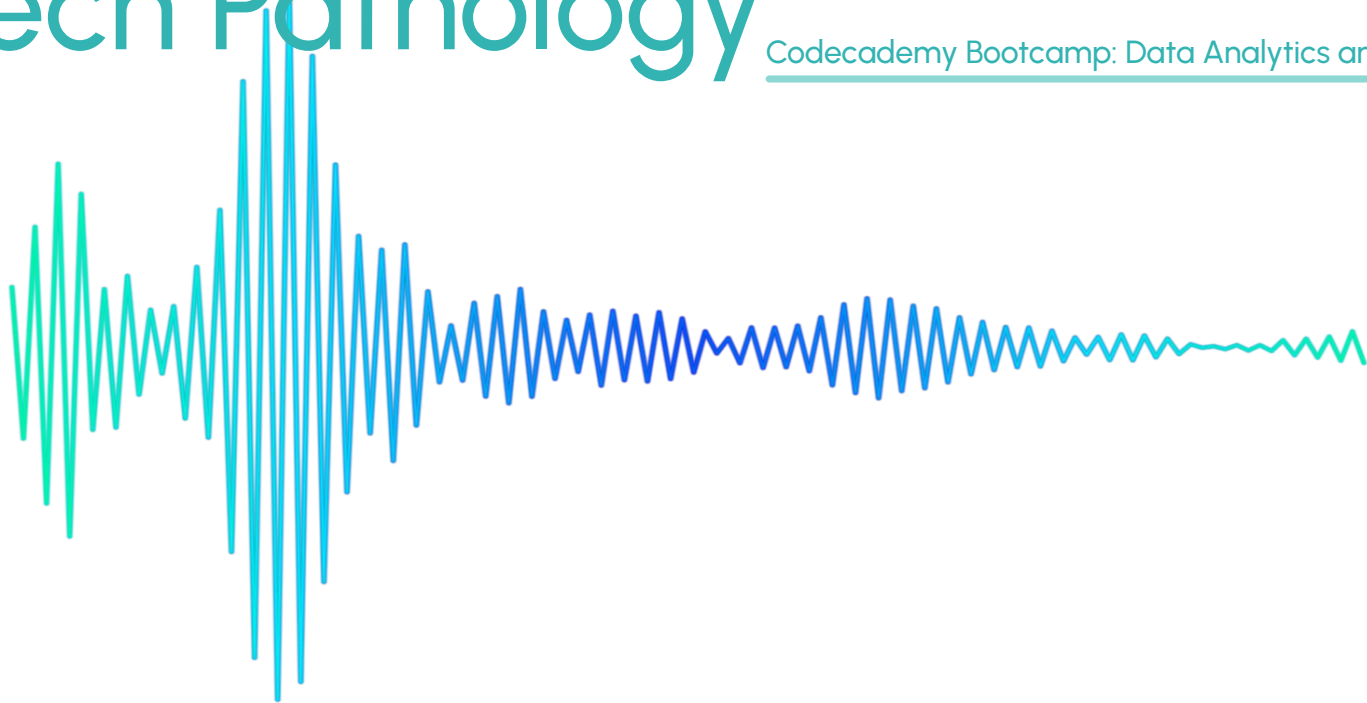


# Capstone Project: Speech Pathology

Codecademy Bootcamp: Data Analytics and AI



# Contents:

1 Project Overview >

2 Data Overview >

3 Approach >

4 Project Workflow >

5 Data Cleaning >

6 Machine Learning >

7 Data Analyses >

8 Conclusion & Recommendations >

# Project Overview:

In this capstone project, you will analyze a dataset containing audio-derived features (e.g., Mel filters, Fourier transformation) extracted from voice recordings of patients labeled as "Healthy" or "Unhealthy." The goal is to perform exploratory data analysis (EDA), clean the data, visualize insights in a Power BI dashboard, and build a predictive model to classify patients based on their health status.

