# Nkululeko
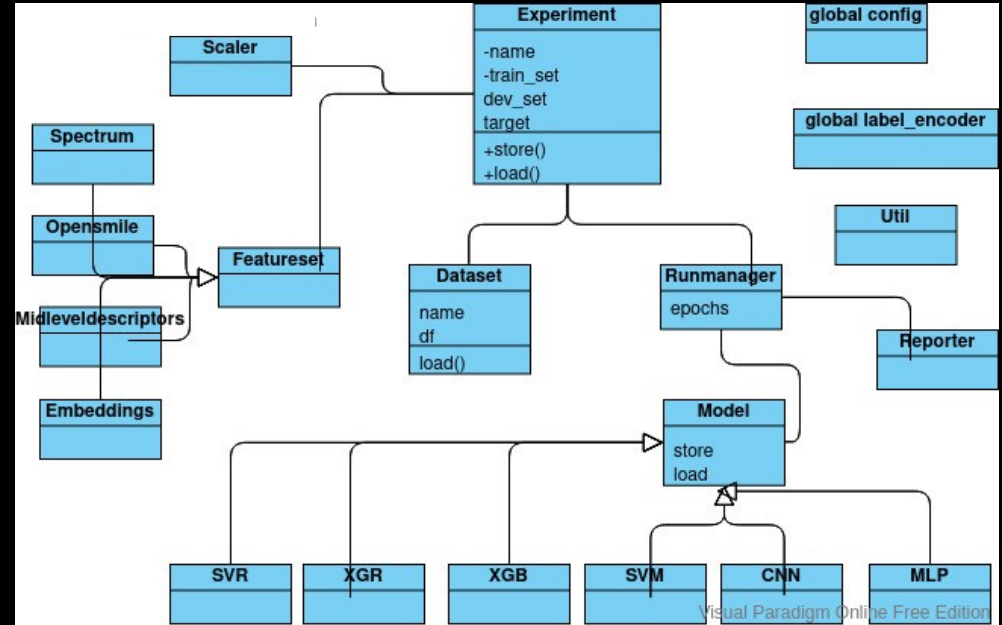## a form based speech machine learning tool

Felix Burkhardt

# outline

- what is Nkululeko

- how to use it

- example experiments

# what is Nkululeko?



- a software written in Python hosted on github*

- a tool to do machine learning (ML) experiments WITHOUT the need to program yourself

- focused on combinations of features and machine learners

- uses configuration file templates described in a blog**

* https://github.com/felixbur/nkululeko

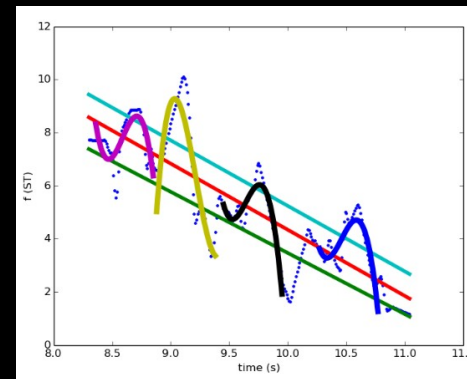** blog.syntheticspeech.de/?s=nkululeko

# motivation

- with the success of Deep Learning, machine learning dominates science

- empiricists sometimes struggle with programming

- teaching students

- re-use of code
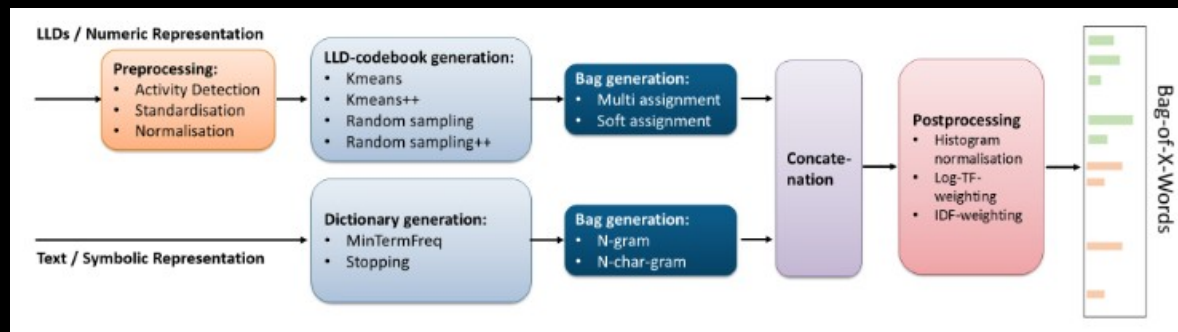
# features

Three kinds:

- expert features

- brute force (needs feature selection)

- learned

  - autoencoders

  - embeddings, latent space

# features

- opensmile [1]
  - (e)GeMAPS – 62/88
  - Compare16 - 6,373
- mid level descriptors [2]
- openXbow [3]
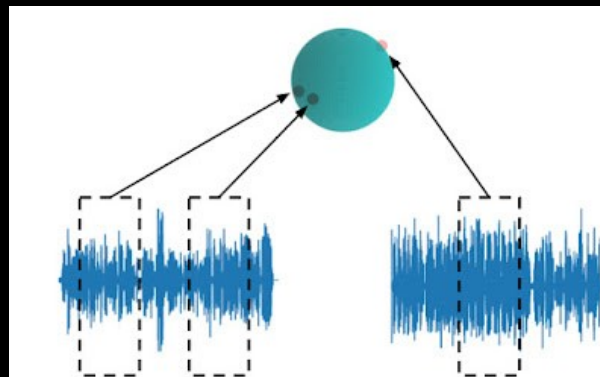


img src: [2]



img src: [3]

[1] Eyben, F., Wöllmer, M. and Schuller, B. (2010). opensmile – the munich versatile and fast open-source audio feature extractor.
[2] Reichel, U., Triantafyllopoulos, A., Oates, C., Huber,S., and Schuller, B. (2020). Spoken language iden-tification by means of acoustic mid-level descriptors
[3] Schmitt, M. and Schuller, B. (2017). openxbow - introducing the passau open-source crossmodal bag-of-words toolkit.
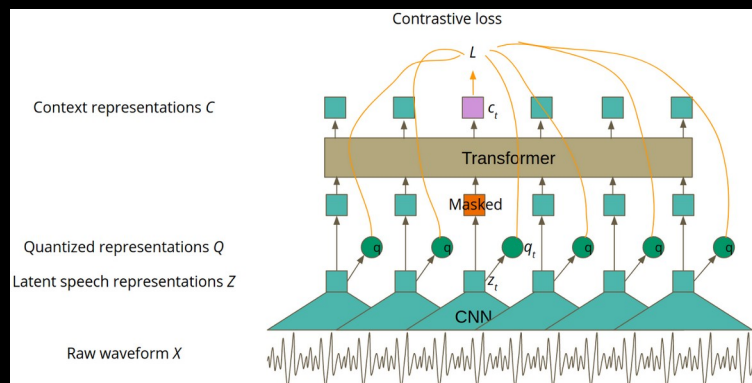
# features cont.

- Logmel spectrograms

- TRILL [1]

- Wav2vec 2.0 [2]

TRILL

img src: https://ai.googleblog.com/2020/06/improving-speech-representations-and.html

Wav2Vec 2.0

img src: https://towardsdatascience.com/wav2vec-2-0-a-framework-for-self-supervised-learning-of-speech-representations-7d3728688cae
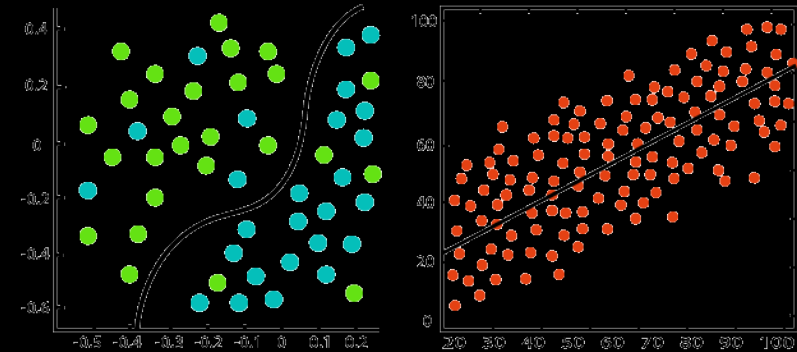
[1] Shor, J., Jansen, A., Maor, R., Lang, O., Quitry, F.,Tagliasacchi, M., Tuval, O., Shavitt, I., Emanuel, D.,and Haviv, Y. (2020).
        Towards learning a universalnon-semantic representation of speech
[2] Baevski, A., Zhou, Y., Mohamed, A., and Auli,M. (2020). wav2vec 2.0: A framework for self-supervised learning of speech representation

