

# Covariate Shift Adaptation for Semi-supervised Speaker Identification

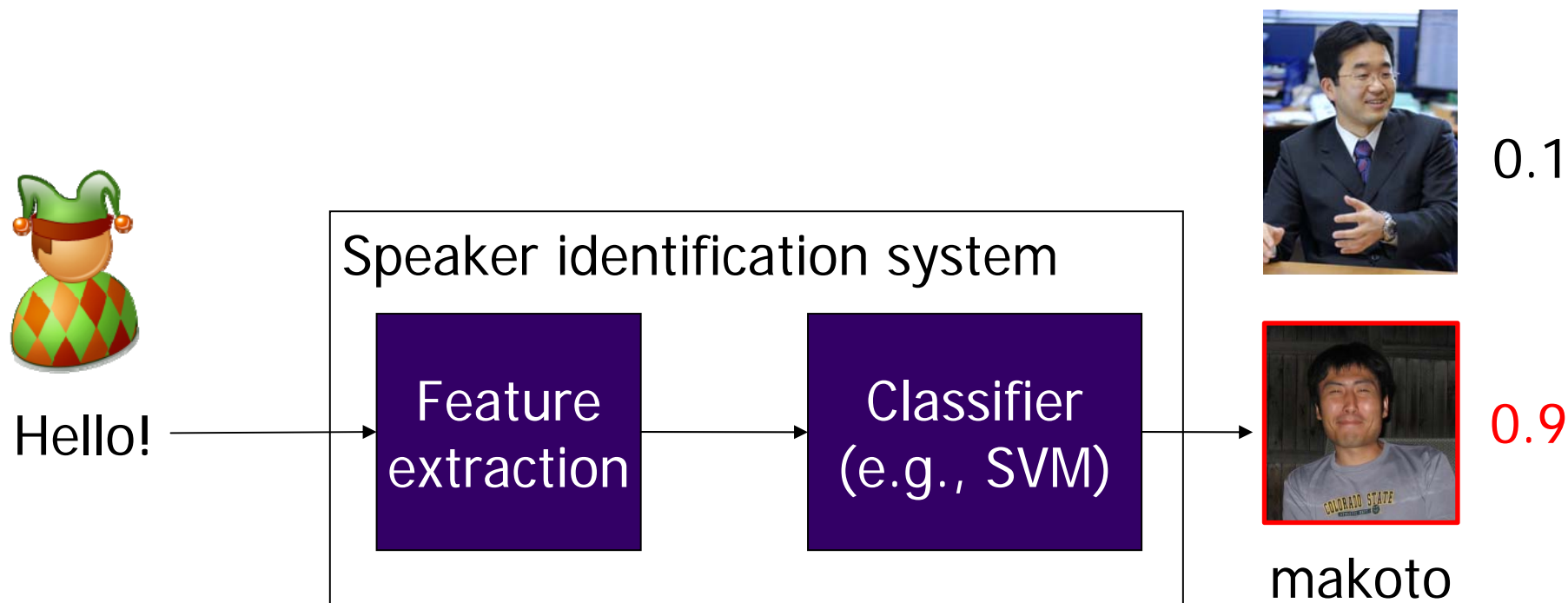
2009/4/22

M. Yamada, M. Sugiyama, and T. Matsui



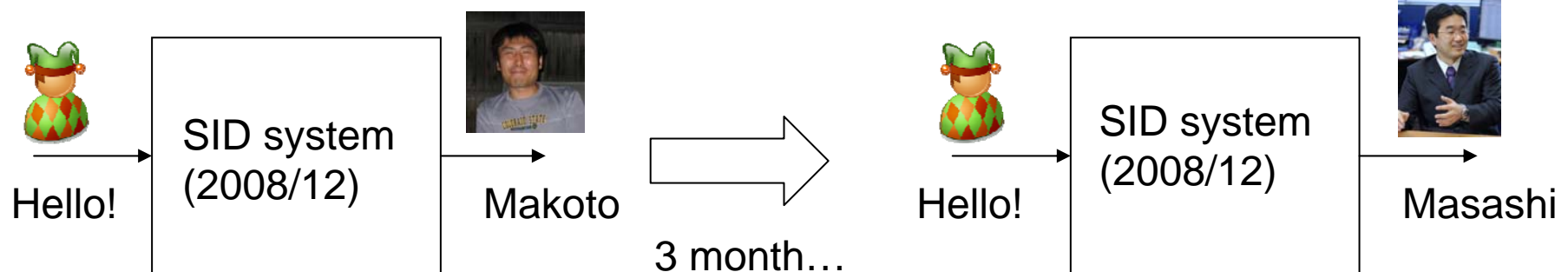
# Speaker Identification

**Task:** Identify speaker using voice.



# Problems in speaker identification

- Feature variation
  - Sound recording environment change
  - Physical condition/emotion
  - Noise
  - Session dependent variation



Can we use the same system trained in 2008/12 for future? **NO!**  
 ⇒ **Speech feature (e.g., MFCC) changes in 3 months.**

1. Recording several sessions of speeches
  - Labeling is required.
  - ⇒ **Very expensive!**
2. Semi-supervised learning
  - We use **unlabeled data** for training.
  - No labeling process required.
  - ⇒ **Reasonable.**

