

CSE343: Machine Learning

Assignment-2

REPORT

Anushka Srivastava (2022086)

SECTION A

- a. The likelihood is 75.18%.

a) Let D be event that company issues a dividend

$$P(D) = 0.8, P(\bar{D}) = 0.2$$

If a company issues dividend

$$\mu = 10\%, \sigma^2 = 36\%$$

Let X be profit increase

$$P(X|D, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X-\mu)^2}{2\sigma^2}}$$

If a company does not issue dividend

$$\mu = 0\%, \sigma^2 = 36\%$$
$$P(X|\bar{D}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X-\mu)^2}{2\sigma^2}} = 0.0403$$
$$P(X=4|\bar{D}) = P(\bar{D}) \cdot P(X=4|\bar{D}) + P(D) \cdot P(X=4|D)$$
$$\Rightarrow 0.0429$$
$$P(D|X=4) = P(D) \cdot P(X=4|D) / P(X=4)$$

The likelihood is 75.18%.

- b. The decision trees are as follows:

b.) Step 1

$P(\text{Yes}) = \frac{7}{12}$, $P(\text{No}) = \frac{5}{12}$

$$H(Y) = -\frac{7}{12} \log_2 \frac{7}{12} - \frac{5}{12} \log_2 \frac{5}{12}$$

$$= 0.979$$

Class Timing

$$H(\text{Morning}) = -\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \left(\frac{1}{4}\right)$$

$$= 0.811$$

$$H(\text{Noon}) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4}$$

$$\Rightarrow H(\text{Noon}) = 0.811$$

$$H(\text{Afternoon}) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4}$$

$$(X)_{4 \times 2} \times (Y)_{2 \times 2} = (X)_{4 \times 2} \times (A)_{2 \times 2}$$

$$I(X) = H(X) = -\log_2 H(Y|X)$$

$$= 0.979 - \frac{4}{12} \times 0.811 = \frac{4}{12} \times 0.158$$

$$= 0.042$$

$$[8P] = \frac{8 \times 0.8}{8 \times 0.8 + 3 \times 0.2} = \frac{8 \times 0.8}{11} = 0.727$$

Sleep

$$H = -p \log_2 p + q \log_2 q$$

$$H(\text{Sleep} = \text{Yes}) = \frac{6}{12} \log_2 \frac{6}{12} + \frac{6}{12} \log_2 \frac{6}{12}$$

$$H(\text{Sleep} = \text{No}) = \frac{1}{12} \log_2 \frac{1}{12} + \frac{5}{12} \log_2 \frac{5}{12}$$

$$= 0.650$$

$$H = -p \log_2 p + q \log_2 q$$

$$24 = 0.979 - \frac{6}{12} \times 0 + \frac{6}{12} \times 0.650$$

$$= 0.654$$

Weather

$$H = -p \log_2 p + q \log_2 q$$

$$H(\text{Cloudy}) = -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} = 0.721$$

$$H(\text{Rainy}) = 0 - \frac{2}{2} \log_2 \frac{2}{2} = 0$$

$$H(\text{Hot}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.970$$

$$H = -p \log_2 p + q \log_2 q$$

$$24 = 0.979 - \frac{5}{12} (0.721) - 0 - \frac{5}{12} \times 0.970$$

$$\Rightarrow 0.274$$

Since sleep has maximum information gain, we use it as root node.

$$H(\text{Sleep} = \text{Yes}) = 0 ; H(\text{Sleep} = \text{No}) = 0.650$$

* Time | Sleep = No

$$H(\text{Morning} | \text{Sleep} = \text{No}) = -(p_1 \log_2 p_1 + p_2 \log_2 p_2) H$$
$$\rightarrow -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) = 1$$

$$H(\text{Noon} | \text{Sleep} = \text{No}) = -(p_1 \log_2 p_1 + p_2 \log_2 p_2) H$$
$$\rightarrow -0 - 0 = 0$$

0.25.0 =

H(Afternoon | Sleep = No)

$$H(\text{Afternoon} | \text{Sleep} = \text{No}) = -(p_1 \log_2 p_1 + p_2 \log_2 p_2) H$$

$$IG = 0.654 - \frac{2}{6} = 0.321 M 25.0 =$$

0.25.0 =

Weather | Sleep = No

$$H(\text{Cloudy} | \text{Sleep} = \text{No}) = -(p_1 \log_2 p_1 + p_2 \log_2 p_2) H$$

$$H(\text{Cloudy} | \text{Sleep} = \text{No}) = -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) H$$

$\rightarrow 1 - 1 = 0$

$$H(\text{Rainy} | \text{Sleep} = \text{No}) = 0 M 25.0 = (0.0) H$$

0.25.0 =

H(Hot | Sleep = No) = 0

$$H(\text{Hot} | \text{Sleep} = \text{No}) = -(p_1 \log_2 p_1 + p_2 \log_2 p_2) H$$

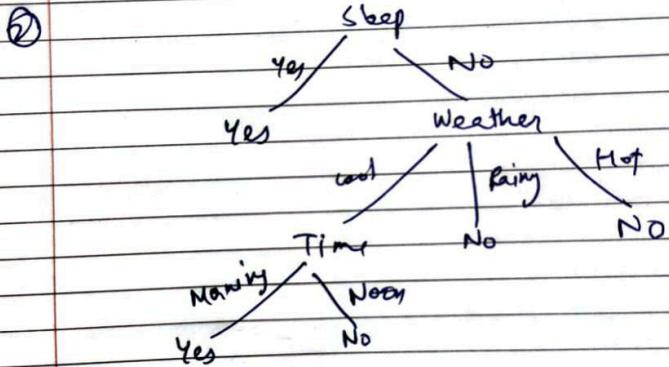
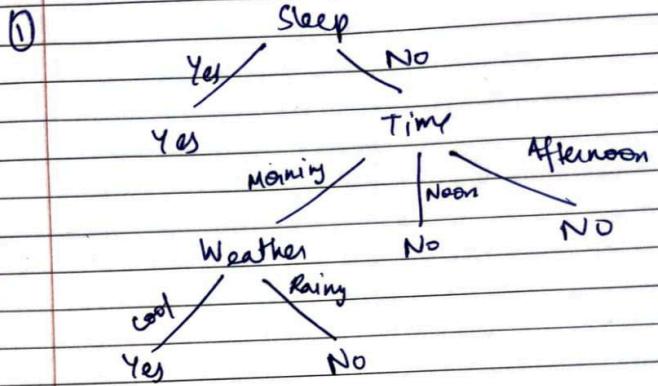
$$IG = 0.654 - \frac{2}{6} = 0.321$$

0.25.0 =

The information gain for both Time and Weather is same and hence we can choose any parameter.

$$0.25.0 = (0.654 - \frac{2}{6}) H = 0.321 H$$

Possible Decision Trees



c. Reference: <https://www.cs.cmu.edu/~avrim/ML10/lect0125.pdf>

C.

