



Hochschule Konstanz
Technik, Wirtschaft und Gestaltung

Signale, Systeme und Sensoren

Aufbau eines einfachen Spracherkenners

Gina Kokoska, Simon Winter

Konstanz, 9. Januar 2022

Zusammenfassung (Abstract)

Thema:	Aufbau eines einfachen Spracherkenners	
Autoren:	Gina Kokoska	gi161kok@htwg-konstanz.de
	Simon Winter	si411win@htwg-konstanz.de
Betreuer:	Prof. Dr. Matthias O. Franz	mfranz@htwg-konstanz.de
	Jürgen Keppler	juergen.keppler@htwg-konstanz.de
	Daniel Castelo	dacastel@htwg-konstanz.de

In diesem Versuch wird ein einfacher Spracherkenners erstellt, der z.B. zur Steuerung eines Staplers in einem Hochregallager dienen könnte. Hierfür werden die Steuerungsbefehle "Hoch", "Tief", "Links" und "Rechts" aufgenommen. Zur Erkennung der Befehlswörter werden die jeweiligen Spektren mit der Windowing-Methode berechnet.

Inhaltsverzeichnis

Abbildungsverzeichnis	III
Tabellenverzeichnis	IV
Listingverzeichnis	V
1 Einleitung	1
2 Versuch 1	2
2.1 Fragestellung, Messprinzip, Aufbau, Messmittel	2
2.2 Messwerte	3
2.3 Auswertung	4
2.4 Interpretation	4
3 Versuch 2	5
3.1 Fragestellung, Messprinzip, Aufbau, Messmittel	5
3.2 Messwerte	6
3.3 Auswertung	7
3.4 Interpretation	7
Anhang	9
A.1 Quellcode	9
A.1.1 Quellcode Versuch 1a und 1b	9
A.1.2 Quellcode Versuch 1c und 1d	11
A.1.3 Quellcode Versuch 2	14
Literaturverzeichnis	19

Abbildungsverzeichnis

Tabellenverzeichnis

Listingverzeichnis

LaTeX_Template/src/aufnahmetriggered.txt	9
LaTeX_Template/src/PlotA1.txt	11
LaTeX_Template/src/rechnungreferenzspektrum.txt	14
LaTeX_Template/src/plotreferenzspek.txt	16
LaTeX_Template/src/korrel.txt	17

Kapitel 1

Einleitung

Diese Versuchsreihe stellt mit Hilfe der Windowing-Methode Referenzspektren der aufgenommenen Steuerbefehle dar. Anschliessend wird der Datensatz mit Steurbefehlen eines anderen Sprechers verglichen. Zur Vergleichsauswertung wird die Korrelation beider Spektren berechnet. Eine starke Korrelation ist als erfolgreiche Spracherkennung zu deuten.

[1] [2]