# Data Overview:

	filename	chroma	st_rmsse	spectral	sc_spectral	ba_rolloff	zero_cross	mfcc1	mfcc2	mfcc3	mfcc4	mfcc5	mfcc6	mfcc7	mfcc8	mfcc9	mfcc10	mfcc11	mfcc12	mfcc13	mfcc14	mfcc15	mfcc16	mfcc17	mfcc18	mfcc19	mfcc20	label	ID	Gender	Age
0	1-a-h_hvay	0.127341	0.355538	1587.507	191.209	2737.972	0.067698	233.566	48.50442	-34.5181	-36.2723	-42.9844	-34.279	23.74783	-11.7224	-32.8372	20.4176	-6.4832	4.146173	-12.5348	-13.1078	-28.1906	1.143172	52.79483	39.0924	49.0706	Healthy	1	w	20	
1	1-b-a_hvay	0.215317	0.250342	1273.801	1850.303	1592.32	0.043325	-17.42	119.7057	0.822231	25.1678	-11.8211	-6.27188	-14.578	31.81473	-12.762	-5.6368	1.126529	7.275009	6.55109	23.2741	7.79023	7.03787	-9.1938	11.16256	-2.17351	-9.42094	Healthy	10	w	22
2	100-a-h_hvay	0.212512	0.250149	1005.208	1014.718	1756.967	0.047079	209.185	167.1357	-37.3882	-12.4357	-13.8215	-15.8407	-21.5113	10.1176	0.882222	-10.4021	4.140114	24.70295	5.14845	-6.0718	-20.4007	-4.15346	-2.72515	-10.6697	-2.39894	74.11264	Healthy	100	m	66
3	1000-a-h_hvay	0.198861	0.157396	1151.971	1529.491	1501.658	0.045487	-210.244	139.8195	-7.7401	13.63091	-52.2359	-10.8682	-6.9796	-4.7747	13.60316	-40.3458	-7.15746	9.966608	-2.41478	12.74221	-10.16355	7.797343	-26.315	-6.22205	0.53512	4.076993	Healthy	1000	m	31
4	1002-a-hv	0.370113	0.143538	1060.676	1346.3	1524.942	0.050109	-175.04	181.9052	-33.1816	-0.09176	-12.4857	-26.2217	12.85998	9.708848	-12.7372	-9.50608	5.919818	1.367481	11.55102	-2.58331	-20.3398	13.30665	-26.5975	-8.32453	10.76817	Healthy	1002	m	25	
5	15003-a-hv	0.232889	0.215562	1073.863	1041.309	1112.077	0.046057	221.869	162.185	-18.4184	-24.692	-36.8539	-13.127	11.10304	-7.8843	-9.21666	-18.7584	11.29885	7.064672	-20.801	-15.1804	8.48992	-29.6265	4.34001	-10.34266	-0.00415	Healthy	15003	m	35	
6	1604-a-hv	0.333185	0.231105	1077.814	1300.157	1360.275	0.045667	-145.841	174.8212	-33.6609	-10.9292	-24.0389	-21.1333	18.6081	16.20833	-20.515	5.81416	10.59185	-5.38708	11.02586	-6.16115	-24.7236	6.25876	-31.9161	17.1853	-4.90731	-7.95104	Healthy	1604	m	43
7	71005-a-hv	0.242202	0.18162	873.3722	1282.35	1068.378	0.040371	-189.152	186.335	12.51912	6.082908	-36.0106	-10.7136	-0.00256	0.80074	-6.3637	-1.41271	-31.1273	23.17631	15.61975	-5.01031	13.58097	-4.13191	-18.0889	-6.0785	-9.18721	14.0434	Healthy	71005	m	33
8	81006-a-hv	0.153282	0.18433	1066.099	1705.306	1183.004	0.043123	-256.187	127.9539	13.82222	-31.0978	-28.4544	12.86928	20.5197	13.28549	-13.7187	7.644126	-4.6204	-10.6621	-0.2249	-11.9769	-18.854	-0.96298	-9.12626	-11.1967	-3.24692	-7.80025	Healthy	81006	w	32
9	91007-a-hv	0.226716	0.203722	1085.429	1284.242	1526.536	0.043089	-206.332	131.8645	-37.6106	-0.29648	-38.6851	-4.26765	13.09653	-15.2264	-6.30196	19.41811	-16.1083	23.30435	-1.67457	-18.9593	-16.5312	0.63583	-4.51473	-4.65408	2.379562	-11.3416	Healthy	91007	m	42
10	1008-a-hv	0.191215	0.17789	820.3114	1130.498	707.5374	0.023662	-335.177	190.2773	26.24766	5.36874	-16.8598	-21.1406	14.13922	2.935519	-14.9148	-1.6841	-10.7849	6.398447	9.302208	-11.0170	0.092198	-12.9934	-11.7686	-12.0347	-1.61468	2.07742	Healthy	1008	w	50
11	11009-a-hv	0.184429	0.212725	1180.663	1647.716	1362.872	0.053674	-179.649	135.8172	0.320124	-12.2397	-44.7117	-9.91046	-1.58335	11.41033	-28.1011	-1.79053	9.74223	14.78432	0.292139	-3.95393	-13.284	-9.33448	-28.5737	17.84802	-14.166	-8.21463	Healthy	11009	m	40
12	1210-a-hv	0.210474	0.247952	1051.921	1019.606	1016.547	0.055105	-232.805	162.8002	-0.26342	1.913021	-56.6432	-14.0908	-16.0168	1.200328	-13.09997	-24.8426	4.546464	13.9295	1.284007	3.050909	-0.5261	-2.78932	-14.5265	-22.5254	3.275426	10.7330	Healthy	1210	m	31
13	13111-a-hv	0.272853	0.12768	993.7004	1137.781	1330.034	0.049121	-225.825	176.1395	-43.3519	-15.373	-24.0438	-12.3213	5.659922	2.796884	-5.1028	7.729187	-32.5297	17.6589	8.462935	-19.7108	-4.03039	2.248461	-13.2598	-3.37933	-6.0358	2.347671	Healthy	13111	m	28
14	141012-a-hv	0.261689	0.17907	1063.826	1252.11	1324.609	0.055161	-236.711	133.0197	-34.0601	25.86259	-49.4118	-14.034	13.65634	-12.6108	-23.145	-19.5293	-5.77098	15.25098	13.59159	-3.58402	13.18977	-23.0774	-15.6380	7.67544	-14.406	-5.7613	Healthy	141012	m	22
15	11009-a-hv	0.184429	0.212725	1180.663	1647.716	1362.872	0.053674	-179.649	135.8172	0.320124	-12.2397	-44.7117	-9.91046	-1.58335	11.41033	-28.1011	-1.79053	9.74223	14.78432	0.292139	-3.95393	-13.284	-9.33448	-28.5737	17.84802	-14.166	-8.21463	Healthy	11009	m	40