The number of mistakes M on S made by the perceptron algorithm is at most $(\frac{1}{\gamma})^{2d}$, which is function of polynomial $(\frac{1}{\gamma})^{2d}$, where

$$\gamma = \min_{x \in S} \frac{|w^* \cdot x|}{\|x\|}$$

where w^* is a unit-length vector. If we scale examples to have Euclidean length 1, then γ is the minimum distance of any example to the plane $w^* \cdot x = 0$.

Proof:

Claim 1: $w_{t+1} \cdot w^* > w_t \cdot w^* + \gamma$

Every time, a mistake is made, the length

Every time, a mistake is made, the dot product of our weight vector with the target increases by at least γ .

If x was a positive example, then we get

$$(w_{t+1} \cdot w^*) - (w_t \cdot w^*) = w_t \cdot w^* + \gamma \cdot w^* \geq w_t \cdot w^* + \gamma.$$

Why, if x was a negative example, we get

$$(w_{t+1} \cdot w^*) - (w_t \cdot w^*) = w_t \cdot w^* - \gamma \cdot w^* \geq w_t \cdot w^* + \gamma.$$

Claim 2: $\|w_{t+1}\|^2 \leq \|w_t\|^2 + 1$.

every time, a mistake is made, the length squared of our weight vector increases by at most 1.

If x was a positive example, we get $\|w_{t+1}\|^2 =$

$$\|w_t\|^2 + 2w_t \cdot x + \|x\|^2. This is less than$$

$$\|w_t\|^2 + 1 because w_t \cdot x is negative.$$

Claim 1 implies that after M mistakes, $w_{M+1} \cdot w^* \geq \gamma M$. On other hand, Claim 2 implies that after M mistakes, $\|w_{M+1}\| \leq \sqrt{M}$. Now, all we need to do is use the fact that $w_t \cdot w^* \leq \|w_t\|$, since w^* is a unit vector. So, this means we must have $\gamma M \leq \sqrt{M}$, and thus $M \leq \gamma^2$.

The number of mistakes (including margin mistakes) made by Margin Perceptron (γ) on S is at most $8/\gamma^2$.

As before, each update increased $w_t \cdot w^*$ by atleast γ . For the original algorithm, we had $\|w_{t+1}\|^2 \leq \|w_t\|^2 + 1$, which implies $\|w_{t+1}\| \leq \|w_t\| + \frac{1}{2\|w_t\|}$.

For the new algorithm, we get

$$\|w_{t+1}\| \leq \|w_t\| + \frac{1}{2\|w_t\|} + \frac{\gamma}{2}$$

which we can see by breaking each x into its orthogonal part and its parallel part. We can now solve this directly, but just to get a simple upper bound, just notice that if $\|w_t\| > 2/\gamma$, then $\|w_{t+1}\| \leq \|w_t\| + 3\gamma/4$.

So, after M updates, we have:

$$\|w_{M+1}\| \leq \frac{2}{\gamma} + \frac{3M\gamma}{4}$$

Solving $4\gamma + M\gamma \leq 2/\gamma + 3M\gamma/4$, we get

$$M \leq 8/\gamma^2, \text{ as desired.}$$

d. Answer

Date _____

d.)

$$a) P(\text{spam}) = \frac{1}{2}, P(\text{not spam}) = \frac{1}{2}$$

$$P(\text{cheap} | \text{spam}) = \frac{1}{2}$$

$$P(\text{buy} | \text{spam}) = 1$$

$$P(\text{cheap} | \text{not spam}) = \frac{1}{2}$$

$$P(\text{buy} | \text{not spam}) = \frac{1}{2}$$

$$b) P(\text{spam} | \text{cheap, not buy}) = P(\text{cheap} | \text{spam}) \times P(\text{not buy} | \text{spam}) \\ \times P(\text{spam})$$

$\rightarrow 0$

$$P(\text{Not spam} | \text{cheap, not buy}) = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8}$$

since $P(\text{not spam} | \text{cheap, not buy}) > P(\text{spam} | \text{cheap, not buy})$, the email is not spam.

c.) While using Naive Bayes estimation, if any one of the conditional probability is 0, then the entire probability evaluates to be zero.
To address this, we use the following:

$$i) \text{ Laplace estimate: } P(A_i | C) = \frac{N_{iC} + 1}{N_C + C}$$

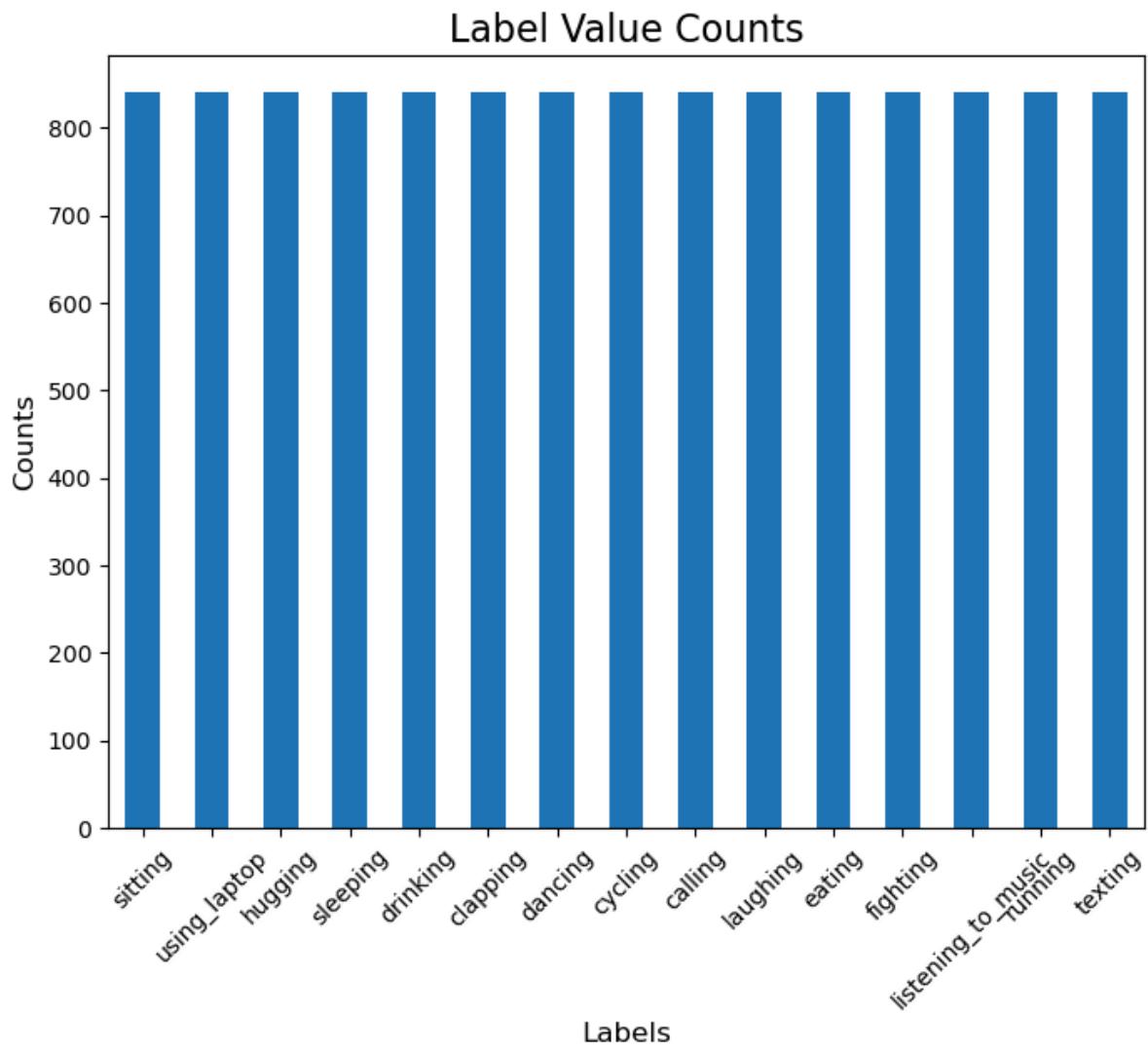
$$ii) m\text{-estimate: } P(A_i | C) = \frac{N_{iC} + mp}{N_C + M}$$

