



Gender Recognition Using Voice

B.MAHESWAR

Gender Recognition



Gender recognition is a technology that uses audio features to predict the biological sex of a speaker. It analyzes speech characteristics like pitch, formants, and vocal tract acoustics to classify the speaker as male or female.

Applications of Gender Recognition:

Smart assistants: Voice-activated assistants like Siri and Alexa can use gender recognition to personalize their responses and interactions.

Security systems: Gender recognition can be integrated with security cameras and access control systems to identify individuals based on their voice.

Voice biometrics: This technology uses voice as a unique identifier like fingerprints or facial recognition. Gender recognition can be used as one of the factors for user authentication.

Why TensorFlow is a Powerful Tool:



TensorFlow is a free and open-source software library for numerical computation and large-scale machine learning. It is particularly well-suited for tasks involving deep neural networks (DNNs) and other types of artificial intelligence (AI).

Its key strengths for gender recognition include:

Flexibility: TensorFlow allows building various deep learning architectures, from CNNs and RNNs to LSTMs, effectively capturing relevant features from speech data.

Scalability: TensorFlow can be easily scaled to handle large datasets and complex models, crucial for efficient training and accurate results.

Community and Resources: TensorFlow has a vast and active community offering extensive resources, tutorials, and pre-trained models, simplifying the development process.

Data Acquisition and Preprocessing



Dataset Description:

- Name: Mozilla's Common Voice Dataset.

Mozilla's Common Voice Dataset Properties:

Common Voice is a massive publicly available dataset of spoken language, designed to help train speech recognition systems. It is a valuable resource for researchers and developers working on speech recognition technologies.

Properties:

Size:

Over 28,750 hours of recorded speech

Over 19,160 hours of validated speech

Data collected in over 100 languages



Speakers:

Over 1 million registered contributors

Diverse demographics across age, sex, accent, and geographic location

Contributions from all over the world

Data Format:

Each entry consists of a unique MP3 audio file and a corresponding text file

Text files contain the spoken sentences

Many entries include metadata like speaker age, sex, and accent

Features:

Offers a wide range of accents and dialects

Provides data for both formal and informal speech styles

Includes recordings from various environments and noise levels

Benefits:

Open-source and freely available

Supports research and development of speech recognition technologies

Encourages the development of more inclusive and diverse language models

Feature extraction

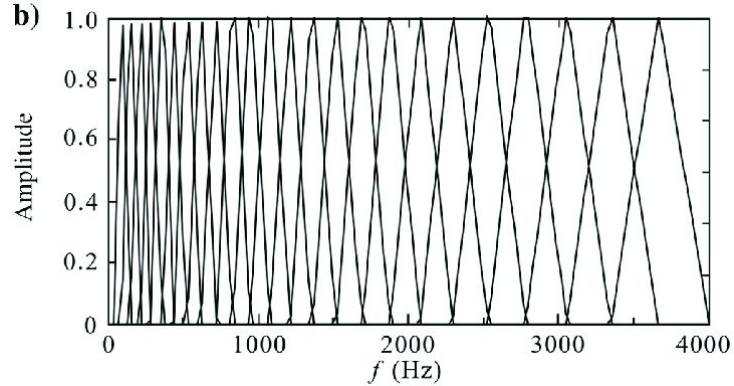
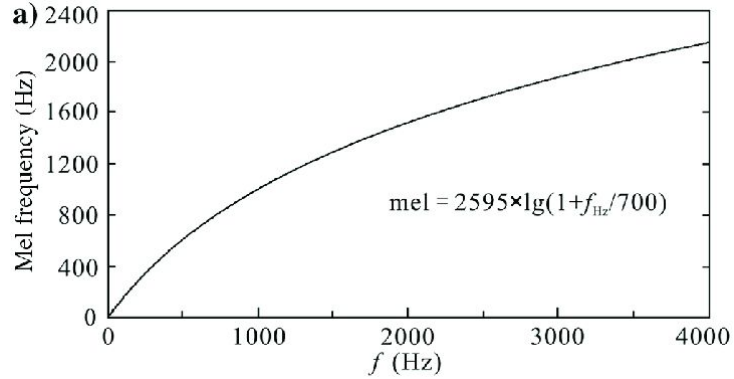


Mel-Frequency Cepstral Coefficients (MFCCs) are a set of features that represent the spectral characteristics of a sound. They are commonly used in audio processing tasks, such as speech recognition, speaker identification, and music genre classification.

