

Visualisierung kontinuierlicher, multimodaler Schmerz Scores am Beispiel akustischer Signale

Masterarbeit

Franz Anders
HTWK Leipzig

Januar 2017

Abstract

Inhaltsverzeichnis

1	Einleitung	6
2	Grundlagen der medizinischen Schrei-Forschung	7
2.1	Schmerz Scores	7
2.2	Schmerz-Schrei aus medizinischer Sicht	9
2.3	Physio-Akustische Modellierung des Weinens	10
3	Grundlagen der Signalverarbeitung	13
3.1	Statistische Merkmale von Signalen	13
3.2	Faltung	13
3.3	Transformationen	13
3.3.1	Fourier-Transformation	13
3.3.2	Laplace-Transformation	13
3.3.3	Z-Transformation	13
3.4	Filter	13
3.5	akustische Modellierung der menschlichen Stimme	13
3.6	Feststellung von Periodizität in Signalen	13
3.6.1	Zero-Crossing-Rate	13
3.6.2	Methoden des Frequenzbereiches	13
3.6.3	Autokorrelation	13
3.6.4	Cepstrum	13
4	Grundlagen der Klassifikation	14
4.1	Gütemaße binärer Klassifikatoren	14
4.2	Klassifikationsalgorithmen	15
4.2.1	CARD	15
4.2.2	Andere	15
5	System zur Visualisierung akustischer Schmerz-Scores	16
5.1	Architektur	16
5.2	Preprocessing	18
5.3	Voice Activity Detection	20
5.3.1	Windowing	20
5.3.2	Feature Extraction	21
5.3.3	Thresholding	29
5.3.4	Markierung der Cry-Units	33
5.3.5	Decision Smoothing	35
5.4	Segmentierung	38
5.5	Feature-Extraction	41
5.6	Ableitung der Schmerz-Scores	41

Inhaltsverzeichnis

5.7 Visualisierung	41
6 Zusammenfassung	42
Appendices	46

Abbildungsverzeichnis

2.1	Veranschaulichung des Grundvokabulars	11
4.1	Confusion-Matrix	14
5.1	Architektur des Systems	17
5.2	Parameter eines Audio-Kompressors	19
5.3	Ergebnis des Preprocessings	19
5.4	Markierung von Schreigeräuschen im Audiosignal	20
5.5	Features des Zeitbereiches	21
5.6	Features des Frequenz-Bereiches	23
5.7	Entstehung des Cepstrums	24
5.8	Cepstrum	25
5.9	Features des Cepstrums	25
5.10	Autokorrelation eines Signals	26
5.11	Features der Autokorrelation	27
5.12	Differenz-Signal des RMS-Wertes	28
5.13	Thresholding eines Feature-Signales	29
5.14	Zusammenfassung klassifizierter Signalfenster zu Cry-Units	34
5.15	Beziehung zwischen agrenzenden Segmenten	35
5.16	Klassifizierung vor dem Decision Smoothing	36
5.17	Klassifizierung vor und nach dem Decision Smoothing	37
5.18	Ergebnis der Segmentierung	38
5.19	Mögliche Segmentierungen eines Signals	38
.1	Boxplot-Auswertung über Sensitivity, Specificity und Accuracy der beiden VAD-Modelle	48

1 Einleitung

2 Grundlagen der medizinischen Schrei-Forschung

2.1 Schmerz Scores

Bei erwachsenen Menschen wird der Schmerzgrad typischerweise durch eine Selbsteinschätzung des Patienten unter der Leitung gezielter Fragen des Arztes vorgenommen. Bei Kindern unter 3 Jahren ist diese Selbsteinschätzung nicht möglich. Schmerz drückt sich in Veränderungen des psychologischen, körperlichen und biochemischen Verhaltens des Säuglings aus. Die für den Arzt am leichtesten feststellbaren Verhaltensänderungen sind von außen wahrnehmbare Merkmale, wie zum Beispiel ein Verkrampfen des Gesichtsausdrucks, erhöhte Körperbewegungen oder lang anhaltendes Weinen. Um eine weitestgehend objektive Schmerzfeststellung zu ermöglichen, wurden sogenannte *Pain-Scores* entwickelt, die durch ein Punktesystem den gesamten Schmerzgrad des Babies quantifizieren.[23] Es existieren *eindimensionale* Pain-Scores, die den Schmerz nur Aufgrund der Beobachtung eines Merkmals beurteilen, so wie beispielsweise die reine Beurteilung des Gesichtsausdrucks. *Mehrdimensionale* (auch *multimodale*) Pain-Scores beziehen mehrere Faktoren in das Scoring mit ein.[1]. Tabelle 2.1 zeigt das Scoring-System „Neonatal Infant Pain Scale“ (NIPS) als Beispiel für eine multimodale Pain-Score. Der Säugling wird anhand der aufgeführten Kategorien bewertet und alle vergebenen Punkte aufsummiert. Ein insgesamt Wert von > 3 zeigt Schmerz an, ein Wert von > 4 großen Schmerz.[10]

Tabelle 2.1: NIPS-Scoring

NIPS	0 points	1 point	2 points
Facial Expr.	Relaxed	Contracted	-
Cry	Absent	Mumbling	Vigorous
Breathing	Relaxed	Different than basal	-
Arms	Relaxed	flexed/stretched	-
Legs	Relaxed	flexed/stretched	-
Alertness	Sleeping	uncomfortable	-

In den meisten mehrdimensionalen Scoring-Systeme werden die Schreieräusche mit einbezogen. Tabelle 2.2 zeigt eine Übersicht über eine ausgewählte Menge an multimodalen Pain-Scores. Alle Pain-Scores sind für Kleinkinder bis 3 Jahren gedacht. In der Übersicht wird nicht wiedergegeben, welche weiteren Merkmale jeweils in das Scoring mit einbezogen werden, oder welche Ingesamtpunktzahlen auf welche Schmerzintensität hinweisen. Es soll an dieser Stelle nur verdeutlicht werden, welche unterschiedlichen Ansätze zur Bewertung des Schreiens aus medizinischer Sicht im Zusammenhang mit Pain-Scores existieren. Folgende Beobachtungen lassen sich aus der Übersicht ziehen:

1. Die zu beobachtenden Eigenschaften des Weinens werden mit subjektiv behafteten Werten charakterisiert. Beispielsweise wird im N-PASS-System ist ein Schmerz-Schrei als „High-pitched or silent-continuous crying“ beschrieben. Es wird nicht fest definiert, was als „crying“ gilt oder welche Tonhöhe als „high-pitched“ ist. Auch die Erstquellen geben keine festen Definitionen.
2. Es gibt verschiedene Ansätze zur Bewertung des Weinens. Bei CRIE ist die Tonhöhe, bei BIIP die Länge und bei COMFORT die Art des Weinens entscheidend.
3. Die Beschreibungen sind kurz und prägnant gehalten, der Arzt hat in keinem der Modelle auf mehr als drei Parameter des Schreiens zu achten. Die Begründung liegt darin, dass bei allen Modellen a.) das Schreien nur eines von mehreren Faktoren ist, und b.) Die Schmerzbestimmung in einem vorgegebenen Zeitrahmen durchführbare sein muss.

System	P.	Description
FLACC[30]	0	No cry (awake or asleep)
	1	Moans or whimpers; occasional complaint
	2	Crying steadily, screams or sobs, frequent complaints
N-PASS[24]	-2	No cry with painful stimul
	-1	Moans or cries minimally with painful stimuli
	0	Appropriate Crying
	1	Irritable or Crying at Intervals. Consolable
	2	High-pitched or silent-continuous crying. Not consolable
BIIP[8]	0	No Crying
	1	Crying <2 minutes
	2	Crying >2 minutes
	3	Shrill Crying >2 minutes
CRIES[3]	0	If no cry or cry which is not high pitched
	1	If cry high pitched but baby is easily consoled
	2	If cry is high pitched and baby is inconsolable
COVERS[14]	0	No Cry
	1	High-Pitched or visibly crying
	2	Inconsolable or difficult to soothe
PAT[11]	0	No Cry
	1	Cry
DAN[4]	0	Moans Briefly
	1	Intermittent Crying
	2	Long-Lasting Crying, Continuous howl

COMFORT[21]	0	No crying
	1	Sobbing or gasping
	2	Moaning
	3	Crying
	4	Screaming
MBPS[22]	0	Laughing or giggling
	1	Not Crying
	2	Moaning quiet vocalizing gentle or whimpering cry
	3	Full lunged cry or sobbing
	4	Full lunged cry more than baseline cry

Tabelle 2.2: Übersicht über Pain-Scores

2.2 Schmerz-Schrei aus medizinischer Sicht

Die Frage ist: Woher kommen diese unterschiedlichen Bewertungen des Weinens in Tabelle 2.2? Gibt es eine Pain-Score, die aus wissenschaftlicher Sicht „recht hat“? Dieser Fragestellung unterliegen unterliegen zwei grundlegendere Fragen: 1.) Ist es überhaupt möglich, anhand der akustischen Eigenschaften den Grund für den Schrei abzuleiten, also beispielsweise Hunger, Einsamkeit oder Schmerz? Anders formuliert: Gibt es überhaupt so etwas wie einen Schmerz-Schrei? 2.) Ist es möglich, anhand der akustischen Eigenschaften den Schweregrad des Unwohlseins abzuleiten (also beispielsweise den Grad des Schrei-Versursachenden Schmerzes)?

Die Annahme, dass es möglich ist, aus dem Schreien den Grund abzuleiten, wird als „Cry-Types Hypothesis“ bezeichnet. Die berühmtesten Befürworter dieser Hypothese ist eine skandinavische Forschungsgruppe, auch bezeichnet als „Scandinavian Cry-Group“, die diese Idee in dem Buch „Infant Crying: Theoretical and Research Perspectives“[2] publik machte. Die Annahme ist, dass die verschiedenen Ursachen *Hunger, Freude, Schmerz, Geburt und Anderes* klare Unterschiede hinsichtlich ihrer akustischen Merkmale aufweisen, welche an einem Spektogramm ablesbar seien. Entsprechende Beispiele werden in dem Buch gegeben. Nur einige Jahre Später zeigte Müller et al [6] in einem Paper, dass bei leichter Veränderung der Bedingungen der Experimente die Unterscheidung nicht möglich ist. Die Gegenhypothese ist, dass Weinen „nichts als undifferenziertes Rauschen“ sei. 50 Jahre später liegt kein anerkannter Beweis für die eine oder andere Hypothese vor. Es gibt nur starke Hinweise dafür, dass die Plötzlichkeit des Eintretens des Schreigrundes hörbar ist. Ein plötzliches Ereignis, wie ein Nadelstich oder ein lautes Geräusch, führen auch zu einem plötzlich beginnenden Schreien. Ein langsam einretendes Ereignis, wie ein langsam immer stärker werdender physischer Schmerz oder langsam eintretender Hunger führen auch zu einem langsam eintretenden Weinen. Da keine Einigung herrscht, wird empfohlen, den Grund aus dem Kontext abzuleiten.[29]

Die Zweite Frage nach der Ableitung der Stärke des Unwohlseins aus den akustischen Eigenschaften des Geschreis wird in der Fachsprache unter dem Begriff *Cry as a graded*

Signal subsumiert. Je „stärker“ das Weinen, desto höher das Unwohlsein (*Level of Distress (LoD)*) des Säuglings. Tatsächlich bemessen wird dabei der von dem Beobachter vermutete Grad des Unwohlsein des Babies, und nicht der tatsächliche Grad, da dieser ohne die Möglichkeit der direkten Befragung des Kindes nie mit absoluter Sicherheit bestimmt werden kann. Dieser vermutete LoD wird entweder durch das subjektive Empfinden der Beobachter oder durch Pain-Scores festgestellt. Ein hohes Level of Distress hat vor allem eine schnelle Reaktion der Aufsichtspersonen zur Beruhigung des Babies zur Folge, womit dem Geschrei eine Art Alarm-Funktion zukommt. Es gibt starke Hinweise darauf, dass das Level of Distress anhand objektiv messbarer Eigenschaften des Audiosignals bestimmt werden kann. So herrscht beispielsweise weitestgehend Einigung darüber, dass ein „lang“ anhaltendes Geschrei auf einen hohen Level of Distress hinweist. Insofern aus dem Kontext des Schreiens Schmerz als wahrscheinlichste Ursache eingegrenzt werden kann, kann aus einem hohen Level of Distress ein hoher Schmerz abgeleitet werden. [29] und [27]

Es herrscht wiederum keine Einigung darüber, welche akustischen Eigenschaften im Detail ein hohes Level of Distress anzeigen. Carlo V Bellieni et al [4] haben festgestellt, dass bei sehr hohem Schmerz in Bezug auf die DAN-Scala (siehe Tabelle 2.2) die Tonhöhe des Schreis steigt. Qiaobing Xie et al [27] haben festgestellt, dass häufiges und „verzerrtes“ Schreien (ohne feststellbares Grundfrequenz, da der Ton stimmlos erzeugt wird) auf einen hohen Level of Distress hinweist.[29] Diese Uneinigkeit hat wahrscheinlich zu den verschiedenen Bewertungen in den Pain-Scores geführt. 2.2.

2.3 Physio-Akustische Modellierung des Weinens

Das Ziel dieses Kapitels ist die Schaffung eines einheitlichen Vokabulares, auf den sich bezogen wird, um das Schreien eines Babys zu beschreiben. Die hier vorgestellten Begriffe stammen sowohl aus dem Buch „A Physioacoustic Model of the Infant Cry “ H Golub und M Corwin [12] als auch aus dem Paper „Rythmic organization of the Sound of Infant Cry “ von Zeskind et al.[25]

Die Lautäußerung eines Neugeborenen, umgangssprachlich auch als „Weinen“ oder „Schreien“ bezeichnet, lässt sich im allgemeinen beschreiben als das „rythmische Wiederholen eines beim ausatmen erzeugen Geräusches, einer kurzen Pause, einem Einatmungs-Geräusch, einer zweiten Pause, und dem erneuten Beginnen des Ausatmungs-Geräusches.“[34].

Das Vokabular, welches insbesondere von H Golub und M Corwin geschaffen wurde, ist sehr umfassend. An dieser Stelle wird eine Auswahl grundlegender Begrifflichkeiten vorgestellt, die in dieser Arbeit gebraucht werden. Sie werden in Abbildung 2.1 veranschaulicht.

Expiration beschreibt den Klang, der bei einem einzelnen, ununterbrochenem Ausatmen mit Aktivierung der Stimmbänder durch das Baby erzeugt wird. [25]. Der von Golub et al [12] verwendete Begriff **Cry-Unit** wird in dieser Arbeit synonym verwendet. Umgangssprachlich ist handelt es sich um einen einzelnen, ununterbrochenen *Schrei*.

