

Don't over think,  
I know your drink

Suzanne Bardelmeijer, Karen Beckers,  
Sanne Eggengoor, Job van Gerwen & Roos Slingerland

The front page image is cover art from Marshall McLuhan

# Don't over think, I know your drink

Suzanne Bardelmeijer 10716971  
Karen Beckers 10811958  
Sanne Eggengoor 10729895  
Job van Gerwen 10378863  
Roos Slingerland 10775935

Project Report  
Course Media Understanding  
Credits: 6 EC

Bachelor Artificial Intelligence

College of Science  
University of Amsterdam  
Faculty of Science  
Science Park 904  
1098 XH Amsterdam

*Course coördinator*  
Dr. A. Visser

Informatics Institute  
Faculty of Science  
University of Amsterdam  
Science Park 904  
1098 XH Amsterdam

April 2th, 2017

# 1 Summary

This research developed a system Sjeel 2.0 that can recognise the personal wishes of the two favourite drinks of students: beer and coffee. It uses two different media types namely speech recognition and speaker identification. Using Google's API for speech recognition the difference between coffee and beer a result was reached of only 3 mistakes in 199 samples. Using multiple features such as pitch, MFCC and RMSE for speaker identification multiple algorithms with 10-fold crossvalidation were applied. Naive Bayes, Decision Tree, Random Forest all reached an accuracy around 70% while Gradient Boosted Trees after parameter finetuning reached an accuracy of 81%. As expected, the level of noise does influence the accuracy of the system. When recording the trainingsset in a quiet environment, an accuracy of 95% was reached. This system could be used in situations where speaking for a long time is an issue, such as with stutters or emergency situations. Future work could focus on extending the number of beverages and users. Also other media forms for speaker identification could be added such as face recognition and user-dependent vocabulary.

## Contents

<b>1 Summary</b>	<b>4</b>
<b>2 Introduction</b>	<b>5</b>
<b>3 Theoretical Foundation</b>	<b>5</b>
3.1 Phonetics . . . . .	5
3.2 Speech recognition . . . . .	6
3.3 Speaker identification . . . . .	6
3.4 Speech recognition and speaker identification combined . . . . .	7
<b>4 Research Method</b>	<b>7</b>
4.1 Data . . . . .	7
4.2 Feature extraction . . . . .	8
4.3 Model training . . . . .	8
4.4 Output: Sjeel in action . . . . .	9
4.5 Evaluation . . . . .	10
<b>5 Results</b>	<b>10</b>
<b>6 Conclusion</b>	<b>13</b>
<b>7 Discussion</b>	<b>13</b>
<b>References</b>	<b>16</b>

## 2 Introduction

'Goodmorning, a large cappuccino I assume?', asks Sjeel. The woman of the canteen of Science Park knows the preferences of her regular customers almost flawless. How great would it be to copy the skills of Sjeel and use them in a system? A system that can recognise the preferences of the users by analysing their voices. And above that, can recognise the kind of beverage you want. To achieve this, 2 different kinds of media are used: speech recognition and speaker identification. This would not only save time, but is also a great solution for people who have troubles with speech.

Imagine a residential care home with many people requiring help. How efficient would it be that in case of emergency a person in help only has to say 'help' and our system would know who it is. And above that, knows the medication history and can therefore provide help more efficiently.

Another example is with people suffering from stuttering. Instead of saying: *I would like to have a cappuccino with almond milk and a dash of cacao, say: Coffee please.* This would save even more time than it would with regular speakers. Secondly, it might decrease the social anxiety of speaking in a busy queue under pressure.

Even though this system could be used very widely, there is decided to define the two-folded problem as follows. First, the system recognises which beverage is ordered by using speech recognition. Secondly, by using speaker identification, the system analyses the personal preferences of this specific user. In this study the number of users is equal to the number of groupmembers, so 5 in total.

For example if a user says: 'Beer please', the system should know that the user ordered the beverage beer. Secondly it should recognise that it was user 3 that ordered the beer and would therefore know that it should be a Heineken. The beverages coffee and beer are chosen because they are widely known as the most popular student drinks.

The system is expected to prove the following hypothesis:

1. Distinct the word beer from coffee perfectly.
2. Distinct the man of the group from the 4 women perfectly.
3. Reach an accuracy of at least 80% for recognising the speaker in a noisy environment.
4. Reach a higher accuracy for recognising the speaker in a quiet environment.

This research tends to bridge the semantic gap by giving meaning to the values of audio signals. What to a computer just a collection of wavelengths is, is to humans the difference between 'beer' and 'coffee' or between the voice of person A and person B.