# machine learners

## Two distinctions

– classifiers versus regressors
  - classifier: prob. of a specific category
  - regressor: predict a scalar

– approach
  - geometric
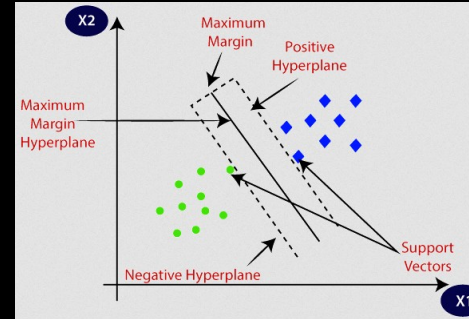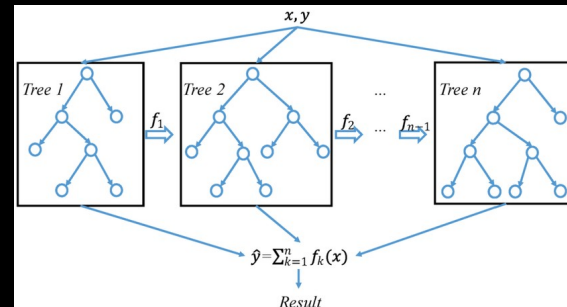  - decision trees
  - ANNs
  - ...



Img. Src: https://www.javatpoint.com/regression-vs-classification-in-machine-learning

# learners

- SVM: support vector machine

- SVR

- XGB: XG-boost

- XGR



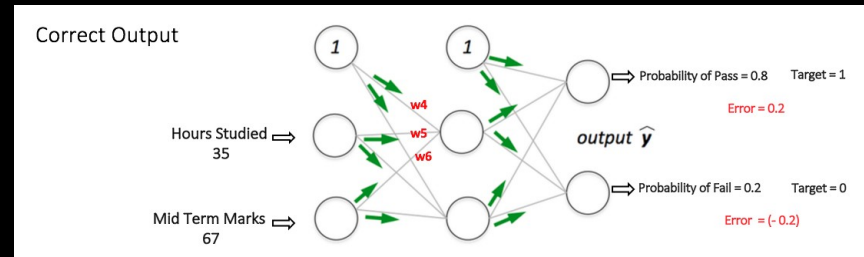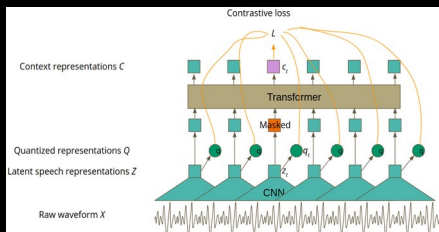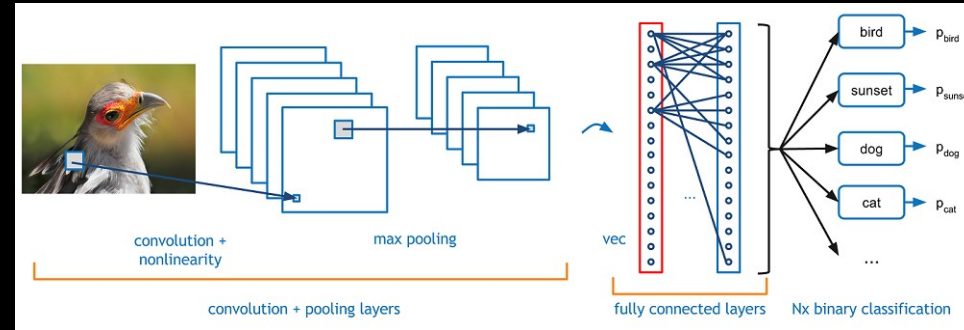img src: https://medium.com/@skilltohire/support-vector-machines-4d28a427ebd



img src: Wang, Yuanchao & Pan, Z. & Zheng, J. & Qian, L. & Mingtao, Li. (2019).
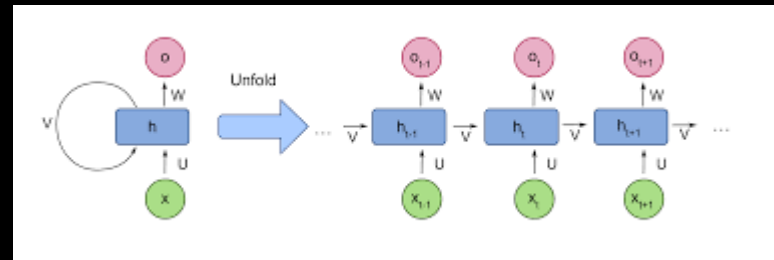A hybrid ensemble method for pulsar candidate classification.