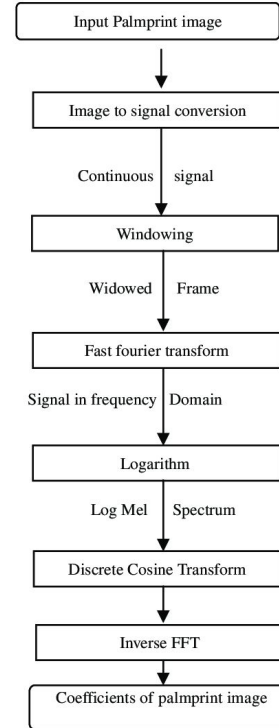
Mel-frequency Scale and Human Hearing:

Human ears perceive sound differently across different frequencies. We are more sensitive to changes in the lower frequencies than the higher ones.

The Mel-frequency scale is a non-linear scale that approximates the human perception of sound. It compresses the high-frequency region and expands the low-frequency region.



MFCC feature extraction process.



Model Building: Introduction to TensorFlow



TensorFlow is an open-source machine learning library widely used for building and training deep learning models. It provides a powerful and flexible platform for researchers and developers to solve complex problems in various domains, including speech recognition, natural language processing, and computer vision.

Here are some key advantages of using TensorFlow for building and training deep learning models:

1. Flexibility and Expressiveness:

TensorFlow offers a high degree of flexibility and expressiveness, allowing you to build and experiment with various deep learning architectures, including CNNs, RNNs, LSTMs, and custom architectures.

It provides pre-built components and tools that can be readily integrated into your models, saving you development time.



2. Scalability and Performance:

TensorFlow is highly scalable, allowing you to train your models on large datasets and distribute the training process across multiple GPUs or TPUs for faster training times.

It offers various optimization algorithms and techniques to improve training efficiency and achieve optimal model performance.

3. Open-source and Community:

TensorFlow is an open-source project with a large and active community of developers and researchers. This means you have access to extensive documentation, tutorials, code examples, and community support.

The open-source nature allows you to freely use and modify the library to suit your specific needs and research interests.



4. Production-ready and Deployment:

TensorFlow provides tools and libraries for deploying trained models into production environments. This includes serving models for real-time inference and integrating them with other applications and systems.

TensorFlow offers various deployment options, ranging from cloud platforms to mobile devices and edge computing devices.

5. Rich Ecosystem and Integrations:

TensorFlow has a vast ecosystem of tools and libraries that can be integrated into your workflow. This includes libraries for data preprocessing, visualization, and model analysis, making it a comprehensive platform for deep learning development.

TensorFlow integrates with other popular frameworks and libraries, such as NumPy, Scikit-learn, and Keras, allowing you to leverage existing tools and libraries for data analysis and model building.



Highlight specific features of TensorFlow relevant to your gender recognition project, such as:

- TensorFlow's support for various audio processing libraries and tools.
- Pre-built Keras layers and models for CNN and RNN architectures commonly used in gender recognition.
- TensorFlow's ability to handle large datasets and train complex models efficiently.

Deep Learning Architectures



Convolutional Neural Networks (CNNs):

CNNs are powerful deep learning architectures capable of automatically extracting spatial features from data like images and audio signals.

They consist of layers of convolutional filters that learn to identify relevant patterns in the data.

In gender recognition, CNNs can effectively learn features related to pitch, formants, and vocal tract characteristics from MFCCs or other spectral representations of speech.

Benefits of using CNNs for gender recognition:

Automatic feature extraction: CNNs automatically learn relevant features from the data, eliminating the need for manual feature engineering.

Robustness to noise: CNNs are robust to noise and variations in speech data, making them suitable for real-world applications.

High accuracy: CNNs have achieved state-of-the-art accuracy in gender recognition tasks.



Recurrent Neural Networks (RNNs):

RNNs are specifically designed to handle sequential data like speech. They contain loops that allow them to process information across different time steps.

RNNs can capture the temporal dynamics of speech signals, which can be helpful for gender recognition.

Long Short-Term Memory (LSTM) networks:

LSTMs are a type of RNN that can learn long-term dependencies in data. This is particularly beneficial for analyzing speech, where long-term dependencies often exist between different parts of an utterance.