## 3 Theoretical Foundation

### 3.1 Phonetics

To understand the conversion from pure sound waveform to text and to the speaker identity, it is necessary to have a better understanding of speech itself.

In this section a short introduction in the field of phonetics is given. In a strict sense, speech is produced by the vocal tract, which is a collaboration between the lungs, larynx and oral cavities. The larynx contains the vocal cords, two small pieces of tissue that vibrate when they are closed and air is pressed through. But with this only the pitch of the voice can be altered. Different sounds can be made by changing the shape of the oral cavities. For example, the difference between an [i] in 'keen' and [o] in 'hope' can be explained by the change in the location of the tongue, the height of the lower jawbone and the tension in the lips. Also the difference between [m] and [n] in 'man' can be explained by the difference in location of the tongue and lips. The difference in sound comes from the resonance in the cavities.

This is of course also visible in the waveform of the sound, which is a combination of the frequency of the vocal cords, and the frequency of resonances in the several cavities. These frequencies are called formants, with the lowest formant (formant 0,  $f_0$ ) as the frequency of the vocal cords. The next two formants contain information about the resonance in the pharynx (cavity in the throat) and mouth. These formants are easily changeable by moving the tongue and jaw and thus contain the most information about what is said. The next formant describes the resonance in the nasal cavity and is sometimes used to identify nasal consonants (such as [m] and [n]). The fourth and fifth formant contain speaker specific information, because this describes the resonance in other cavities, which cannot be changed when speaking. Therefore the first, second and third formant are the most useful in defining what is said and the fundamental frequency ( $f_0$ ) and the fourth and fifth formant contain the most useful information about who spoke.

### 3.2 Speech recognition

There is an overall agreement in the literature that Hidden Markov Models (HMM) is the best approach for Automatic Speech Recognition (ASR) (Baker et al., 2009). In almost all applications a HMM forms the basis of the system. The models are trained for each word in a large dataset that contains labeled speech data. When a new (non-labeled) input is given, the sentence is recognised by the most likely model sequence (Foote, 1999). For large-vocabulary recognition datasets a different approach is applied. Instead of making a HMM for every of the thousands of words in the dataset, a much smaller number of sub-word models are constructed that can be combined to form a word model. Since the that will be build in this research only has to recognise two keywords, "beer" and "coffee", the first approach will probably be best for this problem. Even though HMM is the standard method for building these kind of systems, new research focusses on deep neural networks. Google has build a neural network based keyword spotting system that is 45% more accurate than Hidden Markov Models (Chen, Parada, & Heigold, 2014).

### 3.3 Speaker identification

As stated in the section about Phonetics, the most useful data in speech are the fourth and fifth formants and the fundamental frequency. To be able to obtain these data from the waveform sounds there needs to be done a feature extraction (Kinnunen & Li, 2010). These features are typically the Mel Frequency Cepstral

Coefficients for the previous 20-30 ms. However, it might be possible to use only the pitch of the voices, because the choice is limited to 5 speakers (Weenink, 2017).

*Mel Frequency Cepstral Coefficients (MFCC)* are based on the nonlinear and logarithmic perception of pitch in the human ear (Hasan, Jamil, Rahman, et al., 2004). The frequencies on the Mel scale are below 1000 Hz linearly spaced, above 1000 Hz the frequencies are logarithmically spaced. This has an analogy with the human ear, where the basilar membrane (located in the cochlea) vibrates at specific points (depending on the frequency of the sound). The spacing of these specific points is the same as the Mel Frequency Scale, which makes this Scale very suitable for speaker identification.

*Gaussian mixture models* are widely used in unsupervised speaker recognition systems (Reynolds, F., & Dunn, 2000). The task of such a system is to determine by who a segment of speech was spoken. The system described by (Reynolds et al., 2000) was the *Gaussian Mixture Model-Universal Background Model* (GMM-UBM). This system is based on the optimal likelihood ratio test. *GMM-UBM* uses *Gaussian mixture models* for likelihood functions. For the representation of alternative speakers in a specific segment of speech a universal background model is used and Bayesian adaptation obtains a hypothesized speaker model.

