

# COVARIATE SHIFT ADAPTATION FOR SEMI-SUPERVISED SPEAKER IDENTIFICATION

<sup>1</sup>Makoto Yamada, <sup>1</sup>Masashi Sugiyama, and <sup>2</sup>Tomoko Matsui

<sup>1</sup>Department of Computer Science, Tokyo Institute of Technology

<sup>2</sup>Department of Statistical Modeling, The Institute of Statistical Mathematics

e-mail: yamada@sg.cs.titech.ac.jp sugi@cs.titech.ac.jp tmatsui@ism.ac.jp

## ABSTRACT

In this paper, we propose a novel semi-supervised speaker identification method that can alleviate the influence of non-stationarity such as session dependent variation, the recording environment change, and physical condition/emotion. We assume that the utterance variation follows the *covariate shift* model, where only the utterance sample distribution changes in the training and test phases. Our method consists of weighted versions of kernel logistic regression and cross-validation and is theoretically shown to have the capability of alleviating the influence of covariate shift. We experimentally show through text-independent speaker identification simulations that the proposed method is promising in dealing with variations in session dependent utterance variation.

**Index Terms**— Speaker identification, covariate shift, semi-supervised learning, kernel logistic regression, importance estimation.

## 1. INTRODUCTION

Speaker identification methods are widely used in various real-world situations such as access control of information service systems and speaker detection in speech dialog and speaker indexing problems with large audio archives. Recently, the speaker identification and indexing problems in meetings attracted a great deal of attention.

Standard methods of text-independent speaker identification includes the *Gaussian mixture model (GMM)* [1] or kernel methods such as the *support vector machine (SVM)* [2]. In these supervised learning methods, it is implicitly assumed that training and test data follow the same distribution. However, the training and test distributions are not necessarily the same in practice since the utterance features vary over time due to session dependent variation, the recording environment change, and physical condition/emotion.

To alleviate the influence of session dependent variation, it is common to use speech samples recorded in several different sessions [3]. However, gathering many speech samples and labeling the speaker ID to the collected data are expensive both in time and cost and therefore not realistic in practice.

A more practical setup would be *semi-supervised learning*, where unlabeled samples are additionally given from the test environment. In semi-supervised learning, it is required that the training and test distributions are related to each other in some sense; otherwise we may not be able to learn anything about the test distribution from the training samples. A popular modeling is called *covariate shift* [4], where the input distributions are different in the training and test phases but the conditional distribution of labels remains unchanged.

In this paper, we formulate the semi-supervised speaker identification problem in the covariate shift framework and propose a method that can cope with utterance variation. Under covariate shift, standard maximum likelihood estimation is no longer consistent—the influence of covariate shift can be asymptotically canceled by weighting the log-likelihood terms according to the *importance* [4]:  $w(X) = p_{te}(X)/p_{tr}(X)$ , where  $p_{te}(X)$  and  $p_{tr}(X)$  are test and training input densities. The importance weight  $w(X)$  is unknown in practice and needs to be estimated from data. For weight estimation, we utilize the *Kullback-Leibler importance estimation procedure (KLIEP)* [5] due to its superior performance. The regularized kernel logistic regression model contains two tuning parameters: the kernel width and the regularization parameter [3]. Usually these tuning parameters are optimized based on *cross validation (CV)*. However, CV is no longer unbiased due to covariate shift and therefore is not reliable as a model selection method. To cope with this problem, we use an importance-weighted version of CV (IWCV) [6] for unbiased model selection. The validity of our approach is experimentally shown through text-independent speaker identification simulations.

## 2. PROBLEM FORMULATION

In this section, we formulate the speaker identification problem based on the kernel logistic regression (KLR) model.

**Text-independent Speaker Identification:** An utterance sample  $X$  pronounced by a speaker is expressed as a set of  $N$  *mel-frequency cepstrum coefficient (MFCC)* [7] vectors of

$d$  dimensions:

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{d \times N}.$$

For training, we are given  $n$  labeled utterance samples

$$\mathcal{Z} = \{(X_i, y_i)\}_{i=1}^n,$$

where  $y_i \in \{1, \dots, K\}$  denotes the index of the speaker who pronounced  $X_i$ . The goal of speaker identification is to predict the speaker index of a test utterance sample  $\mathbf{X}$  based on the training samples. We predict the speaker index  $c$  of the test sample  $\mathbf{X}$  following the *Bayes decision rule*:

$$\max_c p(y = c | \mathbf{X}).$$

For approximating the class-posterior probability, we use

$$p(y = c | \mathbf{X}; \mathbf{V}) = \frac{\exp f_{v_c}(\mathbf{X})}{\sum_{l=1}^K \exp f_{v_l}(\mathbf{X})},$$

where  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_K]^\top \in \mathbb{R}^{K \times n}$  is the parameter,  $^\top$  denotes the transpose, and  $f_{v_l}$  is a discriminant function corresponding to speaker  $l$ . This form is known as the *softmax* function and widely used in multiclass logistic regression. We use the following kernel regression model as the discriminant function  $f_{v_l}$ :

$$f_{v_l}(\mathbf{X}) = \sum_{i=1}^n v_{l,i} \mathcal{K}(\mathbf{X}, \mathbf{X}_i) \quad l = 1, \dots, K,$$

where  $\mathbf{v}_l = (v_{l,1}, \dots, v_{l,n})^\top \in \mathbb{R}^n$  are parameters corresponding to speaker  $l$  and  $\mathcal{K}(\mathbf{X}, \mathbf{X}')$  is a kernel function. In this paper, we use the *sequence kernel* [2] as the kernel function since it allows us to handle features with different size; for two utterance samples  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{d \times N}$  and  $\mathbf{X}' = [\mathbf{x}'_1, \dots, \mathbf{x}'_{N'}] \in \mathbb{R}^{d \times N'}$  (generally  $N \neq N'$ ), the sequence kernel is defined as

$$\mathcal{K}(\mathbf{X}, \mathbf{X}') = \frac{1}{NN'} \sum_{i=1}^N \sum_{i'=1}^{N'} \exp \left( \frac{-\|\mathbf{x}_i - \mathbf{x}'_{i'}\|^2}{2\sigma^2} \right).$$

Note that kernel logistic regression is a modeling assumption; thus the true class-conditional probability may not be exactly realized by the kernel logistic regression model. This implies that there exists some *model error*, i.e., even when the parameter is chosen optimally, there remains an approximation error. This setup is not of course preferable, but more or less there exists a model error in practice since it is not generally possible to prepare an exactly correct model. Traditional machine learning theories often assume that the model at hand is correct (i.e., no model error exists). However, this is not realistic and not useful in practice, so in this paper we

explicitly take into account *model misspecification* within the covariate shift framework.

**Kernel Logistic Regression [3, 8]:** We employ maximum likelihood estimation for learning the parameter  $\mathbf{V}$ . The negative regularized log-likelihood function  $\mathcal{P}_\delta^{\log}(\mathbf{V}; \mathcal{Z})$  for the kernel logistic regression model is given by

$$\mathcal{P}_\delta^{\log}(\mathbf{V}; \mathcal{Z}) = - \sum_{i=1}^n \log P(y_i | \mathbf{X}_i; \mathbf{V}) + \frac{\delta}{2} \text{tr}(\mathbf{V} \mathbf{K} \mathbf{V}^\top),$$

where  $\frac{\delta}{2} \text{tr}(\mathbf{V} \mathbf{K} \mathbf{V}^\top)$  is a regularizer introduced for avoiding overfitting and  $\mathbf{K} = [\mathcal{K}(\mathbf{X}_i, \mathbf{X}_j)]_{i,j=1}^n$  is the kernel Gram matrix.  $\mathcal{P}_\delta^{\log}(\mathbf{V}; \mathcal{Z})$  is a convex function with respect to  $\mathbf{V}$  and therefore its unique minimizer can be obtained by, e.g., the Newton method.

**Model Selection in KLR:** KLR includes two tuning parameters—the Gaussian width  $\sigma$  and the regularization parameter  $\delta$ . One of the popular approaches to model selection is *cross validation (CV)*.

Let us divide the training set  $\mathcal{Z} = \{(X_i, y_i)\}_{i=1}^n$  into  $k$  disjoint non-empty subsets  $\{\mathcal{Z}_i\}_{i=1}^k$  of (approximately) the same size. Let  $\hat{y}(\mathbf{X}; \mathcal{Z}_j)$  be an estimate of a speaker of a test utterance sample  $\mathbf{X}$  obtained from  $\{\mathcal{Z}_i\}_{i \neq j}$  (i.e., without  $\mathcal{Z}_j$ ). Then the  $k$ -fold CV (kCV) score is given by

$$\hat{R}_{kCV}^{\mathcal{Z}} = \frac{1}{k} \sum_{j=1}^k \frac{1}{|\mathcal{Z}_j|} \sum_{(\mathbf{x}, y) \in \mathcal{Z}_j} I(y = \hat{y}(\mathbf{X}; \mathcal{Z}_j)),$$

where  $|\mathcal{Z}_j|$  denotes the number of samples in the subset  $\mathcal{Z}_j$  and  $I(\cdot)$  denotes the indicator function.