16	161014-a-hv	0.169917	0.250966	840.8565	1293.724	1037.47	0.049189	-197.065	166.4753	9.887601	-10.7365	-48.1591	-7.32388	-15.5345	-8.52953	19.48509	-37.1951	8.00126	17.07967	-1.06497	6.281257	3.924834	2.404845	-18.3032	-13.8327	-7.72009	-1.84695	Healthy	161014	m	28
17	171015-a-hv	0.190823	0.189345	1295.454	1781.977	1516.218	0.046419	-254.541	100.8923	-4.8488	-0.0274	-34.6001	17.8183	-32.5624	-3.93736	-9.27433	-19.0388	11.00425	18.33415	7.01702	-13.493	-4.29769	-9.37422	6.973831	-21.8699	4.040941	6.528244	Healthy	171015	m	37
18	181016-a-hv	0.229435	0.143123	879.8641	1126.972	1101.053	0.045475	-257.978	164.5786	-18.1533	-12.1705	-18.1586	-15.1097	-12.5038	-0.45713	10.10017	-11.6412	-18.1618	6.989252	9.242693	-12.14851	-22.3267	4.880014	-10.2933	-23.5521	2.54314	5.061826	Healthy	181016	m	38
19	191017-a-hv	0.254731	0.242359	878.3044	1394.277	976.1719	0.038974	-229.321	148.4775	-25.892	-2.367888	-48.2009	-0.27344	-12.3562	-12.788	-2.33135	-13.3278	10.4264	8.923584	18.85686	-12.1371	-17.2474	6.235973	-20.6057	-24.7925	21.2742	-15.417	Healthy	191017	m	40
20	201810-a-hv	0.247134	0.233589	964.3654	1106.205	1028.305	0.05485	-194.77	185.265	-15.299	-0.62522	-41.7592	-47.2965	23.2203	14.5826	-15.4471	-27.4441	-3.45169	-6.257	8.900938	7.262973	-6.61984	22.08159	-27.9281	-11.6575	2.577745	-1.73373	Healthy	201810	m	39
21	211019-a-hv	0.200396	0.185512	882.968	1294.231	1072.319	0.043390	-250.989	145.9424	-31.00002	5.693005	-52.673	-27.5874	10.57889	-1.82043	-27.9604	-10.5383	11.85348	-5.95109	-10.7766	18.43419	-12.6854	-6.09332	-28.7678	-1.03378	0.379338	-16.5848	Healthy	211019	m	48
22	22102-a-hv	0.142181	0.272746	1257.555	1398.59	1321.941	0.076514	-238.868	90.16582	-32.0026	-41.2494	-43.3551	-18.967	9.642206	0.045936	-29.2259	-12.6706	-6.40118	-9.9711	-12.9075	-6.10106	9.20323	27.05626	-40.4719	42.00658	9.747732	-16.3897	Healthy	22102	w	39
23	231026-a-hv	0.201341	0.215174	856.8432	1438.306	1122.5	0.025484	-212.237	174.5252	5.367576	-23.0866	-12.1742	-14.8659	14.14931	-7.60075	-7.50047	1.742459	-24.8648	33.14514	-11.7611	-6.0053	-3.29222	-18.1827	-6.17507	12.7316	-13.87	Healthy	231026	m	25	
24	241021-a-hv	0.235878	0.336685	970.978	1173.716	1061.143	0.03393	-185.126	178.0196	7.024777	7.230892	-53.281	-15.5793	-1.14908	-3.301	-4.23589	-24.166	-15.2145	14.50845	0.262472	4.801905	11.60072	1.124244	-17.4257	-4.78706	0.959997	-8.72108	Healthy	241021	m	45
25	251022-a-hv	0.191324	0.181218	1896.835	1349.616	0.048994	-245.765	102.385	8.678334	-17.8823	-33.349	9.419088	29.3332	3.165472	-1.4471	4.781082	-19.6813	-3.69641	-19.5549	-16.8053	-12.4545	-19.9085	-5.96716	-7.85587	-0.66018	3.276174	Healthy	251022	w	45	
26	261023-a-hv	0.248022	0.24855	1013.142	1232.537	1177.443	0.055152	-199.762	146.3228	-6.77592	1.888419	-33.1431	-21.1373	10.27794	5.797043	3.932112	-2.5485	-7.67138	26.62555	-13.8686	0.682614	-11.2853	-5.51002	-21.7131	2.23833	6.165217	-4.38755	Healthy	261023	m	41
27	271024-a-hv	0.147948	0.22476	892.112	1150.728	1133.528	0.03292	-274.996	151.8642	-15.3651	-15.4538	-14.1412	-5.837359	22.0709	14.64664	-21.1819	10.7658	-4.53367	-22.9127	-3.8076	-15.2564	-12.9628	-15.5208	-2.67843	-13.1534	-8.19123	-5.0984	Healthy	271024	m	37
28	281025-a-hv	0.29972	0.242206	1152.941	1320.318	2210.645	0.043167	-154.793	167.3794	-35.7273	22.1703	-32.9656	-54.556	5.068939	21.97803	-12.578	16.71859	5.793391	3.916778	-12.7065	11.01613	-12.649	-20.796	-10.5844	6.785952	Healthy	281025	m	42		
29	291026-a-hv	0.338367	0.212788	1099.158	1228.026	1243.742	0.046097	-197.174	175.6335	-18.930	-0.742	-38.0512	-19.5037	15.90117	1.536258	-35.3672	14.18629	8.661491	-1.8429	-20.1189	-1.66257	-4.39172	-5.4821	6.19047	-5.4604	Healthy	291026	m	42		
30	301027-a-hv	0.196824	0.232659	1712.219	1059.829	1033.369	0.044281	-231.947	168.9327	1.448849	-9.25401	-27.1973	-14.8941	-5.55498	2.99549	-2.5737	-16.2638	-2.59925	3.42211	10.0404	-0.58738	-6.66813	-0.73017	-14.1097	-11.7302	-3.94905	0.80511	Healthy	301027	m	37
31	311028-a-hv	0.226819	0.160989	1240.619	1437.617	1654.788	0.07732	-234.711	113.8309	-4.717	-17.829	-15.1081	-8.5337	12.1091	4.07663	-14.1573	2.79961	-1.29861	2.1864	3.74881	-17.3843	15.34601	-1.7945	-15.2754	9.937485	-5.76761	Healthy	311028	m	41	
32	321029-a-hv	0.195103	0.165637	1066.025	1535.501	1478.295	0.059719	-260.582	136.4148	-30.9684	-22.5926	-19.444	-21.0282	-26.85797	-11.1487	-9.36114	10.26789	6.87454	-1.9998	-14.0744	-16.1845	8.260423	-1								