# learners cont.

- MLP: multi layer perceptron

- CNN: convolutional neural net

- RNN: recurrent neural net

- Transformers



img src: https://ujjwalkarn.me/2016/08/09/quick-intro-neural-networks/



img src: https://medium.datadriveninvestor.com/convolution-neural-network-22565e6d8156



img src: Wikimedia

Felix Burkhardt – Nkululeko - ESSV 2022

# the configuration file

- Key-value pairs
- Organized in sections
  - EXP
  - DATA
  - FEAT
  - MODEL
  - PLOT

```
[EXP]
root = ./tests/
name = exp_syntact
runs = 1
epochs = 1
save = True
[DATA]
databases = ['syntact']
syntact = /home/felix/data/research/syntAct/syntact/
syntact.split_strategy = speaker_split
syntact.testsplit = 50
syntact.value_counts = True
target = emotion
labels = ['angry', 'happy', 'neutral', 'sad']
[FEATS]
#type = trill
type = os
scale = standard
[MODEL]
type = svm
save = True
[PLOT]
value_counts = True
tsne = True
```

# conf. file EXP section

- name

- type [classification | regression]

- #epochs

- #runs

https://github.com/felixbur/nkululeko/blob/main/ini_file.md
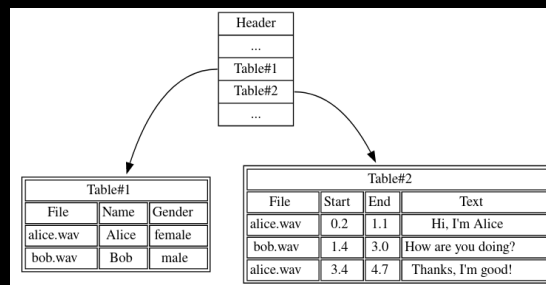
# conf. file DATA section

- databases
  - type: [audformat* | CSV]
  - table specifics
  - train/test splits
- type [cross corpus | train-test]
- trains, tests
- label mapping
- binning (scalar → classes)
- sex: data filter



```
x/sample.wav, s1, female, happy
...
```
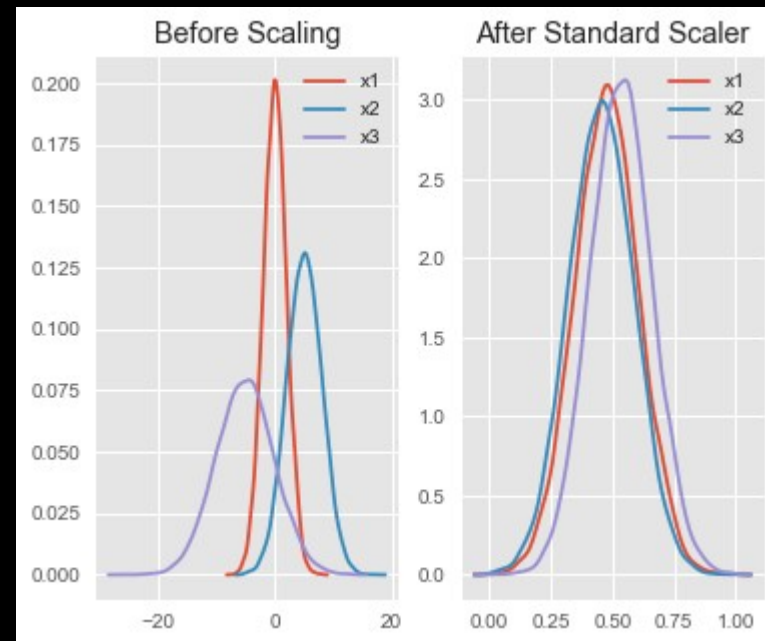or with age:
```
x/sample.wav, roger, male, 45
```



*https://audeering.github.io/audformat/

https://github.com/felixbur/nkululeko/blob/main/ini_file.md

# conf. file FEATS section

- type: [os | spectra | mld | xbow | trill | wav2vec]

