

Biometrics and Security Speaker Recognition

Jonas Marquardt und Maik Riestock

Otto von Guericke University Magdeburg - Advanced Multimedia and Security Lab (AMSL)

speaker reco
orientiert am hyke(database)
Closed set speaker authentication
projektion auf 'Doddingtons Zoo'

Categories and Subject Descriptors:

Additional Key Words and Phrases: Speaker Recognition, Audio Feature Extraction, Doddingtons Zoo

1. MOTIVATION

Um Personen eindeutig zu identifizieren gibt es verschiedene Verfahren. Eines davon ist die Erkennung der Stimme. Als Aufnahmegerät ist ein handelsübliches Mikrophone ausreichend. In dieser Übung ging es darum, herauszufinden wie eine Stimmenerkennung umgesetzt wird und welche Eigenheiten dieses Verfahren mit sich bringt.

2. HYKE-SYSTEM

Da sich unter projekt von dem des Hyke ableiten wollen wir zuerst einmal im Folgenden das Hyke Projekt vorstellen.

Das Hyke-System hat seinen Ursprung in der Region Rajasthan in Nordwest Indien. In dieser indischen Region sind ausschließlich kleine Schulen anzufinden die meist aus 1 bis 3 Klassen bestehen. Dennoch unterstehen sie der behördlichen Bildungseinrichtung der Region die sicherstellen möchte, dass in den Schulen ein regelmäßiger Unterricht stattfindet.

Das bisherige System basierte auf eine visuelle Überprüfung der Anwesenheit. Dazu waren die Lehrer angehalten zwei mal täglich ein Bild von sich aufzunehmen und an die zentrale Einrichtung zu schicken. Dort wurde dieses Bild von Mitarbeitern verwendet um eine manuelle Authentifizierung vorzunehmen. Der Nachteil dieses Systems bestand in den hohen Kosten die es erzeugte, zum einen durch das Vorhandensein einer Kamera an jedem Standort und zum anderen durch die manuelle Überprüfung der Aufnahmen.

Aus diesem Grund entschloss man sich für die Entwicklung von Hyke. Ein System welches die Lehrer anhand ihrer Stimme authentifiziert. Dies hat den Vorteil, dass kaum Anschaffungen der Aufnahmegeräte von Nuten waren, da bereits 75% der Schulen bereits ein Telefon verfügten. Außerdem ist nun eine automatisierte Authentifizierung der Lehrer möglich. Lehrer möglich sein.

This report was created in the context of the course Biometrics and Security [BIOSEC] winter term 2014/15. This course was held by: Prof. Dr.-Ing. Jana Dittmann and Prof. Dr.-Ing. Claus Viehauer; Research group Multimedia and Security, Otto-von-Guericke-University of Magdeburg, Germany. The course was supported by: Dr.-Ing. Christian Krtzer, M.Sc. Kun Qian

3. UNSER ANSATZ

Das Thema unseres Projektes ist die **Speaker Recognition**, wobei wir anhand von Aufnahmen der Stimme von Personen versuchen diese zu authentifizieren. Genauer definiert ist die Speaker Recognition folgendermaßen:

Speaker recognition, sometimes referred to as speaker biometrics, includes identification, verification (authentication), classification, and by extension, segmentation, tracking and detection of speakers. It is a generic term used for any procedure which involves knowledge of the identity of a person based on his/her voice. [Beigi 2011]

In dem Rahmen des Projektes sollten die folgenden Aufgabenstellungen erfolgreich bearbeitet werden:

- Closed set speaker authentication on the Hyke speech database
- Compare the results achieved (in terms of authentication performance) to the results presented in [Azarias Reda 2011]
- A projection of the samples in your data set to the characters of 'Doddingtons Zoo' [George Doddington 1998]

Bei den Arbeitsschritten in unserem Projekt haben wir uns an das allgemeine Modell für das Authentifizieren von Benutzern anhand von Biometrischen Daten orientiert, welches in [?] vorgestellt wird. Dieses Modell haben wir an unser Thema, der *Speaker Recognition*, und der Aufgabenstellung angepasst:

- Data Acquisition: Verwendung der Hyke Databasis
- Pre-processing: Einteilung der Datenbasis in geschlechtsspezifischen Sets
- Feature Extraction: Merkmals-Extraktion der Audio-Daten mittels AAFE.
- Post-processing: Aufarbeitung der Feature-Matrizen
- Comparison and Classification: Klassifizierung der Samples und Authentifizierung der Sprecher

Zur Bearbeitung des Projektes wurden uns zwei Programme bereitgestellt:

- **AAFE** (AMSL Audio Feature Extractor), ist ein Tool für die Extraktion von Merkmalen in Audio-Dateien und entstammt dem AMSL Audio Steganalysis Toolset (AASST). [?] Anwendung fand das AAFE-Tool in dem Kapitel: 6.
- **WEKA**, ist eine Sammlung von Algorithmen des Maschinellen Lernens für Aufgaben im Bereich des Data-Mining. [?] Anwendung fand dieses Tool in den Kapiteln: 7 und 8.

4. DATENBASIS

Die Datenbasis wurde dem Hyke-Projekt entnommen. Sie kann unter folgender URL heruntergeladen werden: Sie umfasst Aufnahmen von 83 verschiedenen Sprechern, davon 48 männlich und 35 weiblich. Von jeder Person gibt es fünf Aufnahmen in denen Abfolgen verschiedener Ziffern gesprochen werden. Die Sprache dabei

ist Englisch. Die Länge der Aufnahmen liegt zwischen 5 und 35 Sekunden. Es gibt auch Aufnahmen, die keine Stimme enthalten. Die Stimmen wurden über das Telefon aufgenommen und bieten daher eine geringere Bandbreite als die menschliche Stimme hat. Bei den Sprechern handelt es sich um Kinder mit verschiedenen Hintergründen. Die Aufnahmen enthalten teilweise Hintergrundgeräusche, vom leisen Rauschen bis zu Gesprächen und Musik.

5. VORVERARBEITUNG

In diesem Kapitel geht es um die Daten unserer Datenbank auf die folgenden Schritte vorzubereiten.

Da wir in unseren Ergebnissen am Ende einen möglichen Unterschied zwischen den Ergebnissen der Authentifizierung beider Geschlechter beobachten zu können, wurden die Datenbank in sechs Sets unterteilt. Hierfür wurde das Set mit Sprechern von beiden Geschlechtern, *mixed set*, aufgeteilt in zwei Sets mit ausschließlich Stimmen von weiblichen Sprechern, *female set*, und mit ausschließlich männlichen Sprechern, *male set*.

Zusätzlich benötigen die Klassifikatoren zwei verschiedene Sets von Daten. Mit dem einen Set wird das Modell trainiert, hier *train set*, und mit dem anderen Set evaluiert, hier *test set*.

Die Datenbasis aufgeteilt in folgende Sets:

- mixed train set
- mixed test set
- female train set
- female test set
- male train set
- male test set

6. FEATURE EXTRACTION

Um aus den Aufnahmen die Features zu extrahieren wurde der *AMSL Feature Extractor* verwendet. Dieser zerlegt eine Audiodatei in sehr kurze Samples und berechnet aus diesen verschiedene Features. Die Länge der Samples kann man frei wählen. Wir haben eine Länge von 1024 gewählt und eine Überlappung von Null. Die Hanning-Fenster-Funktion wurde aktiviert. Der Feature Extractor berechnet aus jedem Sample 593 verschiedene Features berechnet.

7. NACHTVERARBEITUNG

Die extrahierten Daten wurden mit Hilfe von Weka aufbereitet. Dadurch sollten bessere Ergebnisse bei der Klassifikation erzeugt werden. Dazu wurde die Features *lbs flipping ratio* (in allen Instanzen 922337203685477.6000) und *lbs flipping rate* (in allen Instanzen 0) entfernt. Weil sie in allen Fällen gleich sind lassen sich an ihnen keine Unterschiede in den Aufnahmen feststellen.