# Approach

## Phase 1: Understanding the Project Overview

- Define the objectives outlined in the project statement to ensure clarity with the end-goal of this assignment.
- Critically assess the provided dataset and deconstruct it into digestible components to improve comprehension.

## Phase 2: Importing the data for Exploratory Data Analysis (EDA)

- Import the provided .csv file into Google Colab, assign it to a dataframe, and evaluate the schema of the dataset.
- Ensure each column is of the correct data type, impute any null values, check for any duplicated data, drop irrelevant columns, standardise column names to lowercase, and perform feature encoding where appropriate (e.g. with the gender column).
- Create visualisations for the dataset to uncover any trends, patterns, and relationships between the features.

## Phase 3: Machine Learning

- Assign audio extractions as features (independent variables) and 'Status' as the target (dependent variable). Segregate the data into train-test split, scale the features where needed, and initialise several machine learning models to see which one yields the highest %.

## Phase 4: Power BI

- Export the cleaned up dataset from Google Colab, so it can be imported into Power BI to transform the data and create interactive visualisations.

## Phase 5: Power Presentation

- Compile the entire process into a power presentation to outline any notable findings.

# Project Workflow:

Understand  
Objectives & Data



Perform  
Exploratory Data  
Analysis



Build a predictive  
model (machine  
learning)



Create an  
interactive results  
dashboard



Design a  
presentation slide



Define the objectives and analyse the provided dataset to ensure that the goal of the project is clear.

Perform EDA on the .csv file, after assigning it to a dataframe, to extract insights and create visualisations.

Create and compare several predictive models to analyse the data and predict the status outcome of participants.

Export the .ipynb file and upload it into Power BI to create an interactive report that outlines the significant findings.

Build a presentation deck to illustrate the process of the project and the key insights which have been extracted from the data set.

# Data Cleaning

The dataset was cleaned by:

1. checking duplicate or irrelevant entries
2. executing a comprehensive age validation to ensure that there are no suspicious age values (e.g. -1 years old, or 125 years old)
3. standardising column names (converting column names into lower case) for consistency
4. dropping irrelevant columns
5. renaming columns for better visibility (e.g. g → gender)
6. feature engineering to transform the continuous, numerical values of the age column into a categorical column comprised of age bins (e.g. 0-12, 12-24, 24-36, etc.)
7. converting categorical variables like "gender" into numerical format (e.g., "gender\_encoded") for analysis.
8. standardising "age\_group" and "status" columns to ensure uniform labeling.

This process resulted in a structured, analysis-ready dataset suitable for further modeling or statistical evaluation.

# Machine Learning:

- Several predictive models were made to classify audio samples into two health status categories (Healthy and Unhealthy) and compared against each other to identify the most effective model.

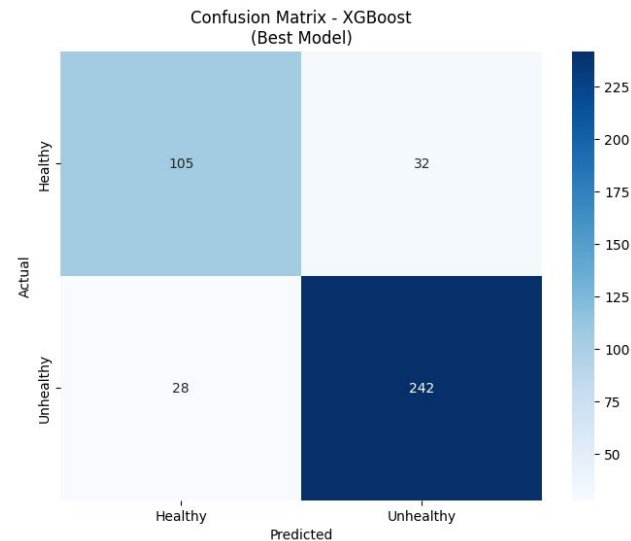
## Data Preparation

- This model utilises a comprehensive set of features extracted from the audio files. These include:
  - Spectral features (chroma stft, spectral centroid, spectral bandwidth, rolloff)
  - Temporal features: rmse and zero crossing rate
  - Cepstral features: 20 Mel-Frequency Cepstral Coefficients (MFCC1-MFCC2)
- The target variable was the 'Status' (Healthy/Unhealthy) column.
- Data splitting: the dataset of 2035 inputs was split into training (80%) and testing (20%) sets, using stratification, to ensure the class distribution was preserved in both splits. A fixed random\_state = 42 was used for reproducibility.
- Feature Scaling: a StandardScaler was applied to normalise the features, which is important for models that are sensitive to the scale of data (e.g. K-Nearest Neighbors). Tree-based models (XGBoost, Random Forest Tree, Decision Tree) were trained on unscaled data.
- Label encoding: the string labels ('Healthy', 'Unhealthy') in the target variable were encoded to numerical values (0,1), using LabelEncoder to be compatible with scikit-learn models.
- 5 candidate models were selected to evaluate different algorithmic approaches: XGBoost, Random Forest, Logistic Regression, K-Nearest Neighbors, Decision Tree.
-



# Machine Learning:

Model	Accuracy	Precision	Recall	F1-Score
XGBoost	85.3%	85.2%	85.3%	85.2%
Random Forest	84.3%	84.0%	84.3%	84.1%
Logistic Regression	82.3%	82.7%	82.3%	82.5%
Decision Tree	76.9%	76.8%	76.9%	76.9%
K-Nearest Neighbors	74.4%	74.4%	74.4%	74.4%



## Performance Metrics and Evaluation Strategy

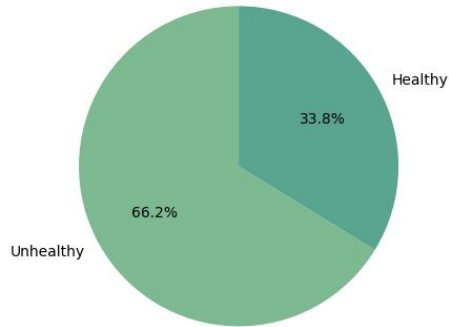
The primary metric for evaluating model performance was accuracy, defined as the proportion of correct predictions among the total predictions.