Inspiration beschreibt den Klang, der beim Einatmen durch das Baby erzeugt wird.

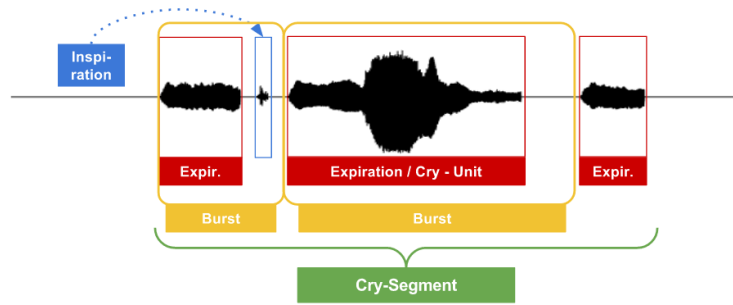


Abbildung 2.1: Veranschaulichung des Grundvokabulars

Burst beschreibt die Einheit von einer Expiration und der darauf folgenden Inspiration. Das heisst, dass die zeitliche Dauer eines Bursts sowohl das Expiration-Geräusch, das Inspiration-Geräusch als auch die beiden Pausen zwischen diesen Geräuschen umfasst. Praktisch ergibt sich das Problem, dass vor allem bei stärkerem Hintergrundrauschen die Inspiration-Geräusche häufig weder hörbar noch auf dem Spektrogramm erkennbar sind. Daher wird die Zeitdauer eines Bursts oder Cry-Unit vom Beginn einer Expiration bis zum Beginn der darauf folgenden Expiration definiert und somit allein von den Expirations auf die Bursts geschlossen. Implizit wird somit eine Inspiration zwischen zwei Expirations angenommen.

Cry die insgesamt klangliche Antwort zu einem spezifischen Stimulus. Eine Gruppe mehrerer Cry-Units.[12] In dieser Arbeit wird ein *Cry* als **Cry-Segment** bezeichnet, um Verwechslungen zu vermeiden.

Weiterhin wurden von H Golub und M Corwin [12] Cry-Units in eine der folgenden drei Kategorien eingeführt:

Phonation beschreibt eine Cry-Unit mit einer „vollen Vibration der Stimmbänder“ mit einer Grundfrequenz zwischen 250 und 700 Hz. Entspricht umgangssprachlich einem Weinen mit einem „klaren, hörbaren Ton“.

Hyper-Phonation beschreibt eine Cry-Unit mit einer „falsetto-artigem Vibration der Stimmbänder“ mit einer Grundfrequenz zwischen 1000 und 2000 Hz. Entspricht umgangssprachlich einem Weinen mit einem „sehr hohen, aber klaren, hörbaren Ton“.

Dysphonation beschreibt eine Cry-Unit ohne klar feststellbare Tonhöhe, produziert durch Turbulenzen an den Stimmbändern. Entspricht umgangssprachlichen dem „Brüllen oder Krächzen“.

Eine Cry-Unit gehört dabei mindestens einer dieser Kategorien an, kann aber auch in seinem zeitlichen Verlauf die Kategorie wechseln. H Golub und M Corwin [12] stellen weiterhin eine Reihe an charakteristischen Eigenschaften vor, die in Bezug auf ein Cry-Segment berechnet werden.

Latency-Period beschreibt die Dauer zwischen dem zufügen eines Schmerz-Stimulus und

dem beginn des ersten Cry-Bursts des Segmentes

Duration beschreibt die insgesamt Zeitdauer des Cry-Segmentes. Es wird keine genaue Definition gegeben, wodurch Beginn und Ende definiert werden. Das Segment endet dort, wo es „scheint, aufzuhören“.

Maximum-Pitch beschreibt die höchste festgetellte Grunfrequenz des Segmentes.

... und viele weitere, die in [12] nachgelesen werden können, aus Platzgründen an dieser Stelle jedoch nicht vollständig genannt werden.

3 Grundlagen der Signalverarbeitung

3.1 Statistische Merkmale von Signalen

3.2 Faltung

3.3 Transformationen

3.3.1 Fourier-Transformation

3.3.2 Laplace-Transformation

3.3.3 Z-Transformation

3.4 Filter

3.5 akustische Modellierung der menschlichen Stimme

3.6 Feststellung von Periodizität in Signalen

3.6.1 Zero-Crossing-Rate

3.6.2 Methoden des Frequenzbereiches

3.6.3 Autokorrelation

3.6.4 Cepstrum

4 Grundlagen der Klassifikation

4.1 Gütemaße binärer Klassifikatoren

Bei einer binären Klassifikation gibt es vier mögliche Klassifikationsergebnisse. Die beiden Klassen werden im allgemeinen als *Positive* und *Negative* beschrieben, was in diesem speziellen Fall *Stimme* und *Stille* entspricht. Eine Klassifikation, bei der ein tatsächliches Positive richtig als Positive vorhergesagt wird, spricht man von einem *True Positive* [TP]. Wird hingegen ein tatsächliches Positive fälschlicherweise als Negative vorhergesagt, spricht man von einem *False-Negative* [FN]. Das System wird entsprechend für die Klassifikation tatsächlicher Negatives angewandt und ergibt. *True-Negatives* [TN] und *False-Positives* [FP]. Die *Confusion Matrix* in Abbildung 4.1 gibt eine Übersicht über die vier möglichen Klassifikations-Ergebnisse.

		Predicted Class	
		Positive	Negative
Real Class	Positive	True-Positive	False-Negative
	Negative	False-Positive	True-Negative

Abbildung 4.1: Confusion-Matrix

Die insgesamt Güte einer Klassifikation wird durch die *Accuracy* nach Formel 4.1 bestimmt. Eine Accuracy von 100% bedeutet, dass *alle* Instanzen richtig klassifiziert werden, eine Accuracy von 50% bedeutet, dass die Hälfte aller Instanzen richtig klassifiziert werden, was der Güte einer rein zufälligen Wahl entspricht.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FN + FP} \quad (4.1)$$

Die Accuracy beziffert die insgesamt Performance des Klassifikators, gibt jedoch keinen Aufschluss darüber, ob der Klassifikator eher eine Tendenz zur falschen Klassifizierung von Positives oder Negatives hat. Bei einer Datenbank mit der selben Anzahl an Positives und Negatives kann eine Accuracy von 50% beispielsweise dadurch entstehen, dass *alle* Instanzen als Positives markiert werden, also sowohl die Positives richtigerweise als Positives, aber die Negatives fälschlicherweise ebenfalls als Positives. Im Umgedrehten Fall ergibt die Klassifizierung aller Instanzen als Negatives ebenfalls eine Accuracy von 50%. In einem dritten Fall irrt sich die Klassifikator gleich oft bei der Einordnung der Negatives und Positives. Die Maße *Sensitivity* und *Specificity* geben Aufschluss über die Güte der Klassifikation hinsichtlich der Positives und Negatives. Die *Sensitivity*, auch bezeichnet als

True-Positive-Rate, bemisst den Anteil tatsächlicher Positives, die auch als solche erkannt wurden, nach Formel 4.2. Eine Sensitivity von 100% bedeutet, dass alle Positives durch den Klassifikator erkannt wurden. Die Erkennungsrate der Negatives hat keinen Einfluss auf die Sensitivity. Eine hohe Sensitivity lässt sich somit „einfach“ erzielen, in dem man *alle* Instanzen immer als Positives klassifiziert. Die Specificity nach Formel 4.3 bestimmt analog zur Sensitivity den Anteil der korrekt als Negatives bestimmten Instanzen. Ein Klassifikator, der alle Instanzen als Positives markiert, hat zwar eine Sensitivity von 100%, aber eine Specificity von 0%. Ergeben zwei verschiedene Klassifikationsmodelle sehr ähnliche Accuracies, hilft die Bestimmung der Sensitivity und Specificity bei der Auswahl des für den Anwendungsfall Adequäteren Klassifikators. So ist beispielsweise bei der Bestimmung von schweren Krankheiten eventuell ein Klassifikator mit höherer Sensitivity wünschbar, um die Wahrscheinlichkeit zu minimieren, dass die entsprechende Krankheit nicht erkannt wird. [19] [26]

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (4.2)$$

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (4.3)$$

4.2 Klassifikationsalgorithmen

4.2.1 CARD

4.2.2 Andere

5 System zur Visualisierung akustischer Schmerz-Scores

5.1 Architektur

Das Ziel dieser Arbeit ist die Ableitung des Schmerz-Scores aus einem Audiosignal sowie die darauf folgende Visualisierung dieser Schmerz-Scores. Um diese komplexe Aufgabenstellung zu lösen, wird sie in einzelne Arbeitspakete zerlegt. Welche Arbeitspakete formuliert werden, ist davon abhängig, aus welcher Perspektive die Aufgabenstellung grundlegend betrachtet wird. So kann sie beispielsweise als Unterproblem der Schrei-Forschung betrachtet werden, woraufhin sie mit Hilfe des in Kapitel 2.3 vorgestellten Vokabulars modelliert wird. Ebenso kann die Aufgabe grundlegend als Problem der automatisierten Spracherkennung verstanden werden, woraufhin ein in diesem Wissenschaftsgebiet üblicher Verarbeitungsprozess entworfen werden würde. Dieser Ansatz wurde beispielsweise in dem von Cohen et al vorgestellten System zur Kategorisierung von Audioaufnahmen von Babys in (0) Nicht-Weinen und (1) Weinen gewählt.[28]

In dieser Arbeit wurde sich für den Ansatz entschieden, die Arbeitspakete aus Sicht der Schreiforschung zu modellieren und dabei das in 2.3 vorgestellte Vokabular zu verwenden. Die einzelnen Arbeitspakete werden daraufhin mit Methoden der Signalverarbeitung gelöst. Hintergrund ist, dass die Arbeitsweise des Systems in Zukunft einfacher Mediziner*innen vermittelbar ist, da die verwendeten Begrifflichkeiten leichter zugänglich sind. Das folgende Vorgehen wird dazu vorgeschlagen:

1. Kontinuierliche Betrachtung des Eingangssignals
2. Markierung der Cry-Units / Expirations in diesem Eingangssignal (und eventuelle Kategorisierung in Phonation, Hyper-Phonation und Dysphonation)
3. Gruppierung der markierten Cry-Units zu Cry-Segments
4. Errechnung der charakteristischen Eigenschaften der Segmente
5. Ableitung der Pain-Score aus den Eigenschaften der Segmente.
6. Visualisierung dieser Pain-Score.

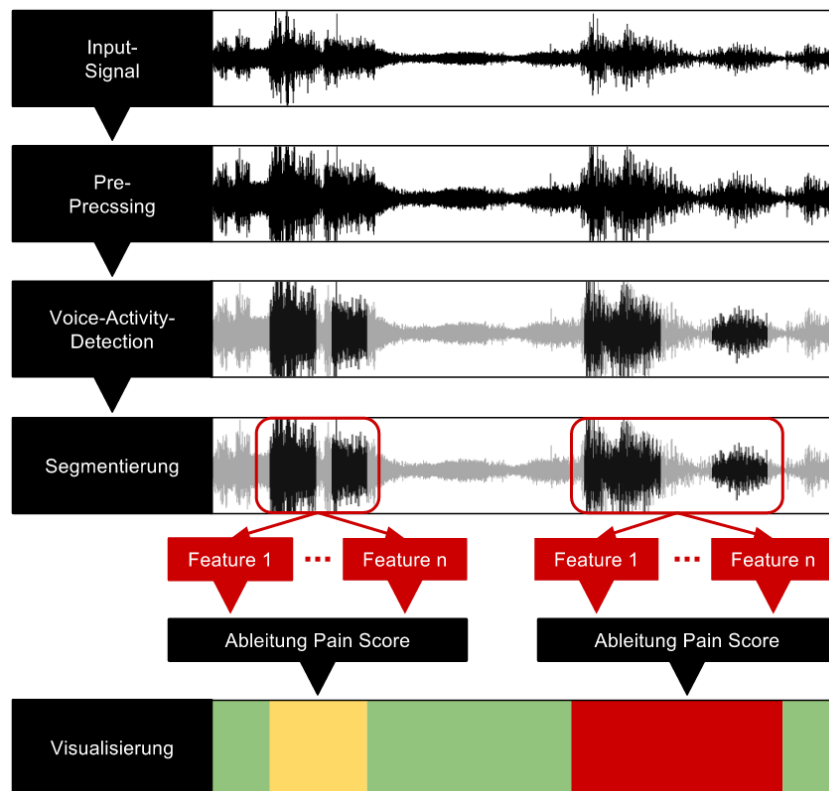


Abbildung 5.1: Architektur des Systems

Diese Arbeitspakete werden nun auf Probleme der automatisierten Spracherkennung abgebildet. Abbildung 5.1 veranschaulicht die Architektur, die sich daraus ergibt. Eingang in das System ist das Audiosignal, welches kontinuierlich in das System gegeben wird. Das Signal passiert daraufhin die folgenden Verarbeitungsschritte:

1. **Pre-Processing.** Das Signal wird vorverarbeitet. Durch die Anwendung von Filtern (linearen wie nicht-linearen) werden die darauf folgenden Verarbeitungsschritte vereinfacht. Das Ergebnis ist ein vorverarbeitetes Signal. Dieser Verarbeitungsschritt wurde zu den Zeiten der manuellen Schrei-Analyse noch nicht mit einbezogen, ist jedoch in der automatisierten Spracherkennung üblich.
2. **Voice-Activity-Detection.** In dem Signal müssen diejenigen Stellen erkannt werden, die die Stimme des Babys enthalten. Dazu wird das Signal in kurze Fenster weniger Millisekunden Länge zerlegt und jedes Fenster als *Stille* (0) oder *Stimme* (1) markiert. Zusammenhängende, als stimmhaft erkannte Signalfenster werden daraufhin zu *Cry-Units* gruppiert.
3. **Segmentierung**, das heißt die Gruppierung der markierten Cry-Units zu *Cry-Segments*.
4. **Feature-Extraction**, das heißt die Berechnung von Eigenschaften für jedes Segment, die für die Ableitung der Pain-Scores von Interesse sind, wie beispielsweise die Duration. Die errechneten Features müssen sich dabei nicht auf diejenigen beschränken, die von H Golub

und M Corwin [12] vorgestellt wurden, da diese zur (vermeintlichen) Unterscheidung der Wein-Ursache, und nicht zur Ableitung einer Pain Score gedacht waren.

