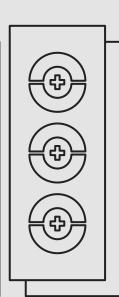
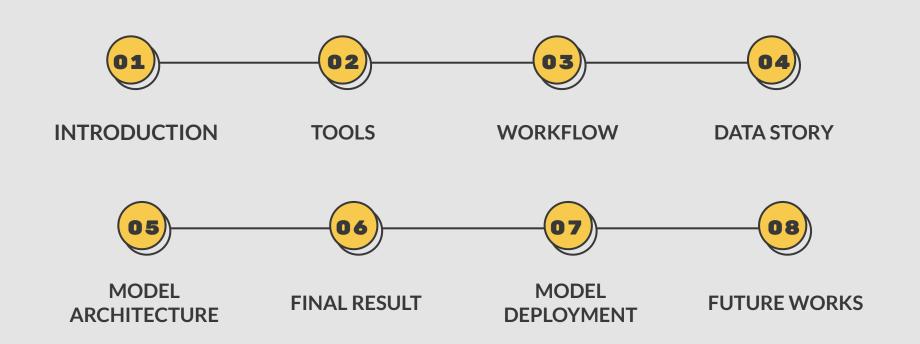
# CLASSIFYING NOISE SOUNDS

Presented by: Alanoud Alosaimi, Raghad Althunayan, Shaikha bin Ateeq



## TABLE OF CONTENTS



#### INTRODUCTION

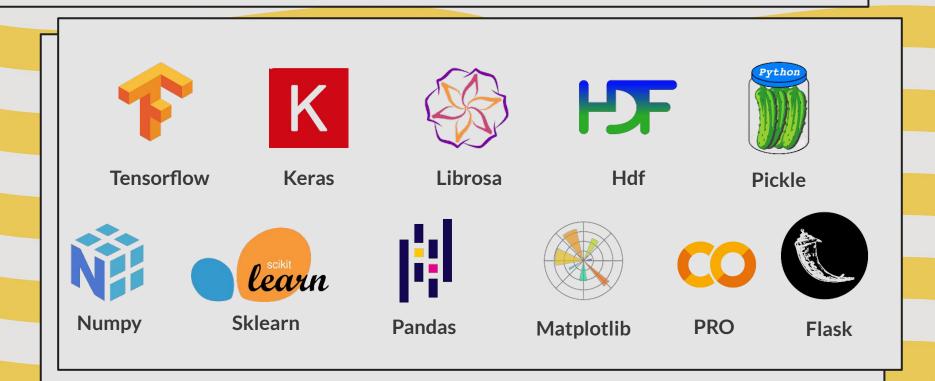
**Sounds are all around us.** Whether directly or indirectly. Sounds outline the context of our daily activities, conversations, music, noise. The human brain is continuously processing and understanding this audio data, so how the machine can understand it?

#### GOAL:

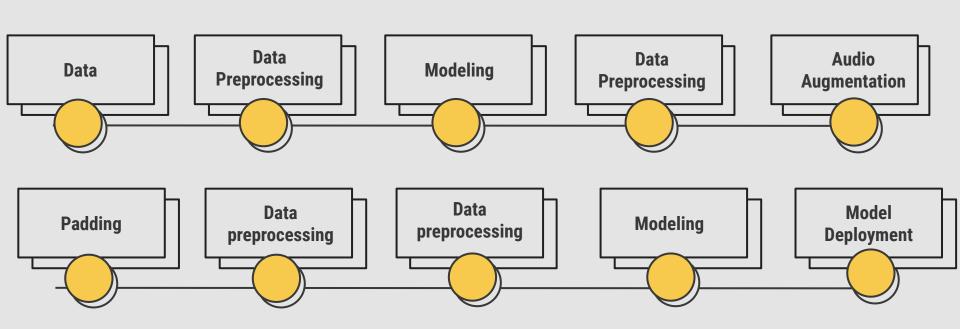
Apply Deep Learning techniques to the classification of environmental sounds:

- Assisting deaf individuals in their daily activities
- Safety and security capabilities
- Smart home use