- scale: [std | spkr | sex]

- model: path to model



Img src: https://shauryauppal.medium.com/how-and-where-to-apply-feature-scaling-machine-learning-93316663cd63

https://github.com/felixbur/nkululeko/blob/main/ini_file.md
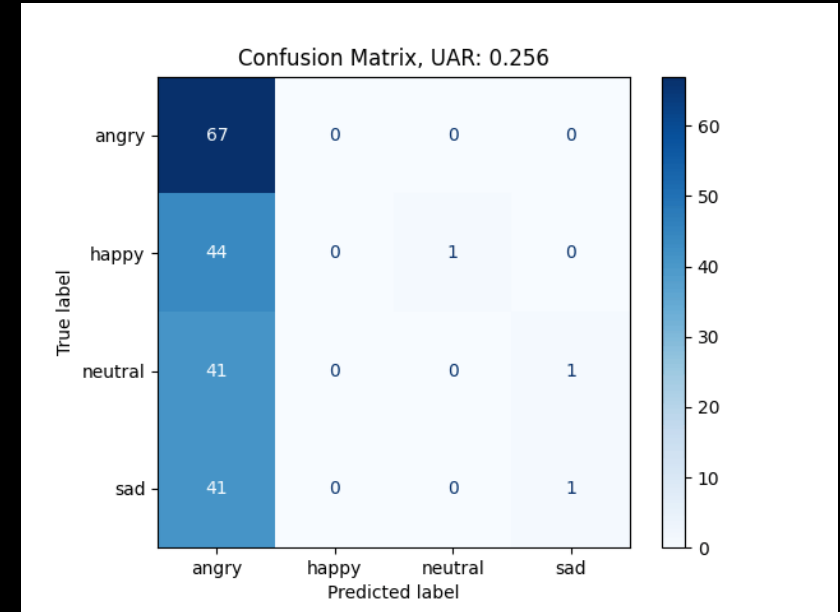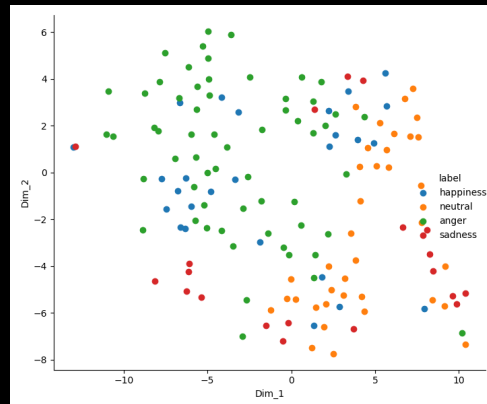
# conf. file MODEL  section

- type: [svm | svr | xgb | xgr
  | mlp | mlp-reg | cnn]

