

SVLDL: IMPROVED SPEAKER AGE ESTIMATION USING SELECTIVE VARIANCE LABEL DISTRIBUTION LEARNING

Zuheng Kang, Jianzong Wang*, Junqing Peng, Jing Xiao

Ping An Technology (Shenzhen) Co., Ltd.

ABSTRACT

Estimating age from a single speech is a classic and challenging topic. Although Label Distribution Learning (LDL) can represent adjacent indistinguishable ages well, the uncertainty of the age estimate for each utterance varies from person to person, i.e., the variance of the age distribution is different. To address this issue, we propose selective variance label distribution learning (SVLDL) method to adapt the variance of different age distributions. Furthermore, the model uses WavLM as the speech feature extractor and adds the auxiliary task of gender recognition to further improve the performance. Two tricks are applied on the loss function to enhance the robustness of the age estimation and improve the quality of the fitted age distribution. Extensive experiments show that the model achieves state-of-the-art performance on all aspects of the NIST SRE08-10 and a real-world datasets.

Index Terms— speaker age estimation, label distribution learning, multi-task learning, gender recognition

1. INTRODUCTION

Speech is the sound produced by the accurate coordinated movement of multiple organs in the human body. Hence, the acoustic characteristics of speech can transmit information about the physical characteristics of the speaker. The rapid development of new speech applications requires techniques capable of estimating information on various biological attributes of such speakers. Recently, deep-learning-based approaches show great performance in extracting hidden speech information, including facial expression [1] and emotion [2], and age [3], etc. If such speech features can be used to automatically estimate a speaker's age, it could be widely used for human-computer interaction, forensics, and other purposes.

Many researchers have studied the performance of human and artificial intelligence systems in estimating age from speech. The results show that the average error of humans judging the age of adults is about 10 years old, and the judgment of the age of children is about 1-year old [4]. The performance of age estimates may also have implications for human development. [5] collected the speech of children. It can be seen that, as children gradually enter puberty, changes in the vocal cords

can affect age estimates and increase uncertainty. In adulthood, the vocal cords are fully developed and the change tends to be slow. However, as we age, various organs experience regular aging: the voice changes from bright to hoarse, and articulation from clear to vague [6, 7]. Judgments at different ages also have different uncertainties, and these uncertainties may vary from age to age, from utterance to utterance.

Traditional methods for speaker age estimation can be generally classified into classification-based and regression-based methods. Most researchers mainly focus on the exploration of backbone model structures, such as deep neural network (DNN) [8], i-vector [9], x-vector [10, 11] or adding attention mechanism [12]. Some researchers have tried different machine learning features, such as the OpenSmile toolbox [13] to study this problem [14, 15]. As manipulated acoustic features, such as mel-filter banks, encounter performance bottlenecks, some researchers use other speech features for modeling, which can capture acoustic features that are imperceptible to the human ear, such as SincNet [16] take full advantage of acoustic information, resulting in improved performance. However, these features are only direct translations of speech signals, not language models for understanding human speech. Self-supervised learning (SSL) generates high-quality speech features with language model (such as wav2vec [17] and WavLM [18]) by learning from a large amount of data [19]. By injecting this prior knowledge, speech age estimation achieves better performance [20]. Although these methods have achieved great results, they ignored the fact that it rarely considers the relationship between labels, such as order and adjacent correlations, which are important clues for speaker age estimation. Since speaker age labels form an ordered set of numbers, significant ordinal relationships and adjacencies between labels should be fully exploited to achieve higher performance.

Label distribution learning (LDL) [21] addresses the above problems by transforming the classification problem into a distribution learning task that minimizes the difference between the predicted and constructed Gaussian distributions of labels. In the field of computer vision, impressive progress has been made in facial age estimation, where LDL shows great potential [22]. Framework [3] applied this method to the speaker age recognition task and achieved good performance. Since the uncertainty of each person is different, i.e.,

*Corresponding author: Jianzong Wang, jzwang@188.com

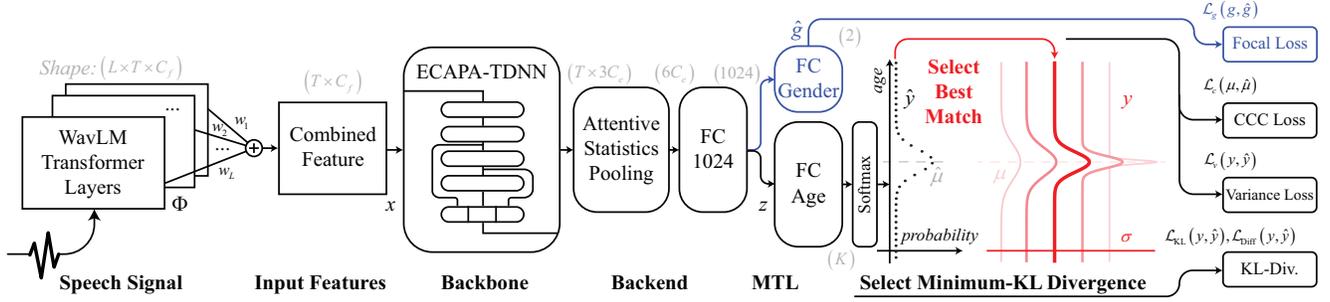


Fig. 1. Network topology of the SVLDDL framework. “FC” denotes a fully connected layer. \oplus denotes element-wise addition.

the variance of the Gaussian distribution varies from person to person, adaptive-based LDL methods have been proposed successively [23, 24, 25]. However, loss functions that measure regression error often use simple metrics, such as L1 distance, which are not dynamically adjusted for a specific distribution at training time. This method does not achieve optimal regression performance. Meanwhile, these algorithms do not get the correct shape of the learned distribution, which may lead to multimodal problems (multiple peaks in the fitted distribution).