## **TOOLS**



## WORKFLOW



## DATA STORY

#### **Urban sounds classification**

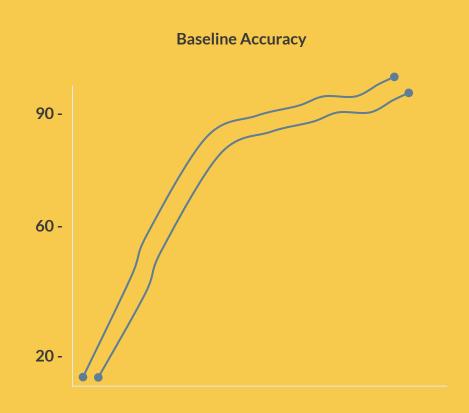


From kaggel

Row = 8733

#### Class label [10]=

- Air Conditioner
- Car Horn
- Children Playing
- Dog bark
- Drilling
- Engine Idling
- Gun Shot
- Jackhammer
- Siren
- Street Music



## DATA STORY

#### **Urban sounds classification**





Row = 8733

#### Class label [10]=

- Air Conditioner
- Car Horn
- Children Playing
- Dog bark
- Drilling
- Engine Idling
- Gun Shot
- Jackhammer
- Siren
- Street Music

Noise sounds classification



From zenodo

Row = 2171

Class label [11]=

- Applause
- Keys\_jangling
- Telephone
- Cough
- Microwave oven
- Laughter
- Tearing
- Fireworks
- Bus
- Scissors
- Computer\_keyboard'

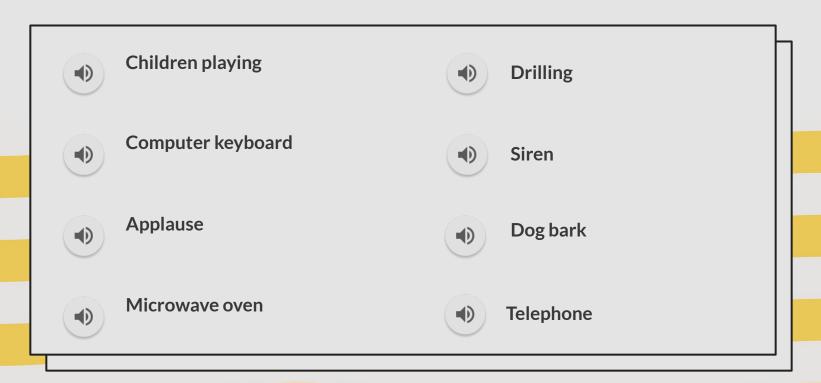
Before = (8733, 10)

After = (10904, 21)

One dataset not enough COMPLEXITY NEEDED!

## DATA STORY

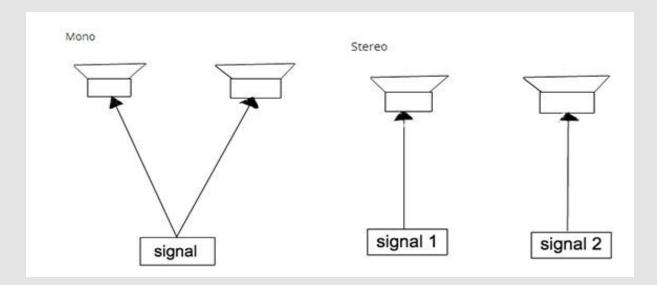
## **Audio Samples**



We heard the audio samples, Do we know now its **PROPERTIES**?