The results show two superior models: XGBoost and Random Forest, having the highest scores across all evaluation metrics. XGBoost's precision (85.2%) and recall (85.3%) indicates a well-balanced model that performs consistently across different classes. XGBoost is slightly more accurate but can be sensitive to hyperparameters. Random Forest is generally more robust and less prone to overfitting on noisy data.

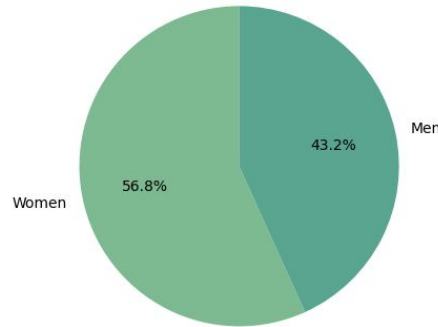
- On the confusion matrix, we see that XGBoost had 105 true negatives and 242 true positives.
- XGBoost is better at identifying unhealthy cases (91% recall) than healthy cases (79% recall), indicating that this model is better at identifying unhealthy voices.

# Data Analysis:

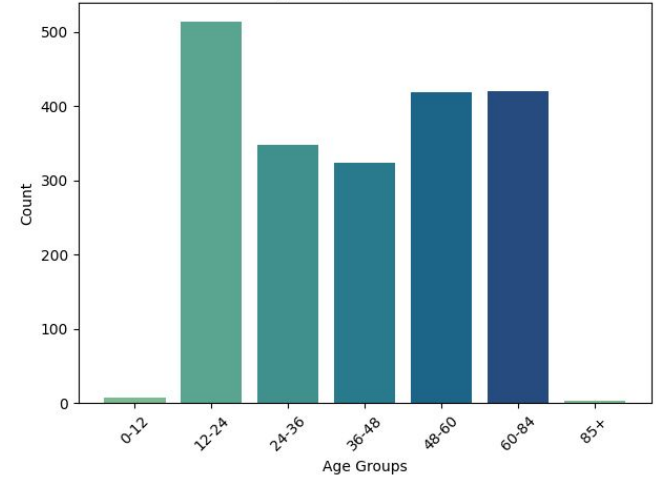
Healthy vs Unhealthy Distribution



Gender Distribution



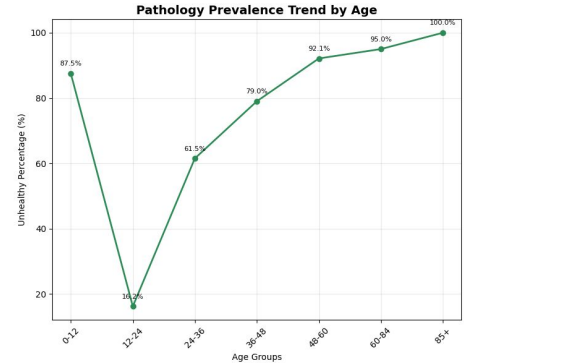
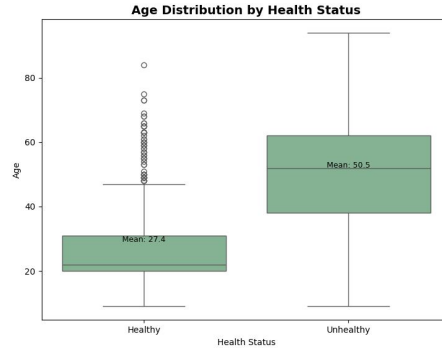
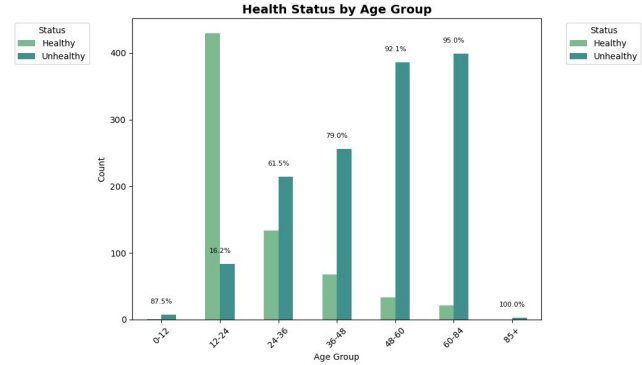
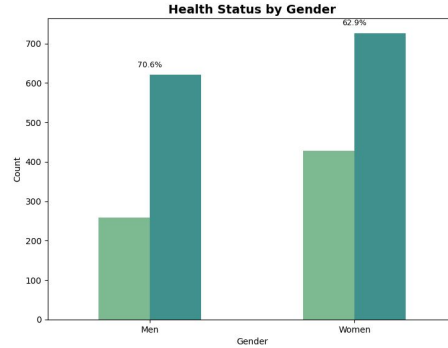
Age Group Distribution



The sample was an inclusive sample, with participants ranging from children, teens, adults, and elderly. The demographic was also consistent of women (1150) and men (880). Each individual's health status also varied between 'Healthy' (687) and 'Unhealthy' (1348), with the latter being the more prominent class.