Kapitel 2

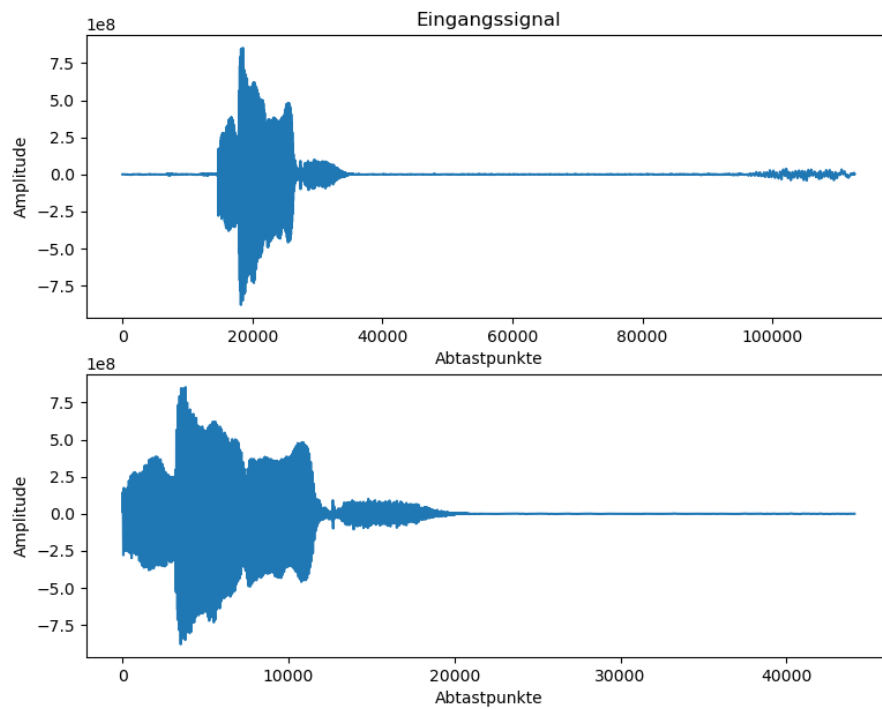
Versuch 1

2.1 Fragestellung, Messprinzip, Aufbau, Messmittel

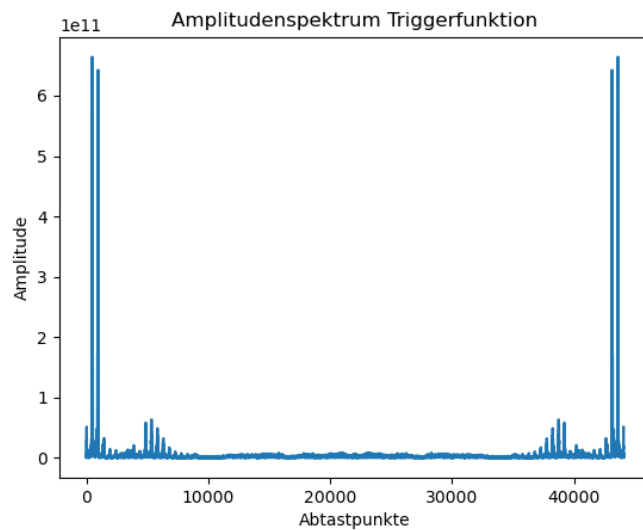
Es wird ein Python-Programm zur Aufnahme akustischer Signale verfasst, welches die aufgenommene Sequenz als wave-Datei und als Plot speichert. Damit die Aufnahmen im Folgeversuch verglichen werden können, wird die Aufnahmezeit auf eine Sekunde beschränkt, mit Bedingungen einer Amplitudenüberschreitung von 150000000. Das getriggerte Signal wird auch als wave-Datei gespeichert und graphisch dargestellt. Anschliessend wird mit Hilfe der Fouriertransformation das Amplitudenspektrum des getriggerten Signals berechnet und dargestellt. Zuletzt wird das Signal in Samples der Länge 512 zerlegt, die sich zur Hälfte überlappen sollen. Jedes Fenster wird mit einer gaußschen Fensterfunktion gewichtet, mit entsprechender Fensterbreite von 4 Standardabweichungen. Erst wird eine lokale Fouriertransformation auf die einzelnen Fenster angewendet. Anschliessend wird die gemittelte Fouriertransformation aller Fenster berechnet. Das Amplitudenspektrum wird berechnet und mit dem getriggertem Spektrum verglichen. Das verwendete Kleinmembran-Kondensatormikrofon ist in ein Callcenter-Headset eingebaut, welches Noise-cancelling verwendet. Es verwendet einen elektrostatischen Wandler, um Luftschall in eine elektrische Spannung um zu wandeln. Es nutzt das Prinzip der Kapazität, in welchem der Abstand einer leitfähigen Membran zur Gegenelektrode gemessen wird. Hierbei wird eine polarisierte Kapsel bestehend aus einer elastischen Membran und einer Gegenelektrode verwendet. Diese bilden die Platten eines Kondensators der Gleichspannung als elektrische Kapazität speichern kann. Auftreffender Schall bringt die Membran zum Schwingen, wodurch sich der Abstand zwischen Membran und Gegenelektrode und damit auch die Kapazität des Kondensators verändert. Je geringer der Abstand desto grösser der eintreffende Schall (antiproportionales Signal). Als Messmittel wird ein Plantronics EncorePro HW725 Headset verwendet.

