

# Audio Language Classifier

Kirsten Regier

December 20, 2020

# Why Audio Classifiers?

**Goal:** Using the VGGish model as a feature extractor, build and train a classifier to identify speaker characteristics (gender or native language) from audio files.

Automatic speech recognition (ASR) and voice assistants are expanding.

ASR systems don't work well with non-standard or non-native accents [Sheng & Edmund, 2017].

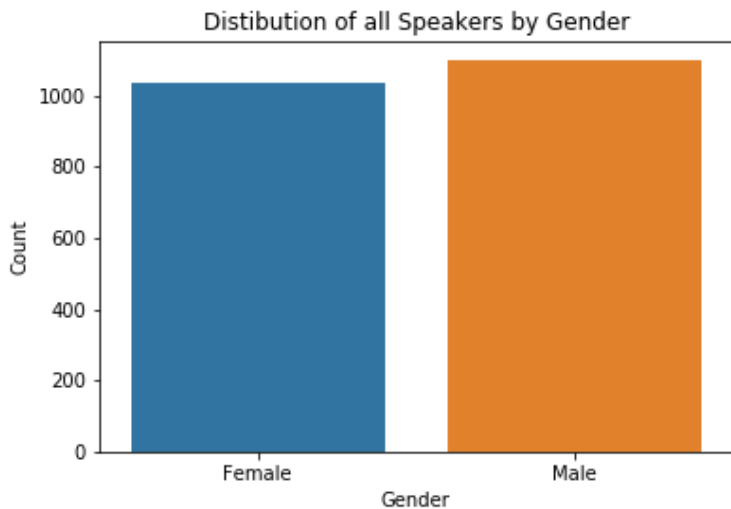
Better speech technology may improve customer experiences and customer service.

# The Data

## Speech Accent Archive [Weinberger, 2013]

- ▶ 2000+ recordings of speech samples.
- ▶ Speakers read a fixed passage in English.
- ▶ Demographic information about speakers: age, sex, birthplace, country, native language.
  - ▶ Speakers from 177 countries.
  - ▶ Native speakers of 199 native languages.
  - ▶ 78 languages represented by a single speaker.
  - ▶ 121 languages with multiple speakers.

# Distribution of Speakers by Gender



# Top 10 Languages in Speech Accent Archive

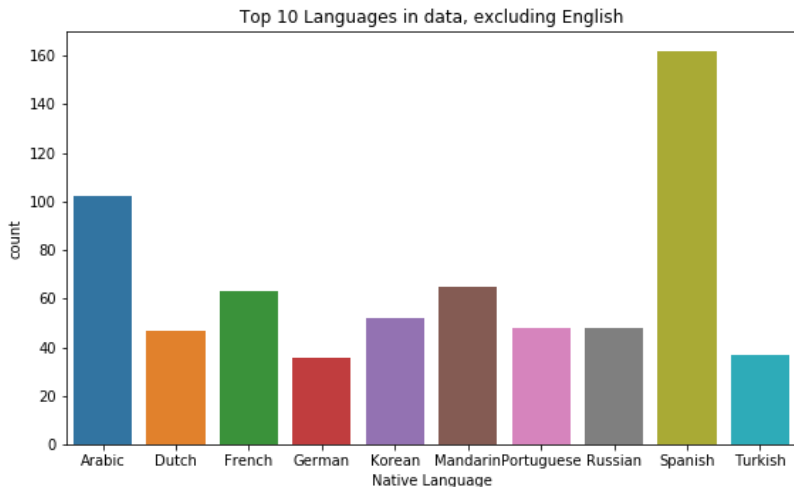


Figure: Number of speakers for top 10 languages, excluding English

# Transfer Learning

Transfer learning from image classification models has been an effective starting point for audio classification.

[[Hershey, et. al., 2017](#)]

[VGGish](#) is an audio embedding model based on the architecture of the VGG image classifier, that has been trained on the [YouTube-8M Segments Dataset](#) for Acoustic Event Detection Classification (AED). [[Hershey, et. al., 2017](#)]

Features extracted by VGGish for AED may be useful for training models for (acoustic) speech classification, with relatively short training times and small datasets.

# Model Architecture

Audio files fed through the VGGish model.

- ▶ Audio segment (10s) converted to mel spectrogram.
- ▶ Mel spectrogram fed through VGGish convolutional layers.
- ▶ VGGish output is a 128-D array of features.

VGGish features fed to a Gender Classifier or a Language Classifier.

Classifier models vary in:

- ▶ Number of (Dense + Dropout) layer repeats.
- ▶ Number of nodes in Dense layer(s).
- ▶ Shape of Flatten layer (based on output of previous layer)

# Gender Classifier Architecture

Table: Structure of the Gender Classifier models.