LSTM networks have shown promising results for gender recognition, particularly when dealing with longer speech segments or complex speech patterns.

Comparison of CNNs, RNNs, and LSTMs



1. Strengths:

CNNs: Efficient feature extraction, robust to noise, high accuracy, suitable for short segments and spatial features.

RNNs: Handles sequential data, captures temporal dynamics, suitable for continuous speech and temporal dynamics.

LSTMs: Captures long-term dependencies, good for longer segments and complex patterns.

2. Weaknesses:

CNNs: May ignore temporal relationships, struggles with long dependencies.

RNNs: Sensitive to noise, vanishing gradient problem.

LSTMs: Computationally expensive, prone to overfitting.



3. Suitable for:

CNNs: Short segments, spatial features.

RNNs: Continuous speech, temporal dynamics.

LSTMs: Longer segments, complex patterns.

4. Typical Architectures:

CNNs: VGGNet, ResNet.

RNNs: Simple RNN, GRU.

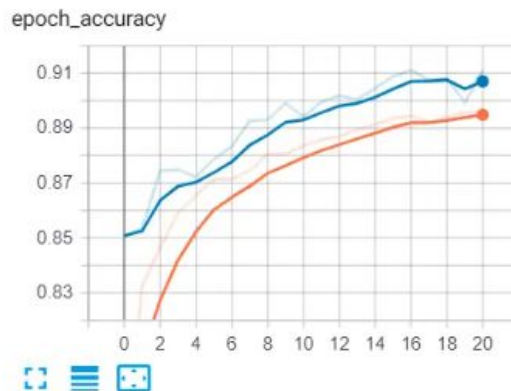
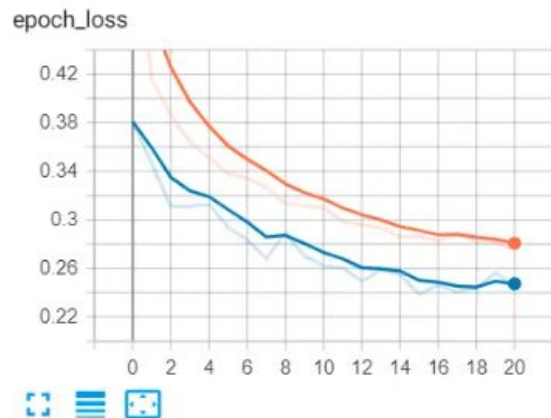
LSTMs: LSTM, Bi-LSTM.

Testing the model:

Evaluating the model using 6694 samples...

Loss: 0.2305

Accuracy: 92.95%



The blue curve is the validation set, whereas the orange is the training set. You can see the loss is decreasing over time, and the accuracy is increasing.

Testing the model with a sample file:

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 256)	33024
dropout (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 256)	65792
dropout_1 (Dropout)	(None, 256)	0
dense_2 (Dense)	(None, 128)	32896
dropout_2 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 128)	16512
dropout_3 (Dropout)	(None, 128)	0
dense_4 (Dense)	(None, 64)	8256
dropout_4 (Dropout)	(None, 64)	0
dense_5 (Dense)	(None, 1)	65

=====
Total params: 156545 (611.50 KB)
Trainable params: 156545 (611.50 KB)
Non-trainable params: 0 (0.00 Byte)

1/1 [=====] - ETA: 0s
1/1 [=====] - 0s 412ms/step
Result: male
Probabilities: Male: 96.57% Female: 3.43%

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 256)	33024
dropout (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 256)	65792
dropout_1 (Dropout)	(None, 256)	0
dense_2 (Dense)	(None, 128)	32896
dropout_2 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 128)	16512
dropout_3 (Dropout)	(None, 128)	0
dense_4 (Dense)	(None, 64)	8256
dropout_4 (Dropout)	(None, 64)	0
dense_5 (Dense)	(None, 1)	65

=====
Total params: 156545 (611.50 KB)
Trainable params: 156545 (611.50 KB)
Non-trainable params: 0 (0.00 Byte)

1/1 [=====] - ETA: 0s
1/1 [=====] - 0s 162ms/step
Result: female
Probabilities: Male: 7.31% Female: 92.69%



Let see the code for a test run.....