In den Aufnahmen gibt es Bereiche, die keine Stimme enthalten. Diese konzentrieren sich auf Anfang und Ende der Datei. Es gibt auch Pausen zwischen den gesprochenen Ziffern. Die "stillen" Bereiche enthalten keine Information über die Stimme und somit den Sprecher. Dadurch wird die spätere Klassifikation erschwert. Um die "Stille" herauszufiltern wurden alle Samples mit einer geringen Amplitude gelöscht. Dazu wurde das Feature *rms amplitude* genutzt und alle Samples mit einem Wert unter 10 gefiltert. Es wurde der *RemoveWithValues* Filter von Weka mit den Parametern -S 10.0 -C 5 -L first-last verwendet. Dadurch wurden von 50.424 Samples 32.026 entfernt. Das heißt es wurden rund 64 Prozent der Datenbasis entfernt.

8. KLASSIFIKATION

Diesem Kapitel beschäftigt sich mit der Klassifikation unserer Daten, also die richtige Zuordnung der Samples zu den Sprechern. Für diese Aufgabe haben wir das Tool WEKA verwendet.

Dabei sind unsere Ausgangsdaten für die Klassifikatoren die aufbereiteten Features-Matrizen, welche in Kapitel 7 vorgestellt wurden und die wir in den folgenden Sets unterscheiden:

- mixed train set as feature matrix
- mixed test set as feature matrix
- female train set as feature matrix
- female test set as feature matrix
- male train set as feature matrix
- male test set as feature matrix

Zur Bestimmung des besten Klassifikators haben wir die Methode *try and error* verwendet. Das heißt wir haben alle anwendbaren Klassifikatoren in der Standardeinstellung auf unsere Datenbasis angewandt und danach die Ergebnisse verglichen.

Ein gutes Ergebnis bestand darin, dass möglichst viele Samples eines Sprechers dem richtigen Sprecher zugeordnet wurden. Also der Klassifikator unter Verwendung des *female/male/mixed test set* eine gute Vorhersagegenauigkeit aufwies.

Dabei hat sich ein Klassifikator als besonders gut erwiesen, der IBK. Dieser Klassifikator erzielte ein Ergebnis von 54.94 % Vorhersagegenauigkeit bei dem *mixed test set*. Als Vergleich haben wir den Klassifikator mit dem zweitbesten Ergebnis mit aufgeführt, der RandomForest. Die Tabelle 1 zeigt das Ergebnis beider Klassifikationen mit den dazugehörigen Konfigurationen des Klassifikators. Die vollständigen Ergebnisse sind im Anhang zu finden.

Table 1. Ergebnisse der Klassifikation des IBK und RandomForest

Datenset	IBK	RandomForest
female test set	54.94%	39.14%
male test set	58.50%	41.7679%
mixed test set	53.16%	33.86%
Konfiguration	-K 1 -W 0 -A	-I 10 -K 0 -S 1

Aus den Ergebnissen ist zu entnehmen, dass es keinen signifikanten Unterschied zwischen den Ergebnissen der geschlechtsspezifischen Sets *female test set* und *male test set* existiert. Der bestehende Unterschied lässt sich aus der geringen Größe der Datenbasis erklären.

Außerdem ist zu beobachten, dass sich die Ergebnisse beider Klassifikatoren verschlechtert haben bei Erhöhung der Anzahl von Sprechern. Dies ist jedoch ein zu erwartendes Ergebnis, da der Klassifikator nun das Sample eines Sprechers mit 82 anderen Samples statt mit 47 bzw. 34 anderen Samples vergleichen muss.

8.1 Authentifizierung

In diesem Kapitel geht es nun um die Aufgabe der *Closed Set Speaker Authentication*. Wobei die akustische Aufnahme eines Sprechers mit der aller anderen möglichen Sprechern verglichen wird und die beste Übereinstimmung als Ergebnis ausgegeben wird. [Beigi 2011] Zu beachten ist, dass hier im Gegensatz zu

der *Open Set Speaker Authentication* es in jedem Fall zu einem Ergebnis kommt.