**KLR, CV, and Covariate Shift:** The use of KLR and CV could be theoretically justified when the training utterance features and the test utterance features independently follow the *same* probability distribution with density  $p(\mathbf{X})$  and the class label  $y$  follows the *common* conditional probability distribution  $p(y | \mathbf{X})$  in the training and test phases. Indeed, if these conditions are fulfilled, KLR is shown to be *consistent*, i.e., the learned parameter converges to the optimal value:

$$\lim_{n \rightarrow \infty} \hat{\mathbf{V}} = \mathbf{V}^*,$$

where  $\hat{\mathbf{V}}$  is the parameter learned by KLR and  $\mathbf{V}^*$  is the optimal parameter that minimizes the expected prediction error for test samples:

$$\mathbf{V}^* = \arg \min_{\mathbf{V}} \iint I(y = \hat{y}(\mathbf{X}; \mathbf{V})) p(y | \mathbf{X}) p(\mathbf{X}) d\mathbf{y} d\mathbf{X}.$$

$\hat{y}(\mathbf{X}; \mathbf{V})$  is an estimate of speaker of an utterance feature  $\mathbf{X}$  for parameter  $\mathbf{V}$ . Also, when  $p(\mathbf{X})$  and  $p(y | \mathbf{X})$  are common in the training and test phases, kCV is (almost) *unbiased*:

$$\mathbb{E}_{\mathcal{Z}} [\hat{R}_{kCV}^{\mathcal{Z}} - R^{\mathcal{Z}}] \approx 0,$$

where  $E_{\mathcal{Z}}$  is the expectation over the training set  $\mathcal{Z}$  and  $R^{\mathcal{Z}}$  is the expected prediction error defined by

$$R^{\mathcal{Z}} = \iint I(y = \hat{y}(X; \mathcal{Z}))p(y|X)p(X)dYdX.$$

However, in practical speaker identification, speech features are not stationary due to utterance variation, the recording environment change, and speaker feeling. Thus, the training and test feature distributions are not necessarily the same. Then the above good theoretical properties are no longer true.

If the training and test feature distributions share nothing in common, we may not be able to learn anything about the test distribution from the training samples. In this paper, we explicitly deal with such changing environment via the *covariate shift* model [4]—the input distributions change between the training and test phases,  $p_{tr}(X) \neq p_{te}(X)$ , but the conditional distribution  $p(y|X)$  remains unchanged.

### 3. IMPORTANCE WEIGHTING TECHNIQUES FOR COVARIATE SHIFT ADAPTATION

In this section, we show how to cope with covariate shift.

**Importance Sampling:** In the absence of covariate shift, the expectation over test samples can be computed by the expectation over training samples since they are drawn from the same distribution. However, under covariate shift, the difference of input distributions should be explicitly taken into account. A basic technique for compensating for the distribution change is *importance sampling*, i.e., the expectation over training samples is weighted according to their importance in the test distribution. Indeed, based on the importance weight

$$w(X) = \frac{p_{te}(X)}{p_{tr}(X)},$$

the expectation of some function  $F(X)$  over the probability density  $p_{te}(X)$  can be computed by

$$E_{p_{te}(X)}[F(X)] = E_{p_{tr}(X)}[F(X)w(X)].$$

**Importance Weighted Kernel Logistic Regression:** If the importance sampling technique is applied in KLR, we have the following *importance weighted KLR (IWKLR)*:

$$\tilde{\mathcal{P}}_{\delta}^{\log}(V; \mathcal{Z}) = - \sum_{i=1}^n w(X_i) \log P(y_i | X_i; V) + \frac{\delta}{2} \text{tr}(VKV^{\top}).$$

IWKLR is consistent even under covariate shift [4]:

$$\lim_{n \rightarrow \infty} \tilde{V} = V^*,$$

where  $\tilde{V}$  is the parameter learned by IWKLR and  $V^*$  is the optimal parameter given by

$$V^* = \underset{V}{\operatorname{argmin}} \iint I(y = \hat{y}(X; V))p(y|X)p_{te}(X)dYdX.$$

Note that  $\tilde{\mathcal{P}}_{\delta}^{\log}(V; \mathcal{Z})$  is still convex and thus the global solution can be obtained by the Newton method.