### 3.4 Speech recognition and speaker identification combined

This research aims at simultaneously recognizing the person who is speaking and what is being said. For the system to work it is therefore necessary that the components mentioned above are combined. To this, though, there is little research done. In 2010 a self-learning speech controlled system was presented by Herbig, Gerl en Wolfgang. (Herbig, Gerl, & Wolfgang, 2010). Their research tried to improve speech recognition by first identifying the person who is speaking. Speaker identification with high accuracy was obtained by first retrieving personal speech characteristics such as the Mel Frequency Cepstral Coefficient. Speech recognition was improved by taking into account how the identified person speaks. Knowing how this person pronounces specific words or formulates sentences improves the recognition of what is being said.

This project will not focus on improving speech recognition, but uses the identification of a speaker to fulfill his or her needs by knowing what his or her preferences are.

## 4 Research Method

### 4.1 Data

Two datasets of both 200 audio files are created. One dataset was situated in a noisy environment and the other in a quiet environment. For each dataset all five users have recorded the dutch sentence "Mag ik koffie?" (Can I have a coffee?) and the sentence "Mag ik een biertje?" (Can I have a beer?) 20 times. A screenshot (Figure 5) of the created dataset can be found in Appendix A.

## 4.2 Feature extraction

To create an application that is able to define the preferences of a specific user, a Python program is created that can capture sound of the input using the *PyAudio* toolkit (Pham, 2017). This toolkit can store an audio input into a Waveform Audio File Format (WAV).

After the audio has been captured and stored, the program calls the *SpeechRecognition* library (Zang, 2015), which makes use of the *Google Speech Recognition API* in order to translate the audio into text. This library can be used for both direct audio input and previously recorded audio. The approach that uses previously recorded audio will be applied in the training phase and the approach using direct audio input will be applied in the test phase.

*RMS Energy* is a metric that measures the continuity of the audio signal level (Weihs, Jannach, Vatolkin, & Rudolph, 2016). A piece of audio that contains many portions of silence will have a large RMSE value whereas audio containing continuous signals will have a small RMS Energy value. The RMS Energy can be computed by taking the root-mean-square (RMS) energy of a signal (Weihs et al., 2016). This project used the RMS Energy feature to trim the audio segments because in this way the silences can be cut off. This results in an audio segment that only contains speech. The differences in RMS Energy between segments of speech of different users will also help by the identification of a speaker. This is because people tend to speak with different dynamics.

In order to identify the speaker, several timbre features will be extracted from the audio files. These features are going to be the data for the classification algorithm. In order to extract these features, the *LibROSA* package (McFee et al., 2017) is employed. This package is commonly used for music and audio analysis and music information retrieval, but is also for this research a helpful tool, since it enables the extraction of a large number of features that can be used for speaker identification. As mentioned before the interview with an expert in the field of linguistics revealed that the most useful features for this research will be pitch and the Mel-Frequency Cepstral Coefficients (Weenink, 2017). These will be the features used for speaker identification.

## 4.3 Model training

The next step is training the algorithm used for identifying the speaker. Only supervised learning algorithms will be tested and a comparison between these algorithms will be made. The literature showed that Naive Bayes classification is commonly used in this area of research (Reynolds et al., 2000). This algorithm will therefore be applied in this study. Besides that, also a Decision Tree, Random Forest and Gradient Boosted Trees algorithm will be tested.

The Naive Bayes algorithm is based on Bayes theorem. Here the assumption of independence between every pair of features is made (Reynolds et al., 2000).

Decision Trees build a tree-like structure based on the characteristics of the data. Every node is a feature and the branches are the different values that a feature can have (Mitchell, 1997).

Random Forest learns multiple Decision Tree models simultaneously. The output is the most frequent label of all the predicted values of the models. Hereby the problem of overfitting that Decision Trees are prone to is solved

(Breiman, 2001).

Gradient Boosted Trees start with a pretty weak model and then tries to enhance that. It calculates the error of the model and fits a new Decision tree to the corresponding cost function. By adding this new model to the previous one, the whole function gets more and more complex. This happens in an iterative way in order to refine the model (Friedman, 2001).

#### 4.4 Output: Sjeel in action

When both speech recognition and speaker identification have been executed, the results will be combined to evaluate what the speaker wants. The preferences of a person are predetermined by a questionnaire for a set of situations. The wish of the user will then be returned as a text message.

For this research the focus will solely be on detecting the preference for a beverage. This will be divided in two categories: beer and coffee. For both categories the users will be given a choice between several brands or flavours. Table 1 shows the preferences of the different users. The system will first tell the user that his/her preferred drink is on its way and will then show a photo of the user with either a cup of coffee or a beer in his/her hand. Figure 1 shows an example of the output.