We assume the speech data follows **covariate shift**

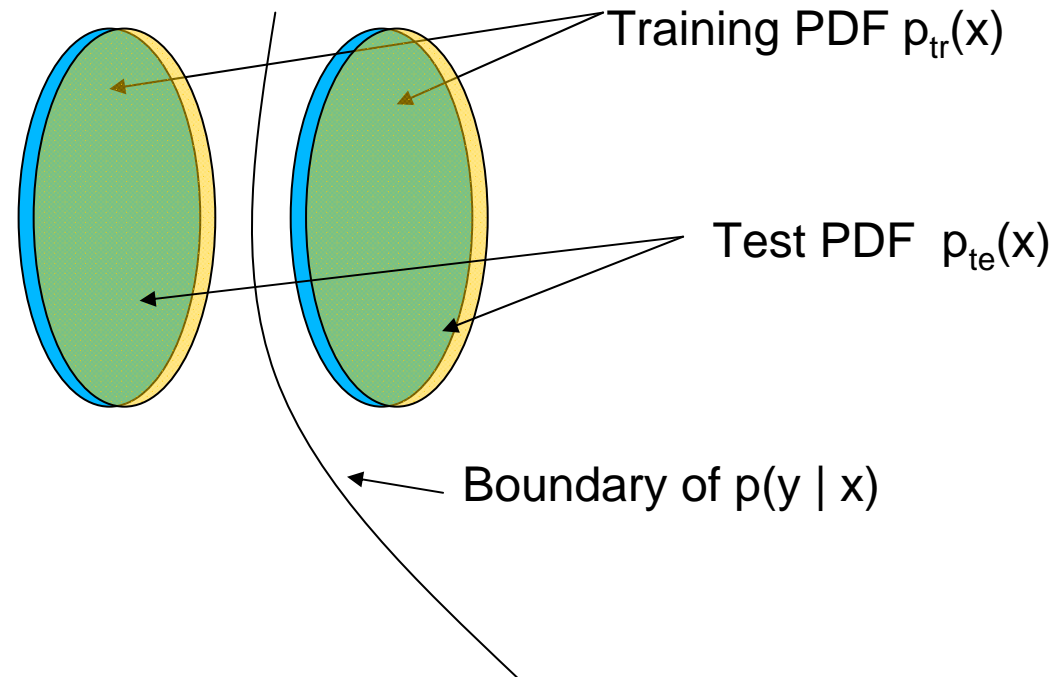
⇒ We model the covariate shift by using

**Importance Weighted Kernel Logistic Regression (IWKLR)**

# Supervised Learning

Assumption in supervised learning

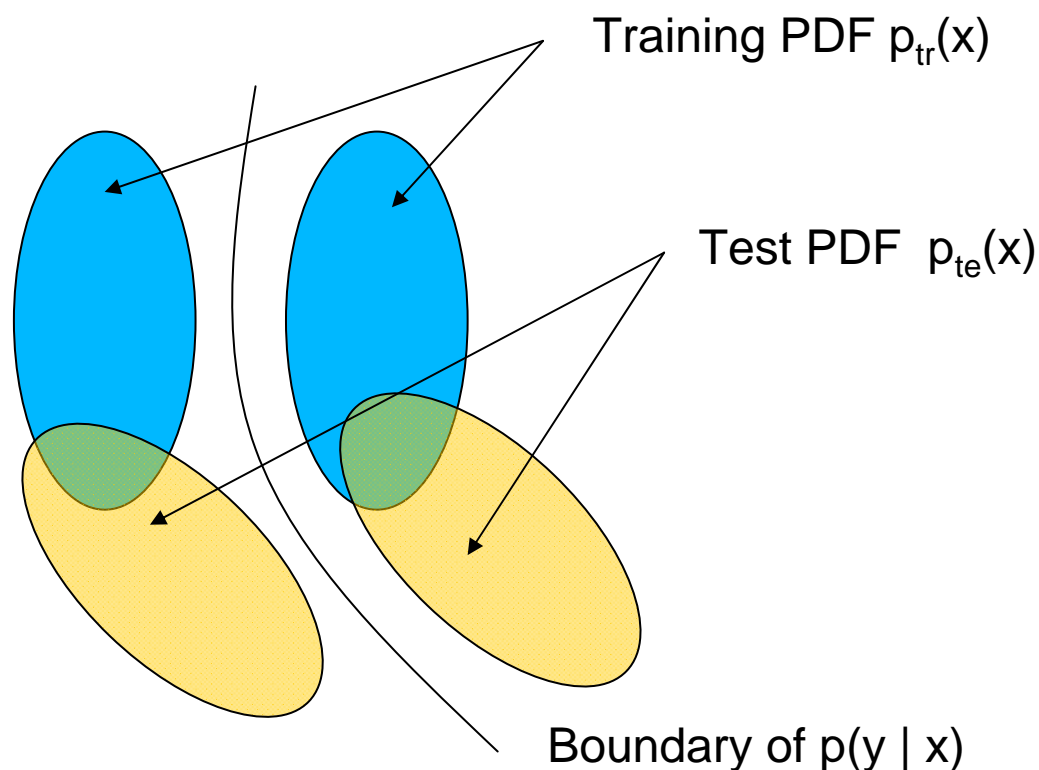
Training and test probability density functions are same.



Is this assumption acceptable in practice? **NO!**

# Covariate shift [1]

- Input probability density changes:  $p_{tr}(x) \neq p_{te}(x)$
- Conditional probability density remains unchanged:  $p(y | x)$



[1] H. Shimodaira, "Improving predictive inference under covariate shift by weighting the log-likelihood function," JSPI, 90, 227-244, 2000

# Covariate shift

Cost function under covariate shift:

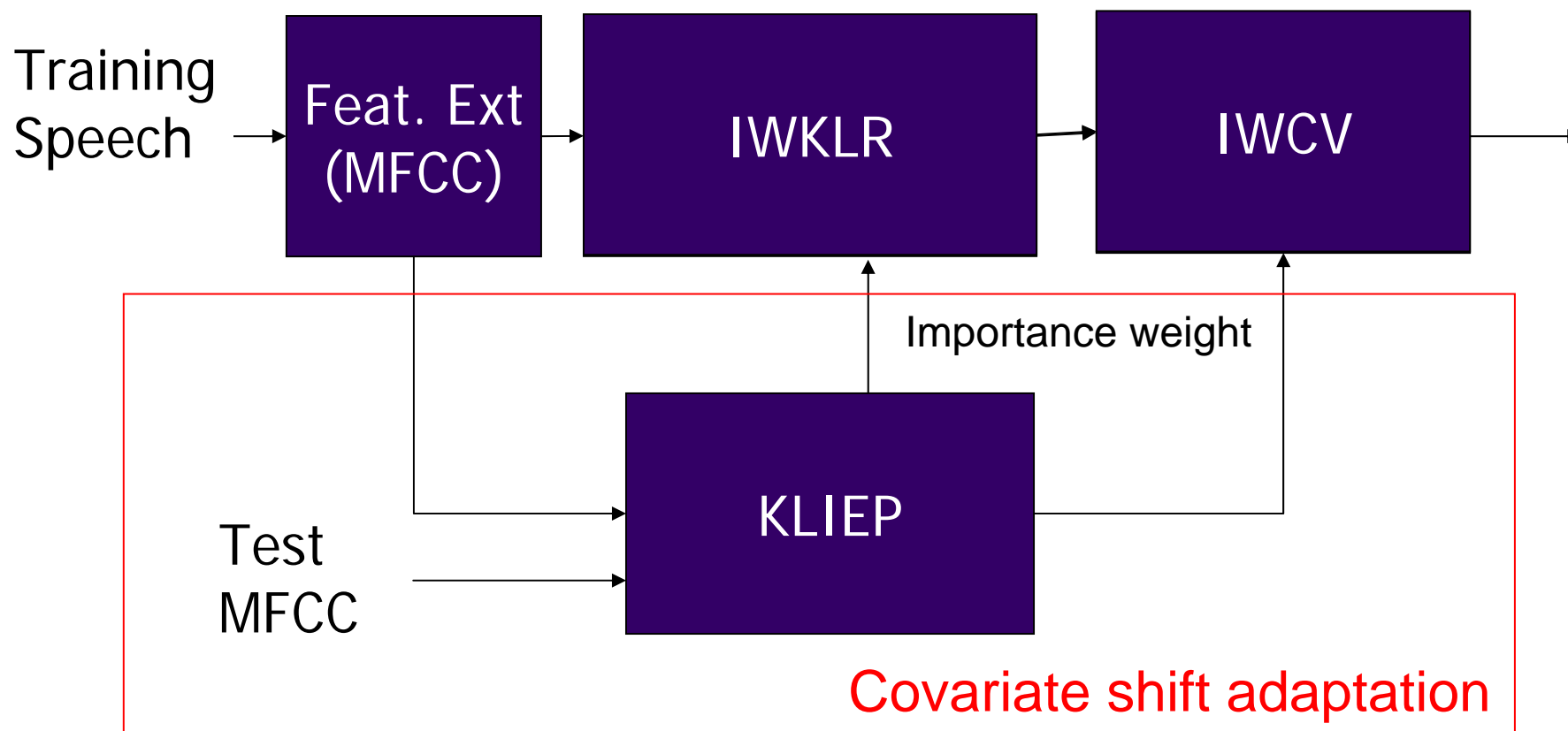
Taking the expectation over **test probability density**.

$$\begin{aligned} E_{p_{te}(\mathbf{X})}[F(\mathbf{X})] &= \int F(\mathbf{X})p_{te}(\mathbf{X})d\mathbf{X} \\ &= \int F(\mathbf{X})w(\mathbf{X})p_{tr}(\mathbf{X})d\mathbf{X} \\ &= E_{p_{tr}(\mathbf{X})}[F(\mathbf{X})w(\mathbf{X})]. \end{aligned}$$

Importance:

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

# Proposed framework



**Proposed method is consistent under covariate shift!**

[2] M. Sugiyama, et.al, ``Covariate shift adaptation by importance weighted cross validation.,'' JMLR, vol. 8 (May), pp.985-1005, 2007

[3] M. Sugiyama, et.al., ``Direct importance estimation for covariate shift adaptation.,'' AISM, vol. 60, no.4, pp.699-746, 2008