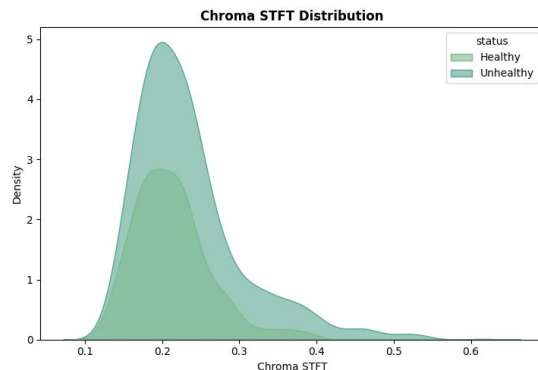
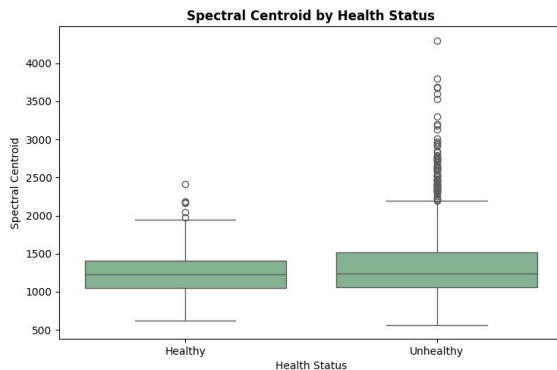
# Data Analysis:

- There is a higher percentage of Unhealthy individuals in men (70.6%) compared to women (62.9%).
- There is a low pathology percentage in the younger demographic, particularly in the 12-24 year range, compared to their older counterparts.
- Mean age distribution shows that, on average, Unhealthy participants are older compared to Healthy participants.
- There is a disproportionately high pathology % on either spectrum (0-12 and 85+ ages), but these groups were underrepresented.
  - We need a bigger sample for these age groups to make any conclusive interpretations.

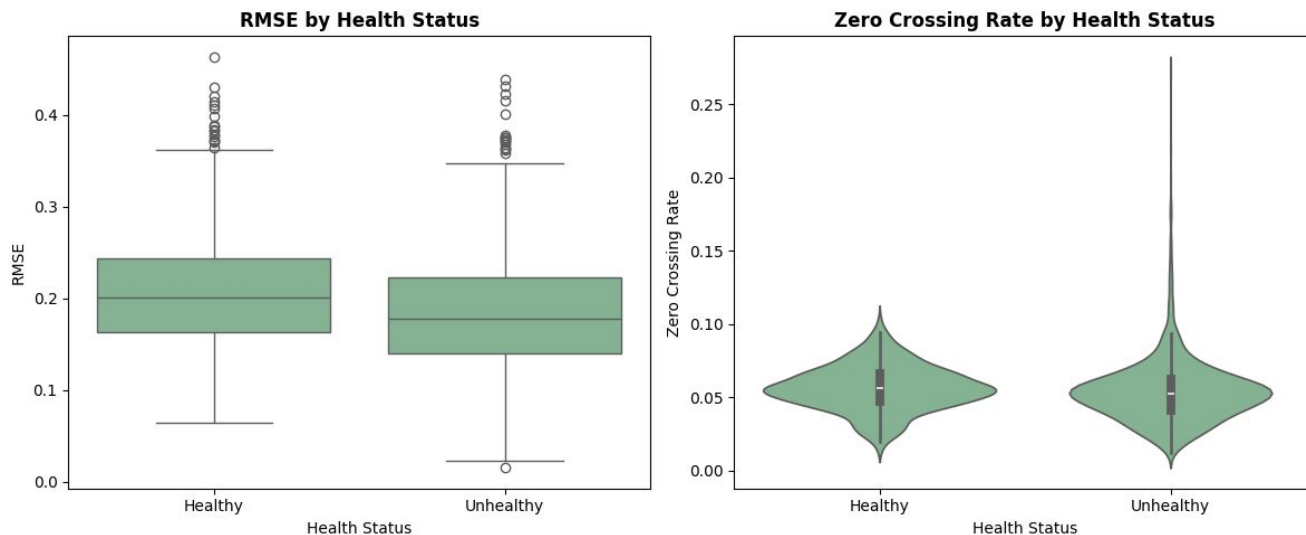


# Data Analysis:

- Spectral Centroid is higher in Unhealthy participants compared to Healthy participants, with the Unhealthy group having outliers up to 4000 Hz.
  - Consequently, the Spectral Roll-Off is also higher in the Unhealthy group as they have more of a high frequency energy.
- The spectral bandwidth is wider, bimodal, in the Unhealthy group compared to the Healthy group.
- The Chroma STFT for the Unhealthy group has a higher density across many pitch classes. Meanwhile the Healthy group has a more concentrated and stable distribution.

