

Multi-Label Speaker Recognition using Recurrent Neural Networks

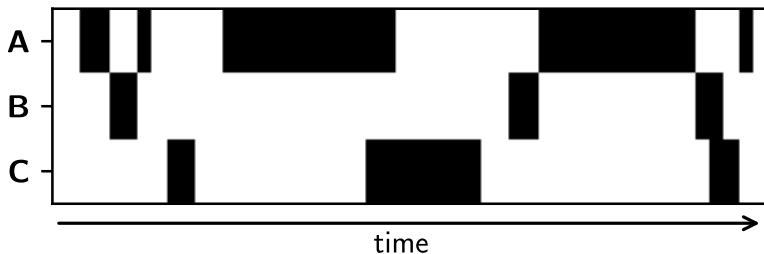
Puhujien tunnistus takaisinkytketyillä neuroverkoilla

Lauri Niskanen

Diplomityöseminaari, 29.10.2018

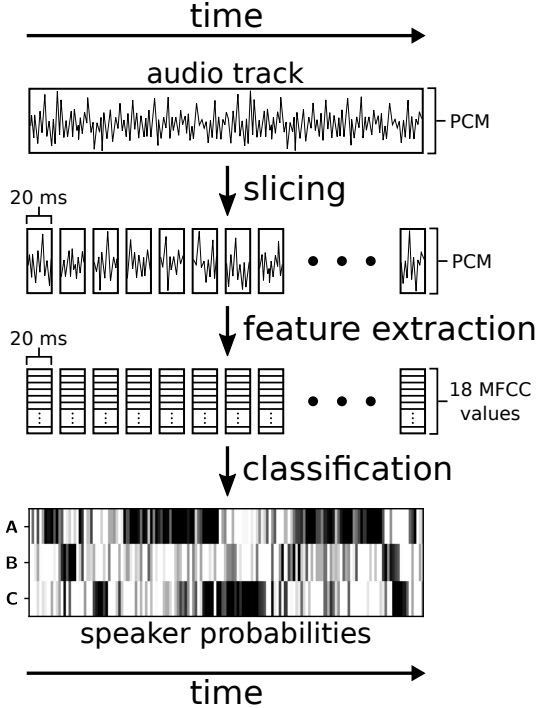
Puhujien tunnistaminen

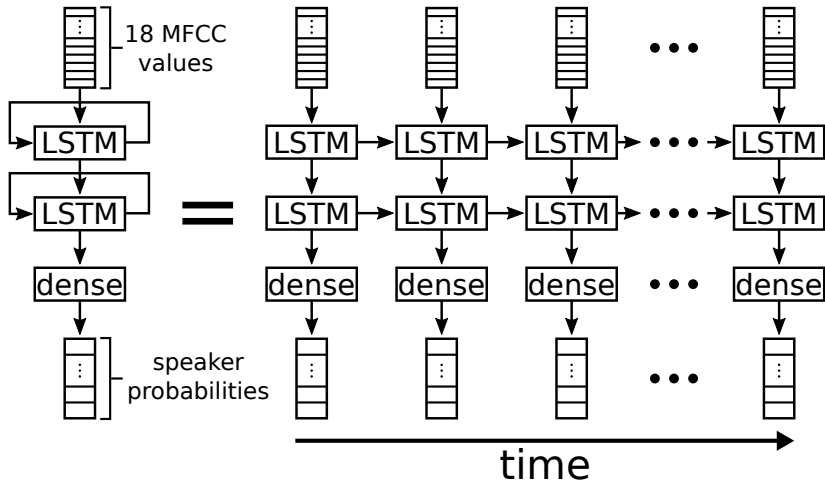
- ▶ Syöteenä on äänitiedosto keskustelusta.
- ▶ Ohjelma tunnistaa milloin kukin puhuja on äänessä.
- ▶ Ei välitetä puheen sisällöstä.