5. **Ableitung der Pain Score** aus den errechneten Eigenschaften.

6. **Visualisierung** der errechneten Pain-Score.

Für jedes dieser Module werden verschiedene Lösungsansätze vorgestellt und zumindest ein Lösungsansatz implementiert. Das Pre-Processing wird in Kapitel 5.2 vorgestellt, die Voice-Activity-Detection in Kapitel 5.3, die Segmentierung in Kapitel

5.2 Preprocessing

Beim Preprocessing wird das Signal so vorverarbeitet, dass Störeinflüsse auf die darauf folgenden Verarbeitungsschritte von vorneherein minimiert werden. Welches Pre-Processing durchgeführt wird, ist Abhängig von der konkreten Aufgabenstellung. So werden beispielsweise bei einigen Algorithmen zur Voice-Activity-Detection, also dem markieren stimmhafter Signalabschnitte, Tiefpass, Hochpass- und Bandpassfilter eingesetzt, um diejenigen Frequenzanteile herauszufiltern, die von der Stimme nicht produziert werden können [18] [31] [17]. Bei einigen Pitch-Detection-Algorithmen wird *Centerclipping* eingesetzt, also das 0-Setzen von Samples mit $x[i] < 0.5 \cdot \text{Maximalaussteuerung}$. [7]

In dieser Arbeit wurde sich für eine Vorverarbeitung entschieden, bei der das Signal hinsichtlich seiner Dynamik im Zeitbereich eingegrenzt wird. Dies ist ein typischer Vorverarbeitungsschritt bei Sprachaufnahmen. Hintergrund ist, dass sehr kurz, aber sehr laute Pegelspitzen weit über dem Durchschnittspegel des Gesamtsignals den Maximalwert des Signals unnötig begrenzen und die Signalenergie so gering halten. Da die Testsignale, die in dieser Arbeit verwendet werden, aus inhomogenen Quellen stammen und sehr unterschiedliche Lautstärken haben, wird so gewährleistet, dass sie zumindest ähnliche Energien haben. An dieser Stelle werden (noch) keine Frequenzanteile herausgefiltert, um keine Frequenzen zu verlieren, die in den späteren Verarbeitungsschritten wieder Voice-Activity-Detection 5.3 oder der Feature-Extraction eventuell noch benötigt werden.

Die Dynamikeinschränkung wird mit Hilfe eines Audiokompressor umgesetzt. Ein Audiokompressor verringert Signalspitzen, die über einen festgelegten *Schwellwert* (*Threshold*) liegen, um ein festgelegtes *Verhältnis* (*Ratio*). Ein Threshold von 0.3 mit Ratio von 0.5 bedeutet beispielsweise, dass alle Signalspitzen, die den Wert 0.3 überschreiten oder -0.3 unterschreiten, um 50% verringert werden. Ein Kompressor kann auf die Überschreitung des Thresholds erst nach einer als *Attack* bezeichneten Verzögerung reagieren, und bei erneuten Verlassen des Thresholdes mit einer als *Release* bezeichneten Verzögerung nachwirken. Signalspitzen werden so verringert und die Lautstärke-Dynamik eingeschränkt. Die tatsächliche Erhöhung der Signalenergie geschieht im Anschluss durch die Anhebung der insgesamt Signallautstärke, wie Beispielsweise der Normalisierung des Signals auf den Maximalpegel. Abbildung 5.2 zeigt die Parameter eines solchen Audio-Kompressors.

Der entwickelte Kompressor automatisiert die Einstellung von Threshold und Ratio auf

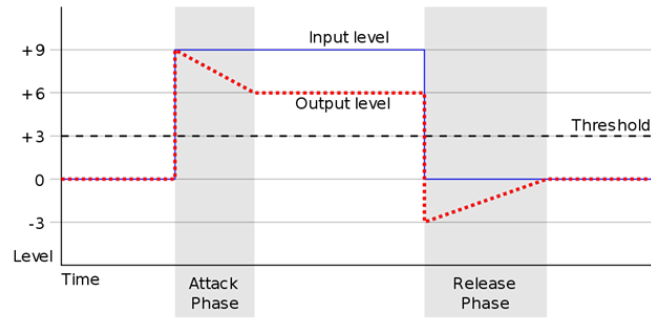


Abbildung 5.2: Parameter eines Audio-Kompressors

Grundlage des Root-Mean-Square (RMS) des Signales x der Länge N . Der RMS-Wert ist ein Maß für die durchschnittliche Signalenergie und wird wie nach Formel 5.1 berechnet. Threshold und Ratio werden nach den Formeln 5.2 und 5.3 berechnet, wobei der Parameter r_a den Ziel-RMS-Wert angibt und mit dem Wert.

$$\text{RMS}(x) = \sqrt{\frac{1}{N} \sum_{n=0}^{N-1} x[n]^2} \quad (5.1)$$

$$\text{THold}(x) = \left[\frac{\text{RMS}(x)}{r_a} \right]^2 \quad (5.2)$$

$$\text{Ratio}(x) = \left[\frac{\text{RMS}(x)}{r_a} \right]^2 \quad (5.3)$$

Abbildung 5.3 zeigt das ein Signal vor und nach dem Preprocessings. Zu sehen ist, dass die Lautstärke der einzelnen Schrei-Einheiten nach der Anpassung einheitlicher ist.

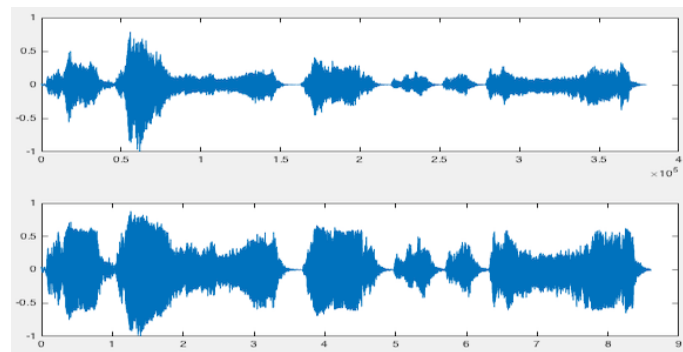


Abbildung 5.3: Ergebnis des Preprocessings

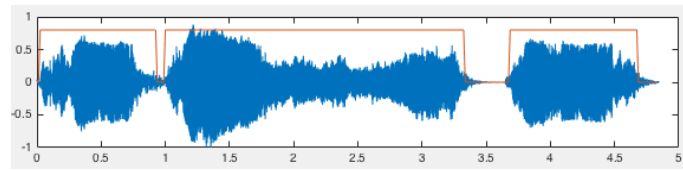


Abbildung 5.4: Markierung von Schreigeräuschen im Audiosignal

5.3 Voice Activity Detection

Das Ziel ist, in einem Audiosignal diejenigen Stellen zu markieren, in denen Schreigeräusche vorliegen. Abbildung 5.4 visualisiert ein Beispiel für eine solche Markierung: Zu sehen ist der Zeitbereich eines Audiosignales mit drei klar erkennbaren Cry-Units. Die Rote Linie, die das Signal überspannt, bildet die Zeiteinheiten des Eingangssignales in die binären Kategorien *Stimmhaft* und *Stille* ab.

Das zu lösende Problem wird als *Voice Activity Detection (VAD)* oder auch *Speech Detection* bezeichnet. Das Ziel ist die Unterscheidung von denjenigen Zeiträumen im Signal, in denen Stimme enthalten ist, von denen, in denen keine Stimme enthalten ist. Zeiträume ohne Stimme werden von nun als „Stille“, und Zeiträume mit Stimme als „Stimmhaft“ bezeichnet. Die größte Herausforderung ist für VAD-Systeme die robuste Erkennung bei Signalen mit Rauschen unbekannter Stärke und Natur.[15]

Der Grundlegende Aufbau eines VAD-Algorithmus ist wie folgt:

1. **Windowing:** Unterteilung des Signals in (einander überlappende) Fenster, typischerweise zweistelligen Millisekunden-Bereich.
2. **Featureextraction** aus den einzelnen Fenstern
3. **Thresholding / Entscheidung** über Präsenz oder Nicht-Präsenz von Stimme für jedes Zeitfenster auf Grundlage der Features mit Hilfe von Entscheidungsregeln wie Grenzwerten.
4. **Decision-Smoothing**, das nachträgliche Hinzufügen oder Entfernen von Entscheidungen mit Hilfe von kontextuellen Informationen der umliegenden Entscheidungen.[15]

Der an dieser Stelle entwickelte Ansatz ist eine Kombination aus den Ideen, die in [20], [33], [18], [31] und [16] vorgestellt wurden.

5.3.1 Windowing

Das Signal x wird in Fenster $x_1 \dots x_n$ unterteilt. Dieser Prozess wird als „Windowing“ bezeichnet. Es wurde sich für die in [18] vorgestellte Fensterlänge von 25 ms entschieden, als Kompromiss zwischen den in [20] empfohlenen 10 ms und den in [31] empfohlenen 40 ms. Die einzelnen Fenster überlappen einander um 50%, das heisst 12.5 ms.

5.3.2 Feature Extraction

Für jedes Signalfenster $x_1 \dots x_n$ à 25 ms werden die folgenden Features aus den Kategorien **Zeitbereich**, **Frequenzbereich**, **Cesptrum** und **Autokorrelation** berechnet.

Zeit-Bereich

Im **Zeit-Bereich** werden die beiden Features *Root-Mean-Square-Wert* [RMS] und *Zero-Crossing-Rate* [ZCR] berechnet.

In [20] wird die Energy eines Zeitfensters für die Voice Activity Detection berechnet. Um den errechneten Wert besser ins Verhältnis zum Signalpegel setzen zu können, wird anstelle der Energie der RMS-Wert nach Formel 5.1 verwendet. Ein höherer RMS-Wert weist auf das Vorhandensein von Stimme im Signalfenster hin, da diese typischerweise eine höhere Energie als das Hintergrundrauschen hat. Als zweites Feature des Zeitbereiches wird die in [31] verwendete *Zero-Crossing-Rate* (ZCR) berechnet. Dabei wird der relative Anteil an Vorzeichenwechseln des Signalfensters x_i nach Formel 5.4 berechnet. Eine höhere ZCR weist auf Stille hin, da Rauschen typischerweise einen höheren ZCR hat als periodische Signale mit tieferer Grundfrequenz.

$$ZCR(x_i[]) = \sum_1^{N-1} |\text{sng}(x_i[n]) - \text{sng}(x_i[n-1])| \quad (5.4)$$

Abbildung 5.5 visualisiert ein Beispielsignal mit Babygeschrei mit einem Hintergrundrauschen mit einem Signal/Rausch-Abstand von 25 dB. Bei dem Signal wurden die stimmhaften Segmente magenta markiert. Für jedes Signalfenster wurden der RMS-Wert und der ZCR-Wert berechnet und als Signale mit im Graphen dargestellt, in dem für den Anfangszeitpunkt eines Zeitfenster der jeweilige berechnete Feature-Wert abgetragen wird. Die beiden Signale wurden so skaliert, dass ihr Maximalwert 1 nicht überschreitet, um ihr Verhalten bezüglich des Vorhandenseins von Stimme klarer erkennbar zu machen.

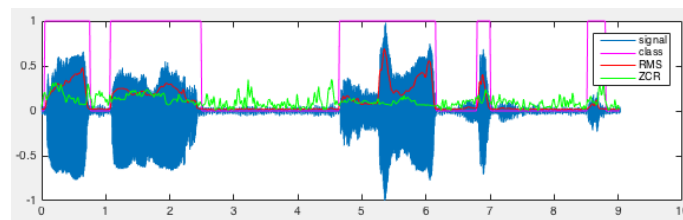


Abbildung 5.5: Features des Zeitbereiches

Frequenz-Bereich

Aus dem **Frequenz-Bereich** werden die drei Features *unnormalisierte spektrale Entropie* [$SEnt_u$], *normalisierte spektrale Entropie* [$SEnt_n$] und *dominanteste frequenzkomponenten*

$[f_{Dom}]$ berechnet.

Als Vorbereitungsschritt wurde die *Kurzzeit-Fourier-Transformation* (engl. *Short-Time-Fourier-Transform* STFT) nach dem Algorithmus von [32] implementiert. So wird jedes Signalfenster x_i mit einem 25 ms-Hamming-Window multipliziert und daraufhin mit einer 2048-Punkte langen *schnellen Fouriertransformation* (FFT) in den Frequenzbereich $F\{x_i\}$ transformiert. Weiterhin wurde für jedes Frequenz-Fenster des Magnitudensignal $X_i = |F\{x_i\}|$ berechnet und das Phasensignal jedes Zeitfensters verworfen, da für keines der Features Phaseninformationen von Interesse waren. Das Ergebnis dieser Transformation aus den Signalfenstern in Zeitbereich $x_1 \dots x_n$ sind die Frequenzfenster $X_1 \dots X_n$.

In [33] wird die *spektrale Entropie* für die Voice-Activity-Detection berechnet. Dabei wird das Spektrum des Frequenzfensters X_i als Wahrscheinlichkeitsverteilung betrachtet. Die *normalisierte spektrale Entropie* wird nach der Formel 5.6 berechnet. Der Wert px_i ergibt sich durch die Normalisierung des N -Punkte langen Spektrums nach Formel 5.5. Neben der in [33] vorgestellten normalisierten spektralen Entropie wird zusätzlich die *unnormierte spektrale Entropie* nach Formel 5.7 berechnet. Bei dieser wird das Spektrum nicht normalisiert, das heißt, es gilt $px_i[f] = X_i[f]$. Somit hat die Energie des Signals einen größeren Einfluss auf den berechneten Wert. Bei der unnormierten spektralen Entropie ist zu erwarten, dass Signalfenster mit Stimme eine höherer Spektrale Entropie haben als Fenster mit Stille. Bei der normalisierten spektralen Entropie ist zu erwarten, dass Fenster mit Stille eine höhere Entropie haben als Fenster mit Stimme.