2.2 Messwerte

Das Programm lieferte folgendes original (oben) und getriggertes (unten) Signal:

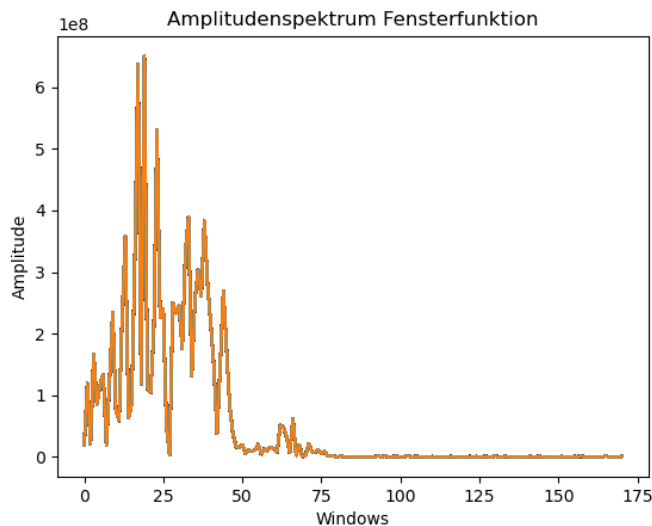


Das Programm lieferte folgende Amplitudenspektrum:



2.3 Auswertung

Das Windowing lieferte folgendes Amplitudenspektrum:



2.4 Interpretation

Schaut man sich die Plots des Eingangssignals sowie des getriggerten Eingangssignals an, so erkennt man die erfolgreiche Zuschneidung des Signals. Anhand der x-Achse kann man die Abtastpunkte auslesen, welche für das getriggerte Signal genau eine Sekunde ergeben. Die Amplitudenüberschreitung ist so eingestellt, das ein eingehendes Signal frühzeitig erkannt wird, ohne das etwas vom gewollten Signal abgeschnitten wird. Vergleicht man das Amplitudenspektrum der Fensterfunktion mit dem graphischen Plot des zugeschnittenen Eingangssignals, so lässt sich eindeutig erkennen, das sich beide Plots in ihrer Form sehr ähnlich sehen. Somit kann man das Amplitudenspektrum der Fensterfunktion eindeutig dem Signal des Eingangssignals zuordnen.