Die Aufgabe besteht nun darin das Ergebnis der Klassifikation der Samples zu interpretieren. Dafür betrachteten wir jeden Sprecher die Verteilung seiner Samples. Hierbei wurde ein Sprecher richtig erkannt, wenn bei ihm die größte Menge an Samples zugeordnet wurden. Dies bedeutet, dass wir auch mit einer geringen Anzahl von richtig klassifizierter Samples einen Sprecher erfolgreich authentifizieren konnten solange die richtigen, falsch klassifizierten, Samples gleichmäßig verteilt waren.

Dieses Verfahren wurde bei allen Sprechern angewandt und das erarbeitete Ergebnis ist zu sehen in Tabelle:II.

Table II. Ergebnisse der Authentifizierung

Datenset	Gesamt	Richtig	Falsch	Anteil
female test set	35	33	2	94.29%
male test set	48	46	2	95.83%
mixed test set	83	79	4	95.18%

Das Ergebnis von 95.18% richtig erkannten Sprechern ist gut und entspricht damit dem *state-of-the-art*. [Beigi 2011] Im Vergleich dazu wurde im HYKE-Projekt ein Ergebnis von 95% erreicht, welches mit unseren nahezu identisch ist. [Azarias Reda 2011]

9. DODDINGTONS ZOO

Bei 'Doddingtons' Zoo geht es um die Beobachtung, dass Sprecher ein unterschiedliches Verhalten bezüglich den Erfolg ihrer Authentifizierung aufzeigen. [George Doddington 1998] Dadurch lassen sich Sprecher in vier Kategorien unterscheiden die jeweils von einem Tier repräsentiert werden.

Beschreibung dieser Kategorien: [Prof. Dr. Jana Dittmann 2014]

- **Sheeps:** außerordentlich *leicht* von dem System erkannt, die Mehrheit der Sprecher gehört dieser Kategorie an
- **Goats:** außerordentlich *schwer* von dem System erkannt
- **Lambs:** außerordentlich *verwundbar* gegenüber Nachahmung
- **Wolves:** außerordentlich *erfolgreich* bei der Nachahmung anderen Sprecher

Nun sollte eine Projektion dieser Kategorien auf die Ergebnisse unserer Klassifizierung vollzogen werden. Als Entscheidungsgrundlage dienen nun nicht nur die erfolgreich klassifizierten Samples sondern auch die Verteilung der falsch klassifizierten Samples. Um die Projektion umzusetzen haben wir folgendes Schema erarbeitet und auf unsere Datenbasis angewandt:

- **Sheeps:** viele richtig klassifizierte Samples
- **Goats:** wenig richtig klassifizierte Samples
- **Lambs:** viele Samples von anderen Sprechern wurden diesem Sprecher zugeordnet
- **Wolves:** viele Samples bei wenigen anderen Sprechern zugeordnet

Dieser Ansatz wurde auf das Ergebnis der Klassifikation mit dem Klassifikator IBK auf das *mixed test set* angewandt. Das Ergebnis dieser Projektion ist zu sehen in Tabelle:III.

Table III. Ergebnisse der Kategorisierung nach Doddingtons Zoo

Animal	female	male	mixed	Anteil
Sheep	32	43	75	90.36%
Goat	2	2	4	4.82%
Lamb	1	1	1	2.41%
Wolf	0	2	2	2.41%

Das Ergebnis der Kategorisierung zeigt, dass der Großteil unserer Sprecher **Sheeps** sind. Dies entspricht den Erwartungen, da es sich dabei um den Standardtypen handelt. Wohingegen die anderen Kategorien eine Minderheit darstellen. Außerdem ist zu beobachten, dass unsere falsch authentifizierten Sprecher zu der Kategorie der **Goat** zugeordnet worden. Dies entspricht ganz ihrer Beschreibung als Kategorie der schwer zu authentifizierenden.

10. ZUSAMMENFASSUNG

Es ist möglich einen Menschen anhand seiner Stimme zu identifizieren. Dies eröffnet Anwendungsbereiche, die mit anderen biometrischen Verfahren nicht möglich sind. Ein Beispiel ist die Identifizierung einer Person über das Telefon. Wir konnten in unseren Experimenten, mit geringem Aufwand, 96,68 Prozent der Personen eindeutig identifizieren. Bei den 3,32 Prozent der nicht identifizierten lag eine schlechte Datenbasis vor. Das heißt der Erfolg bei der Identifizierung hängt signifikant von der Datenbank ab.