Additionally, Multi-task learning (MTL) uses a shared backbone model to simultaneously optimize objectives for different tasks. The advantage comes from adding more useful information while optimizing the original model. In speaker age estimation, adding the task of gender recognition has been shown to improve performance [20, 26]. Meanwhile, in regression problems, Lin’s consistent correlation coefficient loss [27] also achieves a lot of performance gains by replacing L1 or L2 distance-based losses.

Considering the above advantages and disadvantages, we have made the following improvements and contributions:

- We improve the original label distribution learning (LDL) method and propose a new selective variance label distribution learning (SVLDDL) method that adaptively selects the optimal distribution that matches the variance.
- The quality of fitted distributions is improved by fitting additional first-order difference distribution, and a brief theoretical proof is given.
- The age estimation performance is enhanced by using Lin’s concordance correlation coefficient [27] loss.
- The performance was improved by adding an auxiliary task for gender recognition and using WavLM as the speech feature extractor.
- Experimental results on the publicly available NIST SRE08-10 dataset and a real-world dataset show that the improved SVLDDL framework achieves state-of-the-art performance compared to the framework [3].

2. METHODOLOGY

2.1. Network Architecture

Figure 1 outlines the pipeline of the proposed method. Since the structure of ECAPA-TDNN [28] has an efficient design structure, such as Res2Net [29] and squeeze excitation blocks (SE) [30], it is used as the backbone model. All the information on the time dimension is collected through attentive statistics pooling (SP). After the SP, there are two fully connected layers, and finally a softmax layer is connected to obtain the output distribution of the labels, denoted as y ; the output of the middle layer is denoted as z , which is also used as input for the auxiliary task of gender recognition.

2.2. Self-supervised Representation

Motivated by the successful application of self-supervised learning (SSL) in various speech domains, we explore the use of WavLM [18] on the task of speaker age estimation. The WavLM model learns speech representations by solving contrastive tasks in a latent space in a self-supervised manner. It tries to recover the randomly masked part of the encoded audio features. By learning from large amounts of real multilingual, multi-channel unlabeled data, SSL models can deeply understand contextual information and produce high-quality speech representations in the latent space.

In our framework, seen from Figure 1, we utilize all latent output of WavLM transformer layers $\Phi = (\phi_1, \dots, \phi_L)$ and assign a trainable weight $W = (w_1, \dots, w_L)$ to each of them. The weighted sum is then used to generate speech features $x = \sum_{i=1}^L (\phi_i \cdot w_i)$, where $\Phi \in \mathbb{R}^{L \times T \times C_f}$, $x \in \mathbb{R}^{T \times C_f}$, T is number of time frames, C_f is the feature size, L is the number of layers of WavLM. In this way, the model can make full use of speech information from shallow to deep, from concrete to abstract.

2.3. Label Distribution Learning

Before introducing SVLDDL, we need to know how LDL works and understand some parameters, $\hat{\mu}_n$ and $\hat{\sigma}_n$ are the mean and standard deviation of the predicted distribution, and μ_n

and σ_n are for the ground-truth of sample n . Where $\hat{\mu}_n = \sum_{k=1}^K k \cdot \hat{y}_n^k$, and $\hat{\sigma}_n = \frac{1}{N} \sum_{k=1}^K (\hat{y}_n^k - \hat{\mu}_n)^2$. To take advantage of the intrinsic relationship in model outputs, we treat these outputs as a distribution representing the predicted age distribution. The label distribution \hat{y}_n^k is a predicted probability distribution, which satisfy $\hat{y}_n^k \in [0, 1]$ and $\sum_{k=1}^K \hat{y}_n^k = 1$, n is the data sample, k is the age label, $k \in [1, K]$ and K denotes the maximum age. In age estimation, age is usually represented using a Gaussian distribution centered around the ground-truth age μ_n . This ground truth probability distribution y_n^k is represented by a Gaussian distribution function.

$$y_n^k = C_n \cdot e^{-(k-\mu_n)^2/(2\sigma^2)} \quad (1)$$

where σ is a fixed value that is reasonably chosen in LDL, C_n is a constant to make $\sum_k y_n^k = 1$. The difference between the ground truth label distribution \hat{y} and the predicted distribution y is measured using the Kullback-Leibler divergence (KL divergence). Therefore, the loss function \mathcal{L}_{KL} can be defined as,

$$\mathcal{L}_{KL}(y, \hat{y}) = \frac{1}{N} \sum_{n=1}^N \text{D}_{KL}(y_n | \hat{y}_n) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K y_n^k \log \left(\frac{y_n^k}{\hat{y}_n^k} \right) \quad (2)$$

where N denotes number of data samples, and $y_n = (y_n^1, \dots, y_n^K)$ and $\hat{y}_n = (\hat{y}_n^1, \dots, \hat{y}_n^K)$. However, not all predicted distributions need to follow the same variance σ , that is, the value of σ needs to be chosen adaptively for each utterance.