User	Coffee preference	Beer preference
1 (Roos)	Glutenfree coffee	Glutenfree beer
2 (Karen)	Latte macchiato with soy cream	Alfa
3 (Sanne)	Black	Grolsch
4 (Job)	Cappuccino	Hertog Jan
5 (Suzanne)	Flat white	Duvel

Table 1: User preferences

```
Stel je vraag aan Sjeel 2.0  
Is goed mop, een Hertog Jan komt eraan!  
Press Enter to continue...
```



Figure 1: Ouput example

#### 4.5 Evaluation

Evaluation of the algorithms will be done by calculating the percentage of the correctly classified training examples. A 10-fold crossvalidation was applied to obtain reliable results.

### 5 Results

First of all the MFCCs of the audio files have been calculated. Examples of these MFCCs are in figure 2.

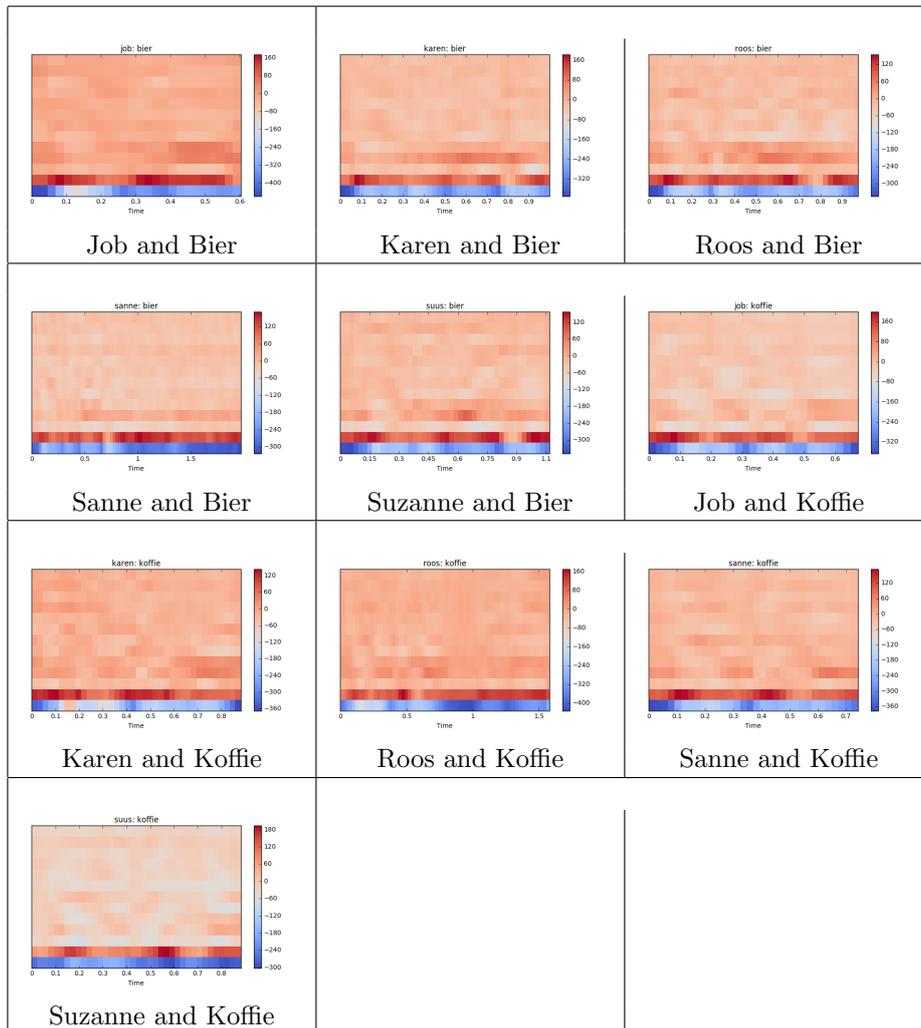


Figure 2: MFCCs of the speakers

Multiple learning algorithms have been used to obtain optimal results for both the quiet and the noisy environment. As shown in Figure 3 *Gradient Boosting Tree* yielded the best results in a noisy environment but *Naive Bayes*, *Decision Tree* and *Random Forest* have been tested as well. For the quiet environment, Figure 4 shows that *Naive Bayes* and *Random Forest* both yielded the same best result.