11. FUTURE WORK

In diesem Kapitel wollen wir mögliche weiterführende Arbeiten an unserem Projekt Ansätze aufzeigen. Da es mehrere interessante Aspekte gibt die wir aus Zeitmangel leider nicht realisieren konnten. Dabei sollte als erstes ein Wechsel der Datenbasis in Betracht gezogen werden, um zu untersuchen ob die erreichten Ergebnisse sich bestätigen lassen.

Ein Ansatzpunkt wäre ein Vergleich der Ergebnisse ohne eine Entfernung der "Stille" in den Aufnahmen, um herauszufinden ob und wenn ja wie signifikant das Ergebnis verändert wurde. Zusätzlich dazu wäre auch interessant welche Auswirkung die Reihenfolge der Filterung der "Stille". Ob das Ergebnis beeinflusst wird, wenn die Filterung nach oder vor der Feature Extraktion vorgenommen wird, da man von der Feature Extraktion noch die Audio-Daten und keine Feature-Matrizen zur Verfügung stehen.

Erweitern könnte man das gesamte Projekt dahingehend, dass man die erarbeitete *Processing-Pipeline* in ein echtzeitfähiges Framework einbettet. Damit wäre eine unmittelbare *user authentication* durch Spracherkennung möglich. Hierfür müsste vorher ein *Reference Storage* erstellt werden mit denen die aktuell eingehenden Authentifizierungs-Daten verglichen werden können.

Ein anderer Aspekt ist die Evaluation der Robustheit unserer Verarbeitungskette. Dies könnte man erreichen indem negativen Einfluss auf die Datenbasis genommen wird. folgende Attacken sind in Betracht zu ziehen:

- schlechtere Aufnahmegeräte (kleinere Bandbreite als die bisherige: 300Hz bis 3400Hz)
- Einfügen von Rauschen in die Aufnahmen
- Einfügen von Hintergrundstimmen in die Aufnahmen
- Cropping, um kurzzeitige Verbindungsabbrüche bei der Aufnahme zu simulieren

REFERENCES

Edward Cutrell Azarias Reda, Saurabh Panjwani. 2011. Hyke: A Low-cost Remote Attendance Tracking System for Developing Regions. *Networked System for Developing Regions* (2011).

Homayoon Beigi. 2011. *Fundamentals of Speaker Recognition*. Springer Science+Business Media.

Alvin Martin-Mark Przybocki Douglas Reynolds George Doddington, Walter Liggett. 1998. SHEEP, GOATS, LAMBS and WOLVES - A Statistical Analysis of Speaker Performance in the NIST 1998 Skeaper Recognition Evaluation. *National Institute of Standards and Technology* (1998).

Prof. Dr.-Ing. Claus Vielhauer Prof. Dr. Jana Dittmann. 2014. *Biometrics and Security - Lecture*. Faculty of Computer Science, Institute of Technical and Business Information Systems, Advanced Multimedia and Security Lab (AMSL).

List of Tables

I	Ergebnisse der Klassifikation des IBK und RandomForest	2
II	Ergebnisse der Authentifizierung	3
III	Ergebnisse der Kategorisierung nach Doddingtons Zoo	3

APPENDIX

A. TASK DESCRIPTION

Run your prototype on the collected data and perform a performance evaluation with your prototype. The evaluation must include:
The evaluation must include:

- Closed set speaker authentication on the Hyke speech database
- Compare the results achieved (in terms of authentication performance) to th results persented in
- A projection of the samples in your data set to the characters of 'Doddingtons Zoo'

```

female ibk test.txt
0.7      0.005      0.7      0.7      0.7      0
0.674    0.015    0.644    0.674    0.659    0
0.462    0.01    0.621    0.462    0.529    0
0.294    0.017    0.208    0.294    0.244    0
0.565    0.005    0.722    0.565    0.634    0
0.457    0.033    0.308    0.457    0.368    0
0.619    0.009    0.565    0.619    0.591    0
Weighted Avg. 0.549 0.014 0.558 0.549 0.548 0

Scheme:weka.classifiers.lazy.IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A "weka.core.EuclideanDistance -R first-last
Relation: audio_steganalysis-weka.filters.unsupervised.instance.RemoveWithValues -S10 -C5 -Lfirst-last
Instances: 6950
Attributes: 591
[list of attributes omitted]
Test mode:user supplied test set: size unknown reading incrementally

=== Confusion Matrix ===

=== Classifier model full training set ===
a b c d e f g h i j k l m n o p q r s t u v w
19 0 0 0 0 1 1 0 0 0 0 0 0 0 2 0 5 3 0 1 0 1 0
IB1 instance-based classifier
using 1 nearest neighbours for classification
0 0 20 0 0 0 2 0 0 0 1 0 0 0 0 0 0 1 1 0 0 0 5
0 0 0 8 0 0 0 0 0 1 1 0 0 1 0 0 0 0 0 0 0 0 0
1 3 1 0 12 2 0 0 5 0 0 1 1 0 0 0 1 0 4 2 0 0
0 0 0 0 0 36 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 1 0 0 2 0 13 0 1 0 0 0 1 0 1 0 0 1 0 0 0 0
2 0 0 0 0 0 3 26 0 0 0 0 0 0 0 1 0 0 2 0 0 0
0 0 0 0 1 0 2 0 11 0 0 0 1 0 3 0 0 1 0 3 0 1
0 0 1 0 0 0 0 0 0 22 0 0 0 0 2 0 0 0 0 0 0 0
Correctly Classified Instances 612 54.9372%
Incorrectly Classified Instances 502 45.0628%
Kappa statistic 0.5351
Mean absolute error 0.0259
Root mean squared error 0.1601
Relative absolute error 46.7266 %
Root relative squared error 96.118 %
Total Number of Instances 1114

=== Detailed Accuracy By Class ===
TP Rate FP Rate Precision Recall F-Measure ROC Area Class
0.475 0.02 0.475 0.475 0.475 0.0728 female_01
0.619 0.017 0.591 0.619 0.605 0.1801 female_02
0.606 0.006 0.741 0.606 0.667 0.0810 female_03
0.615 0.005 0.615 0.615 0.615 0.8050 female_04
0.3 0.029 0.279 0.3 0.289 0.0638 female_05
0.923 0.004 0.9 0.923 0.911 0.0960 female_06
0.5 0.014 0.464 0.5 0.481 0.0743 female_07
0.703 0.008 0.743 0.703 0.722 0.0849 female_08
0.333 0.016 0.393 0.333 0.361 0.0659 female_09
0.88 0.007 0.733 0.88 0.8 0.0938 female_10
0.571 0.016 0.4 0.571 0.471 0.0779 female_11
0.548 0.01 0.607 0.548 0.576 0.0769 female_12
0.516 0.01 0.593 0.516 0.552 0.0753 female_13
0.783 0.004 0.818 0.783 0.8 0.889 female_14
0.25 0.016 0.393 0.25 0.306 0.617 female_15
0.667 0.005 0.783 0.667 0.72 0.831 female_16
0.298 0.01 0.56 0.298 0.389 0.644 female_17
0.605 0.022 0.489 0.605 0.541 0.791 female_18
0.444 0.011 0.5 0.444 0.471 0.717 female_19
0.405 0.03 0.319 0.405 0.357 0.688 female_20
0.286 0.018 0.286 0.286 0.286 0.634 female_21
0.276 0.013 0.364 0.276 0.314 0.631 female_22
0.762 0.019 0.615 0.762 0.681 0.872 female_23
0.343 0.021 0.343 0.343 0.343 0.661 female_24
0.85 0.004 0.81 0.85 0.829 0.923 female_25
0.512 0.023 0.457 0.512 0.483 0.744 female_26
0.762 0.002 0.889 0.762 0.821 0.88 female_27
0.891 0.01 0.788 0.891 0.837 0.941 female_28

```