- tuning_params: 5 fold
  cross optimization

- layers, loss-function,
  learning rate: ANN specs

- class_weight

https://github.com/felixbur/nkululeko/blob/main/ini_file.md

# conf. file PLOT   section

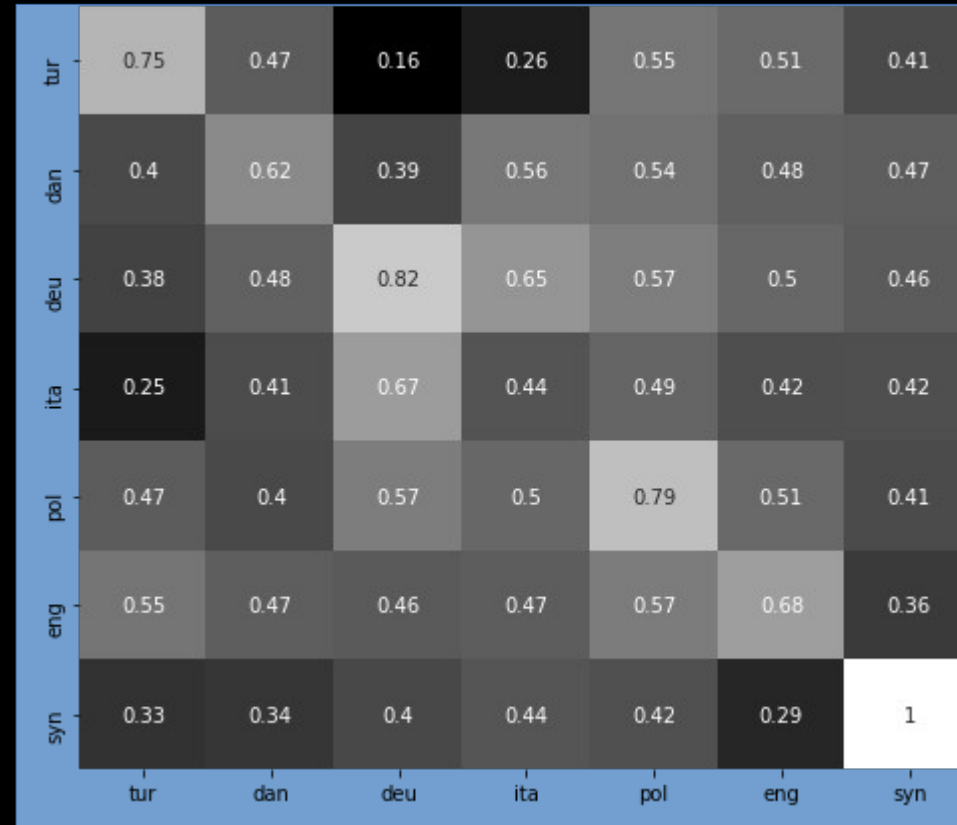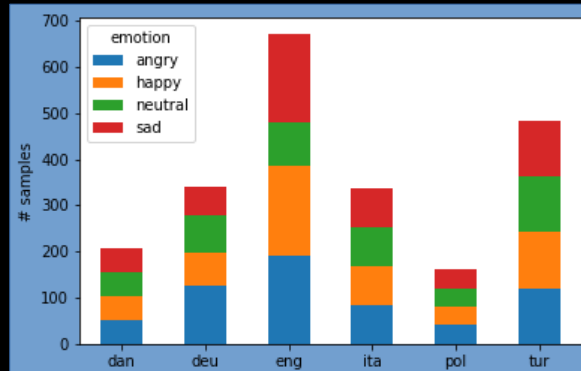- plot_epochs
- plot_anim_progression
- plot_epoch_progression
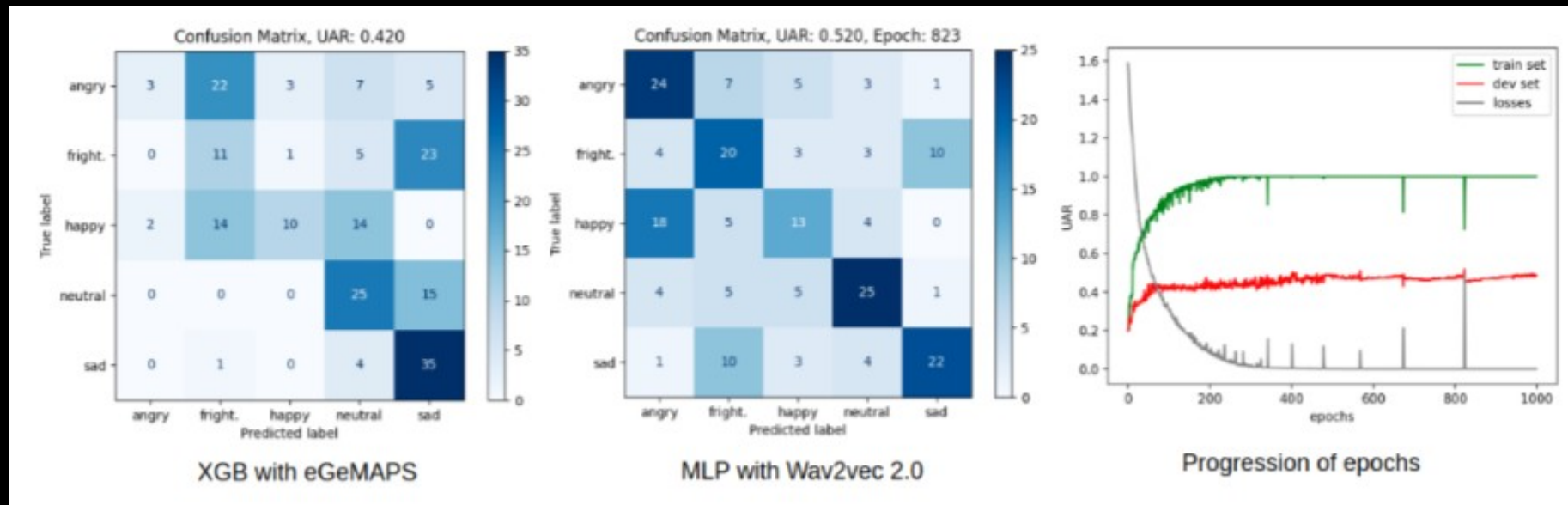- plot_best_model
- t-SNE