In die Berechnungen wurden nur Frequenzen von 250 - 8000 Hz betrachtet, da nach [27] die Grundfrequenz der Stimme eines Babies mindestens 250 Hz beträgt, und nach [16] Stimme keine wichtigen Informationen oberhalb von 8000 Hz überträgt.

$$px_i[f] = \frac{X_i[f]}{\sum_{k=1}^N X_i[k]} \quad (5.5)$$

$$SEnt_n(X_i) = - \sum_{k=1}^N px_i[f] \cdot \log(px_i[k]) \quad (5.6)$$

$$SEnt_u(X_i) = - \sum_{k=1}^N X_i[k] \cdot \log(X_i[k]) \quad (5.7)$$

Weiterhin wird die in [20] vorgestellte *dominanteste Frequenzkomponente* berechnet. Für jedes Frequenzfenster X_i wird diejenige Frequenz nach Formel 5.8 berechnet, welches die höchste Magnitude hat. Ein stimmhaftes Signal hat typischerweise eine höhere f_{Dom} als ein nicht stimmhaftes Signal.

$$f_{Dom}(X_i) = \arg \max_f \{X_i[f]\} \quad (5.8)$$

Abbildung 5.6 visualisiert diese Features für das selbe Eingangssignal aus 5.5. Die Features wurden wie bei Abbildung 5.5 beschrieben für die Darstellung skaliert

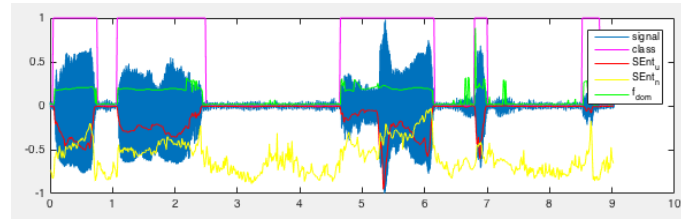


Abbildung 5.6: Features des Frequenz-Bereiches

Cepstrum

Das **Cepstrum** wird definiert als die inverse Fourier-Transformation des Logarithmus des Magnituden-Signals nach Formel 5.9.

$$c[n] = F^{-1}\{\log |F\{x[n]\}|\} \quad (5.9)$$

Die Idee wird anhand von Abbildung 5.7 verdeutlicht. Zu sehen ist ein typisches absolutes Spektrum eines Signalfensters mit Stimme $|F\{x\}|$. Durch die Logarithmisierung des Spektrums wird die Dynamik der Amplitudenunterschiede eingegrenzt. Das entstandene Signal $\log |F\{x[n]\}|$ zeigt klar die Lage der Formanten, dargestellt durch die durchgezogene, schwarze Linie. Die vielen, kleinen Signalspitzen, die sich entlang der Formantenlinie bewegen, sind die harmonischen Obertöne, die von dem periodischen Signal der Stimmbänder erzeugt wurden. Der Grundgedanke ist, dieses Spektrum wiederum als Signal zu interpretieren. Wäre dieses Signal im Zeitbereich (obwohl es eigentlich den Frequenzbereich darstellt), könnte es durch die Amplituden-Modulation eines Quasi-periodischen Signales entstanden. Das höherfrequente Signal (Spektrum der Stimmbänder) wäre dementsprechend der Träger, und das niederfrequente Signal (Spektrum der Formanten) der Modulator. Um diese beiden Signalanteile voneinander zu trennen, würde man eine weitere Fourier-Transformation implementieren (welche an dieser Stelle mathematisch betrachtet einer inversen Fourier-Transformation gleichkommt, da das Phasensignal entfernt wurde). Man erwartet in dem so entstehenden Spektrum $F^{-1}\{\log |F\{x[n]\}|\}$ eine Spitze im niederfrequenten Bereich, bedingt durch den niederfrequenten Modulator, und eine Spitze im höheren Frequenzbereich, bedingt durch den höherfrequenten Träger.

Der Bereich dieser „Fouriertransformation der Fouriertransformation“ wird als *Cepstrum* bezeichnet. Cepstrum ist ein Wortspiel, welches durch die Umkehrung der ersten vier Buchstaben des Wortes *Spectrum* entsteht. Die Unabhängige Variable des Cepstrum folgt dem Wortspiel und wird als *Quefrequency* bezeichnet. Damit wird verdeutlicht, dass die unabhängige Variable des Cepstrum zwar mathematisch betrachtet die Zeit darstellt, jedoch als Frequenz interpretiert wird.

Insbesondere das Vorhandensein einer Spitze im Quefrequency-Bereich ab 3 ms im Cepstrums weist auf das Vorhandensein von harmonischen Obertönen eines periodischen Eingangssignals hin.

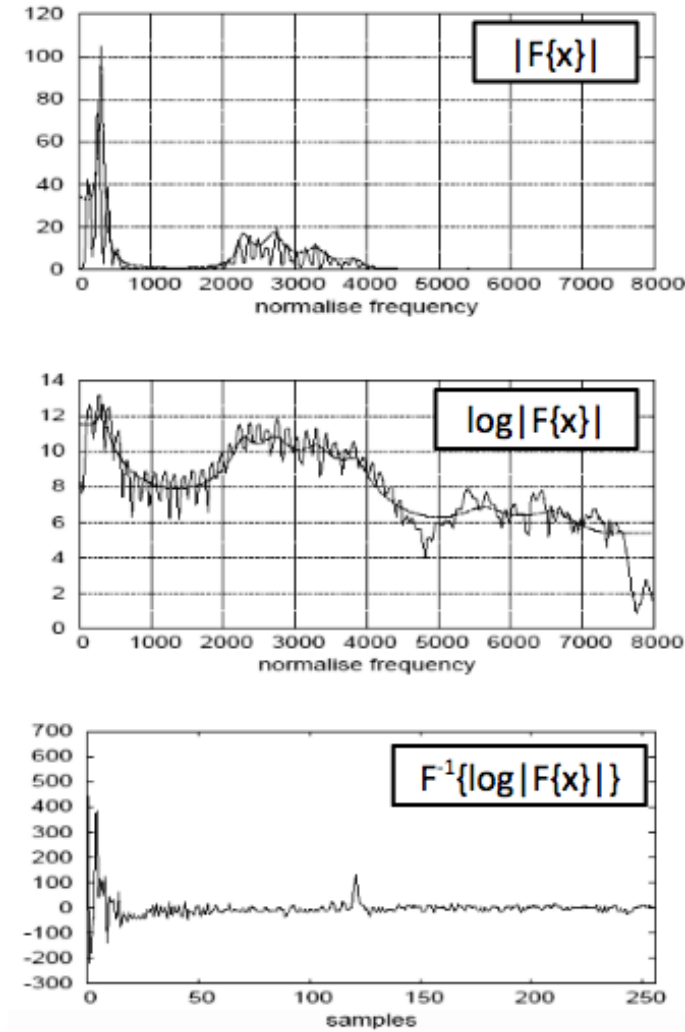


Abbildung 5.7: Entstehung des Cepstrums

hin, wie es von Stimmbändern erzeugt wird. Eine Grundfrequenz von f_0 dieser periodischen Schwingung entspricht einer Signalspitze bei der Queffrequency von $\frac{1}{f_0}$. Abbildung 5.8 zeigt das Log-Spektrum und das Cepstrum eines Sprachsignales, welches in Fenster à 50 ms zerlegt wurde. Die Frames 1-7 sind stimmlos, die Frames 8-15 stimmhaft. Die Stimmhaftigkeit spiegelt sich klar im auftauchen einer Signalspitze bei einer Queffrequency von 10 ms wieder.[9]

Die Informationen des Cepstrum eignen sich ebenfalls zur Voice-Activity-Detection. Es werden die Features *Upper Cepstrum Peak* [$Ceps_{mag}$] und *Upper Cepstrum Peak Location* [$Ceps_{loc}$] berechnet.

Wie in [31] und [33] beschrieben, wurde die *höchste Magnitude im oberen Queffrequency-Bereich* (Upper Cepstrum Peak) nach Formel 5.10 berechnet. Nach [27] liegt die Grundfrequenz eines Kinderweins zwischen 250 und 2000 Hz, was einem Queffrequency-Bereich von 5 - 40 ms entspricht. Folglich werden bei der Berechnung nach Formel 5.10 nur Queffrequency-Werte in diesem Bereich betrachtet werden. Eine hoher $Ceps_{mag}$ -Wert weist auf das Vorhanden-

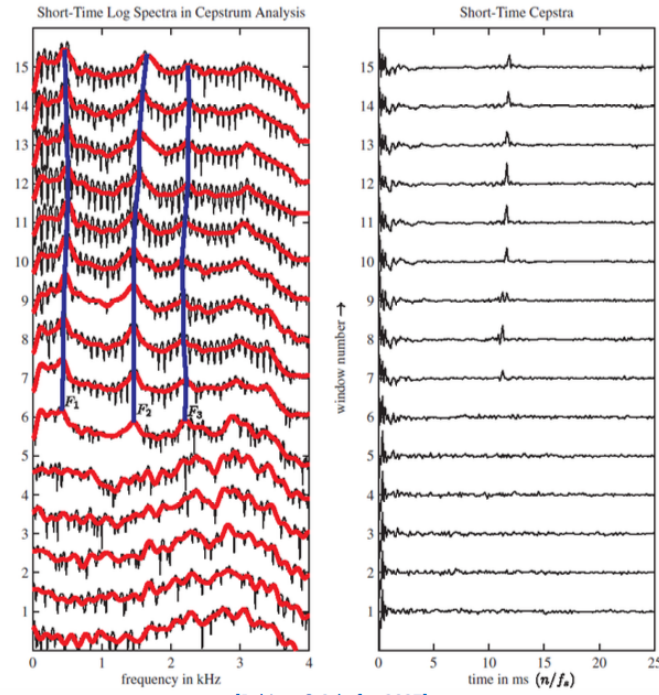


Abbildung 5.8: Cepstrum

sein von Stimme für das aus dem Fenster x_i Berechneten Cepstrum x_i hin. Als zweites Features des Cepstrums wird die Quefrequency der so gefundenen höchsten Magnitude (Upper Cepstrum Peak Location) nach Formel 5.11 berechnet. Bei Signalfenstern mit Stille ist es wahrscheinlicher, dass sich die Spitze am Mindest- oder Maximum-Wert des Cepstrum befindet.

$$Ceps_{mag}(c_i) = \max_k \{\text{mag}(c[k])\} \quad (5.10)$$

$$Ceps_{loc}(c_i) = \arg \max_k \{\text{mag}(c[k])\} \quad (5.11)$$

Abbildung 5.9 visualisiert die beiden Features des Cepstrums auf die selbe Art und Weise, wie es bei Abbildung 5.5

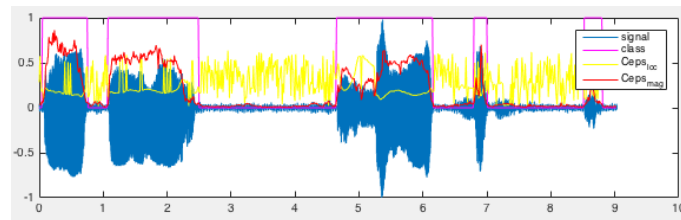


Abbildung 5.9: Features des Cepstrums

Autokorrelation

Die vierte Methode, die zur Voice-Activity-Detection verwendet wird, ist die **Autokorrelation**. Dabei wird das Signalfenster x_i mit einer verzögerten Variante von sich selber nach Formel 5.12 korreliert. Der Parameter k wird als Lag bezeichnet und definiert die Verzögerung des Signals, mit dem die Korrelation durchgeführt wird. Ein hoher Autokorrelationswert bei einem bestimmten Lag k weist darauf hin, dass das Signal eine Periodizität bezüglich diesem Verzögerungswert aufweist. Möchte man beispielsweise berechnen, ob ein bei 16 000 Hz gesampeltes Signal bezüglich 400 Hz periodisch ist, so wird die Autokorrelation mit einem Lag von $k = \frac{\text{Sampling-}f_s}{\text{Detect-}f_s} = \frac{16\,000\text{ Hz}}{400\text{ Hz}} = 40$ durchgeführt. Formel 5.13 zeigt die *normalisierte Autokorrelation*, bei der das Autokorrelations-Signal sowohl bezüglich der Signalenergie als auch der Anzahl der Lags normalisiert wird.

$$a_i[k] = \sum_{n=k}^N x_i[n] \cdot x_i[n - k] \quad (5.12)$$

$$acorr_i[k] = \frac{\sum_{n=k}^N x_i[n] \cdot x_i[n - k]}{\sqrt{(\sum_{n=1}^{N-k} x_i[n]^2)} \cdot \sqrt{(\sum_{n=k}^N x_i[n]^2)}} \quad (5.13)$$

Abbildung 5.10 zeigt das $acorr_i$ -Signal eines Sprachsignals mit den Lags 50 - 320. Die höchste Signalmagnitude bei Lag = 145 ist ein Indikator dafür, dass bei einer Samplingrate von 16 000 Hz die Grundfrequenz $\frac{16000\text{ Hz}}{145} = 110.34\text{ Hz}$ am dominantesten im Signal enthalten ist, was bei einem Stimmsignal die Grundfrequenz f_0 ist.

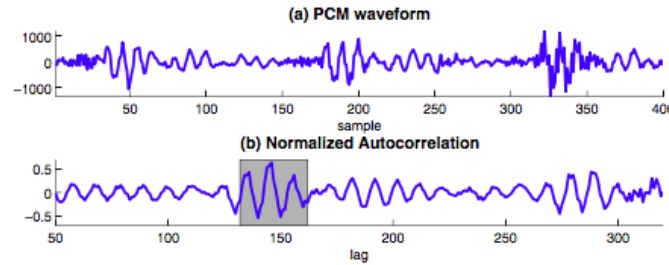


Abbildung 5.10: Autokorrelation eines Signals

Da Stimmsignale eine höhere Periodizität aufweisen als das Hintergrundrauschen, eignet sich die Autokorrelation zur Voice-Activity-Detection. Es werden die Features *Maximum Autocorrelation Peak* [$aMax$] und (*Autocorrelation Peak Count*) [$aCount$] berechnet.

Beide Features werden in [15] zur VAD verwendet. Die *höchste Magnitude der Autokorrelation* (*Maximum Autocorrelation Peak*) wird nach der Formel 5.14 definiert. Eine höherer $aMax$ -Wert weist auf das Vorhandensein von Periodizität im Signal und somit auf Stimme hin. Das zweite Feature ist die *Anzahl an Autokorrelations-Spitzen* nach Formel 5.15. Da nach [27] die Grundfrequenz eines Kindergeschreis zwischen 250 und 2000 Hz liegt, wurden

auch nur Lags für diesen Bereich verwendet.

$$aMax(x_i) = \max_k \{ \text{mag}(acorr_i[k]) \} \quad (5.14)$$

$$aCount(x_i) = \text{count}_k \{ \text{mag}(acorr_i[k]) \} \quad (5.15)$$

Abbildung 5.11 visualisiert die Features der Autokorrelation auf die selbe Art und Weise wie bei Abbildung 5.5.

