentering stelsel te ontwerp wat die spraaksegmentering van die spraakherkenning en die terugvoer van neurale netwerke saamvou. Daarvan uitgaande dat VMM's een goed voorbeeld was om te gebruik om die goedkoopste en mees effektiewe te handhaaf, was die tweede doelwit om die uitgebreide terugvoer neurale netwerke, bidireksionele terugvoer neurale netwerke, ver-skulde Markov modelle en die Viterbi dekoder te kombineer.

Sleutelwoorde: spraaksegmentering, spraakhertkenning, terugvoer neurale netwerke, uitgebreide terugvoer neurale netwerke, bidireksionale terugvoer neurale netwerke, ver-skulde Markov modelle

Spraakherkenning is tans een van die mees populêre intervlakke tussen mens en masjien. Talle probleme word egter nog steeds met hierdie tegnologie ondervind, veral wanneer dit in 'n onbeperkte vorm gebruik word.

Hierdie verhandeling pak twee ambisieuse doelwitte aan. Die eerste is die betroubare, outomatiese segmentering van spraakseine. Die stelsel word beperk tot 'n aflyn taak. Die tweede doelwit is om die foneemherkenning akkuraatheid van 'n hoë verrigting ver-skulde Markov model (VMM) gebaseerde herkenning stelsel te verbeter, deur van die segmentering inligting gebruik te maak. Hier word 'n nuwe tegniek, om die segmentering inligting in die Viterbi dekodering proses in te sluit, voorgestel en daar word gewys dat hierdie tegniek beter resultate lewer as ander metodes wat tans gebruik word.

Die segmenteringstelsel bestaan uit 'n bidireksionale terugvoer neurale netwerk (BTNN), ook 'n uitgebreide terugvoer neurale netwerk (UTNN) genoem. Konvensionele terugvoer neurale netwerke gebruik slegs spraakvektore van die verlede, hede en moontlik vaste venster van die toekoms. In kontras hiermee, gebruik BTNN'e al die inligting van die verlede, hede en die toekoms. Die probleem van die keuse van 'n vaste grootte venster word dus vermy ten koste van 'n moeilike intydse implementasie. Die BTNN kan aferig word met 'n aangepaste weergawe van terugvoering deur tyd (TVDT). Die inset na die

BTNN is die volledige sin met spraak en die twee BTNN uitsette is die waarskynlikheid dat 'n grens tussen foneme in daardie raam van spraak voorkom, asook die waarskynlikheid dat geen grens voorkom nie. Die segmenteringstelsel segmenteer die spraak in foneme deur slegs van die spraaksein gebruik te maak. Geen hoër vlak leksiese kennis, soos bv. die volgorde van foneme, word gebruik nie.

Die herkenningsstelsel kan die segmentering waarskynlikhede in een van twee maniere insluit. Die eerste is om die VMM oorgang waarskynlikhede aan te pas deur die VMM oorgang waarskynlikhede en BTNN uitsette te kombineer. Die tweede metode, ontwikkel in hierdie verhandeling, maak van aanpasbare woord (foneem) oorgang penalisasie gebruik. Voorheen was slegs van vaste oorgang penalisasie tussen woorde (foneme) gebruik gemaak. Deur die oorgang penalisasie term aanpasbaar te maak (gebaseer op segmentering inligting), kan die foneemherkenning akkuraatheid aansienlik verhoog word. Daar word ook gewys dat aanpasbare woord oorgang penalisasie die VMM oorgang waarskynlikheid aanpassing tegniek, soos deur ander gebruik, oortref.

Al die eksperimente wat in hierdie verhandeling gedoen word, word op die TIMIT databasis gedoen, om sodoende 'n geriflike manier daar te stel vir die vergelyking van resultate met dié van ander in die veld. Die stel van verskuilde Markov model programme (HTK) van Cambridge Universiteit word vir alle foneemherkenning eksperimente gebruik.

Acknowledgements

I would like to thank Professor E.C. Botha for allowing me the opportunity to study under her. She provided me with the academic guidance that made the work given here possible. It was an honour to learn from such a great intellectual.

Special thanks goes to Darryl Purnell for all the advice and discussions we had throughout the completion of my studies. His invaluable comments guided me towards achieving something that I am proud of. I also thank Willie Smit, with whom I had many discussions about concepts and ideas.

I am grateful to my family and friends, for their love, encouragement, and patience. I have great appreciation for my parents who coped with me, and always encouraged and supported me in my studies.

Finally, I would like to give my deepest thanks to my Saviour, Jesus Christ, for giving me the talents and gifts to be able to think.

2 Theory

2.1 Special agent of change

2.1.1 Impression

2.1.2	Pronunciation rules	25
2.1.3	Blessing	27
2.1.4	Signal processing	29

Contents

2.2	Speech processing	31
2.3	Speech segmentation and recognition	33
1	Introduction	1
1.1	Problem statement	3
1.2	Summary of related work	4
2	Theory	5
2.1	Speech signal processing	5
2.2	Speech segmentation	8
2.3	Speech recognition	16
1.3	Approach and research hypotheses	18
1.4	Contributions of this study	20
2.4	Dissertation outline	21
2	Theory	22
2.1	Speech signal processing	22
2.2	Speech segmentation	24
2.3	Speech recognition	26

