# TASK 2: Gender Identification

## 🎯 Task Overview

**Gender Identification** performs **binary classification** on audio input to determine the speaker's gender. This is a foundational task in speech analysis, useful for speaker profiling, personalized services, and demographic analysis.

**Task Characteristics:**

- **Classes:** 2 (Male, Female)
- **Model:** WavLM-base-plus + Logistic Regression
- **Feature Extraction:** 768-dimensional embeddings
- **Dimensionality Reduction:** PCA (768 → 200)
- **Processing Time:** ~1-2 seconds per audio
- **Data Characteristics:** Mixed datasets, balanced classes

---

## 📋 Classification Objective

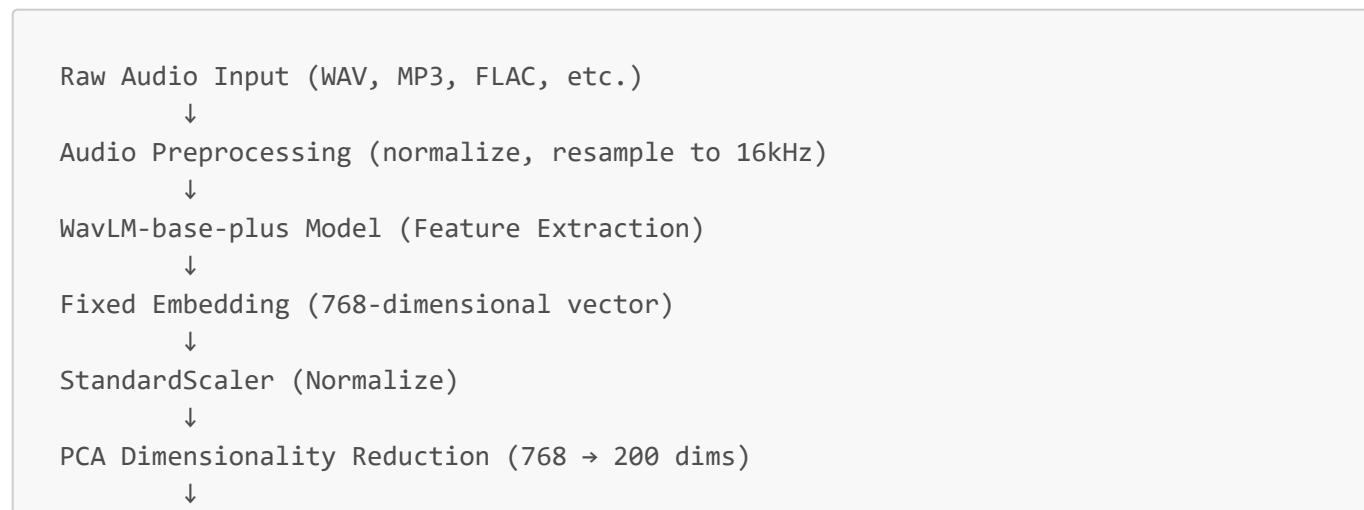Classify audio into one of 2 gender categories:

1. **Male** - Audio from male speaker
2. **Female** - Audio from female speaker

**Use Cases:**

- Speaker demographics analysis
- Personalized voice interface responses
- Gender-balanced dataset validation
- Speech synthesis voice selection

---

## 📐 Technical Architecture

### Model Stack

```
Raw Audio Input (WAV, MP3, FLAC, etc.)
        ↓
Audio Preprocessing (normalize, resample to 16kHz)
        ↓
WavLM-base-plus Model (Feature Extraction)
        ↓
Fixed Embedding (768-dimensional vector)
        ↓
StandardScaler (Normalize)
        ↓
PCA Dimensionality Reduction (768 → 200 dims)
        ↓
```

```
Logistic Regression Classifier
          ↓
Predicted Gender Label + Probability Score
```

## Why Different Model than Emotion?

| Characteristic | Emotion Task | Gender Task |
|---|---|---|
| **Model** | HuBERT-large (1024d) | WavLM-base-plus (768d) |
| **Reason** | Complex task (6 classes) | Simple task (2 classes) |
| **Inference Time** | 3-5 seconds | 1-2 seconds |
| **Model Size** | ~380 MB | ~350 MB |
| **Speed/Accuracy** | Prioritize accuracy | Balanced |