where $C = \text{no. of classes}$, $p = \text{prior probability}$, $m = \text{parameter}$

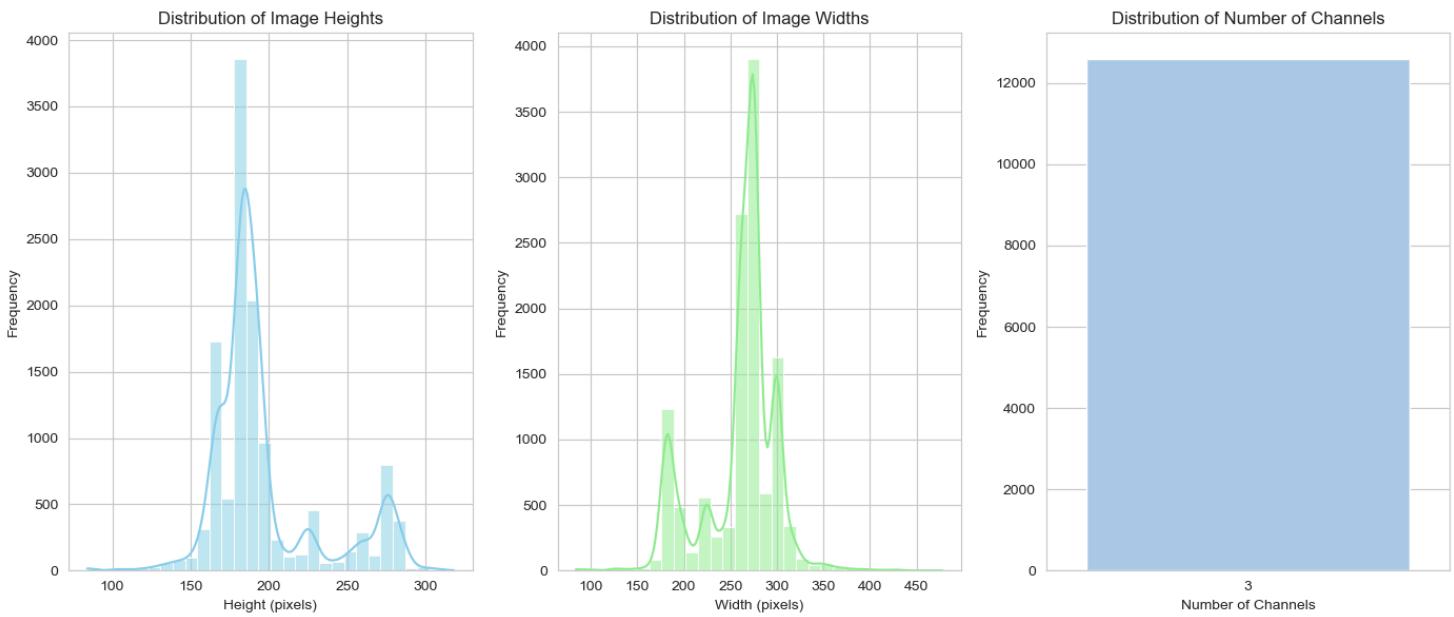
SECTION C

1. EDA

- a. We create a bar plot of Labels vs Counts.

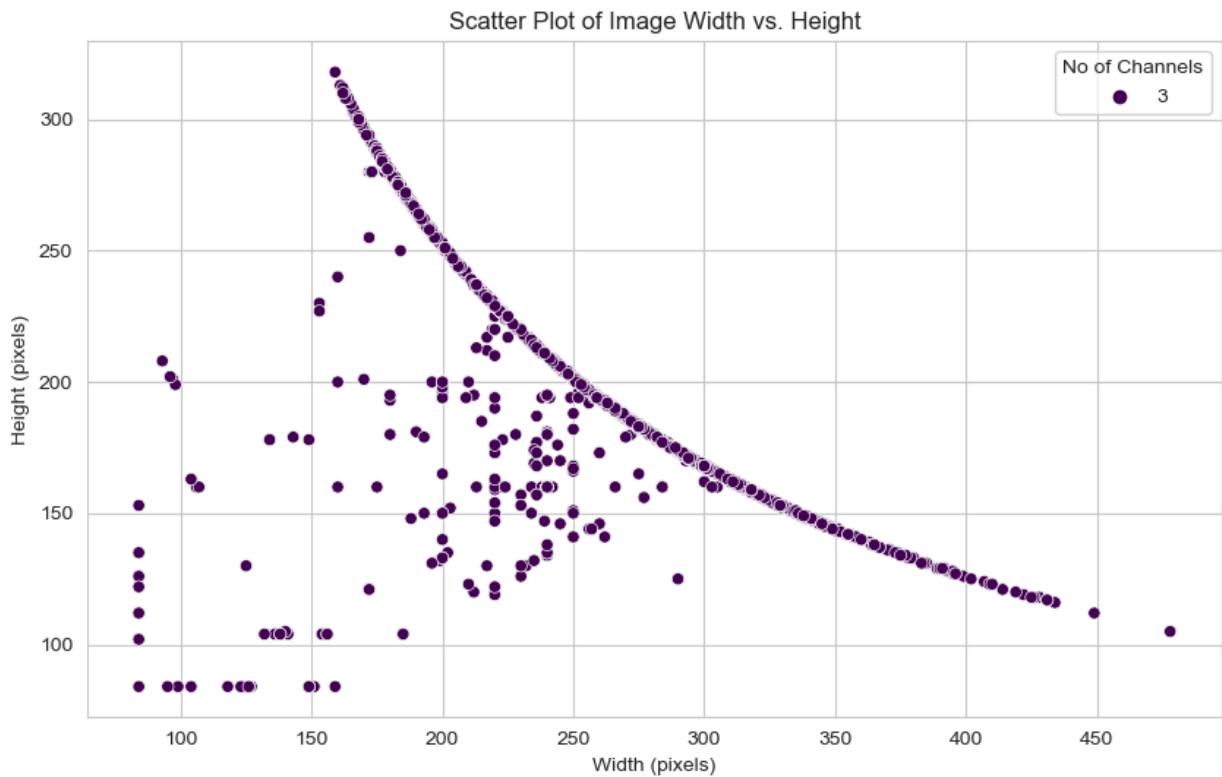


We notice that every class has an equal number of images, that is 840. Hence, we do not notice any class imbalance in our dataset.
We plot graphs for the distribution of image sizes and the number of channels.