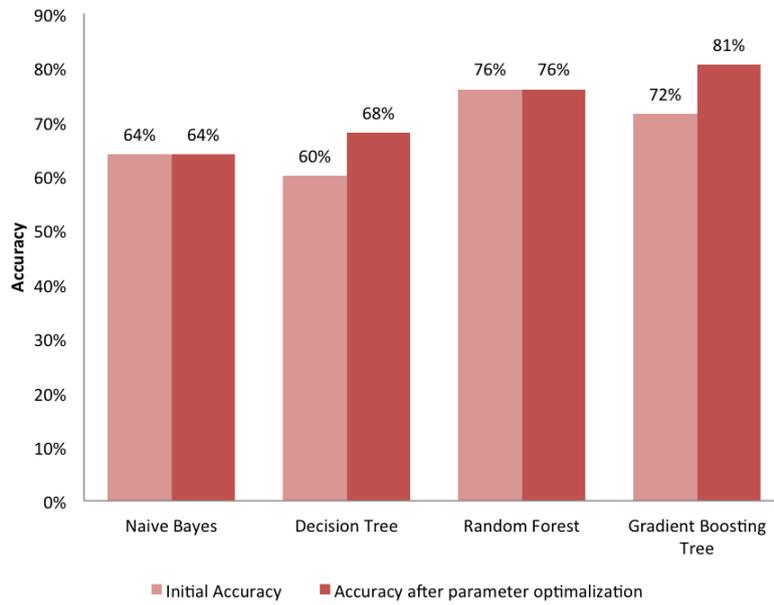


Figure 3: Obtained accuracies (crossvalidation 10) noisy environment

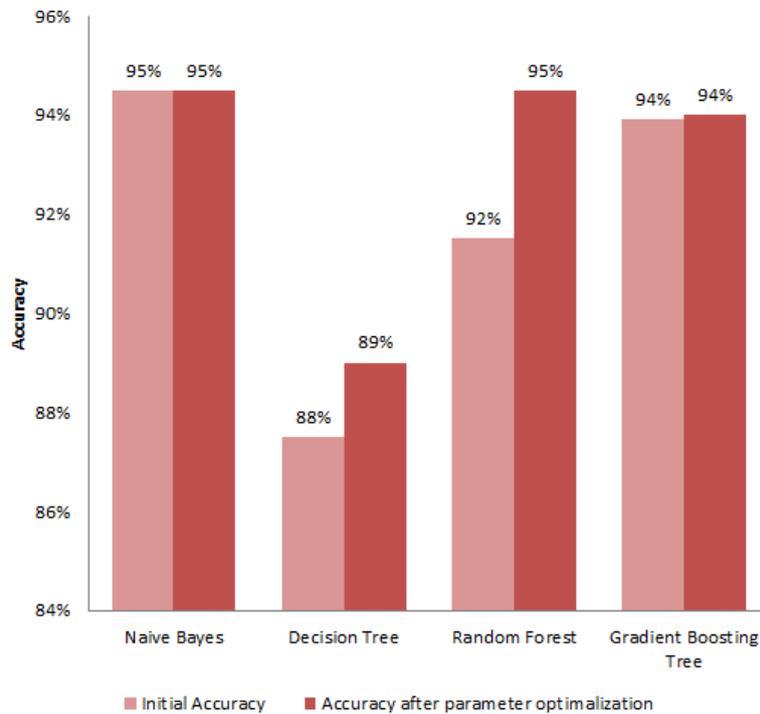


Figure 4: Obtained accuracies (crossvalidation 10) quiet environment

All tests made use of *10-fold-cross-validation*. An accuracy of 80,5% was obtained by using the *GradientBoostingClassifier* in Python. The corresponding parameter values after optimisation are shown in Table 2.

Parameter	Value
n estimators	100
learning rate	0.5
max depth	1
random state	0

Table 2: Parameter values

## 6 Conclusion

In this paper a system was created that copied the skills of coffeelady Sjeel. By using the media speech recognition and speaker identification a twofolded problem was tackled. The system could both distinct between beverages as between users. When looking back at the hypothesis the following conclusions can be drawn.

1. *Distinct the word beer from coffee perfectly.* The system made only 3 mistakes in 199 descriptions, so this hypothesis was correct.
2. *Distinct the man of the group from the 4 women perfectly.* The system made no mistakes by distinguishing user 4 from the other users, so this hypothesis was correct.
3. *Reach an accuracy of at least 80% for recognising the speaker in a noisy environment.* The highest accuracy was reached after parameter finetuning of the Gradient Boosted Trees algorithm namely 81%, so this hypothesis is correct.
4. *Reach a higher accuracy for recognising the speaker in a quiet environment.* The highest accuracy was reached after parameter finetuning of the Gaussian Naive Bayes algorithm namely 95%, so this hypothesis is correct as well.

