# Audio analysis: speech recognition

Carlota Sans and Mariola Monfort
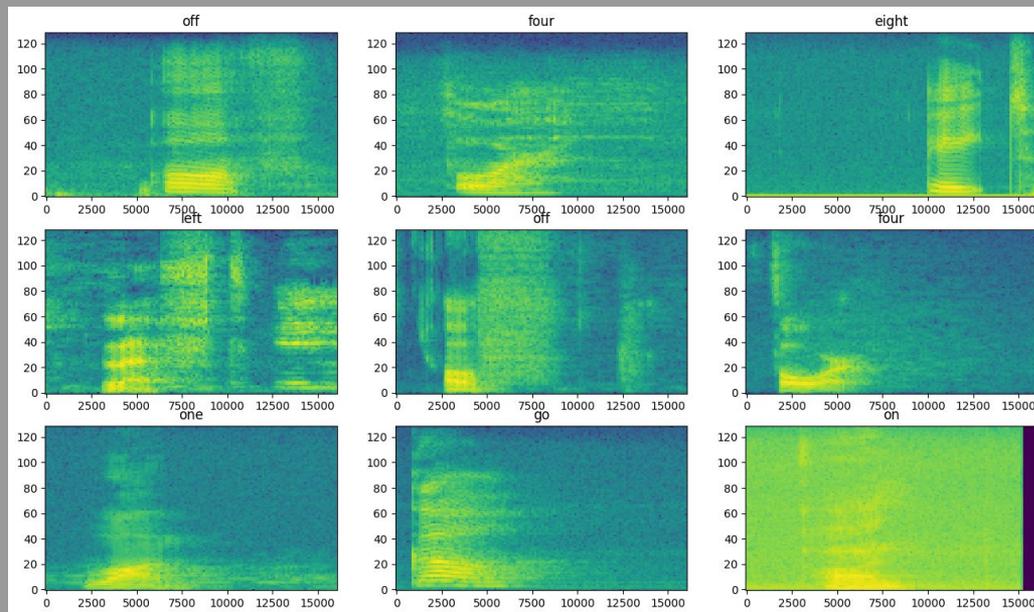
Técnicas Avançadas de Análise de Dados
Supervisor: Filipe Veloso

1 2 9 0

UNIVERSIDADE Ð
COIMBRA

# INDEX

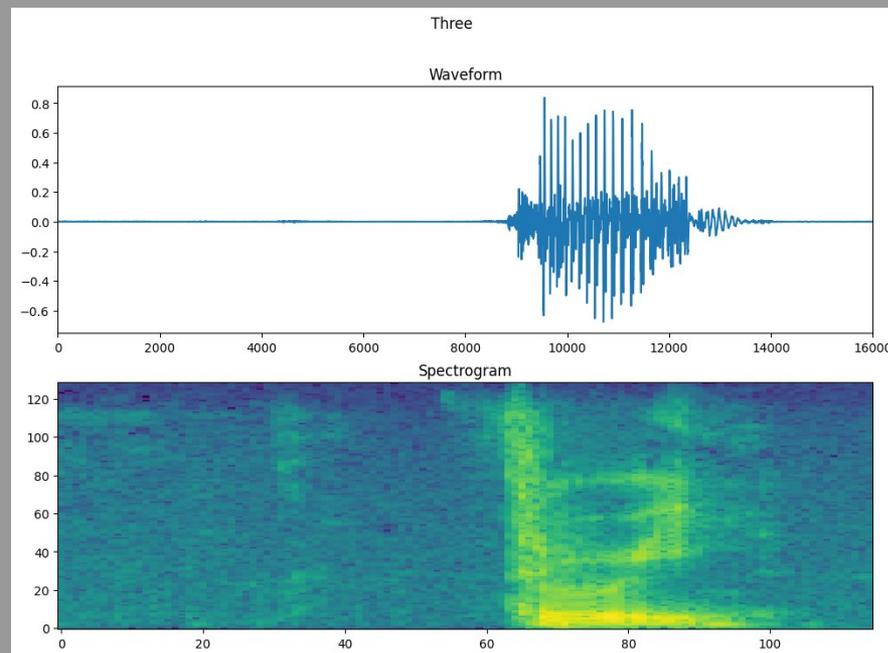- The data
- Our code
- Optimization
- Conclusions

# The Data

The goal of speech recognition is to transform human voice into text or instructions that a computer is able to process. It has many applications and it is very useful nowadays.