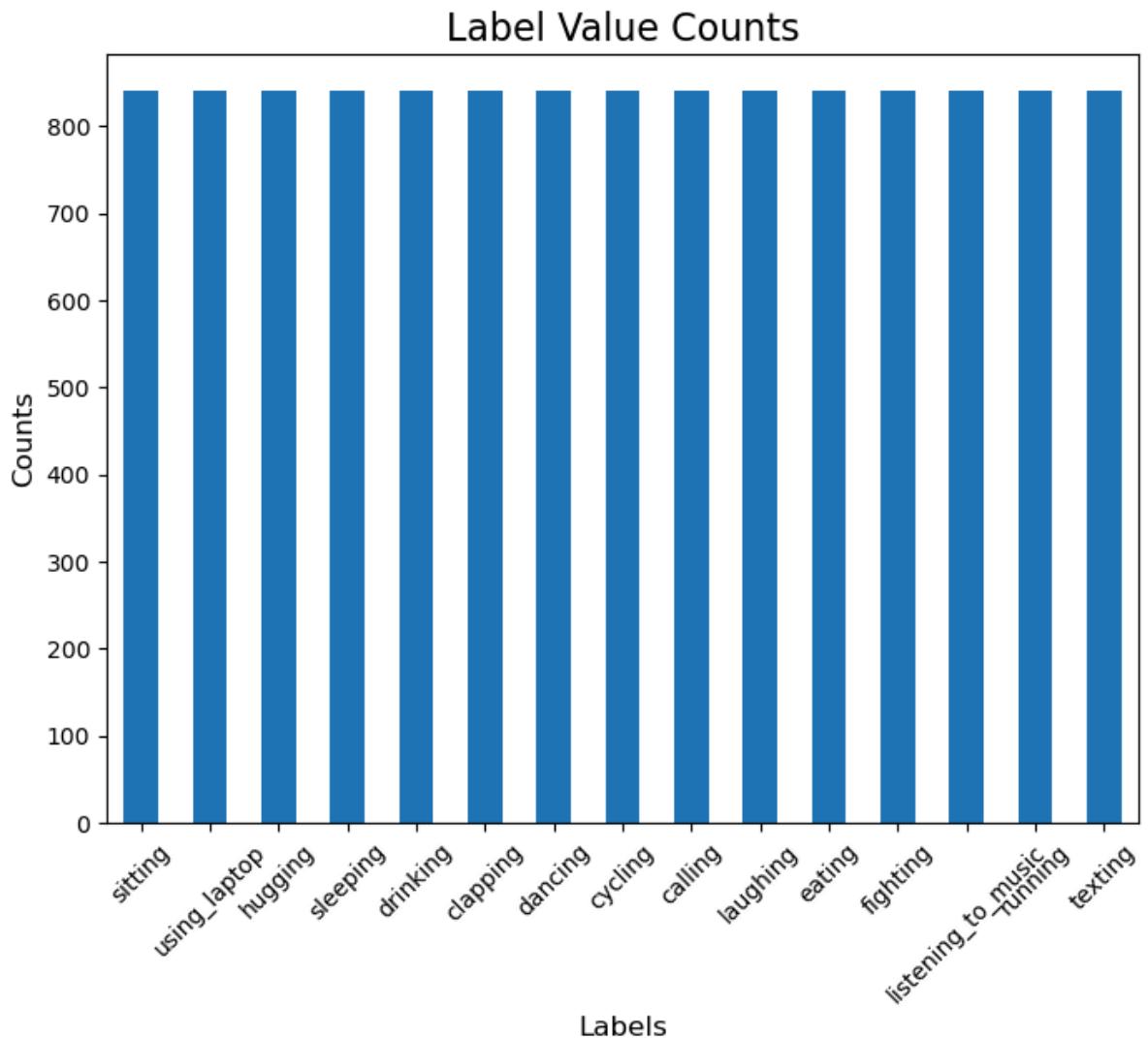
We notice that the height and weight follows a nearly normal distribution. Most images have a height between 150-200 pixels and a width between 250-300 pixels. All images have 3 channels.

We also plot a relationship between height and width.



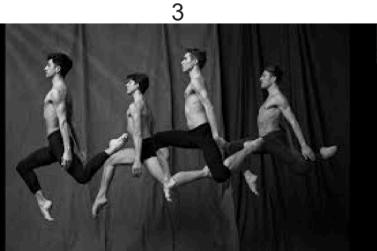
We notice that the height and width have an inverse relationship, that is, when the width increases, the height decreases, and vice versa.

b. Visualization to display the distribution of classes



Each class has an equal number of objects.

We display 1 image from each class.