2.4. Selective Variance Label Distribution Learning

According to the principles discussed earlier, the age distribution of learning should vary by the utterance. To achieve this goal, we propose a novel selective variance label distribution learning method that fully adapts to the variance of each utterance. That is, the process of selecting the best matching distribution from a series of red candidate distributions, shown in Figure 1. These candidate Gaussian distributions can be defined as,

$$y_n^k(s) = C_n \cdot e^{-(k-\mu_n)^2/s} \quad (3)$$

Where $s \in S$, and S is a set of predefined candidate variance values. Among the candidate distributions obtained using these values, there should be one that matches the ground-truth age distribution as closely as possible. Therefore, the problem turns into choosing the smallest difference between a set of candidate label distributions $[\hat{y}]$ and the predicted distribution y . Denote that s^* is the variance of the best matching case, which is related to the ground-truth value of standard deviation with $s^* = \sigma_n^2$. Then the loss function \mathcal{L}_{KL} as follows,

$$\begin{aligned} \mathcal{L}_{KL}(y, \hat{y}) &= \frac{1}{N} \sum_{n=1}^N \left(\arg \min_{s \in S} (\text{D}_{KL}(y_n(s) | \hat{y}_n)) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \text{D}_{KL}(y_n(s^*) | \hat{y}_n) \end{aligned} \quad (4)$$

In this way, the algorithm can adaptively select the best matching variance of the Gaussian distribution for training.

2.5. Unimodal Distribution Constraints

In experiments, we observe that the baseline based on mean-variance learned distributions is multimodal for some instances, in Figure 2. Namely, there will be multiple peaks in the distribution. We propose an approach to overcome this issue by simultaneously learning the first-order differences of the distributions. Suppose $\Delta(\cdot)$ is the first-order difference function of a discrete distribution, in Equation 4 with variance s^* . This method is denoted as Diff. The loss function $\mathcal{L}_{\text{Diff}}$ used to optimize the first-order difference of the distribution is as follows,

$$\begin{aligned} \mathcal{L}_{\text{Diff}}(y, \hat{y}) &= \frac{1}{N} \sum_{n=1}^N (\Delta(y_n(s^*)) - \Delta(\hat{y}_n))^2 \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^{K-1} (\Delta(y_n^k(s^*)) - (\hat{y}_n^{k+1} - \hat{y}_n^k))^2 \end{aligned} \quad (5)$$

Proof: To demonstrate that Equation 5 constrains the distribution to be unimodal, assuming that the first difference of distribution $\Delta(y_n^k)$ is proportional to the first derivative of the distribution $y'_n = dy_n^k/dk$. Since y is a Gaussian distribution, when $k < \mu_n$, $y'_n > 0$, and when $k > \mu_n$, $y'_n < 0$. In order to show how our loss to be unimodal, we take a case of $k < \mu_n$ for illustration, where $\hat{y}_n^{k+1} - \hat{y}_n^k > 0$. If it is not a unimodal case at $\hat{y}_n^{k+1} - \hat{y}_n^k < 0$, to verify that the loss function $\mathcal{L}_{\text{Diff}}$ can constrain the distribution to a single mode, we calculate the gradient of this loss function over \hat{y}_n^k and \hat{y}_n^{k+1} respectively.

$$\frac{\partial \mathcal{L}_{\text{Diff}}}{\partial \hat{y}_n^k} \propto 2 \left(y'_n - (\hat{y}_n^{k+1} - \hat{y}_n^k) \right) > 0 \quad (6)$$

$$\frac{\partial \mathcal{L}_{\text{Diff}}}{\partial \hat{y}_n^{k+1}} \propto -2 \left(y'_n - (\hat{y}_n^{k+1} - \hat{y}_n^k) \right) < 0 \quad (7)$$

According to Equations 6 and 7, \hat{y}_n^k decreases due to its positive gradient and \hat{y}_n^{k+1} increases due to its negative gradient. In addition, the magnitude of this gradient is taken from the first-order difference of the Gaussian distribution, so the loss function can better constrain the distribution to the shape of the Gaussian distribution, thereby improving the quality of the fitted distribution.