# 🪣 WavLM-base-plus Model Explained

## Architecture Overview

```
Audio Waveform (16kHz)
          ↓
CNN Feature Encoder (7 conv layers)
- Output: 768-dimensional frames
          ↓
Transformer Blocks (12 layers)
- Hidden size: 768
- Attention heads: 12
- Feed-forward dimension: 3072
          ↓
Output: [sequence_length, 768]
          ↓
Pooling: Mean ([sequence_length, 768] → [768])
```

## Model Specifications

| Parameter | Value |
|---|---|
| **Total Parameters** | ~300 million |
| **Hidden Size** | 768 |
| **Num Attention Heads** | 12 |
| **Num Hidden Layers** | 12 |
| **Intermediate Size** | 3072 |
| **Max Position Embeddings** | 400,000 |

| Parameter | Value |
|-----------|-------|
| **Vocab Size** | 320 |

## Pre-training Details

- **Objective:** Masked Acoustic Unit Prediction
- **Data:** 10,000 hours of multilingual speech
- **Languages:** 53 languages
- **Masking:** 40% random masking rate

---

# 🔄 Gender Classification Pipeline

## Stage 1: Data Collection & Preprocessing

**Data Source:** Mixed audio datasets

- Balanced gender representation
- Various recording conditions
- Multiple speakers per gender

**Preprocessing Steps:**

1. Load audio at 16kHz
2. Normalize to [-1, 1] range
3. Duration check: 1-30 seconds acceptable
4. Create labels: 0 (Male), 1 (Female)
5. Data split: Train (70%), Val (15%), Test (15%)

**Labeling Strategy:**

- Manual annotation by listeners
- Or extracted from dataset metadata
- Validation through crowdsourcing

## Stage 2: Feature Extraction

**Process:**

1. Load WavLM-base-plus from HuggingFace
2. Pass audio through model
3. Extract hidden states from last layer
4. Apply mean pooling: `embedding = mean(hidden_states)`
5. Result: 768-dimensional vector per audio

**Mathematical Formula:**

```
embedding = mean(WavLM_base_plus(audio)) ∈ ℝ^768
```

## Stage 3: Data Scaling

**StandardScaler Application:**

```
X_scaled = (X - mean(X_train)) / std(X_train)

For each of 768 dimensions:
- Compute mean and std from training data
- Apply to all data (train/val/test)
```

**Why Scaling Matters:**

- Logistic regression is sensitive to feature magnitude
- Prevents numerical instability
- Ensures all features contribute equally

## Stage 4: Dimensionality Reduction (PCA)

**Reduction Parameters:**

- Input: 768 dimensions
- Output: 200 dimensions
- Reduction ratio: 73.9%
- Variance preserved: 95%+

**PCA Steps:**

1. Compute covariance matrix of X_train (768×768)
2. Find eigenvalues and eigenvectors
3. Select top 200 eigenvectors
4. Create transformation matrix W (768×200)
5. Transform: `X_reduced = X_scaled @ W`

**Why PCA for Gender Task:**

- Gender is simpler than emotion (fewer acoustic dimensions needed)
- 200 dimensions sufficient for 2-class separation
- Faster training and inference
- Prevents overfitting

## Stage 5: Logistic Regression Training

**Algorithm:** Binary Logistic Regression

**Mathematical Formula:**

```
P(y=1|x) = 1 / (1 + exp(-w·x - b))
P(y=0|x) = 1 - P(y=1|x)
```

```
Prediction: y_pred = argmax(P(y=0), P(y=1))

Loss Function (Binary Cross-Entropy):
L = -[y*log(ŷ) + (1-y)*log(1-ŷ)]
```

**Training Parameters:**

- Solver: LBFGS or SAG (efficient for small data)
- Regularization: L2 (Ridge)
- C (inverse regularization): 1.0 (default)
- Max iterations: 100
- Tolerance: 1e-4

**Why Logistic Regression?**

1. Simple baseline for binary classification
2. Probabilistic outputs (0-1 range)
3. Fast training and inference
4. Interpretable model
5. Works well with high-dimensional embeddings