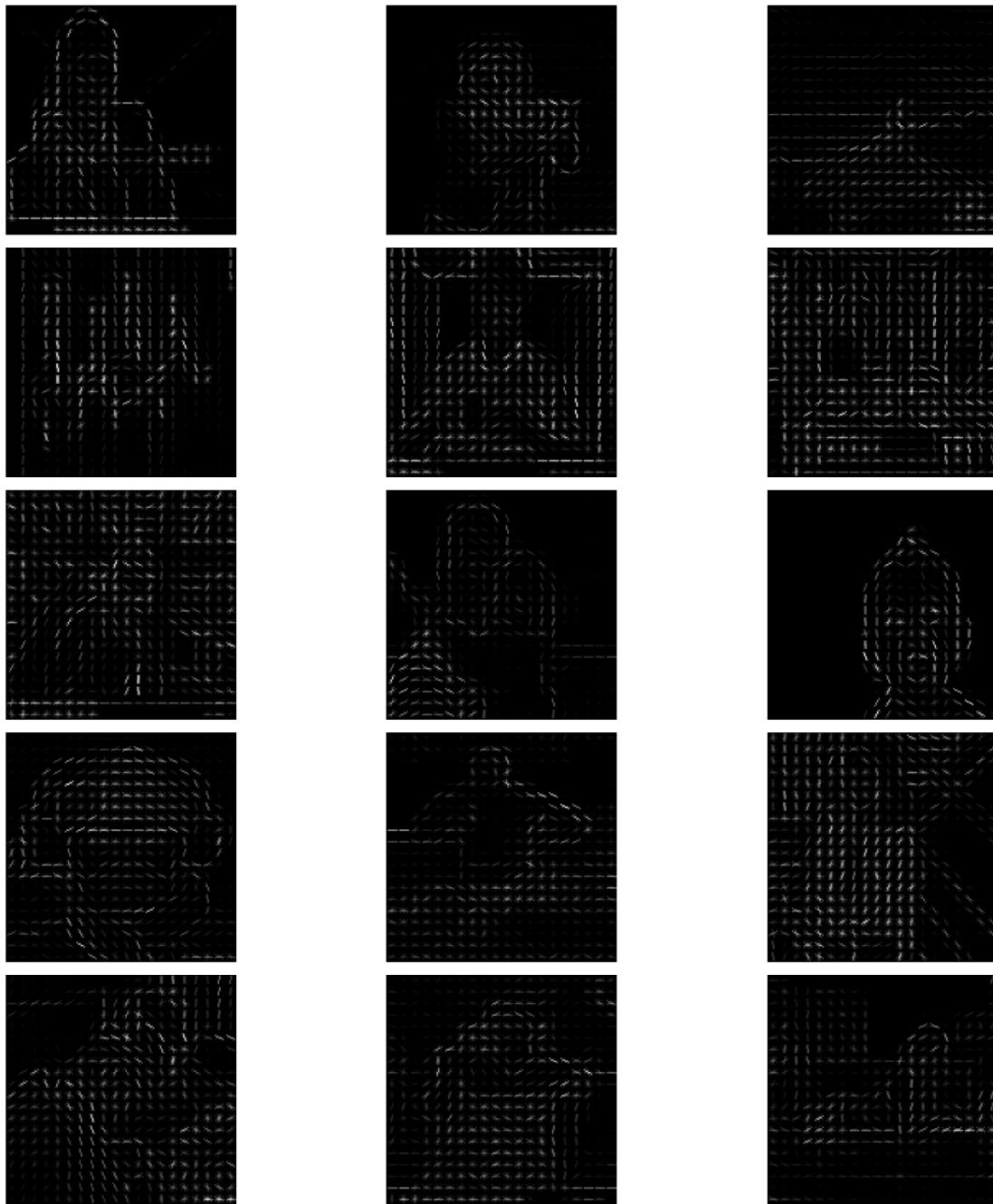
Our findings of an inverse relationship between height and width are verified as we notice from these images that the width decreases as height increases and vice versa.

- c. From our bar plot, we notice that there are no class imbalances in our dataset, and hence, there is no need for any resampling or augmentation techniques.

2. Feature Extraction

We perform Local Binary Pattern, Color Histogram, and SIFT Feature Extraction on our dataset. We also extract the Histogram of Oriented Gradients (HOG) from selected images for visualization purposes, but we do not use the extracted features from HOG to train our dataset.

HOG Image Representations

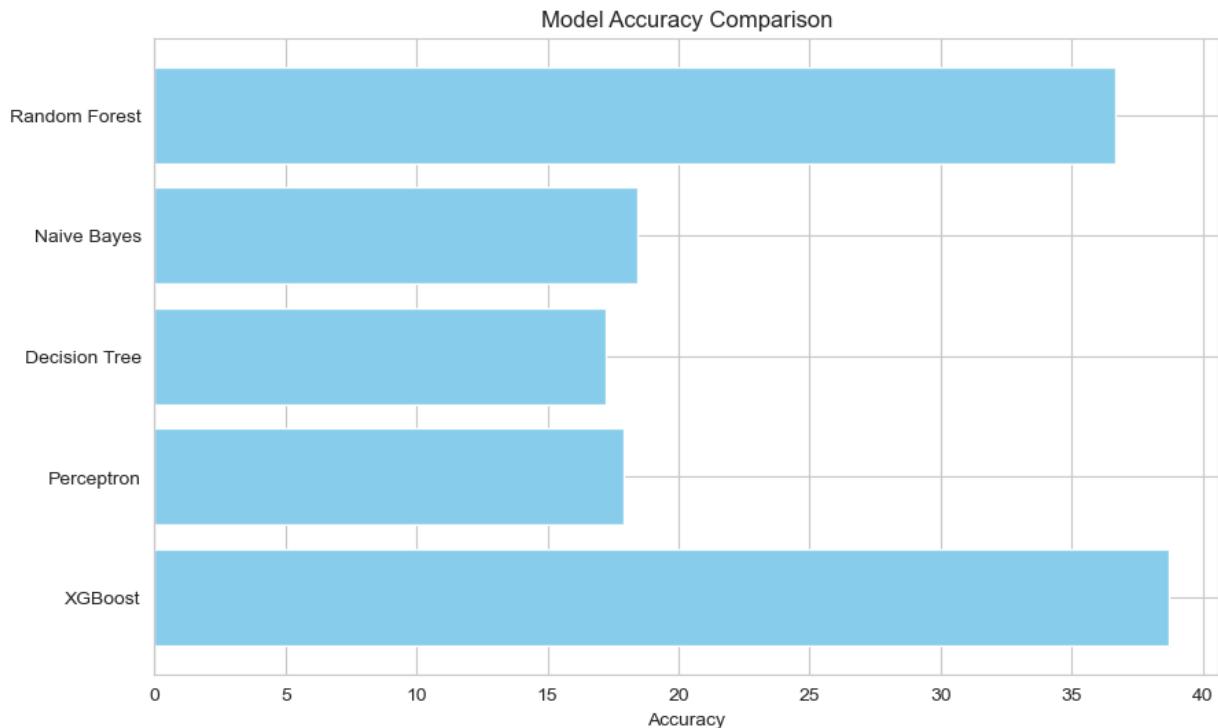


We noticed that the features extracted using LBP, Color Histogram, and SIFT are of constant length. LBP returns a Numpy array of length 26. Color Histogram returns a Numpy array of length 512, and SIFT returns a Numpy array of length 128.

3. Model Training

We train our dataset on Random Forest, Naive Bayes, Decision Tree, Perceptron, and XGBoost. The following accuracies are observed:

- Random Forest: 36.67%
- Naive Bayes: 18.45%
- Decision Tree: 17.22%
- Perceptron: 17.90%
- XGBoost: 38.69%



XGBoost gives the best accuracy, followed by Random Forest.