This research could in the future help people that struggle with speaking such as people who stutter, suffer from anxiety in social situations and in case of emergency where only a few words can be pronounced. Above that, it saves time of explaining all the details of your personal order.

## 7 Discussion

Looking back there are a few critical points to be made about this research. First of all the group of users was small (only 5), for a real-world system this could be expanded for more users. The evaluation of how well the user was identified can then have more meaning.

Secondly, the number of different beverages could be expanded as well. Above drinks, other preferences could be added, such as medication or favorite songs.

When looking at the speaker identification other media could have been used as well. Visual recognition could have been an option as well. This could have been done by recognising the faces or even lengths or postures of the users. This research now chose to let the users speak the exact same words. Literature shows that by using, for example, fillers, speaker identification can be improved. Future work can therefore focus on enlarging the number of users, beverages (or other preferences) and subjoin other user recognition techniques such as visual techniques and user-dependent vocabulary.

# Appendix A

Pitch	RMS	MFCC1	MFCC2	MFCC3	MFCC4	MFCC5	MFCC6	MFCC7	MFCC8	MFCC9	MFCC10	MFCC11	MFCC12	MFCC13	User
22	-0.23	0.0574984865732	-359.274715003	139.7390965	-30.4845936324	38.4005384131	6.8800760301	-49.364002076	-12.357004349	2.8110993328	-14.105893328	-14.165386834	-22.209849239	-12.582295791	job
184	-0.12	0.044482961297	-387.697361636	152.93550861	-17.1216728896	31.00129036	6.674761017	-28.089378635	-2.4291916254	7.88733122642	-16.3170225059	-12.0330971731	-11.051567139	-16.9035655663	saus
130	0.22	0.0655399518874	-300.254475391	-24.7189723515	44.8920238014	30.254475391	2.91592654824	-36.8484141571	3.3225114025	0.944038295144	-19.0395978802	-23.4861893102	-20.862813174	-7.4760072214	same
23	-0.45	0.0582986424792	-342.345775933	139.794747257	-35.718799665	35.4071779865	-0.86615420007	-40.2225025862	13.70262815703	0.04425279387	-13.4835902073	-19.0298248676	-23.0440406118	-12.8497057644	job
99	0.01	0.08087430737	-333.493196356	117.399116703	-42.1947982709	55.906699041	11.8443092599	0.44457864587	-13.70262815703	-22.1843034806	-10.2679477456	-17.19954704331	-15.6327872601	-0.1603050438	roos
66	0.01	0.0780565589666	-346.504626701	124.311860684	-12.7802886946	24.1720019764	-19.5947713042	-8.2026967466	2.11453100106	-9.37676688146	-29.2982328215	-18.2088684962	1.2903339375	-26.8855907845	karen
172	0.1	0.0728292935302	-330.089943021	146.25567477	-22.705512592	63.4018672285	5.59589625551	-51.5923386995	-20.6129280165	2.44423172446	-23.7863205144	-10.848227395	-10.3603541027	-13.3283067064	saus
165	-0.49	0.0532290004194	-362.730951575	141.17757577	-25.6624713158	52.3381195498	12.1395596335	-40.137554778	9.3473248736	10.068993543	-28.0822374183	-9.87733011786	-17.8728332331	-5.94945873472	same
5	0.06	0.0694398194947	-294.630349995	146.44883875	-61.804711011	31.2644796924	21.4848366024	-29.0701953637	-20.0646473615	-21.6681747784	-11.684720177	-15.29191914486	-24.1567792034	-12.8152933505	job
169	0.06	0.0705331414938	-338.710502096	130.391408812	-34.8071255622	62.7186186096	0.73063128556	-46.3250433724	-22.3941524408	0.068336076603	1.62911771546	-23.8116532097	-21.4000326915	-10.979120767	saus
16	0.3	0.0589476265013	-320.496098773	131.97276316	-66.1563129036	44.4247398346	-22.1175718242	-40.6835300142	-17.3865059183	17.5682329389	-8.136646176384	-12.8849346826	-27.4880847277	-11.2643507944	job
88	-0.3	0.0883329981948	-331.32192711	110.224920261	-19.1832520896	61.781741792	-0.00836149859	-48.7571753889	-0.0213723236	-2.7295131287	-19.0629008596	-12.9814848512	-10.6857273075	-10.3288217338	roos
18	0.39	0.0718938112259	-310.868868584	137.757076369	-62.2099213066	43.1217018971	27.5425255699	-38.687023001	-18.2474082527	15.6955147908	-24.3856035978	-12.7594309513	-26.732653541	-18.5608676047	job
20	0.18	0.0598741434515	-324.33981563	138.281834314	-32.5740710442	50.148378294	4.070344396	-52.4698051556	-7.8861646955	1.9570895534	-21.9493082821	-12.4753476844	-17.8152419379	-19.866929329	job
135	0.37	0.0640745684505	-365.963152797	119.188043504	-20.349325039	50.8292832686	6.0086540259	-35.0210693645	7.0057222047	-0.387500663213	-22.448838318	-16.1313400819	-16.0933934929	-12.4131773961	saus
35	0.19	0.0311452847958	-399.164348049	138.602306856	-53.0712174353	68.5564569099	-28.4659280026	-54.8303370482	-32.661937928	-18.1071166037	-15.8607551758	-17.7952688797	-13.9413633481	-21.701955582	saus
