

Extending Linear Dynamical Systems with Dynamic Stream Weights for Audiovisual Speaker Localization

Christopher Schymura, Tobias Isenberg and Dorothea Kolossa

Cognitive Signal Processing Group, Institute of Communication Acoustics, Ruhr-University Bochum, Germany

Introduction

An important aspect of audiovisual speaker localization is the appropriate fusion of acoustic and visual observations based on their time-varying reliability. This study extends speaker localization and tracking frameworks based on linear dynamical systems with dynamic stream weights. A similar strategy has already been utilized for hidden Markov models in the context of audiovisual automatic speech recognition [1]. This work shows that dynamic stream weights can be integrated into the Kalman filter state estimation framework for linear dynamical systems and proposes a computationally efficient method for computing the corresponding acoustic and visual Kalman gains. Additionally, a method for estimating oracle dynamic stream weights from observation sequences with known state trajectories is introduced.

State Estimation

Linear dynamical system with conditionally independent acoustic and visual observations:

$$\begin{aligned} \mathbf{x}_k &= \mathbf{A}_k \mathbf{x}_{k-1} + \mathbf{v}_k, & \mathbf{v}_k &\sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_k) \\ \mathbf{y}_{A,k} &= \mathbf{C}_{A,k} \mathbf{x}_k + \mathbf{w}_{A,k}, & \mathbf{w}_{A,k} &\sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{A,k}) \\ \mathbf{y}_{V,k} &= \mathbf{C}_{V,k} \mathbf{x}_k + \mathbf{w}_{V,k}, & \mathbf{w}_{V,k} &\sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{V,k}) \end{aligned}$$

A recursive state estimation framework similar to the conventional Kalman filter can be derived for this class of linear dynamical systems, using the joint likelihood function at discrete time-step k with dynamic stream weights λ_k

$$p(\mathbf{x}_0, \dots, \mathbf{x}_k, \mathbf{y}_{A,1}, \dots, \mathbf{y}_{A,k}, \mathbf{y}_{V,1}, \dots, \mathbf{y}_{V,k}) \propto p(\mathbf{x}_0) \prod_{k'=1}^k p(\mathbf{x}_{k'} | \mathbf{x}_{k'-1}) p(\mathbf{y}_{A,k'} | \mathbf{x}_{k'})^{\lambda_{k'}} p(\mathbf{y}_{V,k'} | \mathbf{x}_{k'})^{1-\lambda_{k'}} \quad (1)$$

with

$$\begin{aligned} p(\mathbf{x}_0) &= \mathcal{N}(\boldsymbol{\mu}_0 | \boldsymbol{\Sigma}_0), & p(\mathbf{x}_k | \mathbf{x}_{k-1}) &= \mathcal{N}(\mathbf{A}_k \mathbf{x}_{k-1} | \mathbf{Q}_k) \\ p(\mathbf{y}_{A,k} | \mathbf{x}_k) &= \mathcal{N}(\mathbf{C}_{A,k} \mathbf{x}_k | \mathbf{R}_{A,k}), & p(\mathbf{y}_{V,k} | \mathbf{x}_k) &= \mathcal{N}(\mathbf{C}_{V,k} \mathbf{x}_k | \mathbf{R}_{V,k}) \end{aligned}$$

Require: Observations $\{\mathbf{y}_{A,k}\}_{k=1}^K$, $\{\mathbf{y}_{V,k}\}_{k=1}^K$, dynamic stream weights $\{\lambda_k\}_{k=1}^K$, initial state mean and covariance matrix $\hat{\mathbf{x}}_0$, $\hat{\boldsymbol{\Sigma}}_0$.

for $k = 1$ to K **do**

// Prediction step

$$\hat{\mathbf{x}}_{k|k-1} = \mathbf{A}_k \hat{\mathbf{x}}_{k-1}$$

$$\hat{\boldsymbol{\Sigma}}_{k|k-1} = \mathbf{A}_k \hat{\boldsymbol{\Sigma}}_{k-1} \mathbf{A}_k^T + \mathbf{Q}_k$$