SECTION B (BONUS)

1. EDA

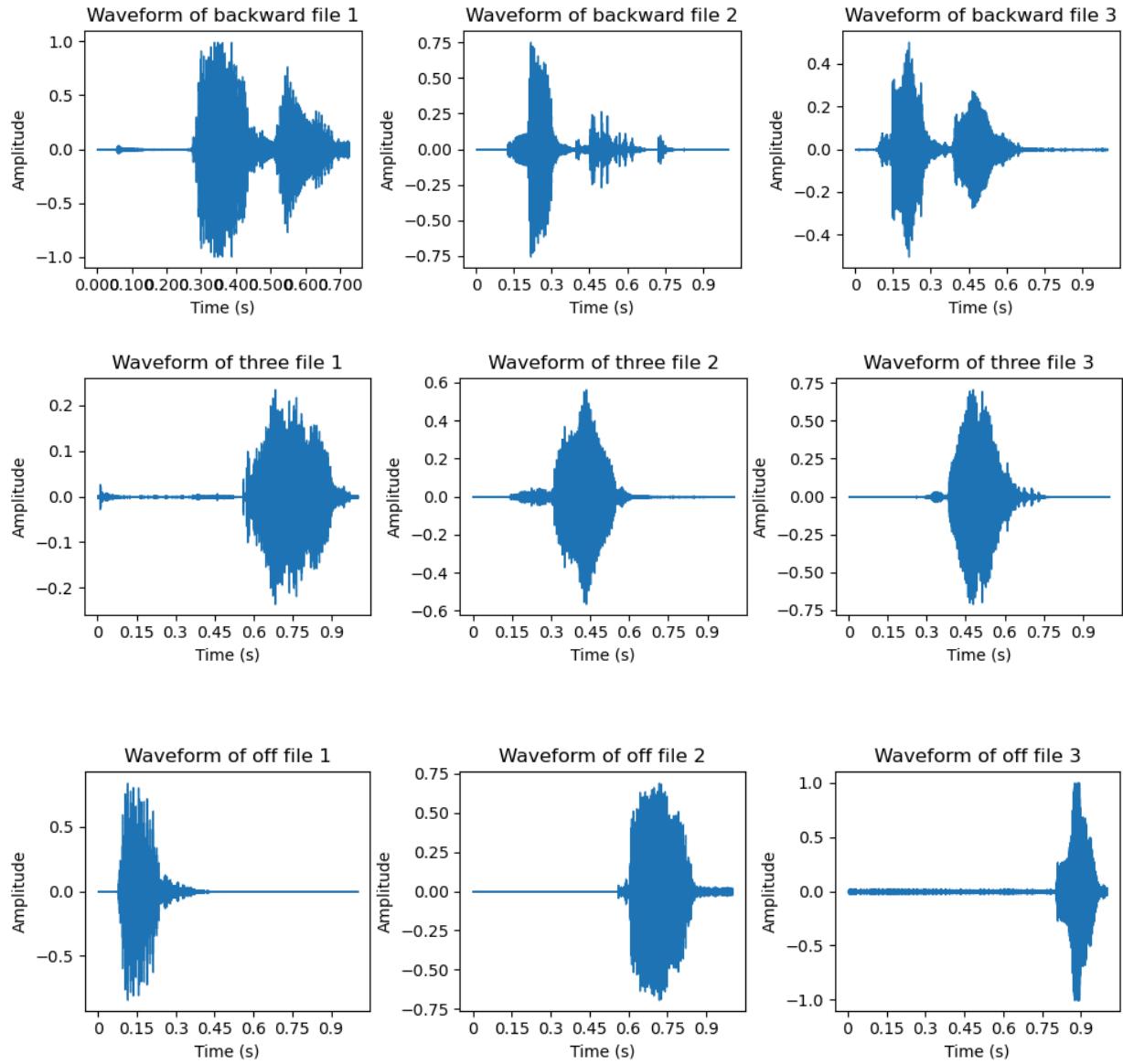
a. Statistical Summary

	mean_amplitude	min_amplitude	max_amplitude	std_amplitude	mean_duration	min_duration	max_duration	std_duration
backward	8.732750e-05	-1.0	0.999969	0.088580	0.986390	0.448000	1.0	0.062522
bed	6.745070e-05	-1.0	0.999969	0.085238	0.970498	0.213312	1.0	0.087145
bird	-1.830991e-04	-1.0	0.999969	0.095966	0.969454	0.325062	1.0	0.092588
cat	-2.848711e-04	-1.0	0.999969	0.075283	0.972065	0.384000	1.0	0.082372
dog	-4.608297e-04	-1.0	0.999969	0.093553	0.972226	0.426625	1.0	0.084865
down	4.458849e-05	-1.0	0.999969	0.087863	0.983559	0.325062	1.0	0.067175
eight	-2.228617e-05	-1.0	0.999969	0.076257	0.980846	0.256000	1.0	0.072684
five	1.695451e-05	-1.0	0.999969	0.092054	0.984320	0.384000	1.0	0.065863
follow	-1.721267e-04	-1.0	0.999969	0.103674	0.981981	0.341313	1.0	0.074247
forward	8.766092e-05	-1.0	0.999969	0.107183	0.984746	0.384000	1.0	0.066149
four	-1.378545e-04	-1.0	0.999969	0.095869	0.982999	0.278625	1.0	0.069495
go	-1.678797e-04	-1.0	0.999969	0.096900	0.978847	0.341313	1.0	0.076190
happy	-1.958642e-04	-1.0	0.999969	0.073446	0.974747	0.384000	1.0	0.079458
house	1.365303e-04	-1.0	0.999969	0.082099	0.974842	0.371500	1.0	0.082296
learn	-1.415776e-05	-1.0	0.999969	0.099211	0.974443	0.371500	1.0	0.091885
left	2.245471e-05	-1.0	0.999969	0.080381	0.984954	0.298625	1.0	0.065198
marvin	-3.048353e-04	-1.0	0.999969	0.095987	0.977728	0.384000	1.0	0.075275
nine	-7.837252e-07	-1.0	0.999969	0.085355	0.984835	0.341313	1.0	0.066458
no	2.396618e-05	-1.0	0.999969	0.088122	0.980128	0.384000	1.0	0.073198
off	3.893072e-05	-1.0	0.999969	0.088514	0.984939	0.394688	1.0	0.060932
on	3.273442e-05	-1.0	0.999969	0.089415	0.981885	0.384000	1.0	0.069206
one	-6.647508e-06	-1.0	0.999969	0.085383	0.978958	0.325062	1.0	0.076265
right	3.956014e-05	-1.0	0.999969	0.083249	0.982247	0.384000	1.0	0.069393
seven	-1.123295e-04	-1.0	0.999969	0.082817	0.984920	0.371563	1.0	0.062049
sheila	-1.687832e-04	-1.0	0.999969	0.083635	0.977648	0.426625	1.0	0.075049
six	-3.445910e-05	-1.0	0.999969	0.069262	0.987542	0.384000	1.0	0.056673
stop	-1.697145e-05	-1.0	0.999969	0.080679	0.984549	0.362687	1.0	0.064572
three	-9.996239e-05	-1.0	0.999969	0.081729	0.983931	0.341313	1.0	0.064122
tree	7.703565e-06	-1.0	0.999969	0.075752	0.970302	0.298625	1.0	0.087430
two	-1.324027e-04	-1.0	0.999969	0.086349	0.981600	0.278625	1.0	0.071753
up	1.073004e-05	-1.0	0.999969	0.074206	0.976156	0.298625	1.0	0.076972
visual	1.372075e-05	-1.0	0.999969	0.091905	0.982301	0.341313	1.0	0.070831
wow	3.483538e-05	-1.0	0.999969	0.096585	0.970225	0.325062	1.0	0.089879
yes	-6.172316e-05	-1.0	0.999969	0.079941	0.983431	0.384000	1.0	0.067037
zero	-7.398344e-05	-1.0	0.999969	0.096609	0.986978	0.298625	1.0	0.057204