196	0.36	0.0831395239838	-350.702823824	117.034865906	-21.9336717551	34.1808463079	4.08883111683	-29.0157871735	-7.4290361618	-15.3250981439	-20.0353963552	-18.2426812916	-10.3139502545	-27.1637419452	saus
35	0.19	0.0531178079545	-349.962059827	136.827079282	-31.8329102243	38.1826715305	-1.90716265042	-45.8714910365	-14.9304508587	2.29029037593	-10.9733737899	-18.311083414	-14.6226837261	-23.8486374246	job
59	-0.4	0.0798298651285	-321.580771276	126.789829328	-28.501726274	45.702658826	6.2727489345	-35.2916174822	-28.616426432	-5.8341418722	-20.150033143	-17.5277085074	-8.2191296368	-18.8178665214	karen
183	-0.42	0.0458607040346	-394.646649877	134.025111811	-12.2013491286	32.9784740768	5.01816104241	-32.663904067	-12.3119091965	3.20629517775	-14.0523021174	-12.800574935	-11.185360887	-19.6823516138	saus
198	-0.19	0.0826717689833	-335.750346376	132.126627996	-17.6660916221	38.2992803379	7.23183574986	-36.878953021	-2.38220856755	-5.93156259284	-24.3219473075	-17.7640803589	-8.6798957082	-23.3786746343	saus
180	0.04	0.040166546789	-392.398538336	132.895321561	-24.8653983698	42.325146557	-6.1337680379	-34.4238549358	-10.932602236	-16.6489217389	-17.3190389248	-19.558461286	-5.84621012575	-7.60399137628	saus
48	0.33	0.077484265512	-313.304954682	131.324554672	-28.6713192821	49.8363378205	-4.3021448501	-42.9873326268	-18.518088842	1.8701595524	-5.64779118453	-30.6964103889	-16.9592813014	-3.01795718011	saus
13	0.35	0.061437437177	-394.100799039	149.429500037	-47.284566627	49.9641615395	-18.702188343	-35.8674815045	-7.34572543695	-13.4121212032	-18.4747819802	-6.20076746282	-9.86650011779	-21.219930084	job
174	0.22	0.0472099930048	-363.576376985	140.847165686	-43.4384792722	56.0726583934	9.68785200888	-55.8504919068	-27.6787519522	-9.84338532807	-2.14268097829	-16.4497937073	-20.450055634	-16.700662452	saus
161	0.04	0.0284143892825	-399.830241415	148.311171444	-38.10964483786	52.5155079562	-20.4215378259	-48.7237789275	-12.2522926504	0.673194248294	-19.3318451841	-17.6183823954	-6.7214889374	-12.226001813	saus
159	-0.19	0.0704978990819	-317.079745238	115.0557226	-23.7413704763	44.2450506652	-2.9295275198	-34.9939817703	-1.0370897792	-3.35084657584	-13.5309069896	-13.4754003034	-23.298585739	-0.466886124072	saus
199	0.07	0.0602126382291	-373.44105234	126.85811591	-16.8526746544	36.5968458007	0.102493382223	-25.6418013696	-2.3172171204	-2.6657009895	-17.3071967398	-21.0740248605	-11.3407853451	-8.214864400137	saus
189	0.23	0.0741381570697	-341.905226857	125.046234549	-14.9214853988	38.5552944625	-6.16877036119	-37.6438191947	-16.655417825	-8.771456267	-18.0077884192	-16.7019155738	-3.75678682653	-28.4019555835	saus
34	0.3	0.0417510650393	-364.981013978	139.452019158	-40.0546438141	46.041043258	2.8940582492	-47.523816571	-9.8396807524	3.48715841303	-13.925431992	-22.0791469266	-14.50696443483	-22.6963181474	saus
182	-0.22	0.0629185272217	-367.457998813	123.054039824	-17.2474232285	31.4041503259	-2.29052373102	-40.1439625727	-9.90212057984	7.52333726397	-13.8822998202	-15.1510463379	-13.1579514036	-18.982872003	saus
4	0.5	0.0759640056485	-288.479342087	135.617149905	-62.5005479829	42.3318813719	19.2036991746	-32.1070517193	-17.2478871343	-17.829859198	-20.0173529941	-10.3332831391	-9.5629490218	-24.1246569607	saus
187	0.37	0.0652193278074	-341.110624612	124.6129398	-23.065599385	41.0135020382	-3.84716074216	-33.5830098587	-6.183361776	-4.5743755155	-13.1494344459	-21.7367284689	-11.3670428686	-27.4167859934	saus
149	0.4	0.0547220119357	-368.179260017	136.897576271	-18.9752840789	37.1292323234	-3.98349487913	-42.9267397659	6.159898058	0.07238071523	-27.1914209585	-21.0654193981	-10.8159477833	-12.021570847	saus
147	0.04	0.0441986585456	-384.924217864	131.298266889	-6.62331843437	38.6147591544	-5.18874895988	-38.0342639912	-6.36906029622	-14.6837492815	-28.2592404204	-20.1593741358	-15.86751541358	-6.2854304536	saus
54	0.04	0.0598438907087	-354.365184646	132.181022163	-28.302041843	54.34416344	7.41414674476	-26.3511173684	-4.17901600307	-4.15471092301	-23.7916137985	-22.008661668	-3.86301127424	-18.1676890572	karen
132	0.28	0.075262823372	-358.678292915	119.89086684	-18.9113514033	47.5024951633	-7.0680357821	-36.4171900707	-7.46842118017	1.32629571607	-22.1267535522	-13.2991592415	-20.2602562599	-12.7176515891	saus
77	0.07	0.0688555853605	-357.938269568	127.765085961	-26.3043800387	26.3043800387	-18.22220278689	-37.5982011726	-4.124567359203	-1.24567379203	-29.17661549218	-13.0444602804	-1.5157118759	-20.7086201103	karen