Kapitel 3

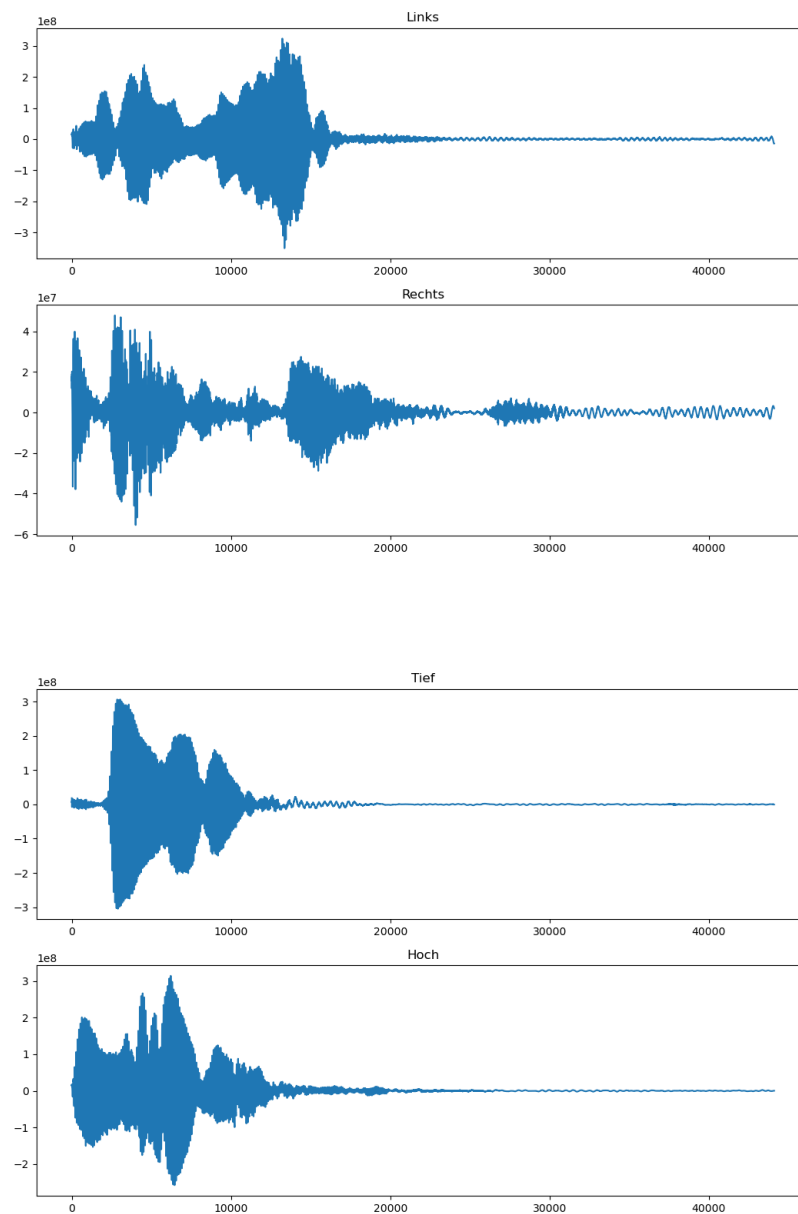
Versuch 2

3.1 Fragestellung, Messprinzip, Aufbau, Messmittel

Für den Spracherkenner werden jeweils 5 Referenzspektren der Steuerungsbeefehle "Hoch", "Tief", "Links" und "Rechts" mit Hilfe des Programms aus Aufgabe 1 aufgenommen, gespeichert und verglichen. Der Sprecher ist für alle Aufnahmen der selbe. Das gemittelte Referenzspektrum jedes Befehls wird graphisch dargestellt. Anschliessend wird ein Testdatensatz eines anderen Sprechers und des original Sprechers erstellt. Beide Datensätze werden durch Berechnung der Korrelationskoeffizienten nach Bravais-Pearson mit einander verglichen. Es ist eine Korrelation nahe 1 bei den selben Befehlswörter zu erwarten, da die identischer Wörter Spektren sehr ähnlich sein sollten. Abweichungen sollten durch Messfehler wie unterschiedliche Sprechweisen erklärbar sein. Zuletzt wird der Spracherkenner anhand der in der Vorlesung behandelten Architektur implementiert. Beide Datensätze werden nochmals mit einander verglichen, um die prozentual erfolgreiche Erkennungsrate und falsch erkannte Fehlerrate zu berechnen. Als Messmittel wird das Headset aus Versuch 1 verwendet.

3.2 Messwerte

Das Programm lieferte folgende Referenzspektren:



3.3 Auswertung

Die Vergleiche zwischen unseren Referenzspektren und unseren Eingabespektren ergeben folgende Werte:

Korrelation Datensatz(Ref) mit original Sprecher G

$\text{korrkoefLLG} = \text{sst.pearsonr}(\text{LRef}, \text{LG})[0] = 0.024579713287452655$

$\text{korrkoefLRG} = \text{sst.pearsonr}(\text{LRef}, \text{RG})[0] = 0.015589028822441266$

$\text{korrkoefLTG} = \text{sst.pearsonr}(\text{LRef}, \text{TG})[0] = 0.011521206726410545$

$\text{korrkoefLHG} = \text{sst.pearsonr}(\text{LRef}, \text{HG})[0] = -0.02717198805765672$

Erkanntes Wort: Links

Korrelation Datensatz(Ref) mit zweitem Sprecher E

$\text{korrkoefLLE} = \text{sst.pearsonr}(\text{LRef}, \text{LE})[0] = 0.017434312307556363$

$\text{korrkoefLRE} = \text{sst.pearsonr}(\text{LRef}, \text{RE})[0] = 0.00036481785038843954$

$\text{korrkoefLTE} = \text{sst.pearsonr}(\text{LRef}, \text{TE})[0] = -0.003291010877873985$

$\text{korrkoefLHE} = \text{sst.pearsonr}(\text{LRef}, \text{HE})[0] = 0.0004718242888293782$

Erkanntes Wort: Links

3.4 Interpretation

Die positive, aber geringe Korrelation lässt auf einen funktionierenden Vergleich zwischen Referenzspektrum und Eingabespektrum schließen. Der höchste, wenn auch immer noch schwach korrelierende Wert, wurde bei dem Vergleich Links - Links erkannt.