Layer	Gender 1	Gender 2
Input	(32, 10, 128)	(32, 10, 128)
Dense Dropout	128 nodes 50%	128 nodes 50 %
Dense Dropout		64 nodes 50%
Flatten	(32, 1280)	(32, 640)
Output	(32, 1)	(32, 1)



# Language Classifier Architecture

Table: Structure of the Language Classifier models.

Layer	Language 1	Language 2	Language 3
Input	(32, 10, 128)	(32, 10, 128)	(32, 10, 128)
Dense	12 nodes	128 nodes	128 nodes
Dropout	50%	50%	50 %
Dense			64 nodes
Dropout			50%
Flatten	(32, 120)	(32, 1280)	(32, 640)
Output	(32, 11)	(32, 11)	(32, 11)

# Results

## Gender Classifier

Identify the gender of the speaker.

- ▶ Convergence after 2-3 epochs.
- ▶ 98% accuracy.

## Language Classifier

Identify the native language of the speaker.

- ▶ 11 language classes.
- ▶ Convergence after 7 epochs.
- ▶ 25% accuracy with majority class distribution of 15%.

# Gender Classifiers - Model Metrics

Table: Summary Metrics

	Gender 1	Gender 2
Loss	0.061640	0.076962
Accuracy	0.981419	0.975507
Precision	0.982818	0.960199
Recall	0.979452	0.991438

Table: Confusion Matrices

		Predicted				Predicted	
		F	M			F	M
Actual	Gender 1	F	M	Gender 2	F	M	M
	M						
	F	590	10		F	576	24
	M	12	572		M	5	579

# Language Classifier - Model Metrics

11 language classes.

Baseline accuracy of 15% based on majority class distribution.

		Language 1	Language 2	Language 3
F1 score	Loss	2.25485	2.323	2.20684
	Accuracy	0.232955	0.238636	0.25
	Precision	0.625	0.5	0.605263
	Recall	0.0142045	0.0738636	0.0653
	Micro	0.232955	0.238636	0.25
	Macro	0.16999	0.195512	0.191103
	Weighted	0.20067	0.228178	0.224565

# Language Classifier - Confusion Matrix

Table: Confusion matrix for Language 3 predictions. The bold numbers on the diagonal represent correct predictions.

lang	R	A	T	K	G	D	S	F	E	P	M	Segments
Russian	<b>0</b>	2	1	4	0	1	7	4	1	1	7	28
Arabic	0	<b>24</b>	0	4	1	7	2	2	4	0	11	55
Turkish	0	1	<b>1</b>	4	0	3	4	0	2	1	5	21
Korean	1	10	0	<b>6</b>	0	3	2	2	4	2	1	31
German	0	0	0	1	<b>0</b>	0	5	1	5	2	1	15
Dutch	0	4	0	0	0	<b>12</b>	3	1	3	1	1	25
Spanish	1	2	0	6	2	1	<b>12</b>	3	8	0	7	42
French	1	3	0	4	1	2	10	<b>5</b>	1	1	4	32
English	0	1	0	3	0	5	2	2	<b>13</b>	0	4	30
Portuguese	0	3	0	3	1	0	7	5	3	<b>0</b>	6	28
Mandarin	0	2	0	8	2	4	9	1	2	2	<b>15</b>	45
Total	3	52	2	43	7	38	63	26	46	10	62	<b>352</b>

# Conclusions

- ▶ Transfer learning is an effective strategy for training speech models.
- ▶ Acoustic features extracted by VGGish for AED are adequate to train a gender classifier with 98% accuracy.
- ▶ Language classifier with 11 classes shows improvement over majority class baseline.

**Future directions:** Optimize language classifier.

- ▶ Add more layers and/or nodes.
- ▶ Use different activation functions.
- ▶ Include additional acoustic features (e.g. tempo/beat tracking).

# References



Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., & Zheng, X. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. <https://www.tensorflow.org/>.



McFee, B., Lostanlen, V., Metsai, A., McVicar, M., Balke, S., Thomé, C. Raffel, C., Zalkow, F., Malek, A., Dana, Lee, K., Nieto, O., Jack Mason, J., Ellis, D., Battenberg, E., Seyfarth, S. Yamamoto, R. Choi, K., viktorandreevichmorozov, Moore, J. Bittner, R., Hidaka, S., Wei, Z., nullmightybofo, Hereñú, D., Stöter, F., Friesch, P., Weiss, A., Vollrath, M., & Kim, T. 2020. librosa/librosa: 0.8.0. Zenodo, <https://doi.org/10.5281/zenodo.3955228>

# References



Hershey, S. et. al., 2017. CNN Architectures for Large-Scale Audio Classification, *ICASSP*.



Sheng, L. M. A., and Edmund, M. W. X. 2017. Deep Learning Approach to Accent Classification. Project Report Stanford University Stanford CA.



Weinberger, S. 2015. Speech Accent Archive. George Mason University.