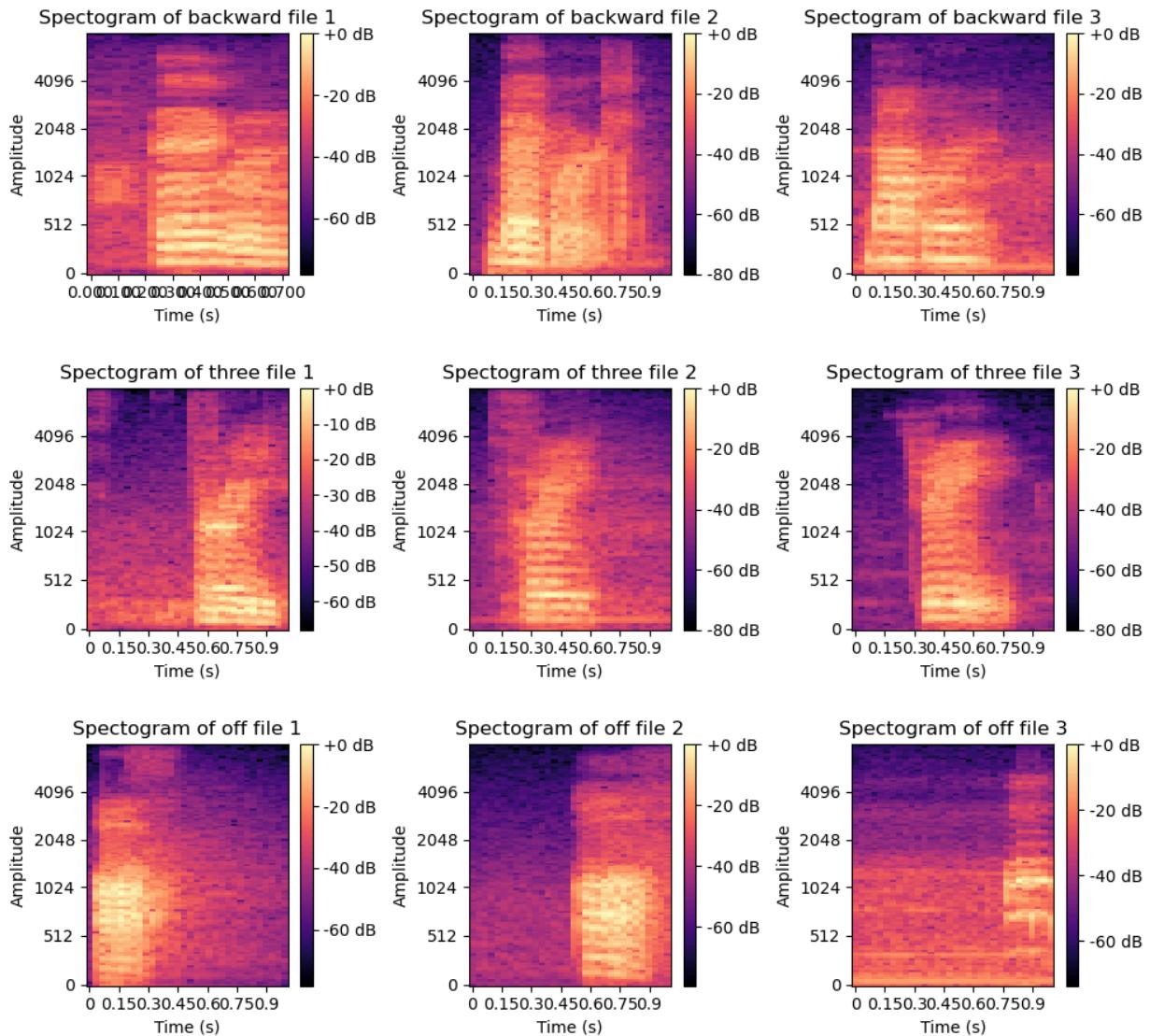
We notice that the minimum and maximum amplitude are the same for all audio classes. Mean and Standard Deviation are very low, indicating that the data seems normalized or already pre-processed.

The maximum duration is 1 s for all classes. The mean durations are consistent.