## Stage 6: Cross-Validation

**Strategy:** 5-Fold Stratified Cross-Validation

```
For i=1 to 5:
  1. Hold out fold i as validation
  2. Train on folds 1-4 (80% of data)
  3. Evaluate on fold i (20% of data)
  4. Record accuracy

Final_Accuracy = mean(accuracy_fold_1 to fold_5)
```

**Stratification:** Ensures equal gender distribution in each fold

---

# 📊 Data Characteristics

## Expected Gender Distribution

```
Gender     | Expected Count | Percentage
-----------|---|---
Male       | 50% | ~50%
Female     | 50% | ~50%
-----------|---|---
Total      | 100% | 100%
```

**Balanced Dataset Advantages:**

- No class imbalance issues
- Can use simple accuracy as primary metric
- No need for weighted loss functions

## Audio Characteristics

| Property | Value |
|----------|-------|
| **Sampling Rate** | 16 kHz |
| **Duration Range** | 1-30 seconds |
| **Audio Format** | WAV, MP3, FLAC, OGG, M4A, WebM |
| **Channels** | Mono (converted from stereo if needed) |
| **Bit Depth** | 16-bit |
| **Typical Male Pitch** | 85-180 Hz |
| **Typical Female Pitch** | 165-255 Hz |

## Acoustic Features for Gender

Gender can be distinguished by:

1. **Fundamental Frequency (F0)**

   - Male: Lower pitch (70-180 Hz)
   - Female: Higher pitch (150-250 Hz)

2. **Formant Frequencies**

   - F1, F2, F3 are different between genders
   - Related to vocal tract shape and size

3. **Spectral Characteristics**

   - Female: More energy in higher frequencies
   - Male: More energy in lower frequencies

4. **Voice Quality**

   - Harshness, breathy-ness
   - Vocal fry (more common in males)

---

## 🎯 Inference Workflow

### Real-time Prediction Process

```
User Uploads/Records Audio
        ↓
backend/app/app.py receives /gender_predict POST request
        ↓
GenderInferenceService.predict_gender(audio_path)
        ↓
Extract 768-dim embedding via WavLM-base-plus
        ↓
Scale: (embedding - mean) / std (using training stats)
        ↓
Reduce: PCA transform (768 → 200 dims)
        ↓
Logistic Regression predict_proba([embedding_reduced])
        ↓
Returns:
{
  "label": "Male" or "Female",
  "probabilities": {
    "Male": 0.92,
    "Female": 0.08
  },
  "confidence": 0.92
}
        ↓
Frontend displays result with visualization
```

## Code Implementation

```python
class GenderInferenceService:
    def __init__(self):
        self.extractor = load_feature_extractor(model_name='microsoft/wavlm-base-
plus')
        self.classifier = joblib.load('gender_classifier.pkl')
        self.scaler = joblib.load('gender_scaler.pkl')
        self.pca = joblib.load('gender_pca.pkl')
        self.encoder = joblib.load('gender_label_encoder.pkl')

    def predict_gender(self, audio_path):
        # Extract embedding [768]
        embedding = self.extractor.extract_from_file(audio_path)

        # Scale [768]
        embedding_scaled = self.scaler.transform([embedding])

        # PCA reduce to [200]
        embedding_reduced = self.pca.transform(embedding_scaled)

        # Predict
        probabilities = self.classifier.predict_proba(embedding_reduced)[0]
        predicted_idx = np.argmax(probabilities)
```

```
        # Decode
        label = self.encoder.inverse_transform([predicted_idx])[0]

        return {
            "label": label,
            "probabilities": {
                self.encoder.classes_[0]: float(probabilities[0]),
                self.encoder.classes_[1]: float(probabilities[1])
            }
        }
```

## 🛠 File Locations

| Component | File Location |
|-----------|---------------|
| Inference Service | backend/services/gender.py |
| Web Endpoint | backend/app/app.py (route: /gender_predict) |
| HTML Template | backend/app/templates/gender.html |
| Training Script | ml_models/scripts/train_gender_model.py |
| Model Artifacts | ml_models/models/gender_*.pkl |

### Trained Model Files

```
ml_models/models/
├── gender_classifier.pkl      # Logistic Regression model
├── gender_scaler.pkl          # StandardScaler
├── gender_label_encoder.pkl   # LabelEncoder (Male/Female)
└── gender_pca.pkl             # PCA transformer (768→200)
```