Eine weitere positive Korrelation zu einem unterschiedlichen Wort wurde festgestellt, diese ist jedoch ca um die Hälfte geringer. Zudem wurde auch eine leicht negative Korrelation errechnet. Der höchste gemessene Korrelationswert wurde bei dem Originalsprecher erreicht, was aufgrund des Referenzspektrums natürlich Sinn macht. Der zweithöchste Wert wurde bei dem Wortvergleich Links - Links beim zweiten Sprecher errechnet, demnach bezeichnen wir unseren Spracherkenner als funktionierend. Die allgemein geringeren Korrelationswerte erklären wir uns mit der Aufnahmequalität des Mikrofons sowie dem eingebau-

ten, automatischen Noise-Cancelling. Wir gehen davon aus, dass bei anderen Mikrofonen stärker positive Korrelationskoeffizienten errechnet werden würden. Es ist anzunehmen, dass der Spracherkenner funktioniert, da der Originalsprecher sowie der Fremdsprecher weiblich waren.

Anhang

A.1 Quellcode

A.1.1 Quellcode Versuch 1a und 1b

```
1 import pyaudio
2 import numpy as np
3
4
5 FORMAT = pyaudio.paInt16
6 SAMPLEFREQ = 44100 # 44 kbps
7 FRAMESIZE = 1024 # each buffer contains 1024 samples
8 NOFFRAMES = 220
9 PATH = "eingangssignal_raw"
10 PATH_T = "referenzdatensatz"
11 PATH_W = "eingangssignal_win"
12 PATH_G = "sprechergina"
13 PATH_L = "sprecherluca"
14
15 p = pyaudio.PyAudio()
16 print('running')
17 # 1a open stream
18 stream = p.open(format=FORMAT, channels=1, rate=SAMPLEFREQ, input=True,
19                 frames_per_buffer=FRAMESIZE)
20 data = stream.read(NOFFRAMES*FRAMESIZE)
21 decoded = np.fromstring(data, np.int)
22 stream.stop_stream()
23 stream.close()
24 p.terminate()
25
26 # 1b Triggerfunktion auf decoded angewendet
27 for i in range(len(decoded)):
28     if decoded[i] > 15000000:
```

```

28     start = i
29     break
30
31 decodedTriggered = decoded[start:start+SAMPLEFREQ]
32
33 def safeTriggeredAsCSV(eingangssignal):
34     csvTrigArray = np.savetxt(PATH_T + "/decodedtrigL1.csv",
35                               eingangssignal, delimiter=",")
36     return csvTrigArray
37
38 safeTriggeredAsCSV(decodedTriggered)
39
40 print('done')

```


A.1.2 Quellcode Versuch 1c und 1d

```
1 import matplotlib.pyplot as plt
2 import pyaudio
3 import numpy as np
4 from scipy import signal
5
6 FORMAT = pyaudio.paInt16
7 SAMPLEFREQ = 44100 # 44 kbps
8 FRAMESIZE = 1024 # each buffer contains 1024 samples
9 NOFFRAMES = 220
10 PATH = "eingangssignal_raw"
11 PATH_T = "referenzdatensatz"
12 PATH_W = "eingangssignal_win"
13
14
15 # oeffne stream
16 p = pyaudio.PyAudio()
17 print('running')
18
19 stream = p.open(format=FORMAT, channels=1, rate=SAMPLEFREQ, input=True,
20                 frames_per_buffer=FRAMESIZE)
21 data = stream.read(NOFFRAMES*FRAMESIZE)
22 decoded = np.fromstring(data, np.int)
23 stream.stop_stream()
24 stream.close()
25 p.terminate()
26
27 # plot decoded und decodedTriggered in einer Abb
28 # plt decodedTriggered in line 75
29 plt.figure(figsize=(9, 7))
30 plt.subplot(211)
31 xvalue = np.arange(0, len(decoded))
32 plt.plot(decoded)
33 plt.title("Grundperiode")
34 plt.xlabel("Zeit in s")
35 plt.ylabel("Amplitude")
36
37 print('done')
38
39
```

```

40 # 1b Triggerfunktion auf decoded angewendet
41 for i in range(len(decoded)):
42     if decoded[i] > 15000000:
43         start = i
44         break
45
46 decodedTriggered = decoded[start:start+SAMPLEFREQ]
47
48
49 # window ueber alle referenz daten
50 windownmr = 171
51
52 windows = np.zeros((windownmr, 512))
53
54 schnitt = 256
55
56 gauss = signal.windows.gaussian(512, 4)
57
58
59 for y in range(0, windownmr):
60     schnitt = schnitt - 256
61     for x in range(0, 512):
62         windows[y, x] = np.mean(np.abs(np.fft.fft(decodedTriggered[
63             schnitt] * gauss)))
64     schnitt = schnitt + 1
65
66 for y in range(0, windownmr):
67     for x in range(0, 512):
68         windows[y] = np.mean(windows[y, x])
69
70
71
72 # plot decodedTriggered
73 # save untriggered vgl mit triggerd
74 deltaT = 1 / SAMPLEFREQ
75 xvalue = np.arange(start, start+SAMPLEFREQ) * deltaT
76 plt.subplot(212)
77 plt.plot(decoded[start:start+SAMPLEFREQ])
78 plt.xlabel("Abtastpunkte")
79 plt.ylabel("Amplitude")
80 plt.savefig('untriggeredANDtriggered.png')