2.6. Hybrid Loss

For regression predicting age, the Lin's Concordance Correlation Coefficient (CCC) [27] is more reliable to use, denoted as ρ_c . CCC is a measure of the agreement between ground-true labels and predicted labels. If the predicted value changes, the score is proportional to its deviation [31]. [32] provides complete proof that CCC outperforms other common regression losses, and we will use experiments to verify that it also holds for speech age estimation. Therefore, the loss function \mathcal{L}_c derived from ρ_c is used as a measure of regression age,

$$\mathcal{L}_c(\mu, \hat{\mu}) = 1 - \rho_c = 1 - \frac{2\sigma_{[pt]}^2}{\sigma_{[p]}^2 + \sigma_{[t]}^2 + (\mu_{[p]} - \mu_{[t]})^2} \quad (8)$$

Where $\sigma_{[pt]}^2 = \text{cov}(\mu, \hat{\mu})$, $\mu_{[p]} = \mathbb{E}(\hat{\mu})$, $\mu_{[t]} = \mathbb{E}(\mu)$, $\sigma_{[p]}^2 = \text{var}(\hat{\mu})$, $\sigma_{[t]}^2 = \text{var}(\mu)$, w.r.t. n .

Since human voice aging is a continuous process, the predicted age should be more likely to be the ground-truth age, and the farther away from this age, the less likely it is. Smaller variance means lower uncertainty in age prediction. The variance loss \mathcal{L}_v reduces the uncertainty in the estimated age distribution,

$$\mathcal{L}_v(y, \hat{y}) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \left(\hat{y}_n^k \cdot (k - \hat{\mu}_n)^2 \right) \quad (9)$$

Due to the physiological differences between men and women, there are distinct differences in speech characteristics – the average formant and fundamental frequency of women’s speech sounds higher than those of men [33, 34]. By using a multi-task learning approach while performing gender recognition tasks, gender information will be implicitly added to the model (blue task in Figure 1). In this task, the gender classification task is trained with a focal loss (FL) [35] with a tunable focus parameter $\gamma \geq 0$. The loss function \mathcal{L}_g is defined as follows,

$$\mathcal{L}_g(g, \hat{g}) = \text{FL}(g, \hat{g}) \quad (10)$$

Where \hat{g} and g are the predicted and the ground-truth gender. The overall loss is that given in Equation 11, where $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ and λ_5 are hyper-parameters.

$$\mathcal{L} = \lambda_1 \cdot \mathcal{L}_c + \lambda_2 \cdot \mathcal{L}_{\text{KL}} + \lambda_3 \cdot \mathcal{L}_v + \lambda_4 \cdot \mathcal{L}_{\text{Diff}} + \lambda_5 \cdot \mathcal{L}_g \quad (11)$$

2.7. Training and Inference

During the training phase, speech activity detection (SAD) preprocesses the audio to remove non-speech frames since speech may contain many silent segments. In our experiments, the rvad model [36] is used for this task. In order to make full use of hardware resources to train models quickly, model training can be divided into two stages.

Offline training: Since the inference speed of the WavLM model is not fast, first convert all the data in the dataset into speech features and save them in Numpy format, and then use these converted data directly to speed up training.

Online training: To improve the robustness of the model, we employ a chain-like augment: (1) Noise was added using MUSAN. (2) The RIR reverb is added. (3) Time stretch augment [37]: time stretching doesn’t change pitch, it simulates a person’s different speech rates.

During the inference phase, the age estimate of the utterance and its uncertainty are the mean age $\hat{\mu}$ and variance $\hat{\sigma}$ of the predicted distribution. At the same time, the auxiliary task of gender recognition will be abandoned. The speech will be processed by SAD first, and then the whole segment will be sent to the model for prediction.

3. EXPERIMENTS

3.1. Datasets

To demonstrate the advantages of the proposed method, we use the same dataset and conduct experimental validation under

the same settings as [3].

NIST SRE08-10 dataset: We use 11,205 utterances (458 male and 769 female speakers) from NIST SRE08 as the training set, and 5,331 telephone-conditioned utterances (236 male and 256 female speakers) from NIST SRE10 as the test set, similar to [38]. The speech in the dataset contains both English and non-English. Neither the speakers nor the recordings in the training and test set overlap. The speech in the dataset contains Chinese and Chinese dialects.

Real-world PA-Age Dataset: This dataset is from the financial insurance domain and contains 69,610 corpora (30,661 male and 28,386 female corpora). The test set used the same 4,000 utterances as [3]. The average duration of effective speech is 28.125 seconds, and the standard deviation of duration is 19.114 seconds.

3.2. Metrics

To evaluate how good our age estimator is, we report regression performance in terms of mean absolute error (MAE) and Pearson’s correlation coefficient (PCC) ρ . It is defined in Equation 12 and 13. The lower MAE and higher PCC, the better.

$$\text{MAE} = \frac{1}{N} \sum_{n=1}^N \|\hat{\mu}_n - \mu_n\|_1 \quad (12)$$

$$\rho = \frac{1}{N-1} \sum_{n=1}^N \left(\left(\frac{\hat{\mu}_n - \mu_{[p]}}{\sigma_{[p]}} \right) \left(\frac{\mu_n - \mu_{[t]}}{\sigma_{[t]}} \right) \right) \quad (13)$$