**Importance Weighted Cross Validation:** IWKLR includes the Gaussian width  $\sigma$  and the regularization parameter  $\delta$  as tuning parameters. Here, we introduce *important weighted cross validation (IWCV)* [6] for model selection: the  $k$ -fold IWCV (kIWCV) score is given by

$$\tilde{R}_{kIWCV}^{\mathcal{Z}} = \frac{1}{k} \sum_{j=1}^k \frac{1}{|\mathcal{Z}_j|} \sum_{(X,y) \in \mathcal{Z}_j} w(X)I(y = \hat{y}(X; \mathcal{Z}_j)).$$

Even under covariate shift, kIWCV is almost unbiased [6]:

$$E_{\mathcal{Z}} [\tilde{R}_{kIWCV}^{\mathcal{Z}} - R^{\mathcal{Z}}] \approx 0,$$

where  $R^{\mathcal{Z}}$  is the expected prediction error defined by

$$R^{\mathcal{Z}} = \iint I(y = \hat{y}(X; \mathcal{Z}))p(y|X)p_{te}(X)dYdX.$$

**Importance Weight Estimation:** As shown above, the importance weight  $w(X)$  plays a central role in covariate shift adaptation. However, the importance weight is usually unknown, so it needs to be estimated from samples. Here, we assume that in addition to the training input samples  $\mathcal{X}^{tr} = \{X_i^{tr}\}_{i=1}^{n_{tr}}$  drawn independently from  $p_{tr}(X)$ , we are given unlabeled test samples  $\mathcal{X}^{te} = \{X_i^{te}\}_{i=1}^{n_{te}}$  drawn independently from  $p_{te}(X)$  (i.e., the semi-supervised setup).

Under the semi-supervised setup, the importance weight may be simply estimated by estimating  $p_{tr}(X)$  and  $p_{te}(X)$  from training and test samples and then taking their ratio. However, density estimation is known to be a hard problem and taking the ratio of estimated quantities tends to magnify the estimation error. Thus such a two-shot process may not be reliable in practice. Below, we introduce a method called the *Kullback Leibler Importance Estimation Procedure (KLIEP)* [5], which allows us to directly learn the importance weight function without going through density estimation.

Let us model the importance function  $w(X)$  by the following linear model:

$$\hat{w}(X) = \sum_{l=1}^b \alpha_l \varphi(X, C_l),$$

where  $\{\alpha_l\}_{l=1}^b$  are parameters to be learned from data samples,  $\{C_l\}_{l=1}^b$  are template points randomly chosen from the test input set  $\{X_i^{te}\}_{i=1}^{n_{te}}$ , and  $\varphi(X, X')$  is a basis function chosen as

$$\varphi(X, X') = \frac{1}{NN'} \sum_{i=1}^N \sum_{i'=1}^{N'} \exp \left( \frac{-\|\mathbf{x}_i - \mathbf{x}_{i'}\|^2}{2\tau^2} \right).$$

We determine the coefficient  $\{\alpha_l\}_{l=1}^b$  by maximum likelihood estimation, which is formulated as

$$\begin{aligned} \max_{\{\alpha_l\}_{l=1}^b} & \left[ \sum_{i=1}^{n_{te}} \log \left( \sum_{l=1}^b \alpha_l \varphi(X_i^{te}, C_l) \right) \right] \\ \text{s.t.} & \sum_{i=1}^{n_{tr}} \sum_{l=1}^b \alpha_l \varphi(X_i^{tr}, C_l) = n_{tr} \quad \text{and} \quad \alpha_1, \dots, \alpha_b \geq 0. \end{aligned}$$

This optimization problem is convex and thus the global solution may be obtained by simply performing gradient ascent and feasibility satisfaction iteratively. Note that the solution  $\{\hat{\alpha}_l\}_{l=1}^b$  tends to be sparse, which contributes to reducing the computational cost in the test phase.

#### 4. EXPERIMENTS

In this section, we report the results of speaker identification in the light of covariate shift adaptation.

Training and test samples were collected from 10 male speakers. The pronounced sentences are common to all speakers, but the test sentences are different from those for training. Moreover, the utterance samples for training were recorded in 1990/12, while the utterance samples for testing were recorded in 1991/3, 1991/6, and 1991/9, respectively. Since the recording time is different between training and test utterance samples, the session dependent variation is expected to be included. So this would be a challenging task. We used three sentences for training and five sentences for the test, where the average duration of the sentences is about 4[s].