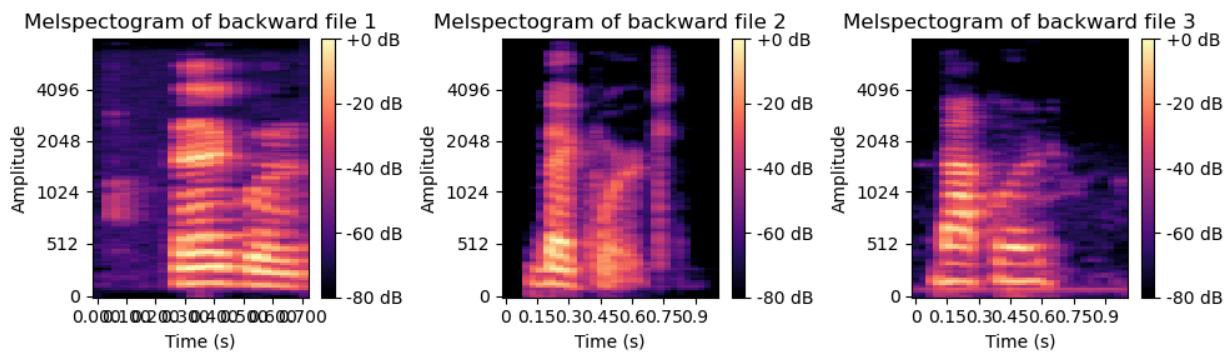
b. Waveforms

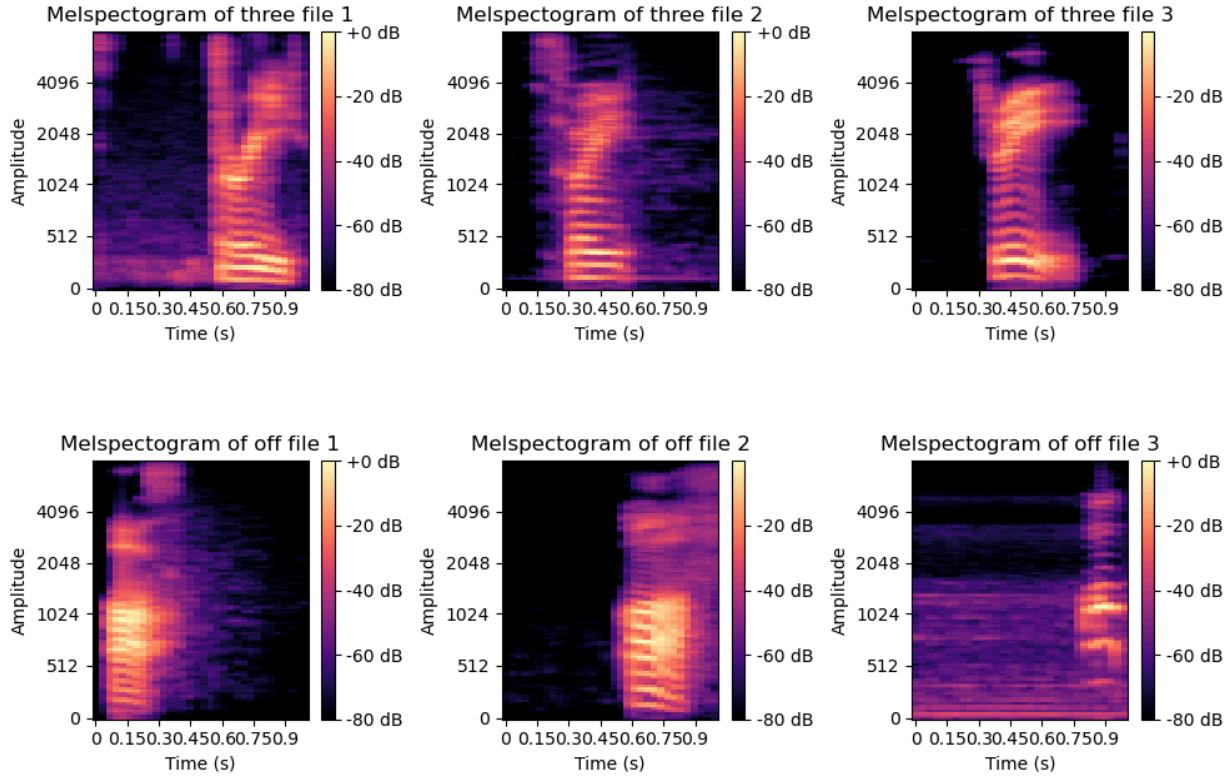


Spectrogram



Melspectrogram



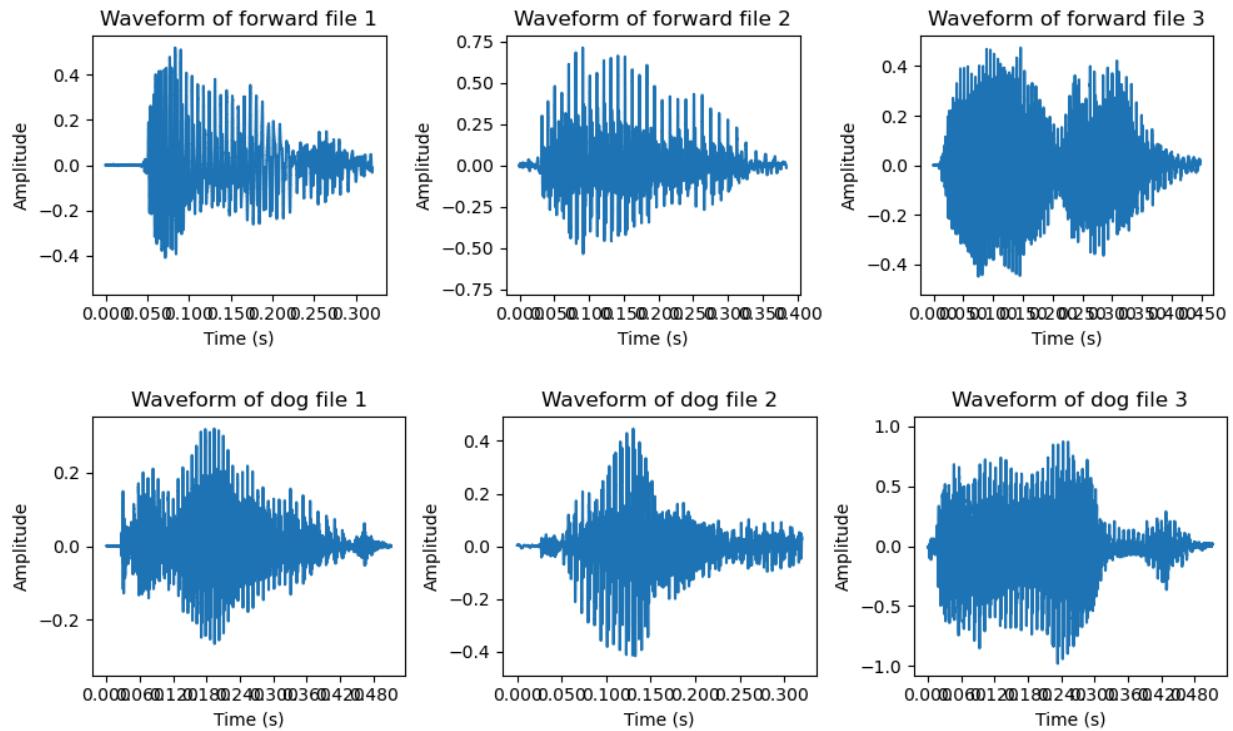


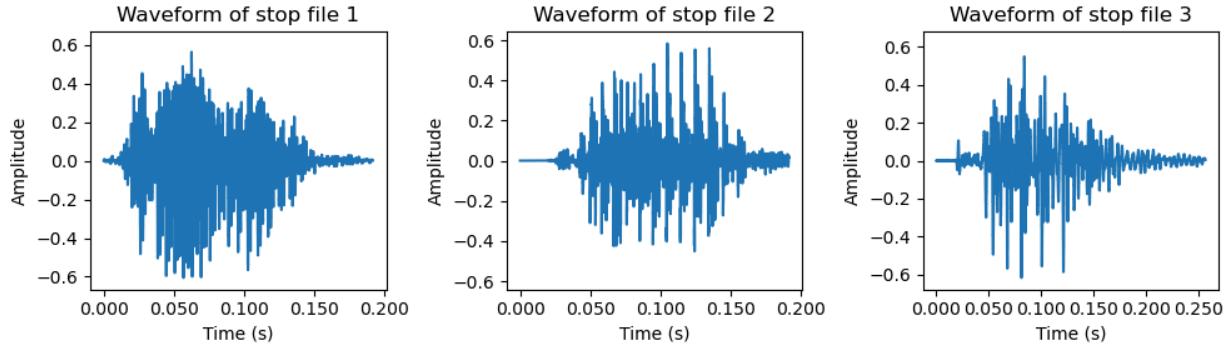
We observe silent segments from the waveforms at the beginning and end of our audio.

- c. Since we notice class imbalance, we select a minimum threshold to filter out from the classes. In this case, we select 1550 as the minimum number and hence, we downsample the dataset to 1550 images per classes. We ignore the _background_noise_ for this purpose.



d. We remove silent segments from our audio and notice the waveforms again.





2. Feature Extraction

We extract MFCC, Chroma, Spectral Centroid, Spectral Bandwidth, Zero Crossing Rate, Spectral Rolloff, and RMS value.

We notice that MFCC Mean returns an array of length 13, and Chroma Mean returns an array of length 12.

This brings our dataset to have 31 columns.

We then apply label encoding to our labels.

3. We split our dataset into 80:20 ratios and apply Min Max Scaling to the training and testing dataset. We get the number of records as 43400 in the training dataset and 10850 in the testing dataset.