To measure whether the resulting distribution is unimodal, we introduce a unimodal coefficient η_q , representing the average number of modes within q standard deviations which can roughly detect the number of peaks in the predicted distribution, in Equation 14. The lower the better. In order to ensure the accuracy of the detection, it is only necessary to consider the age within q standard deviations ($q = 2$ in experiment).

$$\eta_q = \frac{1}{N} \sum_{n=1}^N \left(\sum_{k_{\min} < k < k_{\max}} \mathbb{1}(\text{cond}(n, k)) \right) \quad (14)$$

$$\text{cond}(n, k) = \left(\Delta(\hat{y}_n^k) < 0 \right) \wedge \left(\Delta(\hat{y}_n^{k+1}) > 0 \right) \quad (15)$$

$$\begin{aligned} k_{\min} &= \max(1, \hat{\mu}_n - q \cdot \hat{\sigma}_n) \\ k_{\max} &= \min(\hat{\mu}_n + q \cdot \hat{\sigma}_n, K - 1) \end{aligned} \quad (16)$$

Where $\mathbb{1}(\cdot)$ is a function that converts a boolean value to an integer. Equation 15 is the conditional function to detect peaks in the distribution. Equation 16 defines the age range to be calculated.

3.3. Hyper-parameters

For speech features, the WavLM model uses the “WavLM Base+” setting in our implementation, which has 13 transformer encoder layers [39], 768-dimensional hidden states, and 8 attention heads. The channel parameter C_b for ECAPA-TDNN in the convolutional layers for the proposed network

is 256. The bottleneck size in the SE-Block and Attention modules is set to 128. The scale size in Res2Block is set to 8. The tunable parameter γ in FL is 10. The maximum age label K is set to 100.

The model is first quickly trained by offline training and then fine-tuned by data augmentation using online training. During the offline training phase, the SGD [40] optimizer uses a momentum of 0.9 and weight decay of $1e-3$. The mini-batch is set to 64 and the initial learning rate is $2e-3$ to train all our models. The speech features are segmented into about 3 seconds (150 WavLM frames) to avoid over-fitting and to speed up training. The set of standard deviation candidates S are: from 0.01 to 10 in steps of 0.1, i.e., $S = \{0.1, 0.2, \dots, 10\}^2$. In the fine-tuning stage with online training mode, the speech segment changed to 6 seconds. Due to hardware resource constraints, we adopted a smaller fine-tuning learning rate of $1e-5$, the weight decay becomes $4e-4$, and a batch size of 128. The values of standard deviation candidates S have become finer: from 0.01 to 10 in steps of 0.01, i.e., $S = \{0.01, 0.02, \dots, 10\}^2$.

3.4. Implementation Details

To verify the effectiveness of our proposed method, we compare it with our last results as a baseline method (results directly copied from [3]). The best-performing model uses the ResNet-18 model as the backbone. The age label distribution is optimized with the mean and variance and the KL divergence of the distribution (denoted as MVKL in Table 1).

To demonstrate that the model based on WavLM combined with ECAPA-TDNN outperforms the previous backbone model, we only replace this part and conduct experiments with the same parameters and methods as before based on optimizing LDL and MVKL (with a fixed variance $\sigma = 1$). To show that the SVLDL and CCC are of great help to the age estimation, let $\lambda_1 = 10$, $\lambda_2 = 1$, $\lambda_3 = 0.1$ and $\lambda_4 = \lambda_5 = 0$ for experiments (denoted as CVKL in Table 1). To eliminate the multimodality of the age distribution, we give a weight to the loss function of the KL distance of the first-order difference distribution to eliminate this effect, let $\lambda_1, \lambda_2, \lambda_3, \lambda_5$ are the same, $\lambda_4 = 0.1$. In addition, by adding the auxiliary task of gender recognition, the model can improve the discrimination of gender, thereby improving the performance of the age recognition task, let $\lambda_5 = 0.01$, rest are the same.

3.5. Evaluation Results

Ablation study results on the two datasets are reported in Table 1. When the backbone model is replaced by WavLM+ECAPA-TDNN, the backend model is still LDL+MVKL, and the performance of model regression is slightly improved. If the backend model is replaced by CVKL, the performance of regression will be greatly improved. If SVLDL is replaced, the multimodal problem can be solved to a certain extent. This shows that under SVLDL, different variances are adaptively assigned to each utterance. This results in the models not forcing

Table 1. Ablation study on our proposed system compared with baseline model. + denotes stacking our methods.