Figure 5: Screenshot of the dataset

## References

- Baker, J. M., Deng, L., Glass, J., Khudanpur, S., Lee, C.-H., Morgan, N., & O’Shaughnessy, D. (2009). Developments and directions in speech recognition and understanding, part 1 [dsp education]. *IEEE Signal Processing Magazine*, 26(3).
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5–32.
- Chen, G., Parada, C., & Heigold, G. (2014). Small-footprint keyword spotting using deep neural networks. In *Acoustics, speech and signal processing (icassp), 2014 IEEE international conference on* (pp. 4087–4091).
- Foote, J. (1999). An overview of audio information retrieval. *Multimedia systems*, 7(1), 2–10.
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189–1232.
- Hasan, M. R., Jamil, M., Rahman, M. G. R. M. S., et al. (2004). Speaker identification using mel frequency cepstral coefficients. *variations*, 1(4).
- Herbig, T., Gerl, F., & Wolfgang, M. (2010). Simultaneous speech recognition and speaker identification. *Spoken Language Technology Workshop (SLT)*.
- Kinnunen, T., & Li, H. (2010). An overview of text-independent speaker recognition: From features to supervectors. *Speech communication*, 52(1), 12–40.
- McFee, B., McVicar, M., Nieto, O., Balke, S., Thome, C., Liang, D., ... Lee, H. (2017). *librosa*. Retrieved from <http://librosa.github.io/librosa/>
- Mitchell, T. M. (1997). *Machine learning - decision tree learning* (Vol. 45) (No. 37).
- Pham, H. (2017). *Pyaudio*. Retrieved from <https://people.csail.mit.edu/hubert/pyaudio/>
- Reynolds, D. A., F., Q. T., & Dunn, R. B. (2000). Speaker verification using adapted gaussian mixture models. *Digital Signal Processing*, 10, 19–41.
- Weenink, D. (2017, 02). Personal Communication.
- Weihs, C., Jannach, D., Vatulkin, I., & Rudolph, G. (2016). *Music data analysis: Foundations and applications*. CRC Press.
- Zang, A. (2015). *Speech recognition*. Retrieved from [https://github.com/Uberi/speech\\_recognition#readme](https://github.com/Uberi/speech_recognition#readme)