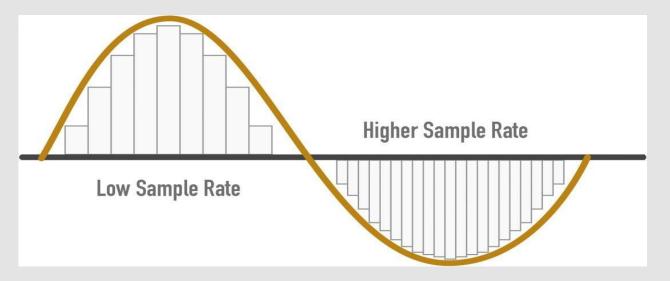
## **AUDIO PROPERTIES:**

Audio Channels



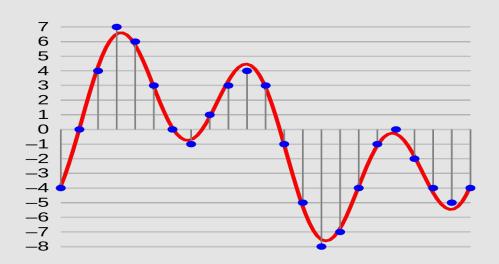
#### **AUDIO PROPERTIES:**

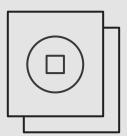
• Sample Rate



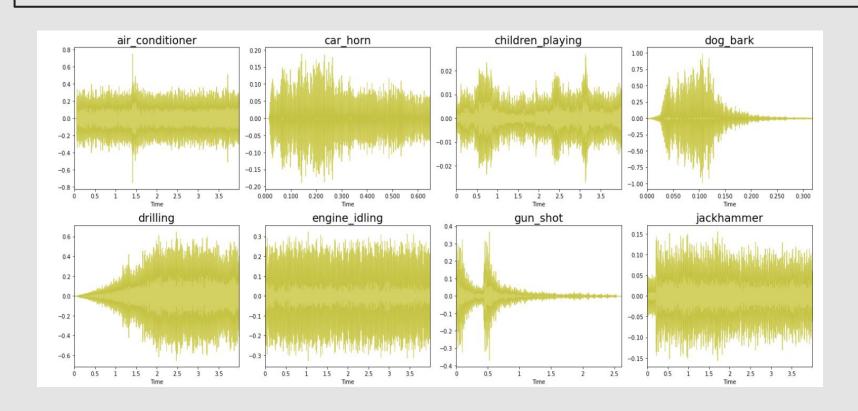
## **AUDIO PROPERTIES:**

• Bit-Depth

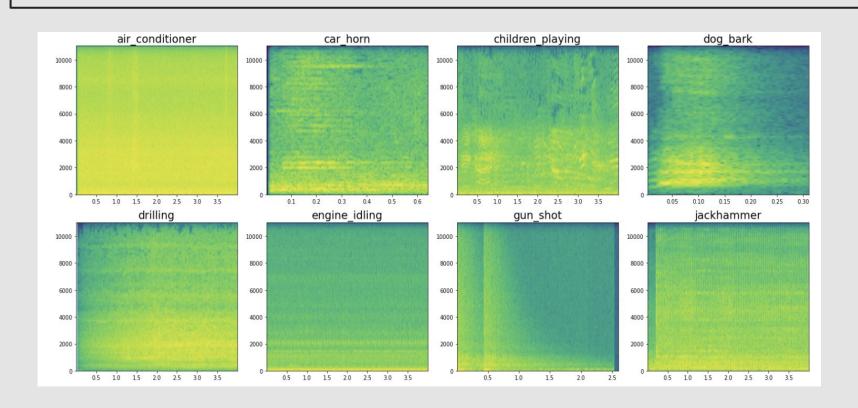




## Waveplot



## **Specgram Plot**



# BUT!!



Can we use the spectrum images as input for our model? or something else?

#### FEATURE EXTRACTION METHOD: MFCC

The MFCC summarises the frequency distribution across the window size, to analyse both the frequency and time characteristics of the sound. These audio representations will allow us to identify features for classification.