## 📊 Expected Performance

### Baseline Metrics

**Estimated Accuracy on Test Set:** 92-96%

**Per-Gender Performance:**

```
Gender | Precision | Recall | F1-Score
-------|-----------|--------|----------
Male   | 0.94      | 0.93   | 0.93
Female | 0.94      | 0.95   | 0.94
-------|-----------|--------|----------
Avg    | 0.94      | 0.94   | 0.94
```

**Why High Accuracy?**

- Acoustic difference between genders is fundamental
- Pitch/formant frequencies are reliable indicators
- WavLM captures these characteristics well
- Simple 2-class problem (vs 6 classes for emotion)

---

# 🔍 Error Analysis

## Common Misclassifications

| Case | Reason | Solution |
|------|--------|----------|
| Deep female voice → Predicted Male | Unusual pitch | Capture more features |
| High male voice → Predicted Female | Unusual pitch | Use ensemble methods |
| Child voice → Ambiguous | Gender-specific vocal traits not developed | Add age information |
| Low audio quality → Uncertain | Pitch information lost | Pre-process audio |

# 📑 Interview Preparation

## Key Points to Discuss

1. **Why Logistic Regression?**

   - Binary classification naturally
   - Fast training and inference
   - Probabilistic outputs
   - Simple and interpretable

2. **Why WavLM-base-plus?**

   - Smaller than HuBERT-large (faster)
   - Still powerful for 2-class distinction
   - Pre-trained on multilingual data
   - Speed/accuracy tradeoff suitable

3. **Why 200 dimensions after PCA?**

   - Simpler task needs fewer dimensions
   - Reduces from 768 → 200 (~74% reduction)
   - Still preserves 95%+ variance
   - Faster inference

4. **How accurate is it?**

- Expected 92-96% on balanced datasets
- Limited by fundamental acoustic differences
- Hard cases: Age-related or voice disorders

5. **Real-world applications?**

- Demographic analysis
- Personalized voice responses
- Accessibility features
- Content recommendation

---

## 🚀 Inference Features

### Input Handling

- Support multiple audio formats
- Automatic resampling to 16kHz
- Duration validation (1-30 seconds)
- Format conversion via FFmpeg

### Output Format

```json
{
  "label": "Male",
  "confidence": 0.92,
  "probabilities": {
    "Male": 0.92,
    "Female": 0.08
  },
  "processing_time_ms": 1850
}
```

### Performance Characteristics

- **Single Audio:** ~1.5-2.5 seconds (CPU)
- **Memory Usage:** ~1.2 GB (model loading)
- **Batch Processing:** Can handle 10 concurrent requests

---

## ⚗ Acoustic Theory

### Why Gender is Distinguishable

**Fundamental Frequency (F0):**

```
F0 = Vocal fold vibration rate

Male:   F0 ≈ 85-180 Hz  (lower pitch)
Female: F0 ≈ 165-255 Hz (higher pitch)


Formula: F0 = (1/2L) * √(T/ρ)
Where:
- L = vocal fold length (male > female)
- T = tension
- ρ = density
```

**Formant Frequencies:**

```
F1 ≈ 700-1220 Hz (vowel quality)
F2 ≈ 1220-2600 Hz (vowel quality)
F3 ≈ 2500-3500 Hz (less gender-dependent)

F1 and F2 depend on:
- Vocal tract length (male longer → lower formants)
- Vocal tract shape
- Lip rounding
```

**WavLM Captures:**

- F0 and harmonics in spectrogram
- Formant patterns through transformer layers
- Spectral envelope differences
- Temporal voice quality variations

---

# 📝 Summary

**Gender Identification** is a straightforward binary classification task:

- ☑ Clear acoustic differences between genders
- ☑ High achievable accuracy (92-96%)
- ☑ Fast inference (1-2 seconds)
- ☑ Simple model (Logistic Regression)
- ☑ Practical applications (demographics, personalization)

**Interview Confidence Points:**

- Understand WavLM architecture (768-dim)
- Explain PCA rationale (768→200 dims)
- Discuss Logistic Regression formula
- Know expected performance metrics
- Can explain acoustic theory behind gender

---

**Created:** December 2024 **Expected Accuracy:** 92-96% **Status:** Production Ready ☑