```

```

81
82
83 #1c fft auf triggered Signal fuer amplitudenspektrum
84 ffttriggerd = np.fft.fft(decodedTriggered)
85 specTriggered = np.abs(ffttriggerd)
86 plt.figure()
87 plt.plot(specTriggered)
88 plt.title('Amplitudenspektrum Triggerfunktion')
89 plt.xlabel('Abtastpunkte')
90 plt.ylabel('Amplitude')
91 plt.savefig('triggeredAmplitudenspektrum.png')
92
93 # fft auf eingangsignal fuer amplitudenspektrum
94 fft = np.fft.fft(decoded)
95 spec = np.abs(fft)
96 plt.figure()
97 plt.plot(spec)
98 plt.title('Amplitudenspektrum')
99 plt.xlabel('Abtastpunkte')
100 plt.ylabel('Amplitude')
101 plt.savefig('ungetriggeredAmplitudenspektrum.png')
102
103
104 # plot fuer windowed funktion
105 plt.figure()
106 plt.plot(windows)
107 plt.title('Amplitudenspektrum Fensterfunktion')
108 plt.xlabel('Abtastpunkte')
109 plt.ylabel('Amplitude')
110 plt.savefig('fensterAmplitudenspektrum.png')

```

A.1.3 Quellcode Versuch 2

```
1 import numpy as np
2 from scipy import signal
3
4 PATH_W = "windowreferenzdatensatz"
5 PATH_R = "refspektren"
6
7
8 # addition aller 5 aufnahmen, mit anschliessender mittelung
9 LRef = (np.genfromtxt(PATH_R + "RefWinL1.csv")
10 + np.genfromtxt(PATH_R + "RefWinL2.csv")
11 + np.genfromtxt(PATH_R + "RefWinL3.csv")
12 + np.genfromtxt(PATH_R + "RefWinL4.csv")
13 + np.genfromtxt(PATH_R + "RefWinL5.csv")) / 5
14
15 # addition aller 5 aufnahmen, mit anschliessender mittelung
16 RRef = (np.genfromtxt(PATH_R + "RefWinR1.csv")
17 + np.genfromtxt(PATH_R + "RefWinR2.csv")
18 + np.genfromtxt(PATH_R + "RefWinR3.csv")
19 + np.genfromtxt(PATH_R + "RefWinR4.csv")
20 + np.genfromtxt(PATH_R + "RefWinR5.csv")) / 5
21
22 # addition aller 5 aufnahmen, mit anschliessender mittelung
23 TRef = (np.genfromtxt(PATH_R + "RefWinT1.csv")
24 + np.genfromtxt(PATH_R + "RefWinT2.csv")
25 + np.genfromtxt(PATH_R + "RefWinT3.csv")
26 + np.genfromtxt(PATH_R + "RefWinT4.csv")
27 + np.genfromtxt(PATH_R + "RefWinT5.csv")) / 5
28
29 # addition aller 5 aufnahmen, mit anschliessender mittelung
30 HRef = (np.genfromtxt(PATH_R + "RefWinH1.csv")
31 + np.genfromtxt(PATH_R + "RefWinH2.csv")
32 + np.genfromtxt(PATH_R + "RefWinH3.csv")
33 + np.genfromtxt(PATH_R + "RefWinH4.csv")
34 + np.genfromtxt(PATH_R + "RefWinH5.csv")) / 5
35
36 # speichern der gemittelten referenzspektren
37 def safeRefWAsCSV(eingangssignal):
38     csvRefArray = np.savetxt(PATH_R + "/decodedtrigR1.csv",
39                               eingangssignal, delimiter=",")
39     return csvRefArray
```