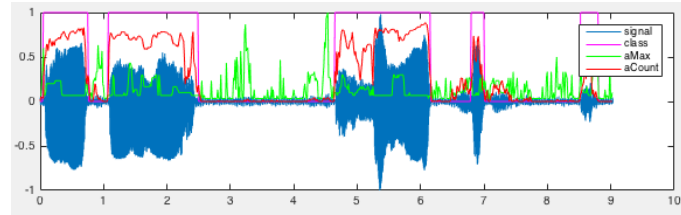


Abbildung 5.11: Features der Autokorrelation

Konstruktion des Feature-Vektors

Für jedes Signalfenster à 25 ms werden die 9 vorgestellten Features RMS , ZCR , $SEnt_u$, $SEnt_n$, f_{Dom} , $Ceps_{mag}$, $Ceps_{loc}$, $aMax$ und $aCount$ berechnet. Wie beschrieben, sollten Signalfenster mit Stimme einen höheren Wert des jeweiligen Features erzeugen als Signalfenster mit Stille (oder, abhängig vom Feature wie der ZCR, einen tieferen).

Abbildung 5.12 zeigt in (a) die RMS-Werte eines Signals mit einem Signal-Rausch-Abstand (SNR) von 50 dB. Die Zeiträume mit Stille haben einen weitaus niedrigeren RMS-Wert als die Signaleile mit Stimme. In (b) ist das selbe Signal mit einem Signal-Rausch-Abstand von 3 dB zu sehen. Nun liegen die RMS-Werte des Rauschens nur noch knapp unter denen des Sprachsignals. Zu sehen ist, dass starkes Hintergrundrauschen ähnlich hohe Feature-Werte erzeugen kann wie die Stimme.

In [20] und [18] wird die Idee präsentiert, den Wert des jeweiligen Features zu messen, der in den stimmlosen Segmenten durch das Hintergrundrauschen erzeugt wird. Jedoch ist für ein Signalfenster x_i zum Zeitpunkt der Berechnung des Features noch nicht bekannt, ob es sich um ein Zeitfenster mit Stille oder Stimme handelt. Genau diese Entscheidung wird später auf Basis der Features getroffen. Man kann jedoch davon ausgehen, dass kein Baby-Schrei länger eine bestimmte Zeitraum t_{max} dauert, bevor das Babie Luft holen muss und somit ein Zeitfenster mit Stille entsteht. Der Mindestwert eines Features in einem Zeitraum t_{max} ist somit der Wert des Features, der mindestens durch das Hintergrundrauschen erzeugt wird. Er wird nach Formel berechnet, wobei t_{xi} die Länge eines Signalfensters x_i ist (in diesem Fall 25 ms). Das *Differenzfeature* wird definiert nach Formel 5.17 als die Differenz des für das aktuelle Signalfenster berechneten Features und dem Mindestwert dieses Feature

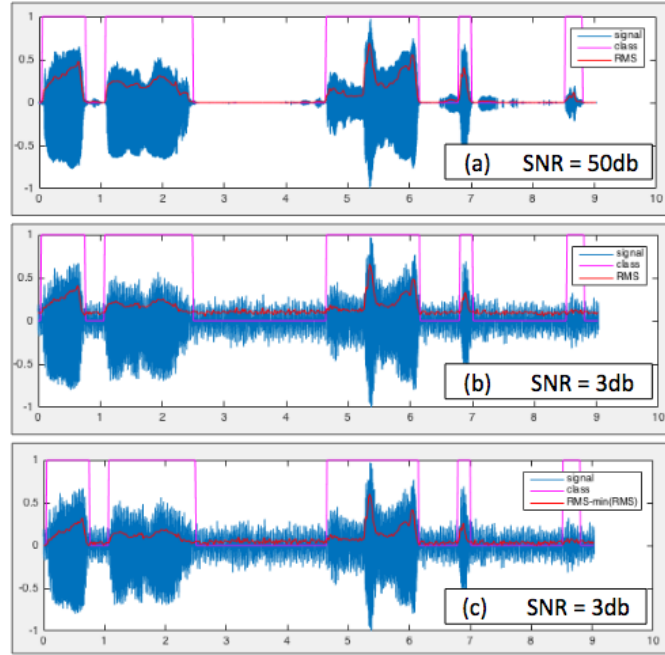


Abbildung 5.12: Differenz-Signal des RMS-Wertes

der letzten t_{max} Sekunden.

$$\text{MinF}(\text{Feat}(x_i)) = \min_{k=i-z\dots i} (\text{Feat}(x_k)), z = \frac{2t_{max}}{t_{xi}} \quad (5.16)$$

$$\text{DiffF}(\text{Feat}(x_i)) = \text{Feat}(x_i) - \text{MinF}(\text{Feat}(x_i)) \quad (5.17)$$

Der Featurevektor, der schlussendlich für jedes Signalfenster x_i berechnet wurde, besteht neben den 9 vorgestellten Features, zusätzlich aus den Differenzfeatures für einen Zeitraum von $t_{max} = 4$ Sekunden. Features, deren Wert bei stimmhaften Signalfenstern geringer ist als bei stimmlosen, werden an der x-Achse gespiegelt, um Formel 5.16 anwenden zu können (das betrifft die Features ZCR, SEnt_u und aCount). Das einzige Feature, welches nicht als Differenzfeature dem Featurevektor beigefügt wurde, ist der *Upper Cepstral Peak Location*-Feature [*Ceps_{loc}*], da es bei Stille sowohl einen höheren als auch einen niedrigeren Wert annehmen kann. Der Feature-Vektor V des Signalfensters x_i wird nach Formel 5.18 berechnet umfasst 17 Features, wobei 9 absolute Features und 8 Differenzfeatures verwendet werden.

$$V(x_i) = \left(\text{RMS}(x_i), \dots, \text{aCount}(x_i), \right. \\ \left. \text{Diff}(\text{RMS}(x_i)) \dots \text{Diff}(-\text{aCount}(x_i)) \right) \quad (5.18)$$

5.3.3 Thresholding

Finden der Grenzwerte

Für jedes Signalfenster $x_1 \dots x_n$ liegt nun ein Featurevektor $v_1 \dots v_n$ vor. Das Ziel ist nun, Grenzwerte für die Features zu finden, bei deren Überschreitung das Signalfenster als *Stimmhaft* kategorisiert wird. Abbildung 5.13 verdeutlicht das Prinzip für das Feature $Ceps_{mag}$. Eine binäre Kategorisierung nach dem Muster $C(x_i) = \{1, \text{wenn } Ceps_{mag}(x_i) \geq 0.2, 0 \text{ sonst}\}$ würde auf den ersten Blick eine weitgehend richtige Kategorisierung vornehmen.

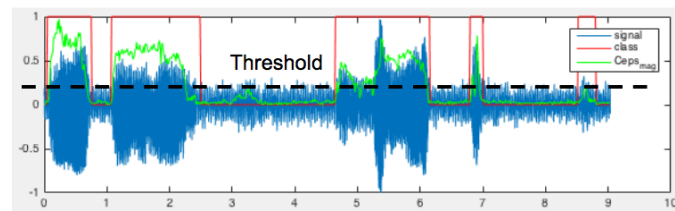


Abbildung 5.13: Thresholding eines Feature-Signales

Es sind Kombinationen von Grenzwerten in Form von Entscheidungsbäumen denkbar, wenn mehrere Features in die Kategorisierung einfließen. Ein Beispiel wird in Listing 5.1 dargestellt, bei dem die Kategorisierung hierarchisch zuerst nach einem Grenzwert für $Ceps_{mag}$ und danach für den RMS-Wert entschieden wird.

Listing 5.1: Beispiel eines CART-Entscheidungsbaums

```

if  $Ceps_{mag}(x_i) > 0.2$ 
|   if  $RMS(x_i) < 0.13$ 
|   |    $C(x_i) = 0$ 
|   |   else
|   |        $C(x_i) = 1$ 
|   else
|        $C(x_i) = 1$ 

```

Zur Festlegung der Grenzwerte wird ein *Classification And Regression Tree* (CART)-Algorithmus verwendet. CART-Algorithmen sind Klassifizierungsalgorithmen, bei denen Entscheidungsbäume wie in Listing 5.1 konstruiert werden. Sie haben den Vorteil, dass die Klassifizierung in Form von Regeln mit Grenzwerten dargestellt werden können, die für den Menschen nachvollziehbar sind (im Gegensatz zu z.B. Neuronalen Netzen). Der durch den CART-Algorithmus konstruierte Entscheidungsbaum wird auch als *Klassifizierungs-Modell* bezeichnet.[5]

Einer der bekanntesten CART-Algorithmen ist der ID3-Algorithmus, entwickelt von J. Ross Quinlan an der University of Sidney. Als Trainingsdatensatz S wird eine Menge an Instanzen verwendet, deren Klassenzugehörigkeit bereits bekannt ist. Der ID-3 Algorithmus funktioniert nur für diskrete Features und Klassen. Ein Beispiel für ein diskretes Feature ist $[Lautstärke = \{laut, leise\}]$. Ein Beispiel für ein numerisches Feature, welches vom ID-3 Algorithmus nicht akzeptiert wird, ist $[Lautstärke = \mathbb{N}]$. Der Algorithmus funktioniert folgendermaßen:

1. Gehören alle Instanzen des Datensatzes S nur einer Klasse an?
 - a) Ja: Markiere diesen Knoten als Blatt und STOP.
 - b) Nein: Erzeuge einen neuen Knoten. Wähle das Feature F , welches den höchsten *Informationsgewinn* bringt. Das heisst, dass dieses Feature die stärkste Unterscheidung zwischen den einzelnen Klassen ermöglicht. Der Informationsgewinn $H(S)$ eines Features F mit den Untermengen X , welche aus den Instanzen der selben Klasse bestehen, wird mit Hilfe der Entropie nach Formel 5.19 berechnet.

$$H(S) = - \sum_{x \in X} F(x) \log_2 F(x) \quad (5.19)$$

2. Unterteilung den Trainingsdatensatz S in die Untermengen $S_1 \dots S_n$, wobei eine Untermenge für jeden möglichen Attributwert des ausgewählten Features gebildet wird.
3. Wiederhole die Schritte 1. und 2. für alle Untermengen $S_1 \dots S_n$ [5]

Der C.45-Algorithmus ist eine Erweiterung des ID-3 Algorithmus, der neben anderen Erweiterungen ebenfalls mit numerischen Attributwerten umgehen kann. Dafür werden für ein Feature alle Instanzen des Datensatzes S nach den Werten dieses Features geordnet und derjenige Grenzwert gesucht, der den höchsten Informationsgewinn bringt. Der Datensatz wird daraufhin an den Knoten in genau 2 Untermengen aufgespalten.[13]

Zum Finden der Grenzwerte der in Kapitel 5.3.2 beschriebenen Features wurde die Implementierung des C.45-Algorithmus *REPTree*¹ der Open Source Data-Mining-Bibliothek *Weka*² verwendet. Diese hat den Vorteil, dass die Tiefe des Entscheidungsbaumes festlegbar ist und somit die Komplexität des Baumes begrenzt werden kann. Ein Nebeneffekt ist, dass so Overfitting, also die Überanpassung des Entscheidungsbaumes auf den Trainingsdatensatz, vermieden wird.

Trainings- und Testdatensätze

Zum Training des REPTree-Algorithmus musste in Trainingsdatensatz S erstellt werden. Dazu wird zunächst eine Menge an Audiosignalen benötigt. Es wurden 6 Audiosignale mit Weinen von Babies von der freien Online-Sound-Bibliothek <https://www.freesound.org/> heruntergeladen und zu Segmenten à 10 Sekunden beschnitten. Diese Audiosignale sind weitestgehend Rausch-frei. Die Segmente der Audiosignale wurden händisch Kategorisiert in die Klassen $\{1 = \text{Stimme}, 0 = \text{Still}\}$. Weiterhin wurden 3 verschiedene Rauschsignale heruntergeladen. Es handelt sich um "realistische" Rauschsignale mit Krankenhausatmosphären. Jedes dieser 3 Rauschsignale wurde mit den 6 Weinsignalen überlagert, einmal mit einem Signal-Rausch-Abstand von 50 dB (fast unhörbares Rauschen), und einmal mit einem

¹Dokumentation von REPTree: <http://weka.sourceforge.net/doc.dev/weka/classifiers/trees/REPTree.html>

²Download von WEKA: <http://www.cs.waikato.ac.nz/ml/weka/>

Signal-Rausch-Abstand von 3 dB (sehr starkes Rauschen). Außerdem wurde ein Testsignal erzeugt, welches eine siebte Tonaufnahme eines Kinderweins enthält, dass mit einem vierten Rauschsignal mit einem SNR von 7 dB überlagert wurde. Dieses Signal spielt eine Sonderrolle, da es nur zur Verifikation verwendet wird und enthält daher nur ein Audiosignal (Siehe Kapitel 5.3.3) So wurden vier Mengen an Audiosignalen erzeugt:

$A_{50\text{dB}}$ enthält $3 \cdot 6 = 18$ Audiosignale, bei dem alle 6 Wein-Signale mit den 3 Rauschsignalen bei einem Signal-Rausch-Abstand von 50 dB überlagert wurden

$A_{3\text{dB}}$ enthält $3 \cdot 6 = 18$ Audiosignale, bei dem alle 6 Wein-Signale mit den 3 Rauschsignalen bei einem Signal-Rausch-Abstand von 3 dB überlagert wurden

$A_{50+3\text{dB}} = \{A_{50\text{dB}} \cup A_{3\text{dB}}\} = 32$ Audiosignale

$A_{7\text{dB}^*}$ enthält 1 Audiosignal, bei dem ein siebtes Wein-Signale mit einem vierten Rauschsignal bei einem Signal-Rausch-Abstand von 7 dB überlagert wurde.

Im nächsten Schritt werden die eigentlichen Trainingsdatensätze $S_{SNR,Feat}$ gebildet, in dem Audiosignale dieser Signalmengen (1) wie in Kapitel 5.2 vorverarbeitet werden, (2) wie in Kapitel 5.3.1 in die Signalfenster à 25 ms zerlegt werden und (3) für jedes Audiosignal die durch Gleichung 5.18 definierte Featurevektoren berechnet werden. Außerdem wird jedem Featurevektor die Klasseninformation *Stimme/Stille* zum Training des REPTree-Algorithmus mitgegeben.

Es ist rechnerisch zu Aufwendig, alle genannten Features in einem kontinuierlichen System zur Voice Activity Detection zu berechnen. Daher werden Untermengen der Features in den Datensätzen gebildet. Das Ziel ist es, diejenige Untermenge an Features zu finden, die sich am besten für die Voice-Activity-Detection sowohl bei niedrigem als auch bei starkem Hintergrundrauschen eignet. Die Untermengen werden in Bezug auf die Methode gebildet, durch die die Features berechnet werden. Das heißt, dass beispielsweise die Untermenge *Zeit* die in Kapitel 5.3.2 beschriebenen Features *RMS* und *ZCR* sowie die dazugehörigen Differenzfeatures $FDiff(RMS)$ und $FDiff(ZCR)$ beinhaltet.