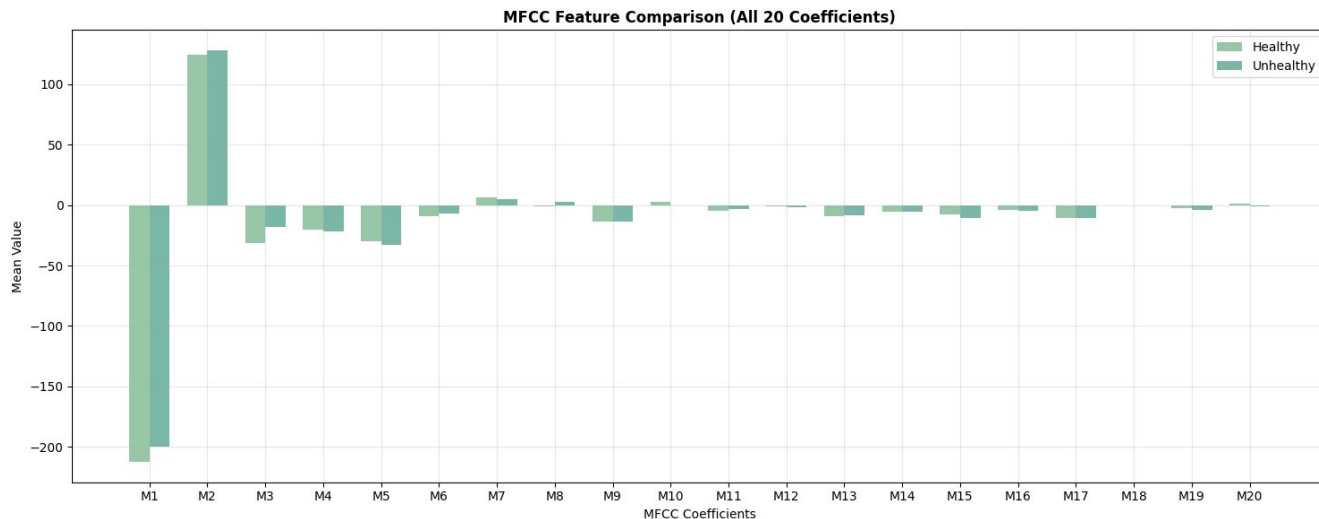


# Data Analysis:



- RMSE (Root Mean Square Energy) is higher for the Healthy group compared to the Unhealthy group, suggesting that the Healthy group produces a stronger, more robust acoustic signal.
- In contrast, the Zero Crossing Rate (ZCR) is lower in the Healthy group, having a distribution tightly clustered at lower ZCR values. This elevated ZCR indicates high-frequency noise (signifying turbulent airflow, breathy phonation, hoarseness).
- Both graphs show that while the Unhealthy group is quieter overall (lower RMSE), their audio files have a noisier and sharper signal, with high-frequency components.

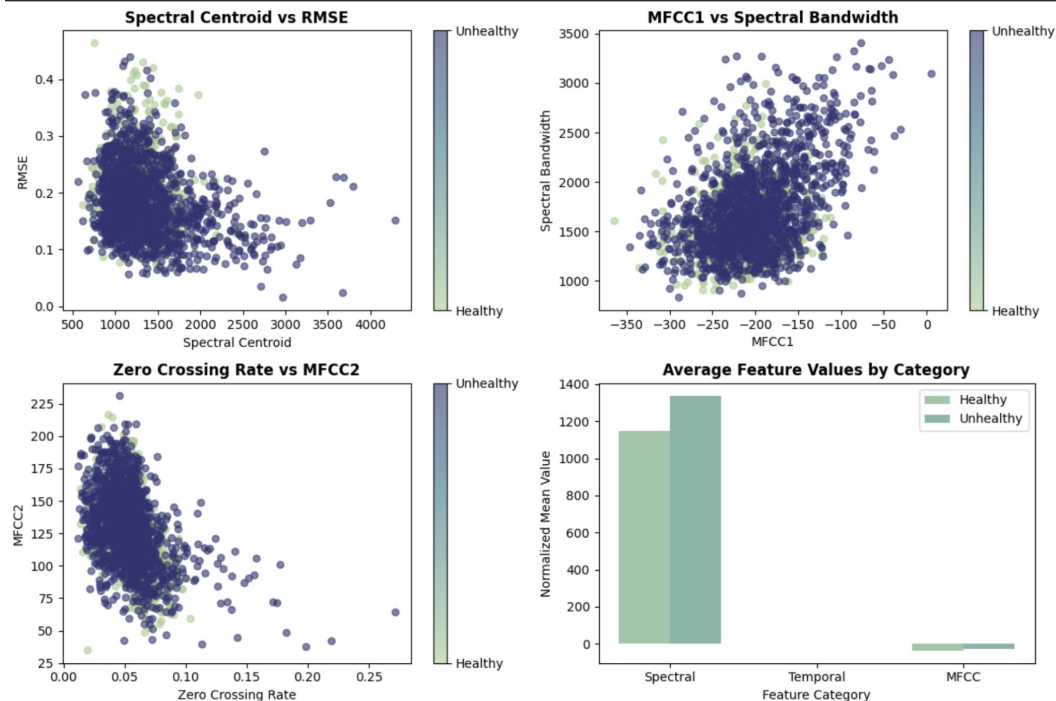
# Data Analysis:



- Shows that for every coefficient, there is a significant difference in mean values for each group.
- MFCC1 shows that there is a difference in spectral tilt between the two groups: the Unhealthy participants have a higher (less negative) MFCC1, thus their voices produce more low frequency energies.
- MFCC2-5 also consistently show that the Unhealthy group have higher mean values compared to the Healthy group, suggesting that the contour of the sound's spectrum varies in between states.
- As expected, the difference in MFCC values show that each group has a distinct acoustic structure. Vowels are defined by the shape of the vocal tract; a healthy individual will have one resonant signature, and an unhealthy individual with a different vocal tract shape (due to asthma, vocal fold nodules, etc) will have a different resonant signature.