	SRE08-10			PA-Age		
	MAE	ρ	η_2	MAE	ρ	η_2
<i>ResNet-18</i> [3]						
+LDL+MVKL	4.62	0.87		6.23	0.82	
<i>WavLM+ECAPA-TDNN (ours)</i>						
+LDL+MVKL	4.48	0.85	1.68	6.17	0.84	1.70
+LDL+CVKL	4.19	0.90	1.59	5.95	0.85	1.67
+SVLDL+MVKL	4.50	0.84	1.32	6.15	0.83	1.41
+SVLDL+CVKL	4.16	0.93	1.34	5.91	0.86	1.43
++Diff	4.19	0.91	1.08	5.93	0.85	1.07
+++Gender	4.14	0.92	1.03	5.82	0.87	1.05

them to optimize with the same variance, which also reduces multimodal problems. Meanwhile, under CVKL, the model uses the CCC loss function, which can be directly optimized for the PCC (because the CCC is the unbiased PCC [32]), and the MAE is further reduced. When the model simultaneously optimizes the L2 distance of the first-order difference of the age distribution, although the regression performance drops slightly, the multimodality problem is greatly solved. When the auxiliary task of gender recognition is added, the model can distinguish gender information, thereby improving the age estimation performance. Meanwhile, the multimodal problem is further solved. Compared to the baseline model, the model achieves an MAE reduction of 10.39% on NIST SRE08-10 and 6.58% on the PA-Age dataset, outperforming the original model and close to the state-of-the-art model results.

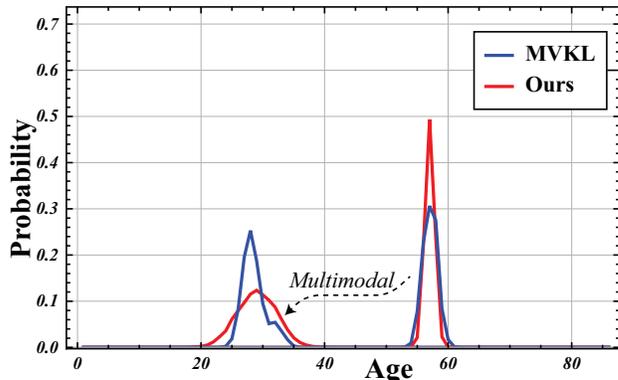


Fig. 2. Distributions predicted by “LDL + MVKL” and “SVLDL + CVKL + Diff + Gender” (ours). The backbone model uses “WavLM + ECAPA-TDNN”. Example from PA-Age dataset.

As seen from Figure 2, our prediction is optimized to be unimodal and adaptively learned based on a specific instance, and has higher prediction performance. That is, SVLDL adaptively selects variance for each instance to predict, Diff ensures

unimodal distribution, and gender recognition improves prediction performance. In contrast, MVKL is optimized for all instances with the same variance and cannot guarantee a unimodal distribution.

Table 2. Ablation study of hyper-parameters for $\lambda_{(2-5)}$ on PA-Age. The first row is when $\lambda_4 = \lambda_5 = 0$, and the second row is when $\lambda_2 = 1$ and $\lambda_3 = 0.1$.

	MAE	η_2	MAE	η_2	MAE	η_2	MAE	η_2	MAE	η_2
$\lambda_{(2,3)}$	(0.1, 0.01)		(0.1, 0.1)		(1, 0.1)		(1, 1)		(10, 1)	
	6.05	1.63	5.97	1.49	5.91	1.43	6.12	1.54	6.47	1.50
$\lambda_{(4,5)}$	(0.01, 0)		(0.1, 0)		(1, 0)		(0.1, 0.01)		(0.1, 0.1)	
	5.98	1.16	5.93	1.07	6.02	1.07	5.82	1.05	5.88	1.06

Table 2 shows the full path in the search for hyperparameter values of $\lambda(s)$ for the loss function, in Equation 11. The first row shows the effect of adjusting λ_2 and λ_3 on predicted age (in MAE) and the multimodal parameter η_2 under $\lambda_4 = \lambda_5 = 0$. When $\lambda_2 = 0.1$, $\lambda_3 = 0.01$, the variance of the age distribution may be too large, resulting in more multimodal problems. Meanwhile, the prediction of age is not very accurate. When λ_3 increases to 0.1, the variance becomes smaller, the age prediction becomes more accurate, and the multimodal problem is slightly solved. The cases of $\lambda_2 = 1$ and $\lambda_2 = 0.1$ are the best-performing combination of λ_2 and λ_3 . But if these two hyperparameters are too large, it will affect the performance of age estimation, and the multimodal problem will appear again.

The second row in Table 2 is an experiment on λ_4 and λ_5 based on the best results from the first row (when $\lambda_2 = 1$ and $\lambda_2 = 0.1$). When $\lambda_5 = 0$, the larger λ_4 is, the better the multimodal problem can be solved. However, this case slightly degrades the age estimation performance. Thus, we choose the best-performing case of $\lambda_4 = 0.1$, and adjust the value of λ_5 . After adding the auxiliary task of gender recognition, the performance of age estimation is further improved and achieves state-of-the-art results. At the same time, the multimodal problem is further solved. Therefore, the optimal hyperparameter combination is $\lambda_1 = 10$, $\lambda_2 = 1$, $\lambda_3 = \lambda_4 = 0.1$ and $\lambda_5 = 0.01$.

Table 3. Effect of duration of test utterance on PA-age estimated by our method compared to the best baseline results.