Die 9 Untermengen sind: { Zeitbereich, Frequenzbereich, Cepstrum, Autokorrelation, Zeit + Frequenzbereich, Zeit + Cepstrum, Zeit + Autokorrelation, Frequenz + Cepstrum, Frequenz + Autokorrelation }.

So enthält beispielsweise der Datensatz $S_{3\text{dB},Zeit}$ die Featurevektoren des Zeitbereiches für die Audiosignale mit einem Signal-Rausch-Abstand von 325 dB. Die Audiosignal-Mengen $[A_{50\text{dB}}]$, $[A_{3\text{dB}}]$, $[A_{50+3\text{dB}}]$ und $[A_{7\text{dB}^*}]$ wurden in Datensätze umgewandelt. Es werden schlussendlich $4 \cdot 9 = 36$ Datensätze gebildet.

Training und Ergebnis

Der REPTree-Algorithmus entwirft einen Entscheidungsbaum, der für den angegebenen Trainingsdatensatz eine möglichst hohe Accuracy gewährleistet (unter den gegebenen Einschränkungen eines CART-Algorithmus). Wird dem REPTree-Algorithmus beispielsweise

der Datensatz $S_{3\text{ dB}, \text{Zeit}}$ (also der Datensatz mit einem SNR von 3 dB unter Verwendung der Zeit-Features, siehe Kapitel 5.3.3) als Input gegeben, entwirft der Algorithmus einen Entscheidungsbaum auf Basis dieses Datensatzes. Wird das gebildete Modell daraufhin für die Klassifikation der $A_{3\text{ dB}}$ -Signalmenge verwendet, kann man aus dem Klassifikationsergebnis die Accuracy des Modells für diesen Signal-Rausch-Abstand berechnen.

Der Datensatz $S_{\text{SNR}, \text{Feat}}$, der als Input für den REPTree-Algorithmus und somit zur Bildung des Entscheidungsbaums verwendet wird, wird als *Trainings-Datensatz* bezeichnet. Die Signalmenge A_{SNR} , auf den das Modell angewandt wird und für den die Accuracy berechnet wird, wird als *Test-Datensatz* bezeichnet. Wird also beispielsweise der Trainings-Datensatz $S_{50\text{ dB}, \text{Zeit}}$ und als Test-Signalmenge $A_{3\text{ dB}}$ verwendet, so kann man berechnen, wie gut sich ein Modell unter Verwendung der Zeit-Features zur Klassifizierung niedriger SNRs eignet, obwohl es für hohe SNRs entworfen wurde.

Das Ziel ist, den Entscheidungsbaum für eine Feature-Untermenge zu finden, die eine möglichst hohe Klassifikations-Accuracy für sowohl hohe als auch niedrige SNR gewährleistet. Die Frage ist, ob ein Entscheidungsbaum, der für einen niedrigen SNR gebildet wird, auch für hohe SNR gut funktioniert, oder ob das Gegenteil zutreffend ist. Daher werden die Entscheidungsbäume sowohl auf Basis verschiedener SNRs als auch verschiedener Feature-Untermengen gebildet. Die Modelle werden daraufhin gegen die Signale mit den hohen und niedrigen SNRs getestet.

Es wurden, wie in Kapitel 5.3.3 beschrieben, $3 \cdot 9 = 27$ Trainings-Datensätze erzeugt ([3 SNR-Werten: 3 dB, 50 dB und 50+3 dB] \times [9 Feature-Untermengen]. Der Datensatz mit einem SNR von 7 dB wird *nicht* zum Training verwendet). Mit diesen 27 Trainingsdatensätzen werden mit Hilfe des REPTree-Algorithmus 27 Klassifikationsbäume erzeugt. Jeder Klassifikationsbaum wurde gegen die 3 Testdatensätze $A_{3\text{ dB}}$, $A_{50\text{ dB}}$ und $A_{7\text{ dB}^*}$ getestet und die Accuracy berechnet. Das Signal $A_{7\text{ dB}^*}$ erfüllt dabei eine Sonderrolle, da weder das Signal des Weinens noch das Rausch-Signal in den Trainingsdatensätzen enthalten sind und somit verifiziert wird, ob der Datensatz nur „auswendig gelernt“ wird oder das Modell auf neue Anwendungsfälle übertragen werden kann. Um Overfitting des Modells zu vermeiden und die Komplexität des Entscheidungsbaumes zu verringern, wurde die maximale Tiefe des REPTree auf 2 gesetzt. Die Ergebnisse sind in Tabelle .1 zu sehen.

Die Features, welche zu den höchsten Accuracy-Werten führten, sind die des *Cepstrum*-Bereiches, genauer gesagt das DiffF(Ceps_{mag})-Feature, da es vom REPTree-Algorithmus als einziges Feature dieses Bereiches für die Entscheidungsbäumen ausgewählt wurde. Die Entscheidungsbäume, die mit dem DiffF(Ceps_{mag})-Feature entworfen wurden, erreichten eine durchschnittliche Accuracy (das heißt, gemittelt über die Testsignale $A_{3\text{ dB}}$, $A_{50\text{ dB}}$ und $A_{7\text{ dB}^*}$) von mindestens 91,45%. Der nächstbeste Entscheidungsbaum mit einer Accuracy von 86,96% wurde unter Verwendung der Features des Zeitbereiches und der rechnerisch aufwendigeren Autokorrelation auf dem Datensatz $S_{50+3\text{ dB}, \text{Zeit}+\text{Correlation}}$ entworfen. Sobald der Cepstrum-Bereich in Verbindung mit den Features der Bereiche *Zeit* und *Frequenz* verwendet wurde, wurde das DiffF(Ceps_{mag})-Feature vom REPTree-Algorithmus bevorzugt, so dass die Features der anderen beiden Bereiche keine Anwendung mehr in den entsprechenden Bäumen fanden.

Auf Basis der Datensätze $S_{3\text{ dB}, \text{Ceps}}$, $S_{3\text{ dB}, \text{Zeit}+\text{Ceps}}$, $S_{3\text{ dB}, \text{Freq}+\text{Ceps}}$, $S_{50+3\text{ dB}, \text{Ceps}}$, $S_{50+3\text{ dB}, \text{Zeit}+\text{Ceps}}$

sowie

$S_{50+3\text{ dB}, \text{Freq}+\text{Ceps}}$ wurde der selbe Entscheidungsbaum erzeugt, der in Listing 5.2 zu sehen ist. Auf Basis der Datensätze $S_{50\text{ dB}, \text{Ceps}}$ und $S_{50\text{ dB}, \text{Zeit}+\text{Ceps}}$ wurde der Entscheidungsbaum in Listing 5.3 erzeugt. Es ist zu sehen, dass (a) beide Entscheidungsbaum einen einfachen Grenzwert für das $\text{Diff}(\text{Ceps}_{\text{mag}})$ -Feature setzen, und zweitens (b) sich die beiden Modelle nur im konkreten Wert des Grenzwertes unterscheiden.

Da das Modell aus Listing 5.2 eine durchschnittliche Accuracy von 92,22% und das Modell aus Listing 5.3 eine unwesentlich geringere Accuracy von 91,45% erreicht, wurden für beide Modelle die Specificity und Sensitivity berechnet, um eine Entscheidung für eines der beiden Modelle fällen zu können. Dazu wurden die Signalmengen $A_{3\text{ dB}}$, $A_{50\text{ dB}}$ und $A_{7\text{ dB}^*}$ in Frames à 100 Windows zerlegt und für jedes Zeitfenster die Sensitivity, Specificity und Accuracy bezüglich der beiden Modelle berechnet. Die Ergebnisse werden als Boxplots in Abbildung .1 dargestellt. Der Unterschied zwischen den Modellen ist am Stärksten beim Testing gegen die Signale mit einem SNR von 3 dB und 7 dB zu sehen. Das Modell mit dem Grenzwert von 0.03 erzielt in beiden Fällen eine höhere Specificity, aber geringere Sensitivity als das Modell mit dem Grenzwert bei 0.02. Es wurde sich für das Modell für mit einem Grenzwert von 0.02 entschieden, da durch die höhere Sensitivity mehr Wein-Signale erkannt werden, die in späteren Verarbeitungsschritten immernoch als False-Positives erkannt und verworfen werden können. Einmal im Prozess der VAD als Stimmlos markierte Fenster werden jedoch nicht weiter verarbeitet und gehen somit „verloren“.

Listing 5.2: Entscheidungsbaum für die VAD mit einem Cepstrum-Grenzwert von 0.02

```

if FDiff(Cepsmag( $x_i$ )) < 0.02
|   C( $x_i$ ) = 0
| else
|   C( $x_i$ ) = 1
    
```

Listing 5.3: Entscheidungsbaum für die VAD mit einem Cepstrum-Grenzwert von 0.03

```

if FDiff(Cepsmag( $x_i$ )) < 0.03
|   C( $x_i$ ) = 0
| else
|   C( $x_i$ ) = 1
    
```

Der Finale Funktion zur Klassifikation eines Signalfensters $C(x)$ in $0 = \text{Stille}$ oder $1 = \text{Stimme}$ ist somit durch Gleichung 5.20 gegeben.

$$C(x) = \begin{cases} 1, & \text{if } \text{Ceps}_{\text{mag}}(x) \geq 0.02 \\ 0, & \text{otherwise} \end{cases} \quad (5.20)$$

5.3.4 Markierung der Cry-Units

Das Ergebnis der Voice-Activity-Detection ist eine Zurdnung aller Signalfenster $x_1 \dots x_n$ zu den Klassen $C(x_i) = 0$ *Stille* oder $C(x_i) = 1$ *Stimme*. In [18] wird die Idee vorgestellt, zu-

sammenhängende und ununterbrochene Ketten als *stimmhaft* klassifizierter Signalfenster zu *Stimm-Segmenten* zusammenzufassen, welche in diesem Zusammenhang eine *Cry-Units* entsprechen. Abbildung 5.14 veranschaulicht diese Gruppierung. Formel 5.21 gibt die Definition des Datentypes *Cry-Unit* [CU]. Eine Cry-Unit definiert sich durch einen Anfangszeitpunkt *start*, einen Endzeitpunkt *end* und der Liste seiner Signalfenster *windows* = $[x_1 \dots x_n]$.

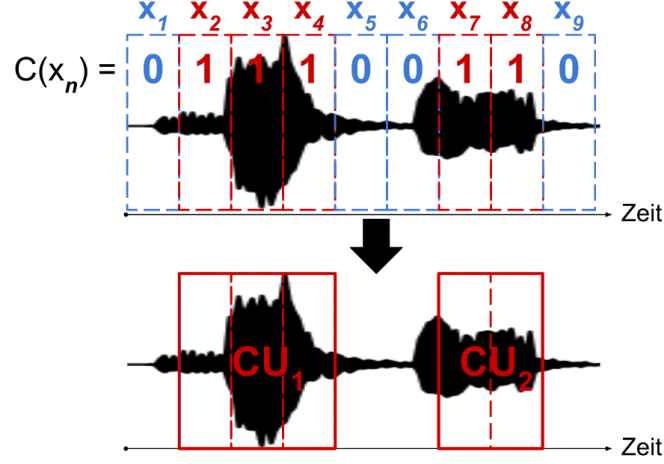


Abbildung 5.14: Zusammenfassung klassifizierter Signalfenster zu Cry-Units

Algorithmus 1 zeigt in Pseudo-Code, wie in der Liste aller Signalfenster $X_{windows} = [x_1 \dots x_n]$ eine Liste von Cry-Units $CU_{all} = [cu_1 \dots cu_m]$ markiert wird. Die Funktion $C(x)$ ist die Klassifikations-Funktion der Signalfenster in Stille/Stimme nach Gleichung 5.20. Die Funktion $getTimeOf(x)$ liefert die Anfangszeitpunkt des Signalfensters x .

$$CU = (windows = [x_1 \dots x_n], start \in Zeit, end \in Zeit) \quad (5.21)$$

$$\lambda(CU) = CU.end - CU.start \quad (5.22)$$

$$d(CU_i, CU_j) = CU_j.start - CU_i.end \quad (5.23)$$

Die Dauer eine Cry-Unit $\lambda(CU)$ wird nach Formel 5.22 berechnet. Der (Stille)-Zeitraum zwischen zwei Cry-Units $d(CU_i, CU_j)$, wird nach Formel 5.23 berechnet. Diese Zusammenhänge werden in Abbildung 5.15 visualisiert.[18]

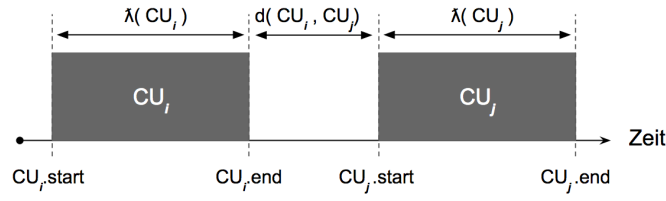


Abbildung 5.15: Beziehung zwischen agrenzenden Segmenten

Algorithm 1 Gruppierung von x-Windows zu Cry-Units

```

1: function TURNWINDOWSINTOCRYUNITS( $X_{windows}$ )
2:    $CU_{all} \leftarrow []$ 
3:    $cu_i \leftarrow ([], 0, 0)$ 
4:   for  $i = 1 \dots \text{length}(X_{windows})$  do
5:      $c_i \leftarrow C(x_i)$ 
6:                                     ▷ Start of Cry-Unit
7:     if  $c_i == 1 \wedge \text{isEmpty}(cu_i.windows)$  then
8:        $cu_i.start \leftarrow \text{getTimeOf}(x_i)$ 
9:        $cu_i.windows \leftarrow [cu_i.windows, x_i]$ 
10:    end if
11:                                     ▷ Inside Cry-Unit
12:    if  $c_i == 1 \wedge ! \text{isEmpty}(cu_i.windows)$  then
13:       $cu_i.windows \leftarrow [cu_i.windows, x_i]$ 
14:    end if
15:                                     ▷ End of Cry-Unit
16:    if  $c_i == 0 \wedge ! \text{isEmpty}(cu_i.windows)$  then
17:       $cu_i.end \leftarrow \text{getTimeOf}(x_i)$ 
18:       $CU \leftarrow [CU, cu_i]$ 
19:       $cu_i.windows \leftarrow []$ 
20:    end if
21:  end for
22:                                     ▷ End last Cry-Unit by force if still open.
23:  if  $! \text{isEmpty}(cu_i.windows) == 0$  then
24:     $cu_i.end \leftarrow \text{getTimeOf}(X_{windows}[\text{end}])$ 
25:     $CU_{all} \leftarrow [CU_{all}, cu_i]$ 
26:  end if
27:  return  $CU_{all}$ 
28: end function
    
```

5.3.5 Decision Smoothing

Abbildung 5.16 zeigt ein Audiosignal mit einem Signal-Rausch-Abstand von 3 dB, bei dem die Klassifizierung nach dem Entscheidungsbaum aus Listing 5.2 durchgeführt wurde. Die rote Linie zeigt die tatsächliche Klassifizierung, und die grüne Linie die gefundene Klassifizierung nach dem vorgestellten Algorithmus. Die tatsächlichen/gefundenen Cry-Units

sind klar zu erkennen als die Bereiche, die von der roten/grünen Linie überspannt werden. Es ist zu sehen, dass False-Negatives und False-Positives in der Klassifizierung enthalten sind. Im folgenden werden drei charakteristische Arten falscher Klassifizierungen näher erläutert:

False Negatives nach (a) : Eine korrekt erkannte, längere Cry-Unit wird zu früh beendet. Oft werden kurz nach dem Ende sehr kurze Cry-Units erkannt, die eigentlich noch zu der längeren, vorhergehenden Cry-Unit gehören.