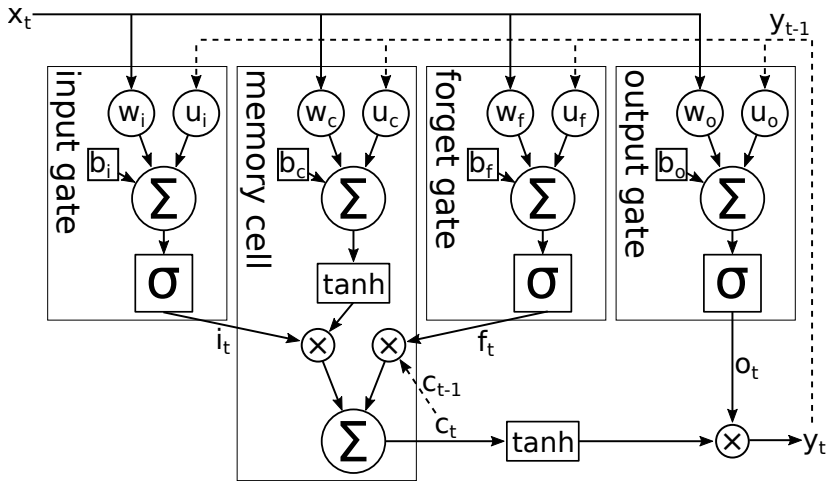


Toteutus

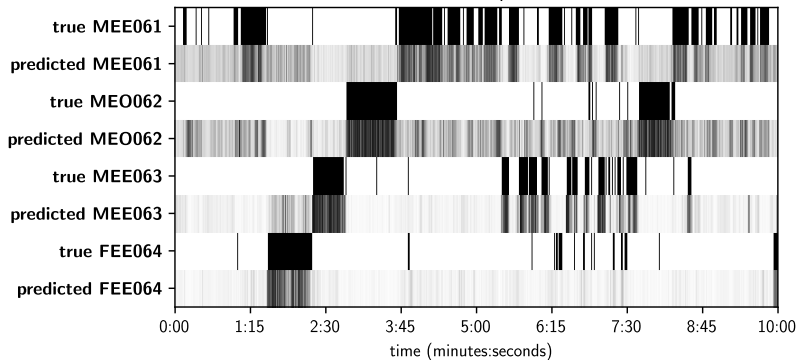
- ▶ Toteutus on tehty koneoppimistekniikoilla.
 - ▶ Luokitin opetetaan datan perusteella tekemään tunnistuksia.
- ▶ Äänen esiprosessointi helpommin käsiteltävään muotoon:
 - ▶ Mel-frequency cepstral coefficients (MFCC) -piirteet
- ▶ Takaisinkytkettyä neuroverkko käy ääntä läpi ja muistaa aiempia ajanhetkiä:
 - ▶ Long short-term memory (LSTM) -neuroverkkoluokitin



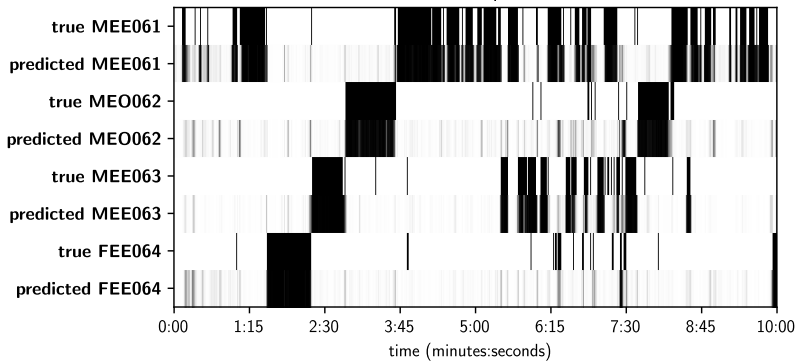




ES2016a - epoch: 1



ES2016a - epoch: 37

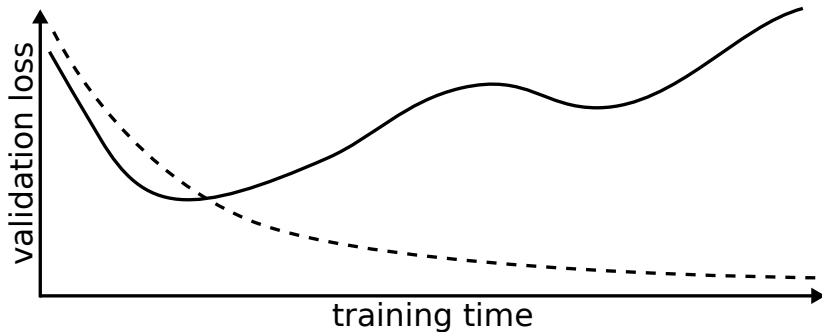


Hyperparametrit

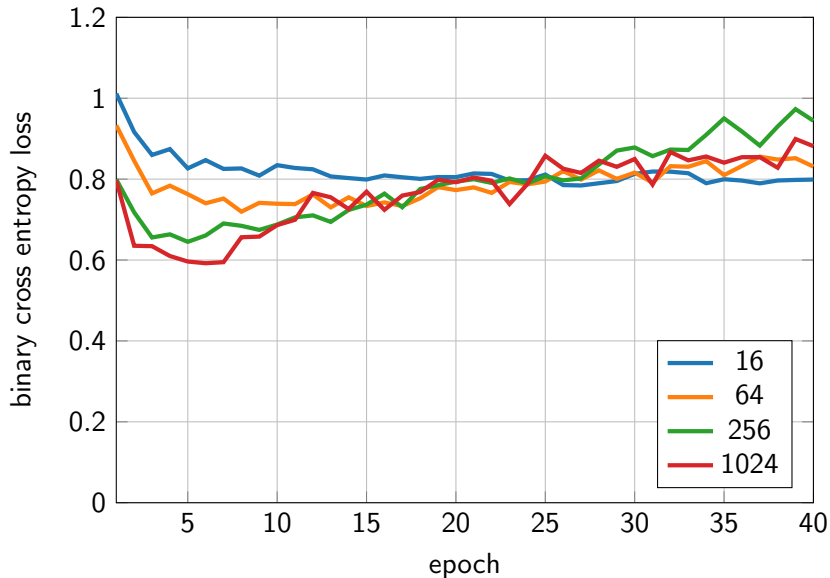
- ▶ Luokitinneuroverkolla on parametrejä, joita verkko ei opi automaattisesti.
- ▶ Säättämällä arvot oikein voi tuloksia saada paremmaksi.
- ▶ Tärkeimmät parametrit säättävät verkon oppimiskykyä.
 - ▶ Liian yksinkertainen verkko ei pysty oppimaan monimutkaisia tehtäviä.
 - ▶ Toisaalta liian suuri verkko ylioppii opetusdatalle, eikä yleisty uudelle datalle.
- ▶ Parametrejä:
 - ▶ LSTM-kerroksien lukumäärä (1, 2)
 - ▶ LSTM-kerrosten koko (16, 64, 256, 1024)
 - ▶ Dropout (0 %, 50 %)