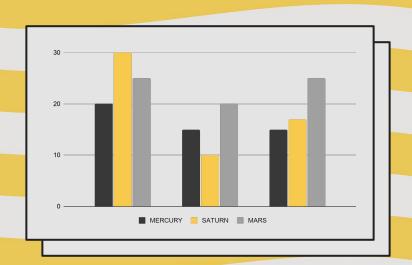


-0.2963816 -0.3560971 -0.27297518 -0.39299247, -0.33859769 -0.4041523, -0.53813744, -0.52863103, -0.23635665, -0.20224652, -0.360925, -0.4068555,	-0.309869 -0.66855148 -0.7531548 -0.7531548 -0.7531548 -0.7531548 -0.9061938 -0.475869574 -0.2672934 -0.474798554 -0.2672934 -0.44798554 -0.2672934 -0.2672934 -0.2672934 -0.2672934 -0.2719564 -0.2672934 -0.2719564 -0.2672934 -0.271956569 -0.271956569 -0.271956569 -0.30666392 -0.306	-0. 29231593 -0. 29231593 -0. 77658434 -0. 778185385 -0. 778185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 78185385 -0. 781853	-0.2699363 -0.22391006, -0.27417266, -0.210177219, -0.17777492, -0.3978389, -0.25736403, -0.3417926, -0.40582865, -0.40582865, -0.40582865, -0.50796515, -0.30450806, -0.19710773, -0.26224333, -0.52654827, -0.48837436, -0.36755478,	-0.36449316, -0.38918048, -0.4472869, -0.545878,
-0.46702236, -0.35711053,	-0.34863254, -0.2876631,	-0.27881584, -0.19479881,	-0.27385667, 0.	-0.3146194 , , dtype=float32)

## All we think data now is PERFECT!

But, the models **didn't predict** some classes (**imbalanced data**)

#### **Classes Imbalanced**



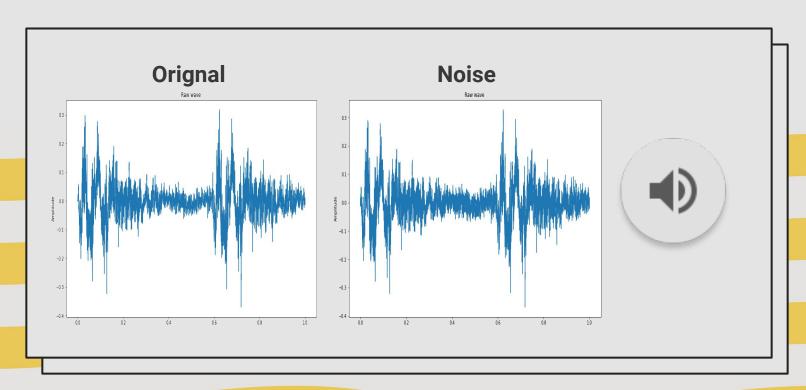
## **Solving Classes Imbalanced**

#### **Data Augmentation**

The objective is to make our model invariant to those perturbations and enhance its ability to generalize.

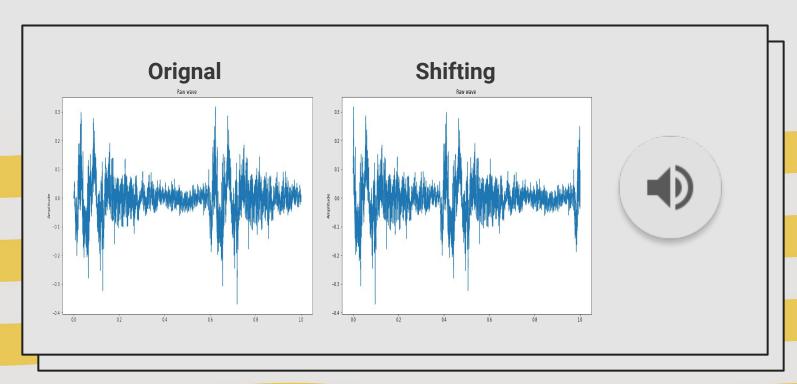
## **Data Augmentation:**

Noise



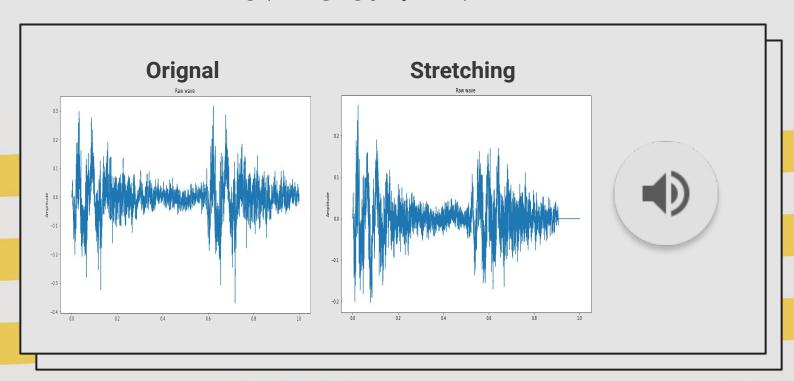
## **Data Augmentation:**

Shifting the Audio



## **Data Augmentation:**

• Time stretching (changing play time)