False Positives nach (b): Kurze Cry-Units werden in eigentlichen Stille-Bereichen erkannt.

False Negatives nach (c): Eine Cry-Unit zerfällt in zwei kürzere Cry-Units, da einige Signalfenster in der Mitte als Stille erkannt wurden.

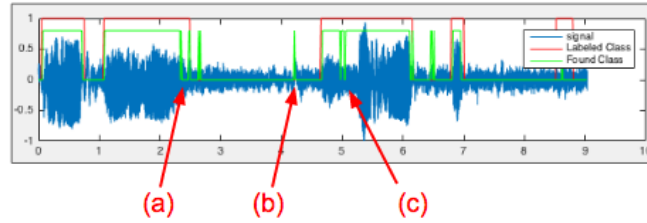


Abbildung 5.16: Klassifizierung vor dem Decision Smoothing

Im Prozess des **Decision Smoothing** werden kontextuelle Informationen genutzt, um nachträglich False-Positives und False-Negatives zu entfernen. Es werden dazu die in [18] präsentierten Ideen verwendet. Es werden zwei Parameter eingeführt: λ_{min} , die Mindestlänge einer akzeptierten Cry-Unit, und d_{min} , die Mindestlänge eines akzeptierten Stille-Segementes. Das Decision Smoothing wird nach den folgenden Entscheidungsregeln durchgeführt:

-
- ist $\lambda(CU_i) \leq \lambda_{min}$?
 - wenn $\lambda(CU_{i-1}) > \lambda_{min}$ und $d(CU_{i-1}, CU_i) \leq d_{min}$, dann vereinige CU_i mit CU_{i-1} . \Rightarrow behebt False-Negatives des Types (a)
 - ansonsten entferne $CU_i \Rightarrow$ behebt False-Negatives des Types (b)
 - wenn $\lambda(CU_i) > \lambda_{min}$ und $d(CU_{i-1}, CU_i) \leq d_{min}$, dann vereinige CU_i mit CU_{i-1} . \Rightarrow behebt False-Negatives des Types (c)
-

Die Entscheidungsregeln greifen Algorithmus greifen nur auf die aktuellen und die letzten bekannte Cry-Unit um, um eine kontinuierliche Analyse zu gewährleisten, weshalb die Entscheidungsregeln jedoch auch komplex sind. Bei einer offline-Analyse können die Entscheidungsregeln vereinfacht werden, da False-Negative Type (a) und (c) mit der selben Regeln abgefragt werden können. Algorithmus 2 zeigt in Pseudo-Code, wie das Decision-

Smoothing durchgeführt wird. Input der Funktion ist die Liste aller Cry-Units CU_{all} , die durch Algorithmus 1 entstanden ist, sowie die Grenzwerte λ_{min} , d_{min} . Ausgang der Funktion ist die Liste aller Cry-Units nach dem Decision-Smoothing $CU_{smoothed}$.

Algorithm 2 Decision-Smoothing of VAD

```

1: function DECISIONSMOOTHING( $CU_{all}$ ,  $\lambda_{min}$ ,  $d_{min}$ )
2:    $CU_{smoothed} \leftarrow [CU_{all}[1]]$ 
3:   for  $i = 2 \dots \text{length}(CU_{all})$  do
4:      $cu_i \leftarrow CU_{all}[i]$ 
5:      $cu_{i-1} \leftarrow CU_{smoothed}[\text{end}]$ 
6:     if  $\lambda(cu_i) > \lambda_{min}$  then
7:       if  $d(cu_{i-1}, cu_i) > d_{min}$  then
8:          $CU_{smoothed} \leftarrow [CU_{smoothed}, cu_i]$ 
9:       else
10:                                     ▷ Erase False-Negative Type (c)
11:          $cu_i \leftarrow \text{vereinige}(cu_i, cu_{i-1})$ 
12:          $CU_{smoothed} \leftarrow [CU_{smoothed}[1 : \text{end} - 1], cu_i]$ 
13:       end if
14:     else
15:                                     ▷ Erase False-Negative Type (a)
16:       if  $\lambda(cu_i) > \lambda_{min} \wedge d(cu_{i-1}, cu_i) \leq d_{min}$  then
17:          $cu_i \leftarrow \text{vereinige}(cu_i, cu_{i-1})$ 
18:          $CU_{smoothed} \leftarrow [CU_{smoothed}[1 : \text{end} - 1], cu_i]$ 
19:       else
20:                                     ▷ Don't accept  $cu_i$ . Erases False-Positives (b)
21:       end if
22:     end if
23:   end for
24:   return  $CU_{smoothed}$ 
25: end function
    
```

Abbildung 5.17 zeigt das Signal vor und nach dem Decision-Smoothing. Die Parameter wurden experimentell mit $\lambda_{min} = 50ms$ und $d_{min} = 50ms$ bestimmt.

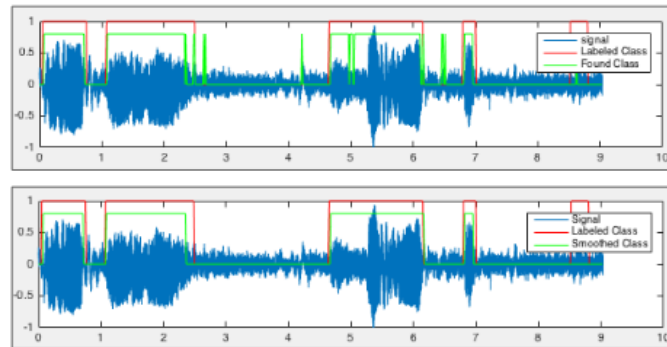


Abbildung 5.17: Klassifizierung vor und nach dem Decision Smoothing

5.4 Segmentierung

Das Ergebnis der Voice-Activiy-Detection ist eine Liste an Cry-Units $CU_1 \dots CU_n$. Das Ziel ist nun, diese Cry-Units zu Cry-Segmenten zu gruppieren. Ein Cry-Segment definiert sich nach Golub et al [12] als „die komplette klangliche Antwort auf einen spezifischen Stimulus. Sie kann mehrere Cry-Units enthalten“. Die Defintion lässt folgende Fragen offen:

- Beginnt das Segment bereits bei Zuführung des Stimulus, oder erst ab der ersten Cry-Unit?
- Wodurch definiert sich der Beginn, wenn ohne Zuführung eines Stimulus das Baby beginnt, zu weinen?
- Endet ein Cry-Segment mit Ende der letzten „Cry-Unit“, oder erstreckt es sich bis zu Beginn des nächsten Cry-Segementes?

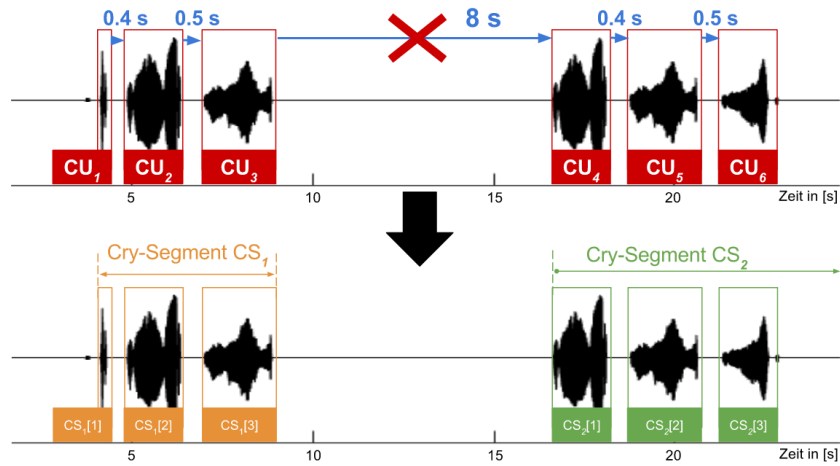


Abbildung 5.18: Ergebnis der Segmentierung

Die Zusammenfassung von Cry-Units zu Cry-Segmentes unterliegt einer gewissen subjektiven Einschätzung, welche Cry-Units als Zusammengehörig angesehen werden, insbesondere, wenn kein erkennbarer Stimulus vorliegt. Abbildung 5.19 verdeutlicht das Problem.

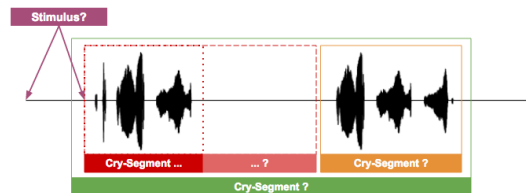


Abbildung 5.19: Mögliche Segmentierungen eines Signals

Um das Problem zu objektivieren, wird es mathematisch formuliert. Eine *Cry-Segment* $[CS]$ wird als Datentyp nach Formel 5.24 definiert. Ein Cry-Segment ist folglich eine

Liste aufeinander folgenden Cry-Units, die gruppiert werden. Der Start-Zeitpunkt eines Cry-Segmentes wird nach Formel 5.25 als der Startzeitpunkt der ersten Cry-Unit des Segmentes definiert. Die Begründung für diese Entscheidung liegt darin, dass rein aus dem Audiomaterial der Zeitpunkt des Stimulus nicht festgestellt werden kann und das Segment mit Sicherheit erst bei der ersten feststellbaren Cry-Unit beginnt. Das Ende eines Segmentes wird definiert als das Ende der letzten Cry-Unit nach Gleichung 5.26. Die Begründung liegt darin, dass das Ende der Reaktion auf den Stimulus ebenfalls rein aus dem Audiosignal abgeleitet werden kann und somit der einzig feststellbare Indikator die letzte Cry-Unit des Segmentes ist.

$$CS = [cu_1, \dots, CU_n] \quad (5.24)$$

$$start(CS) = CS[1].start \quad (5.25)$$

$$end(CS) = CS[end].end \quad (5.26)$$

Wurde bei der kontinuierlichen Analyse des Signals ein Segment geschlossen, führt die Markierung einer neuen Cry-Unit zur Eröffnung eines neuen Segmentes, dessen Start-Zeitpunkt der Start-Punkt dieser Cry-Unit ist. Die Frage ist, welches Kriterium zum schließen dieses Segments führt. Laut Golub et al [12] ein Cry-Segment „die komplette klangliche Antwort auf einen spezifischen Stimulus“. Eine mögliche und objektiv messbare Interepration dieses Endes ist, dass nach dem Auftreten von Cry-Units eine längere Stille mit einer Abwesenheit von Cry-Units festgestellt wird, da das Baby „aufgehört hat, zu weinen“. Übertragen auf die in 5.3.4 vorgestellte Terminologie heißt das, dass ein Segment beendet und ein neues begonnen wird, wenn die Distanz (Zeitraum der Stille) zwischen zwei benachbarten Cry-Units $d(CU_i, CU_{i+1})$ einen gewissen Grenzwert $t_{seg-max}$ überschreitet. Gleichung 5.27 formalisiert diesen Zusammenhang. Daraus lässt sich schlussfolgern, dass die Distanzen zwischen allen benachbarten Cry-Unit eines Segmentes unter diesem Grenzwert $t_{seg-max}$ liegen. Gleichung 5.28 formalisiert diese Nebenbedingung an die Cry-Units eines Segmentes.

$$d(cu_i, cu_{i+1}) > t_{seg-max} \rightarrow CS_n = [CS_n, cu_i] \wedge CS_{n+1} = [cu_{i+1}] \quad (5.27)$$

$$\forall i = 1 \dots \text{length}(CS) - 1 : d(CS[i], CS[i + 1]) \leq t_{seg-max} \quad (5.28)$$

Die einfachste Art, $t_{seg-max}$ festzulegen, ist, einen festen Grenzwert von t s zuzuweisen. Abbildung 5.18 visualisiert die so resultierende Segmentierung an einem Beispiel. Jeder Grenzwert mit $t_{seg-max} > 0.5$ s würde zu der gezeigten Segmentierung führen.

Algorithmus 3 zeigt die Segmentierung nach diesem Prinzip in Pseudo-Code. Input des Algorithmus ist die Liste aller Cry-Units $CU_{all} = [cu_1 \dots cu_n]$, die nach dem Decision-

Smoothing nach Algorithmus 2 entstanden ist. Das Ergebnis des Algorithmus ist die Liste, die alle gefundene Cry-Segmente $[cs_1 \dots cs_n]$ enthält.