# Problem formulation

Mel-frequency cepstrum coefficient (MFCC):

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

$m$  labeled samples and  $(n-m)$  unlabeled samples:

$$\mathcal{Z}_{tr} = \{X_i, y_i\}_{i=1}^m$$

$$\mathcal{Z}_{te} = \{X_i\}_{i=m+1}^n$$

Speaker index:

$$y_i \in \{1, \dots, K\}$$

# Kernel based Speaker Identification

Posterior probability:

$$p(y = c | \mathbf{X}, \mathbf{V}) = \frac{\exp f_{\mathbf{v}_c}(\mathbf{X})}{\sum_{l=1}^K \exp f_{\mathbf{v}_l}(\mathbf{X})},$$

Discriminative function:

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

Sequence kernel[4]:

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

- [4] J. Mariethoz and S. Bengio, "A kernel trick for sequences applied to text-independent speaker verification systems," Pattern Recognition, 40, 2315-2324, 2007

Negative regularized importance weighted log-likelihood:

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

Regularizer :  $\frac{\delta}{2} \text{tr}(\mathbf{V} \mathbf{K} \mathbf{V}^{\top})$

Gram matrix:  $\mathbf{K} = [\mathcal{K}(\mathbf{X}_i, \mathbf{X}_j)]_{i,j=1}^n$

Importance weight:  $w(\mathbf{X})$

Negative log likelihood is **convex**

⇒ Easy to compute via Newton method.

# Importance Weighted Cross Validation (IWCV)

## Model parameters in IWKLR

Kernel width :  $\sigma$   
 Regularization parameter :  $\delta$

k-fold importance weighted cross validation (IWCV):

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

$\{\mathcal{Z}_i\}_{i=1}^k$  : Subset of  $\mathcal{Z} = \{(X_i, y_i)\}_{i=1}^n$

$|\mathcal{Z}_j|$  : Number of samples in the subset

$I(\cdot)$  : Indicator function

$w(X)$  : Importance weight

# Kullback–Leibler Importance Estimation Procedure (KLIEP)[3]

Model:

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

Basis:

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

Cost function:

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

[3] M. Sugiyama, et.al., "Direct importance estimation for covariate shift adaptation," AISM, vol. 60, no.4, pp.699-746, 2008

- Simulation condition
  - 10 speakers
  - Training data (1990/12)
  - Test data (1991/3, 1991/6, 1991/9)
  - 16kHz sampling
  - Speech length 12sec × 10 speakers
  - 12 MFCC +  $\Delta$ MFCC + log power +  $\Delta$  log power
  - Utterance data (300ms)  $X_i \in \mathbb{R}^{26 \times 30}$
  - 5-fold CV (KLR)
  - 5-fold IWCV (IWKLR)

# Evaluation result

Date	IWKLR + IWCV (1.5s)	KLR + CV (1.5s)	IWKLR + IWCV (4.5s)	KLR + CV (4.5s)
1991/3	<b>86.8</b> (1.2, 0.0001)	86.1 (1.2, 0.0001)	<b>92.6</b> (1.2,0.0001)	92.3 (1.2, 0.0001)
1991/6	<b>83.9</b> (1.3,0.0001)	82.0 (1.2, 0.0001)	<b>93.7</b> (1.3,0.0001)	92.7 (1.2, 0.0001)
1991/9	<b>92.0</b> (1.2, 0.0001)	91.7 (1.2, 0.0001)	<b>99.9</b> (1.2,0.0001)	99.7 (1.2, 0.0001)
Average	<b>87.6</b>	86.6	<b>95.4</b>	94.9

\*  $\sigma$  and  $\delta$  are in the bracket.

Model parameters are selected CV/IWCV.

# Conclusion & Future works

---

- Conclusion
  - Propose the semi-supervised speaker identification
  - Session dependent variation was alleviated by using the covariate shift adaptation
- Future works
  - Detection of Covariate shift
  - Modeling the physical condition/emotion using Covariate shift.