The input utterance is sampled at 16kHz with close-talking microphone. A feature vector consists of 26 components: 12 MFCCs, the normalized log energy, and their first derivatives. Feature vectors are derived at every 10[ms] over the Hamming-windowed speech segment of 25.6[ms]. We divide each utterance sequence into 300[ms] disjoint segments, each of which corresponds to a set of features of size  $X_i \in \mathbb{R}^{26 \times 30}$ . We compute the speaker identification rate at every 1.5[s] and judge the speaker ID at time  $t$  based on the average posterior probability  $\frac{1}{5} \sum_{i=1}^5 p(Y_{t-i} | X_{t-i}; V)$ .

We compare KLR and IWKLR in terms of speaker identification for 1991/3, 1991/6, and 1991/9 [9]. For KLR training, we only use the 1990/12 dataset (inputs  $\mathcal{X}^{tr}$  and their labels), where the Gaussian width  $\sigma$  and the regularization parameter  $\delta$  are selected based on 5-fold CV. For IWKLR training, we use unlabeled samples  $\mathcal{X}^{te1}$ ,  $\mathcal{X}^{te2}$ , and  $\mathcal{X}^{te3}$  in addition to the training inputs  $\mathcal{X}^{tr}$  and their labels. We first estimate the importance weight from training and test data pairs  $(\mathcal{X}^{tr}, \mathcal{X}^{te1})$ ,  $(\mathcal{X}^{tr}, \mathcal{X}^{te2})$ , or  $(\mathcal{X}^{tr}, \mathcal{X}^{te3})$  by KLIEP, and then use 5-fold IWCV to decide the Gaussian width  $\sigma$  and regularization parameter  $\delta$ .

Table 1 summarizes the speaker identification rates, showing that IWKLR+IWCV outperforms KLR+CV for all sessions. This result implies that importance weighting is useful in coping with the influence of non-stationarity in practical

**Table 1.** Identification rates in percent. IWKLR+IWCV refers to IWKLR with  $\sigma$  and  $\delta$  chosen by 5-fold IWCV, and KLR+CV refers to KLR with  $\sigma$  and  $\delta$  chosen by 5-fold CV. Values of chosen  $\sigma$  and  $\delta$  are described in the bracket.

Test date	IWKLR+IWCV	KLR+CV
1991/3	<b>86.8</b> (1.2, 0.0001)	86.1 (1.2, 0.0001)
1991/6	<b>83.9</b> (1.3, 0.0001)	82.0 (1.2, 0.0001)
1991/9	<b>92.0</b> (1.2, 0.0001)	91.6 (1.2, 0.0001)
Average	<b>87.6</b>	86.6

speaker identification such as utterance variation, the recording environment change, and physical condition/emotion. Therefore, we conclude that IWKLR+IWCV is a novel promising approach to handling session dependent variation.

#### 5. CONCLUSIONS

In this paper, we proposed a novel semi-supervised speaker identification method that can alleviate the influence of non-stationarity such as session dependent variation, the recording environment change, and physical condition/emotion. We conducted a text-independent speaker identification simulation and experimentally found that the covariate shift formulation is useful in dealing with session dependent variations.

#### 6. REFERENCES

- [1] D.A. Reynolds, T.F. Quatieri, and R.B. Dunn, "Speaker verification using adapted Gaussian mixture models", *Digital Signal Processing* 10(1):19–41, 2000.
- [2] J. Mariethoz and S. Bengio, "A kernel trick for sequences applied to text-independent speaker verification systems", *Pattern Recognition*, 40(8):2315–2324, 2007.
- [3] T. Matsui and K. Tanabe, "Comparative Study of Speaker Identification Methods: dPLRM, SVM, and GMM", *IEICE Transactions on Information and Systems*, E89-D(3):1066–1073, 2006.
- [4] H. Shimodaira, "Improving predictive inference under covariate shift by weighting the log-likelihood function", *Journal of Statistical Planning and Inference*, 90(2):227–244, 2000.
- [5] M. Sugiyama, T. Suzuki, S. Nakajima, H. Kashima, P. von Büna, and M. Kawanabe, "Direct importance estimation for covariate shift adaptation", *Annals of the Institute of Statistical Mathematics*, 60(4):699–746, 2008.
- [6] M. Sugiyama, M. Krauledat, K.-R. Müller, "Covariate shift adaptation by importance weighted cross validation", *Journal of Machine Learning Research*, 8:985–1005, 2007.
- [7] L. Rabiner and B.-H. Juang, *Fundamentals of Speech Recognition*, Prentice Hall, 1993.
- [8] K. Tanabe, "Penalized logistic regression machines: New methods for statistical prediction 1", Technical Report 143, Institute of Statistical Mathematics, 2001.
- [9] T. Matsui and S. Furui, "Concatenated phoneme models for text-variable speaker recognition", *Proceedings of ICASSP'93*, II:391–394, 1993.