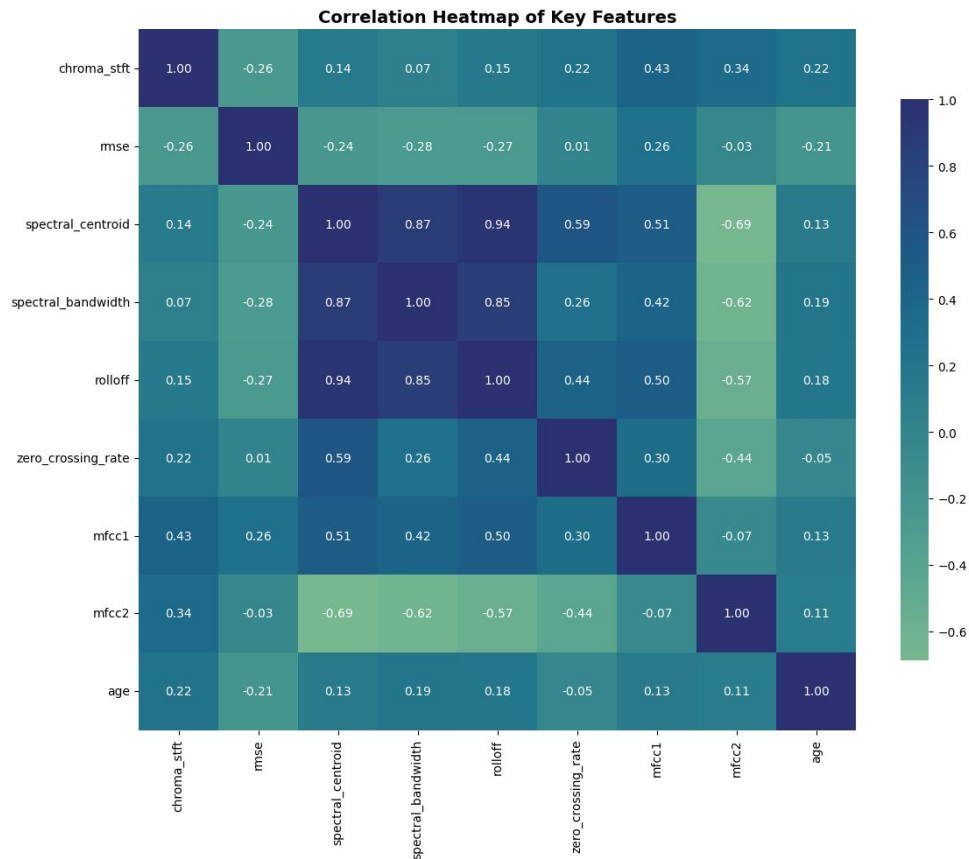
# Data Analysis:

- **Spectral Centroid vs RMSE** shows two distinct clusters with a negative correlation between the features.
  - The Healthy cluster has a high RMSE and a low Spectral Centroid (dull, low frequency sound)
  - The Unhealthy cluster has a low RMSE and a high spectral centroid (bright, high-frequency sound)
- **MFCC1 vs Spectral Bandwidth:**
  - The Healthy cluster predominantly occupies a specific region of MFCC1 (consistent spectral tilt) and has a lower Spectral Bandwidth, indicating that the sound energy is more focused.
  - The Unhealthy cluster is more scattered around MFCC1, and it fills a wider range of the bandwidth (frequencies), consistent with noise and turbulent signals.
- **Zero Crossing Rate vs MFCC2**
  - The Healthy cluster has a low ZCR (tonal, smooth sound) and has a stable spectral shape for the MFCC2.
  - The Unhealthy cluster has a higher ZCR (noisy, chaotic sound) and an altered spectral shape, showing that their audio files has a noisier texture due to changes in timbral quality (MFCC2).
- **Average Feature Values by Category** show you that for every feature recorded, the Healthy group's average differs from the average of the Unhealthy group.
  - The temporal bars are not visible as they are too small to fit on the same scale as the spectral values.



# Data Analysis:

- Strong positive correlations:
  - Spectral Centroid and rolloff (0.94): indicates that “brighter” (higher centroid) sounds have their energies spread out to higher frequencies.
  - Spectral Bandwidth and rolloff (0.85): indicates that there is a strong relation between high frequency cutoff and overall frequency.
- Moderate positive correlations:
  - Zero Crossing Rate and Spectral Centroid (0.59): as the sound becomes more high pitched and ‘brighter’, it also becomes more noisy and less tonal (high ZCR)
  - MFCC1 has moderate correlations with plenty of spectral features (centroid, bandwidth, rolloff), which shows that its spectral energy shape is influenced by broader properties.
- Strong negative correlations:
  - Spectral Centroid and MFCC2 (-0.69): indicates that as the sound becomes brighter, the value of MFCC2 decreases. Inverse relationship shows that the timbral texture (MFCC2) systematically changes as the pitch of the voice increases.





# Conclusions & Recommendations

## Conclusions:

- Men had a higher percentage of Unhealthy participants compared to women. This may be due to behavioural factors, such as men are more likely to engage in risk factors, like smoking. Additionally, there is an over-representation of men in construction and industrial work, which may exacerbate their risk of developing vocal pathologies.
- As expected, vocal pathologies became more prevalent in the older demographic . Presbyphonia (age-related voice age) leads to vocal folds thinning, loss of elasticity and mass, which can lead to tremor and hoarseness. Additionally, neurological decline is often observed in individuals as they age (e.g. stroke, dementia, Parkinson's), resulting in poor control of the larynx.
- The 'audio fingerprint' of Healthy participants is consistently and measurably different from Unhealthy participants, allowing them to produce acoustically distinct and mechanically classifiable sounds.
  - The Unhealthy group had an acoustic signature that was energetically weaker, spectrally busier, noisier, less harmonic, and had a different timbral texture than its Healthy counterpart.
- This study provides a strong foundation for developing an accurate, non-invasive, and automated diagnostic tool based on audio analysis.

# Conclusions & Recommendations

## Recommendations:

- **Have better class balance:** there was a significant imbalance between the number of Healthy and Unhealthy participants. This could lead to the model being biased towards the majority class, learning to predict “Unhealthy” most of the time and still maintaining high accuracy.
- **Have a more diverse dataset:** this dataset had more women than men, overlooking any anatomical and biological differences between the two groups. For example, women have smaller larynxes and shorter vocal tracts; physical aspects which will affect the measurements of their features.
  - A model trained on a female-skewed dataset will perform poorly when deployed on a general population as it will be systematically biased against men.
- **Have better representation of key demographic groups:** this dataset only had 8 young children (ages 0-12) and 4 elderly aged 85 and above, demonstrating paediatric and geriatric neglect. This means that the “Unhealthy” profile may be confounded by age and gender.
- **Add more confounding variables:** other factors, such as body habitus, can affect results. For example, an obese individual may have an impacted speech production, and therefore, will have differing MFCCs compared to a leaner individual.
  - Excess fat deposition in the neck, pharynx, and uvula changes the geometry of the upper airways, which could impact their MFCC values.