## How do we analyze our data???

- Original data: .wav files
- Input data: spectrograms
- STFT

# Our Code

**How does the code work? How does it analyse the data?**
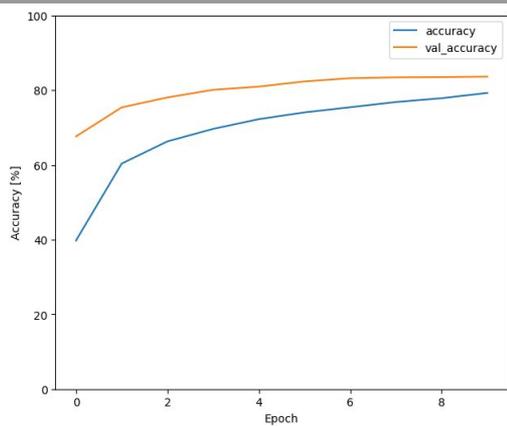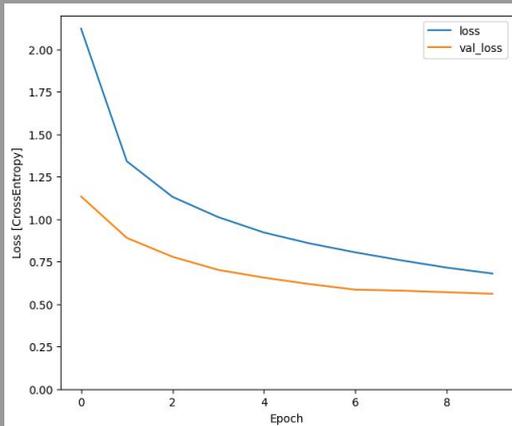


**What were the initial parameters?**

# Our Code

- Data split: 80% train, 10% validation, 10% test
- Transformation of the data
- Neural network: CNN
- Training of the model

```
Input shape: (124, 129, 1)
Model: "sequential"

Layer (type)                Output Shape            Param #
=================================================================
resizing (Resizing)         (None, 32, 32, 1)       0

normalization (Normalizati  (None, 32, 32, 1)       3
on)

conv2d (Conv2D)             (None, 30, 30, 32)      320

conv2d_1 (Conv2D)           (None, 28, 28, 64)      18496

max_pooling2d (MaxPooling2  (None, 14, 14, 64)      0
D)

dropout (Dropout)           (None, 14, 14, 64)      0

flatten (Flatten)           (None, 12544)           0

dense (Dense)               (None, 128)             1605760

dropout_1 (Dropout)         (None, 128)             0

dense_1 (Dense)             (None, 36)              4644

=================================================================
Total params: 1629223 (6.21 MB)
Trainable params: 1629220 (6.21 MB)
Non-trainable params: 3 (16.00 Byte)
_____
```

```
Epoch 1/13
1323/1323 [==============================] - 345s 260ms/step - loss: 2.1112 - accuracy: 0.4034 - val_loss: 1.1466 - val_accuracy: 0.6874
Epoch 2/13
1323/1323 [==============================] - 309s 234ms/step - loss: 1.3272 - accuracy: 0.6084 - val_loss: 0.9220 - val_accuracy: 0.7494
Epoch 3/13
1323/1323 [==============================] - 312s 236ms/step - loss: 1.1374 - accuracy: 0.6637 - val_loss: 0.7768 - val_accuracy: 0.7817
Epoch 4/13
1323/1323 [==============================] - 310s 234ms/step - loss: 1.0212 - accuracy: 0.6935 - val_loss: 0.7079 - val_accuracy: 0.7950
Epoch 5/13
1323/1323 [==============================] - 310s 234ms/step - loss: 0.9446 - accuracy: 0.7173 - val_loss: 0.6902 - val_accuracy: 0.8004
Epoch 6/13
1323/1323 [==============================] - 319s 241ms/step - loss: 0.8882 - accuracy: 0.7312 - val_loss: 0.6240 - val_accuracy: 0.8189
Epoch 7/13
1323/1323 [==============================] - 307s 232ms/step - loss: 0.8464 - accuracy: 0.7434 - val_loss: 0.6078 - val_accuracy: 0.8241
Epoch 8/13
1323/1323 [==============================] - 319s 241ms/step - loss: 0.8074 - accuracy: 0.7541 - val_loss: 0.5998 - val_accuracy: 0.8309
Epoch 9/13
1323/1323 [==============================] - 308s 233ms/step - loss: 0.7734 - accuracy: 0.7620 - val_loss: 0.5809 - val_accuracy: 0.8320
Epoch 10/13
1323/1323 [==============================] - 309s 234ms/step - loss: 0.7506 - accuracy: 0.7706 - val_loss: 0.5597 - val_accuracy: 0.8399
Epoch 11/13
1323/1323 [==============================] - 308s 233ms/step - loss: 0.7227 - accuracy: 0.7778 - val_loss: 0.5438 - val_accuracy: 0.8456
Epoch 12/13
1323/1323 [==============================] - 321s 242ms/step - loss: 0.7044 - accuracy: 0.7846 - val_loss: 0.5555 - val_accuracy: 0.8408
Epoch 13/13
1323/1323 [==============================] - 318s 241ms/step - loss: 0.6865 - accuracy: 0.7883 - val_loss: 0.5374 - val_accuracy: 0.8467
```

```
[25] model.evaluate(test_spectrogram_ds, return_dict=True)

166/166 [==============================] - 26s 156ms/step - loss: 0.5456 - accuracy: 0.8423
{'loss': 0.5455626845359802, 'accuracy': 0.8422739505767822}
```
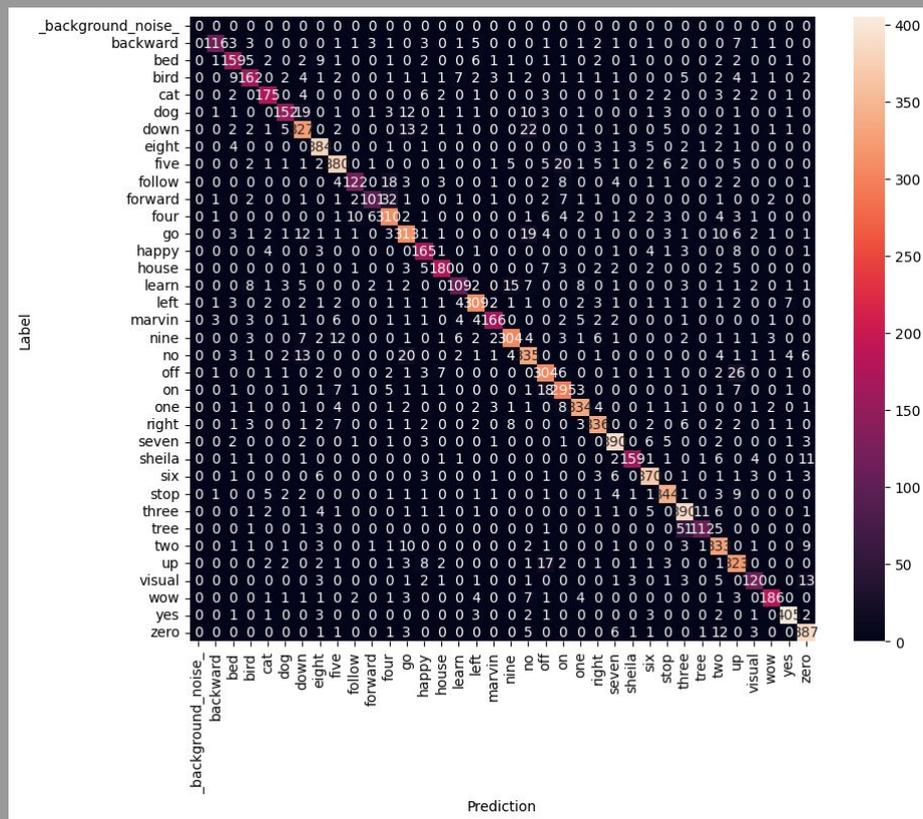
batch size = 64