MAE Model	Test segment lengths(s)			
	10	15	20	full
ResNet-18+LDL+MVKL [3]	13.16	11.04	6.14	6.10
Proposed method (ours)	8.35	6.42	5.86	5.82

Table 3 compares the effects of different test utterance durations on age estimation between the proposed method ‘‘WavLM + ECPAP-TDNN + SVLDL + CVKL + Diff + Gender’’ and the baseline model. Here, only the utterances longer than 10 seconds are cut and selected as test data. Compared

to the baseline model, the proposed method achieves great improvement on short test utterances. It may be because the WavLM speech features are trained based on big data, and the data augmentation method is used during training, which improves the robustness of the model. At the same time, the pooling layer has an attention mechanism, which will make the model pay more attention to the information of speech and reduce over-fitting.

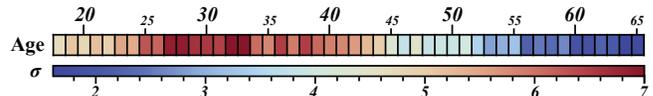


Fig. 3. The heat map to visualize the adaptively learned variances σ corresponding to different ages.

Figure 3 shows the median values of variance for adaptation at different ages. It can be seen that before the age of 27, the voice is in a young and stable state. From the age of 30 to 45, people’s vocal cords gradually age, and the degree of aging varies from person to person, so the uncertainty is large. After the age of 50, almost everyone’s voice becomes older, the variance becomes smaller, and the vocal characteristics become more recognizable.

4. CONCLUSIONS

In this paper, a selective variance labeled distribution learning (SVLDL) method is proposed to accommodate variances of different age distributions. The robustness of the age regression is enhanced by using Lin’s consistent correlation coefficient (CCC) loss compared to the mean-variance-based loss. Since existing methods suffer from multimodality in the age distribution, the quality of the fitted distribution is improved here by optimizing the L2 distance of the first-order difference of the distribution, and a reasonable proof is given. The performance is further improved by using the speech features of WavLM and adding the auxiliary task of gender recognition. Experiments show that the model achieves MAE reduction and multimodal problem-solving on both the NIST SRE08-10 and real-world PA-Age datasets, outperforming the original model to achieve state-of-the-art results.

5. ACKNOWLEDGEMENT

This paper is supported by the Key Research and Development Program of Guangdong Province under grant No.2021B0101400003. Corresponding author is Jianzong Wang from Ping An Technology (Shenzhen) Co., Ltd (jzwang@188.com).

6. REFERENCES

- [1] Shijing Si, Jianzong Wang, Xiaoyang Qu, Ning Cheng, Wenqi Wei, Xinghua Zhu, and Jing Xiao, “Speech2video: Cross-modal distillation for speech to video generation,” in *Proceedings of the Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2021.
- [2] Zuheng Kang, Junqing Peng, Jianzong Wang, and Jing Xiao, “Speecheq: Speech emotion recognition based on multi-scale unified datasets and multitask learning,” in *IEEE Conference of the International Speech Communication Association (INTERSPEECH)*, 2022.
- [3] Shijing Si, Jianzong Wang, Junqing Peng, and Jing Xiao, “Towards speaker age estimation with label distribution learning,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 4618–4622.
- [4] Mark Huckvale and Aimee Webb, “A comparison of human and machine estimation of speaker age,” in *International Conference on Statistical Language and Speech Processing*. Springer, 2015, pp. 111–122.
- [5] Prashanth Gurunath Shivakumar, Somer Bishop, Catherine Lord, and Shrikanth Narayanan, “Phone duration modeling for speaker age estimation in children,” *arXiv preprint arXiv:2109.01568*, 2021.
- [6] Benjamin V Tucker, Catherine Ford, and Stephanie Hedges, “Speech aging: Production and perception,” *Wiley Interdisciplinary Reviews: Cognitive Science*, vol. 12, no. 5, pp. e1557, 2021.
- [7] Yuki Kitagishi, Hosana Kamiyama, Atsushi Ando, Naohiro Tawara, Takeshi Mori, and Satoshi Kobashikawa, “Speaker age estimation using age-dependent insensitive loss,” in *2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2020, pp. 319–324.
- [8] Shareef Babu Kalluri, Deepu Vijayasenan, and Sriram Ganapathy, “A deep neural network based end to end model for joint height and age estimation from short duration speech,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6580–6584.
- [9] Mohamad Hasan Bahari, Mitchell McLaren, David A van Leeuwen, et al., “Speaker age estimation using i-vectors,” *Engineering Applications of Artificial Intelligence*, vol. 34, pp. 99–108, 2014.
- [10] Pegah Ghahremani, Phani Sankar Nidadavolu, Nanxin Chen, Jesús Villalba, Daniel Povey, Sanjeev Khudanpur, and Najim Dehak, “End-to-end deep neural network age estimation,” in *IEEE Conference of the International Speech Communication Association (INTERSPEECH)*, 2018, pp. 277–281.
- [11] Damian Kwaśny, Paweł Jemioło, and Daria Hemmerling, “Explaining predictions of the x-vector speaker age and gender classifier,” in *International Conference on Dependability and Complex Systems*. Springer, 2021, pp. 234–243.
- [12] Manav Kaushik, Tran The Anh, Eng Siong Chng, et al., “End-to-end speaker age and height estimation using attention mechanism and triplet loss,” in *2021 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2021, pp. 1–8.
- [13] Florian Eyben, Martin Wöllmer, and Björn Schuller, “Opensmile: the munich versatile and fast open-source audio feature extractor,” in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 1459–1462.
- [14] Ftoon Abu Shaqra, Rehab Duwairi, and Mahmoud Al-Ayyoub, “Recognizing emotion from speech based on age and gender using hierarchical models,” *Procedia Computer Science*, vol. 151, pp. 37–44, 2019.
- [15] Felix Burkhardt, Markus Brückl, and Björn W Schuller, “Age classification: Comparison of human vs machine performance in prompted and spontaneous speech,” *Studientexte zur Sprachkommunikation: Elektronische Sprachsignalverarbeitung*, vol. 2021, pp. 35–42, 2021.
- [16] Yilin Pan, Venkata Srikanth Nallanthighal, Daniel Blackburn, Heidi Christensen, and Aki Härmä, “Multi-task estimation of age and cognitive decline from speech,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 7258–7262.
- [17] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 12449–12460, 2020.
- [18] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE Journal of Selected Topics in Signal Processing*, 2022.
- [19] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song, “Using self-supervised learning can improve