https://github.com/felixbur/nkululeko/blob/main/ini_file.md

# experiments



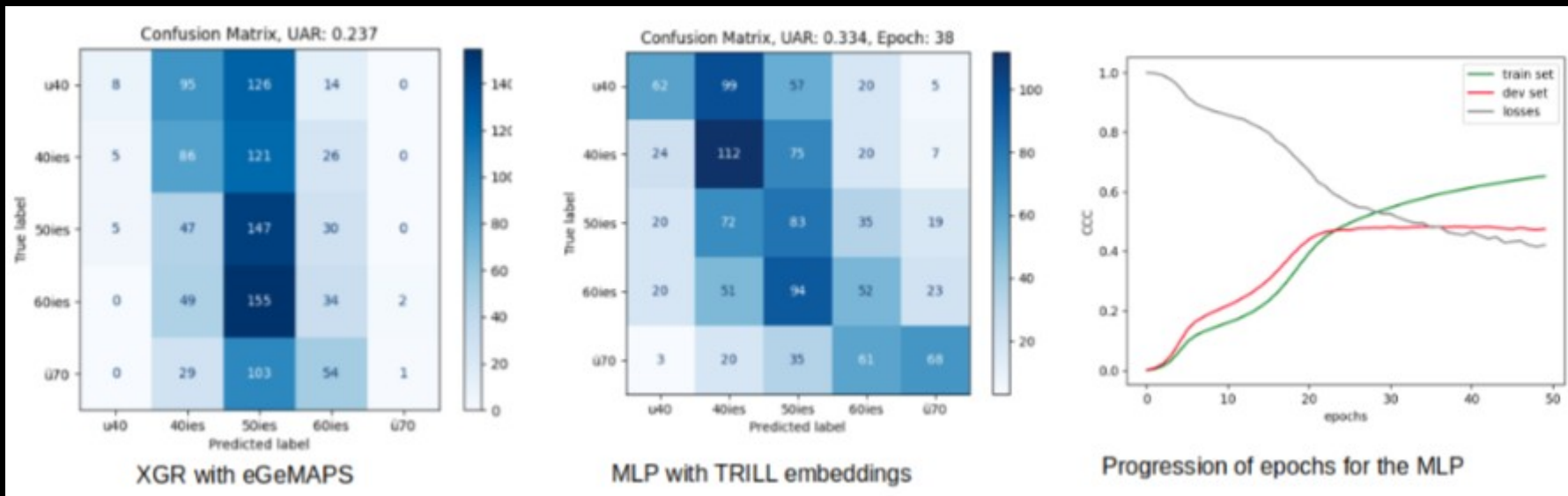- 6 European acted emotional databases + synthesized emotion

# experiments cont.



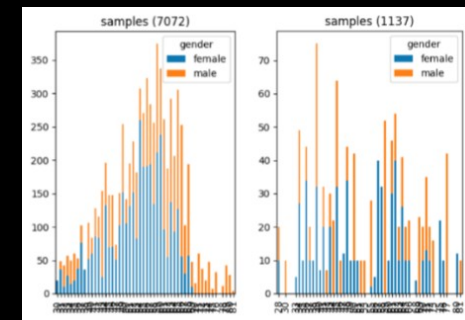Confusion Matrix, UAR: 0.420 — XGB with eGeMAPS

Confusion Matrix, UAR: 0.520, Epoch: 823 — MLP with Wav2vec 2.0

Progression of epochs

- comparing expert with learned features on cross databases acted emotion learning
- Berlin EmoDB vs. Polish data

# experiments cont.



XGR with eGeMAPS — MLP with TRILL embeddings — Progression of epochs for the MLP

- comparing expert with learned features on age regression
- Data: German parliament data



train and test distribution

# wrap up

- introduced Nkululeko

- a software to do machine learning experiments on spoken data without programming

- combines features and learners

- available open source with MIT license