2.1.2	Pre-emphasis filter	25
2.1.3	Blocking	27
2.1.4	Signal bias removal	28
2.1.5	Windowing	29
2.1.6	Spectral analysis	30
2.1.7	Energy and mel cepstrum coefficients	31
2.1.8	Normalisation	34
2.1.9	Delta energy and cepstrum coefficients	36
2.2	Hidden Markov models	37
2.2.1	Basic elements and problems of an HMM	37
2.2.2	The Viterbi decoding process	41
2.2.3	Use of a language model	44
2.2.4	Use of a word transition penalty	46
2.3	Recurrent neural networks	47
2.3.1	Conventional recurrent neural networks	47
2.3.2	Bi-directional recurrent neural networks	52
3	Speech segmentation	56
3.1	Preprocessing	56

3.2 Segmentation	57
3.2.1 Hidden Markov models	59
3.2.2 Recurrent neural networks	61
3.3 Postprocessing	62
3.4 Accuracy measure	63
4 Speech recognition	69
4.1 Recognition (baseline)	69
4.2 Recognition using segmentation information	71
4.2.1 HMM transition probability modification	72
4.2.2 HMM word transition penalty modification	75
4.3 Accuracy measure	79
5 Experiments	80
5.1 Experiment 1: Speech segmentation	80
5.1.1 Hidden Markov models	81
5.1.2 Recurrent neural network	88
5.2 Experiment 2: Speech recognition (baseline)	98
5.2.1 Number of HMM mixtures	99

5.2.2	Language model	99
5.2.3	Word transition penalty	104
5.2.4	Combined language model and word transition penalty	104
5.3	Experiment 3: Speech recognition using segmentation information . . .	105
5.3.1	HMM transition probability modification	106
5.3.2	HMM word transition penalty modification	109
5.3.3	Combined transition probability and word penalty modification	110
6	Summary and conclusion	113
6.1	Summary of results	113
6.2	Statistical significance	119
6.3	Conclusion	121
6.4	Shortcomings and future work	123
A	Recurrent neural network training	125
A.1	Gradient descent training	125
A.2	Backpropagation through time	126
A.3	BRNN training equations	127
A.3.1	Overview of the training process	128

A.3.2	Forward pass	129
A.3.3	Backward pass	130
A.3.4	Gradient calculation	131
A.3.5	Update of weights	132
A.3.6	Generalisation capability	132
and RN – Bidirectional recurrent neural network		
Bibliography	<i>etc.</i> recurrent neural network	133

DC	Direct current
DP	Dynamic programming
DTW	Dynamic time warping
ERNN	Extended recurrent neural network
FBDYN	Forward-backward auditory dynamic experiments
FFT	Fast Fourier transform
HMM	Hidden Markov model
HTK	Hidden Markov Toolkit
KL	Kullback-Leibler
LP	Linear prediction
LPC	Linear predictive coding
MAP	Maximum a posteriori
MFC	Mel frequency cepstrum
MPOC	Mel frequency cepstrum coefficient
MIT	Massachusetts Institute of Technology
MLP	Multi-layer perceptron
NIST	National Institute of Standards and Technology
NN	Neural networks
OGI	Oregon Graduate Institute
PDF	Probability density function

List of Abbreviations

	Pronunciation
ADC	Analogue to digital converter
AM	Auditory modeling
AR	Autoregressive
ASP	Auditory speech processing
ASR	Automatic speech recognition
BCD	Bayesian changepoint detector
BRNN	Bi-directional recurrent neural network
CRNN	Chaotic recurrent neural network
DC	Direct current
DP	Dynamic programming
DTW	Dynamic time warping
ERNN	Extended recurrent neural network
FBDYN	Forward-backward auditory dynamic cepstra
FFT	Fast Fourier transform
HMM	Hidden Markov model
HTK	Hidden Markov toolkit
KL	Kullback Leibler
LP	Linear prediction
LPC	Linear predictive coefficient
MAP	Maximum a posteriori
MFC	Mel frequency cepstrum
MFCC	Mel frequency cepstrum coefficient
MIT	Massachusetts Institute of Technology
MLP	Multilayer perceptron
NIST	National Institute of Standards and Technology
NN	Neural network
OGI	Oregon Graduate Institute
PDF	Probability density function

PSM	Polynomial segment model
PSS	Pseudosyllables
RASTA	Relative spectral
RMLP	Recurrent multilayer perceptron
RMS	Root mean square
RNN	Recurrent neural network
S/M-R	Synchrony/mean-rate
SCVQ	Sequence-constrained vector quantisation
SDU	Semantic dialogue unit
SLAM	Segmentation and labeling automatic module
SRN	Simple recurrent network
SVF	Spectral variation function
TI	Texas Instruments
VCV	Vowel-consonant-vowel, or vocal-consonant-vocal
VQ	Vector quantisation

efficient way of storing and retrieving them.

The automatic recognition of features and the analysis of speech patterns, which difficult problem, is often considered to be the core task of the automatic continuous speech, and the representation of the patterns, from the speaker to the environment, is far more broad and complex than the analysis and perception, and to the realiser that produces the general form of speech patterns of speech, and consequently, applying the technical solution to analyse them. In particular, the use of only one speaker instead of many speakers (referred to as speaker independent systems), the use of a small number of large vocabulary, and a well structured dialogue between user and to solve the task of an easily implemented system. By realising the limitations of the system to facilitate the analysis of speech application specific, such speech recognition becomes an important factor in its performance.

Most speech recognition systems use some form of parametric model. That is, we