- model robustness and uncertainty,” *Advances in neural information processing systems*, vol. 32, 2019.
- [20] Tarun Gupta, Duc-Tuan Truong, Tran The Anh, and Chng Eng Siong, “Estimation of speaker age and height from speech signal using bi-encoder transformer mixture model,” *arXiv preprint arXiv:2203.11774*, 2022.
- [21] Xin Geng, “Label distribution learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 7, pp. 1734–1748, 2016.
- [22] Huiying Zhang, Yu Zhang, and Xin Geng, “Practical age estimation using deep label distribution learning,” *Frontiers of Computer Science*, vol. 15, no. 3, pp. 1–6, 2021.
- [23] Xin Geng, Qin Wang, and Yu Xia, “Facial age estimation by adaptive label distribution learning,” in *22nd International Conference on Pattern Recognition*. IEEE, 2014, pp. 4465–4470.
- [24] Xin Wen, Biying Li, Haiyun Guo, Zhiwei Liu, Guosheng Hu, Ming Tang, and Jinqiao Wang, “Adaptive variance based label distribution learning for facial age estimation,” in *European Conference on Computer Vision*. Springer, 2020, pp. 379–395.
- [25] Qiang Li, Jingjing Wang, Zhaoliang Yao, Yachun Li, Pengju Yang, Jingwei Yan, Chunmao Wang, and Shiliang Pu, “Unimodal-concentrated loss: Fully adaptive label distribution learning for ordinal regression,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 20513–20522.
- [26] Hieu Trung Huynh and Hoang Nguyen, “Joint age estimation and gender classification of asian faces using wide resnet,” *SN computer science*, vol. 1, no. 5, pp. 1–9, 2020.
- [27] I Lawrence and Kuei Lin, “A concordance correlation coefficient to evaluate reproducibility,” *Biometrics*, pp. 255–268, 1989.
- [28] Brecht Desplanques, Jenthe Thienpondt, and Kris Demuyne, “Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification,” *arXiv preprint arXiv:2005.07143*, 2020.
- [29] Shang-Hua Gao, Ming-Ming Cheng, Kai Zhao, Xin-Yu Zhang, Ming-Hsuan Yang, and Philip Torr, “Res2net: A new multi-scale backbone architecture,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 2, pp. 652–662, 2019.
- [30] Jie Hu, Li Shen, and Gang Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [31] Bagus Tris Atmaja and Masato Akagi, “Evaluation of error-and correlation-based loss functions for multitask learning dimensional speech emotion recognition,” in *Journal of Physics: Conference Series*. IOP Publishing, 2021, vol. 1896, p. 012004.
- [32] Vedhas Pandit and Björn Schuller, “The many-to-many mapping between the concordance correlation coefficient and the mean square error,” *arXiv preprint arXiv:1902.05180*, 2019.
- [33] Ke Wu and Donald G Childers, “Gender recognition from speech. part i: Coarse analysis,” *The journal of the Acoustical society of America*, vol. 90, no. 4, pp. 1828–1840, 1991.
- [34] Elvira Mendoza, Nieves Valencia, Juana Muñoz, and Humberto Trujillo, “Differences in voice quality between men and women: Use of the long-term average spectrum (Itas),” *Journal of voice*, vol. 10, no. 1, pp. 59–66, 1996.
- [35] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [36] Zheng-Hua Tan, Najim Dehak, et al., “rvad: An unsupervised segment-based robust voice activity detection method,” *Computer speech & language*, vol. 59, pp. 1–21, 2020.
- [37] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” *arXiv preprint arXiv:1904.08779*, 2019.
- [38] Seyed Omid Sadjadi, Sriram Ganapathy, and Jason W Pelecanos, “Speaker age estimation on conversational telephone speech using senone posterior based i-vectors,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 5040–5044.
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [40] Léon Bottou, “Stochastic gradient descent tricks,” in *Neural networks: Tricks of the trade*, pp. 421–436. Springer, 2012.