40

41 `safeRefWAsCSV(LRef)`

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4
5
6 # plotten der referenz spektren
7 PATH_R = "refspektren"
8
9 LRef = np.genfromtxt(PATH_R + "/LRef.csv")
10 RRef = np.genfromtxt(PATH_R + "/RRef.csv")
11 HRef = np.genfromtxt(PATH_R + "/HRef.csv")
12 TRef = np.genfromtxt(PATH_R + "/TRef.csv")
13
14 fig1, axs = plt.subplots(2)
15 axs[0].plot(LRef)
16 axs[0].set_title("Links")
17 axs[1].plot(RRef)
18 axs[1].set_title("Rechts")
19
20
21 fig2, axs = plt.subplots(2)
22 axs[0].plot(HRef)
23 axs[0].set_title("Hoch")
24 axs[1].plot(TRef)
25 axs[1].set_title("Tief")
26
27 plt.show()

```

```

1 import numpy as np
2 import scipy.stats as sst
3
4
5 PATH_R = "refspektren"
6 PATH_G = "sprechergina"
7 PATH_E = "sprecherelena"
8
9 # Referenzdatensatz
10 LRef = np.genfromtxt(PATH_R + "/LRef.csv")
11 RRef = np.genfromtxt(PATH_R + "/RRef.csv")
12 TRef = np.genfromtxt(PATH_R + "/TRef.csv")
13 HRef = np.genfromtxt(PATH_R + "/HRef.csv")
14
15 # selber sprecher wie referenz
16 LG = np.genfromtxt(PATH_G + "/winL.csv")
17 RG = np.genfromtxt(PATH_G + "/winR.csv")
18 TG = np.genfromtxt(PATH_G + "/winT.csv")
19 HG = np.genfromtxt(PATH_G + "/winH.csv")
20
21 # anderer sprecher wie referenz
22 LE = np.genfromtxt(PATH_E + "/winL.csv")
23 RE = np.genfromtxt(PATH_E + "/winR.csv")
24 TE = np.genfromtxt(PATH_E + "/winT.csv")
25 HE = np.genfromtxt(PATH_E + "/winH.csv")
26
27
28 # korrelation datensatz mit selben sprecher
29 korrkoefLLG = sst.pearsonr(LRef, LG)[0]
30 korrkoefLRG = sst.pearsonr(LRef, RG)[0]
31 korrkoefLTG = sst.pearsonr(LRef, TG)[0]
32 korrkoefLHG = sst.pearsonr(LRef, HG)[0]
33
34 # korrelation datensatz mit anderem sprecher
35 korrkoefLLE = sst.pearsonr(LRef, LE)[0]
36 korrkoefLRE = sst.pearsonr(LRef, RE)[0]
37 korrkoefLTE = sst.pearsonr(LRef, TE)[0]
38 korrkoefLHE = sst.pearsonr(LRef, HE)[0]
39
40
41 # berechne welches wort am besten erkannt wird original sprecher
42 maxValG = max(max(korrkoefLLG, korrkoefLRG), max(korrkoefLTG, korrkoefLHG))

```

```

    )
43
44
45 # berechne welches wort am besten erkannt wird anderer sprecher
46 maxValE = max(max(korrkoefLLE, korrkoefLRE), max(korrkoefLTE, korrkoefLHE)
    )
47
48
49 print(korrkoefLTG)
50 print(korrkoefLLG)
51 print(korrkoefLHG)
52 print(korrkoefLRG)
53 print(maxValG)
54
55 print(korrkoefLTE)
56 print(korrkoefLLE)
57 print(korrkoefLHE)
58 print(korrkoefLRE)
59 print(maxValE)

```


Literaturverzeichnis

- [1] Franz, Prof. Dr. Matthias O.: *Vorlesung 10 - Sprache und Spracherkennung: Kurzzeit-Fouriertransformation, Erzeugung und Wahrnehmung von Sprache, Mustererkennung durch Korrelation*. In: *Vorlesung Technische Grundlagen der angewandten Informatik*, 2016.
- [2] Franz, Prof. Dr. Matthias O.: *Vorlesung 1 - Einführung*. In: *Vorlesung Technische Grundlagen der angewandten Informatik*, 2016.