Oppimiskäyrät ja ylioppiminen

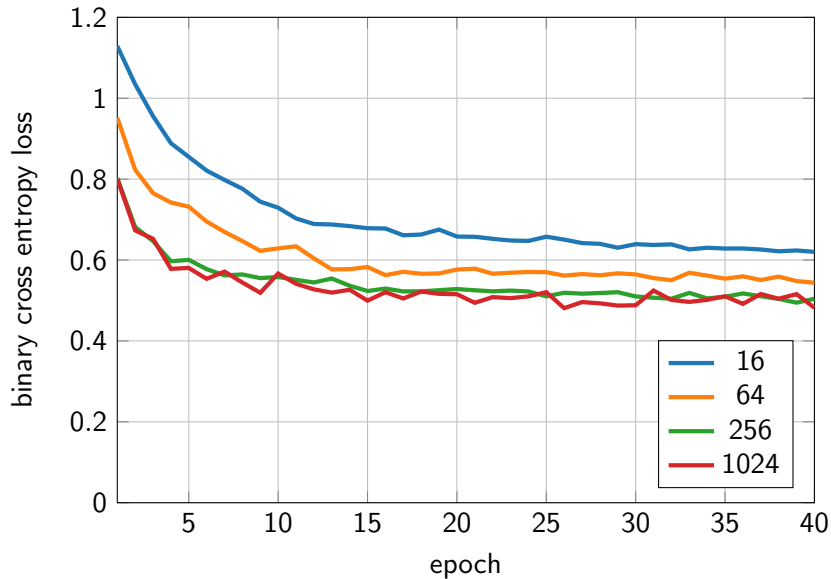
- ▶ Tavoitteena on saada luokittimen virhe mahdollisimman pieneksi.
- ▶ Virhe pienenee kun opetukseen käyttää enemmän aikaa.
- ▶ Jos luokitin alkaa ylioppia, niin virhe alkaakin kasvaa.



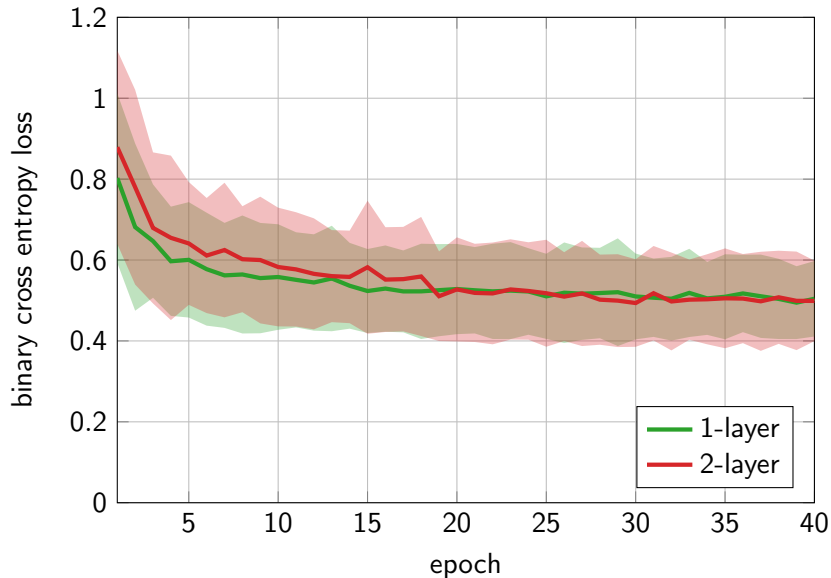
Verkon koon vaikutus



Dropout



Kaksi LSTM-kerrosta



Opetukseen kuluva aika

training time per epoch (s)		
	1 layer	2 layers
layer size 16	35.1 ± 0.78	73.8 ± 3.50
layer size 64	36.7 ± 1.94	72.1 ± 2.68
layer size 256	35.0 ± 1.01	74.0 ± 1.84
layer size 1024	121.9 ± 5.22	340.1 ± 13.89

- ▶ Pienemmät verkot on nopeampi opettaa.

Tuloksia