## **Data Augmentation**

INITIAL DATA NEW DATA MERGED DATA 10903 54515 65418 (10903, 21)(65418, 21)



## **PADDING**

#### **Length sounds problem**



▶ 0:22 / 0:22 **------**

**Least** length (fram\_num = 174)

**Longest** length (fram\_num = 1292)

Least, why?

## SPLIT SIZE

Train := 0.8 Test: 0.2

Train := 0.75 Validations: 0.25



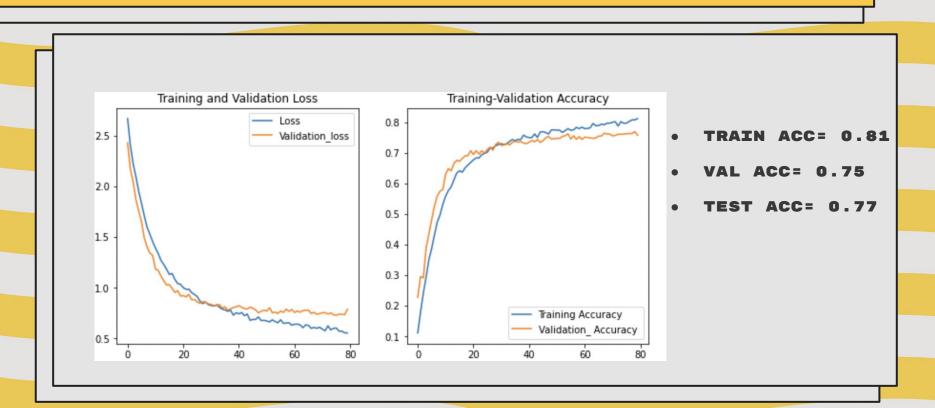
## MODELS RESULT (1)

	TRAIN ACC	VAL ACC	EPOCHS	ватсн
BEASLINE	0.61	0.45	500	50
CNN2 (1)	0.28	0.25	200	50
CNN2 (2)	0.60	0.44	1000	50
CNN2 (3)	0.33	0.21	300	300
CNN2 (4)	0.48	0.44	500	50

## MODELS RESULT (2)

	TRAIN ACC	VAL ACC	EPOCHS	ватсн
BEASLINE	0.69	0.63	500	50
CNN2D (1)	0.80	0.77	300	300
CNN2D (2)	0.89	0.78	300	300
CNN2D (3)	0.81	0.75	300	300
CNN1D	0.27	0.3	500	50
LSTM	0.3	0.34	500	50

## FINAL RESULT



#### **MODEL ARCHITECTURE**

Conv2D (Filter size = 128)

MaxPooling2D

Conv2D (Filter size = 128)

MaxPooling 2D

Droupout (0.8)

Flatten

Dense (512, activation = 'relu') Droupout (0.8) Dense (512, activation = 'relu')

Droupout (0.8)

Dense (18, activation = softmax)

## **CHALLENGES**





**JUPYTER** 

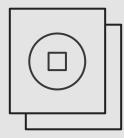
Ram crash

RUN

Long time

## **FUTURE WORK**

- Build an app that helps the deaf in their daily life
- Transfer Learning Model



# MODEL DEPLOYMENT DEMO

## **APPENDIX**

## **DATA SOURCE:**

- 1. <a href="https://urbansounddataset.weebly.com/urbansound8k.html">https://urbansounddataset.weebly.com/urbansound8k.html</a>
- 2. <a href="https://zenodo.org/record/2552860#.Yek-aVhBy3K">https://zenodo.org/record/2552860#.Yek-aVhBy3K</a>

## **THANKS**