frame length= 225

frame step = 128

dropout rate 1 = 0.25

dropout rate 2 = 0.5

# Our Code



batch size = 64

frame length = 225

frame step = 128
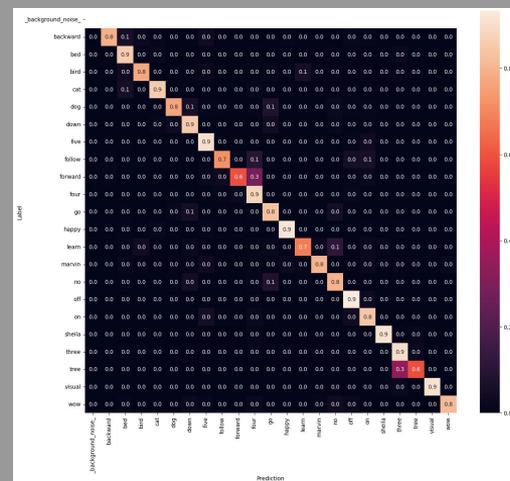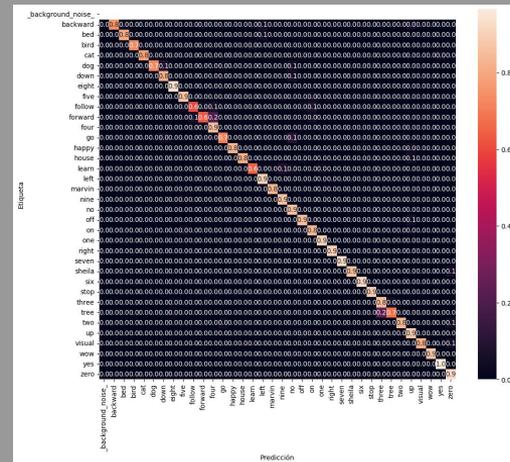
dropout rate 1 = 0.25

dropout rate 2 = 0.5

# Optimization

1. Set the number of epochs to 13
2. Understand the confusion matrix
3. Reduce the data
4. Use Optuna to find the best parameters
   - dropout rate = 0.4315 and 0,4661
   - batch size = 64
   - frame step = 128
5. Apply them to the whole data set



```
[25] model.evaluate(test_spectrogram_ds, return_dict=True)

    166/166 [==============================] - 37s 223ms/step - loss: 0.4891 - accuracy: 0.8613
    {'loss': 0.4890890121459961, 'accuracy': 0.8613179922103882}
```

# Conclusions

- Improvement of the loss (0.06) and accuracy (2%).
- Understanding of how machine learning works and how the parameters within our code affect the final result.