Algorithm 3 Gruppierung von Cry-Units zu Cry-Segments

```

1: function SEGMENTCRYUNITS( $CU_{all}, t_{seg-max}$ )
2:    $CS_{all} \leftarrow []$ 
3:    $cs_i \leftarrow [CU_{all}[1]]$ 
4:   for  $i = 2 \dots \text{length}(CU_{all})$  do
5:      $cu_i \leftarrow CU_{all}[i]$ 
6:      $cu_{i-1} \leftarrow CU_{all}[i-1]$ 
7:     if  $d(cu_{i-1}, cu_i) < t_{seg-max}$  then
8:        $cs_i \leftarrow [cs_i, cu_i]$ 
9:     else
10:       $CS_{all} \leftarrow [CS_{all}, cs_i]$ 
11:       $cs_i \leftarrow [cu_i]$ 
12:    end if
13:  end for return  $CS_{all}$ 
14: end function

```

Algorithmus 3 kann zwar kontinuierlich durchgeführt werden, da er jeweils nur auf die aktuelle gefundene und eine vergangene Cry-Unit zurückgreift, hat in dieser Form jedoch den nachteil, dass das Ende eines Segmentes später als notwendig festgestellt wird. Angenommen, ein Grenzwert von $t_{seg-max} = 20$ s wurde festgelegt

Bei einer kontinuierlichen durchgeführten Segmentierung wird das erste Segment dann eröffnet, sobald die erste Cry-Unit durch die VAD markiert wurde, und diese Cry-Unit dem Segment hinzugefügt. Die Dauer der Stille nach dieser Cry-Unit wird kontinuierlich gemessen. Wird ein nächste Cry-Unit festgestellt, bevor die Stille $d_{seg-max}$ übersteigt, so wird diese Cry-Unit dem Segment hinzugefügt und das Messen der Stille nach dieser Cry-Unit beginnt von vorne. Dieser Prozess wird so lange wiederholt, bis die Dauer der Stille nach einer hinzugefügten Cry-Unit $d_{seg-max}$ übersteigt. Dann wird das Segment beendet und der Endzeitpunkt des Segmentes auf den Endzeitpunkt der letzten Cry-Unit gesetzt. Abbildung 5.18 zeigt die resultierende Segmentierung für Beispielsignal mit $d_{seg-max} = 3$ s. Tatsächlich würde in dem Beispiel jeder Grenzwert $d_{seg-max} > 0.5$ s zur gezeigten Segmentierung führen.

Es gibt verschiedene Möglichkeiten, die höchst mögliche Pause $d_{seg-max}$ zu definieren. Der Einfachste Fall, der auch in Abbildung 5.18 angenommen wurde, ist das Setzen eines global festgelegten Grenzwertes. Weitere Möglichkeiten sind, $d_{seg-max}(CS)$ als Funktion des Segmentes selber zu gestalten. So könnte beispielsweise ein längeres Segment eine höhere maximal-Pause erzeugen.

Schlussendlich konnten in der Fachliteratur keine konkreten Hinweise zur Bestimmung von $d_{seg-max}$ gefunden werden. Daher wurde entschieden, dem Arzt, der das System benutzt, selber einstellen zu lassen. Es wird mit einem Festen

5.5 Feature-Extraction

5.6 Ableitung der Schmerz-Scores

5.7 Visualisierung

6 Zusammenfassung

Literaturverzeichnis

- [1] K J S Anand. *Pain in Neonates and Infants*. Elsevier, 2007.
- [2] Zachariah Boukydis Barry Lester. *Infant Crying: Theoretical and Research Perspectives*. Springer, 1985.
- [3] Judy Bildner. *CRIES Instrument Assessment Tool of Pain in Neonates*. City of Hope Pain, 1997. Online unter <http://prc.coh.org/pdf/CRIES.pdf>.
- [4] R Sisto & Giuseppe Buonocore Carlo Bellieni, Franco Bagnoli. Cry features reflect pain intensity in term newborns: An alarm threshold. *Pediatric Research*, 5:142–146, 1. Online unter https://www.researchgate.net/publication/297827342_Cry_features_reflect_pain_intensity_in_term_newborns_An_alarm_threshold.
- [5] Douglas Dankel. The ID3 Algorithm , 1997. Online unter <http://www.cise.ufl.edu/~ddd/cap6635/Fall-97/Short-pa/2.htm>.
- [6] H. Hollien & T Murry E Müller. Perceptual responses to infant crying: identification of cry types. *Journal of Child Language*, 1(1):89–95, 1974. Online unter <https://www.cambridge.org/core/journals/journal-of-child-language/article/perceptual-responses-to-infant-crying-identification-of-cry-types/4F0F8088116FCE381851D8D560697A5F>.
- [7] B. Simak E. Verteletskaya. Performance Evaluation of Pitch Detection Algorithms, 2009. Online unter <http://access.feld.cvut.cz/view.php?cisloclanku=2009060001>.
- [8] Jan Hamers & Peter Gessler Eva Cignac, Romano Mueller. Pain assessment in the neonate using the Bernese Pain Scale for Neonates. *Early Human Development*, 78(2):125–131, 2004. Online unter <http://www.sciencedirect.com/science/article/pii/S0378378204000337>.
- [9] Ricardo Gutierrez-Osuna. Introduction to Speech Processing - Cepstral analysis. Online unter <http://research.cs.tamu.edu/prism/lectures/sp/19.pdf>.
- [10] Health Facts For You. *Using Pediatric Pain Scales Neonatal Infant Pain Scale (NIPS)*, 2014. Online unter <https://www.uwhealth.org/healthfacts/parenting/7711.pdf>.
- [11] Hodgkinson. Neonatal Pain Assessment Tool , 2012. Online unter http://www.rch.org.au/uploadedFiles/Main/Content/rchcpg/hospital_clinical_guideline_index/PAT%20score%20update.pdf.
- [12] Michael J Corwin Howard L Golub. A Physioacoustic Model of the Infant Cry. In *Infant*

- Crying - Theoretical and Research Perspectives*, chapter 3, pages 59 – 82. Plenum, 1985.
- [13] Giorgio Ingargiola. Building Classification Models: ID3 and C4.5. Online unter <http://cis-linux1.temple.edu/~giorgio/cis587/readings/id3-c45.html>.
 - [14] Donna Geiss Laura Wozniak & Charles Hall Ivan Hand, Lawrence Noble. COVERS Neonatal Pain Scale: Development and Validation. *International Journal of Pediatrics*, 2010, 2010. Online unter <https://www.hindawi.com/journals/ijpedi/2010/496719/>.
 - [15] J Gorriz & J Segura J Ramorez. Voice Activity Detection. Fundamentals and Speech Recognition System Robustness. *Robust Speech Recognition and Understanding*, page 460, 2007. Online unter http://cdn.intechopen.com/pdfs/104/InTech-Voice_activity_detection_fundamentals_and_speech_recognition_system_robustness.pdf.
 - [16] Jieh-weih Hung & Lin-shan Lee Jia-lin Shen. Robust Entropy-based Endpoint Detection for Speech Recognition in Noisy Environments. 1998. Online unter https://www.researchgate.net/publication/221489354_Robust_entropy-based_endpoint_detection_for_speech_recognition_in_noisy_environments.
 - [17] Carol Espy-Wilson & Tarun Pruthi Jonathan Kola. Voice Activity Detection. *MERIT BIEN*, 2011. Online unter http://www.ece.umd.edu/merit/archives/merit2011/merit_fair11_reports/report_Kola.pdf.
 - [18] Kim Weaver & Fathi M. Salam Khurram Waheed. A robust Algorithm for detecting speech segments using an entropic contrast. *IEEE*, 2003. Online unter <http://ieeexplore.ieee.org/document/1187039/>.
 - [19] Tze-Wey Loong. Understanding sensitivity and specificity with the right side of the brain. *BMJ*, 327(7417), 2003. Online unter <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC200804/>.
 - [20] M M Homayounpour M H Moattar. A simple but efficient real-time Voice Activity Detection Algorithm. Signal Processing Conference, IEEE, August 2009. Online unter <http://ieeexplore.ieee.org/document/7077834/?arnumber=7077834&tag=1>.
 - [21] Hans M Koot Dick Tibboel Jan Passchier & Hugo Duivenvoorden Monique van Dijk, Josien de Boer. The reliability and validity of the COMFORT scale as a postoperative pain instrument in 0 to 3-year-old infants. *Pain*, 84(2):367—377, 2000. Online unter <http://www.sciencedirect.com/science/article/pii/S0304395999002390>.
 - [22] Taddio Nulman. A revised measure of acute pain in infants. *J Pain Symptom Manage*, 10:456–463, 1995. Online unter [http://geriatricphysio.yolasite.com/resources/Modified%20Behavioral%20Pain%20Scale%20\(MBPS\)%20in%20infants.pdf](http://geriatricphysio.yolasite.com/resources/Modified%20Behavioral%20Pain%20Scale%20(MBPS)%20in%20infants.pdf).
 - [23] J L Mathew P J Mathew. Assessment and management of pain in infants. *Postgrad Med J*, 79:438–443, 2003. Online unter <http://pmj.bmj.com/content/79/934/438.full>.
 - [24] Steven Creech & Marc Weiss. Patricia Hummel, Mary Puchalski. N-PASS: Neonatal

- Pain, Agitation and Sedation Scale – Reliability and Validity. *Pediatrics/Neonatology*, 2(6), 2004. Online unter <http://www.anestesiarianimazione.com/2004/06c.asp>.
- [25] Susan Parker-Price & Ronald Barr Philip Zeskind. Rythmic organization of the Sound of Infant Cry. *Dev Psychobiol*, 26(6):321–333, 1993. Online unter <https://www.ncbi.nlm.nih.gov/pubmed/8119482>.
- [26] David M W Powers. Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness & Correlation. *Journal of Machine Learning Technologies*, 2(1):37–63, 2011. Online unter https://dl.dropboxusercontent.com/u/27743223/201101-Evaluation_JMLT_Postprint-Colour.pdf.
- [27] R Ward & C Laszlo Qiaobing Xie. Automatic Assessment of Infants’ Levels-of-Distress from the Cry Signals. *IEEE Transanctions on Speech and Audio Processing*, 4(4):253–265, 1996. Online unter <http://ieeexplore.ieee.org/document/506929/>.
- [28] Yizhar Lavner Rami Cohen. Infant Cry Analysis and Detection. IEEE 27-th Convention of Electrical and Electronics Engineers in Israel, 2012. Online unter <http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=6376996>.
- [29] Brian Hopkins & James Green Ronald Barr. *Crying as a Sign, a Symptom, and a Signal*. Mac Keith Press, 2000.
- [30] J R Shayevitz & Shobha Malviya Sandra Merkel, Terri Voepel-Lewis. The FLACC: A Behavioral Scale for Scoring Postoperative Pain in Young Children. *Pediatric Nursing*, 23(3):293–7, 1996. Online unter https://www.researchgate.net/publication/13998379_The_FLACC_A_Behavioral_Scale_for_Scoring_Postoperative_Pain_in_Young_Children.
- [31] Andreas Spanias Sassan Ahmadi. Cepstrum-Based Pitch Detection Using a New Statistical V/UV Classification Algorithm. *IEEE Transactions on Speech and Audio Detection*, 7(3):333–338, 1999. Online unter <http://ieeexplore.ieee.org/document/759042/>.
- [32] Julius Smith. *Spectral Audio Signal Processing*. Center for Computer Research in Music and Acoustics (CCRMA), 1993. Online unter https://www.dsprelated.com/freebooks/sasp/Short_Time_Fourier_Transform.html.
- [33] Sabine Deligne & Peder Olsen Trausti Kristjansson. Voicing Features for Robust Speech Detection. Interspeech Lisboa, September 2005. Online unter <http://papers.traustikristjansson.info/wp-content/uploads/2011/07/KristjanssonRobustVoicingEurospeech2005.pdf>.
- [34] P H Wolff. The role of biological rhythms in early psychological development. *Bulletin of the Menninger Clinic*, 31:197–218, 1967.

Appendices

Tabelle .1: Accuracy-Werte der Grenzwertfindung mit REPTree

$S_{Training}$	3 dB				50 dB				50+3 dB			
A_{Test}	3 dB	50 dB	7 dB*	Mean	3 dB	50 dB	7 dB*	Mean	3 dB	50 dB	7 dB*	Mean
Zeit	77.81%	79.02%	86.04%	80,96%	49.33%	94.70%	48.66%	64,23%	77.54%	92.47%	84.38%	84,80%
Freq	82.05%	89.28%	82.71%	84,68%	70.52%	94.37%	55.06%	73,31%	81.75%	91.22%	74.90%	82,62%
Ceps	88.98%	94.72%	92.96%	92,22%	86.83%	94.68%	92.83%	91,45%	88.98%	94.72%	92.96%	92,22%
Corr	80.45%	73.47%	84.89%	79,60%	73.07%	87.14%	77.98%	79,39%	77.90%	84.88%	82.84%	81,87%
Zeit+Freq	82.05%	89.28%	82.71%	84,68%	70.52%	94.37%	55.06%	73,31%	81.75%	91.22%	74.90%	82,62%
Zeit+Ceps	88.98%	94.72%	92.96%	92,22%	86.83%	94.68%	92.83%	91,45%	88.98%	94.72%	92.96%	92,22%
Zeit+Corr	80.45%	73.47%	84.89%	79,60%	49.33%	94.70%	48.66%	64,23%	80.32%	92.35%	88.22%	86,96%
Freq+Ceps	88.98%	94.72%	92.96%	92,22%	70.65%	94.75%	55.06%	73,49%	88.98%	94.72%	92.96%	92,22%
Freq+Corr	82.05%	89.28%	82.71%	84,68%	70.52%	95.60%	95.60%	87,24%	81.75%	94.42%	74.90%	83,69%

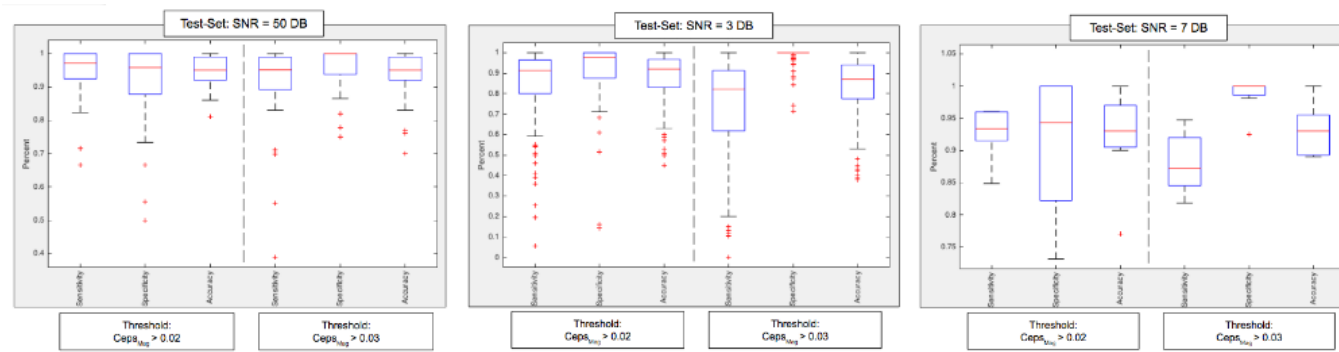


Abbildung .1: Boxplot-Auswertung über Sensitivity, Specificity und Accuracy der beiden VAD-Modelle