// Compute acoustic and visual Kalman gain

$$\mathbf{S}_{A,k} = \mathbf{R}_{A,k}^{-1} \mathbf{C}_{A,k} \hat{\boldsymbol{\Sigma}}_{k|k-1}; \quad \mathbf{S}_{V,k} = \mathbf{R}_{V,k}^{-1} \mathbf{C}_{V,k} \hat{\boldsymbol{\Sigma}}_{k|k-1}$$

$$\begin{bmatrix} \mathbf{K}_{A,k}^T \\ \mathbf{K}_{V,k}^T \end{bmatrix} = \begin{bmatrix} \mathbf{I} + \lambda_k \mathbf{S}_{A,k} \mathbf{C}_{A,k}^T & (1 - \lambda_k) \mathbf{S}_{A,k} \mathbf{C}_{V,k}^T \\ \lambda_k \mathbf{S}_{V,k} \mathbf{C}_{A,k}^T & \mathbf{I} + (1 - \lambda_k) \mathbf{S}_{V,k} \mathbf{C}_{V,k}^T \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{S}_{A,k} \\ \mathbf{S}_{V,k} \end{bmatrix}$$

// Update step

$$\tilde{\mathbf{y}}_{A,k} = \mathbf{y}_{A,k} - \mathbf{C}_{A,k} \hat{\mathbf{x}}_{k|k-1}; \quad \tilde{\mathbf{y}}_{V,k} = \mathbf{y}_{V,k} - \mathbf{C}_{V,k} \hat{\mathbf{x}}_{k|k-1}$$

$$\hat{\mathbf{x}}_k = \hat{\mathbf{x}}_{k|k-1} + \lambda_k \mathbf{K}_{A,k} \tilde{\mathbf{y}}_{A,k} + (1 - \lambda_k) \mathbf{K}_{V,k} \tilde{\mathbf{y}}_{V,k}$$

$$\hat{\boldsymbol{\Sigma}}_k = \left(\mathbf{I} - \lambda_k \mathbf{K}_{A,k} \mathbf{C}_{A,k} - (1 - \lambda_k) \mathbf{K}_{V,k} \mathbf{C}_{V,k} \right) \hat{\boldsymbol{\Sigma}}_{k|k-1}$$

end for

return estimated state trajectory $\{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_K\}$

Algorithm 1: State estimation algorithm based on the conventional Kalman filter recursions, which are extended to cope with dynamic stream weights in linear dynamical systems. The mathematical derivation is based on the MMSE criterion using the joint likelihood function with dynamic stream weights in Eq. (1).

Oracle Dynamic Stream Weights

If a dataset with audiovisual observation sequences $\mathbf{Y}_A = \{\mathbf{y}_{A,k}\}_{k=1}^K$, $\mathbf{Y}_V = \{\mathbf{y}_{V,k}\}_{k=1}^K$ and the corresponding true state sequence of speaker positions $\mathbf{X} = \{\mathbf{x}_k\}_{k=1}^K$ is available, the sequence of oracle dynamic stream weights $\boldsymbol{\lambda}^* = \{\lambda_k\}_{k=1}^K$ can be computed according to

$$\boldsymbol{\lambda}^* = \arg \max_{\lambda_1, \dots, \lambda_K} \prod_{k=1}^K p(\lambda_k) p(\mathbf{y}_{A,k} | \mathbf{x}_k)^{\lambda_k} p(\mathbf{y}_{V,k} | \mathbf{x}_k)^{1-\lambda_k}.$$

By assuming i.i.d. oracle stream weights and imposing a Gaussian prior with mean μ_λ and variance σ_λ^2 , a closed-form solution can be obtained:

$$\lambda_k^* = \mu_\lambda + \sigma_\lambda^2 \log \left\{ \frac{p(\mathbf{y}_{A,k} | \mathbf{x}_k)}{p(\mathbf{y}_{V,k} | \mathbf{x}_k)} \right\}$$

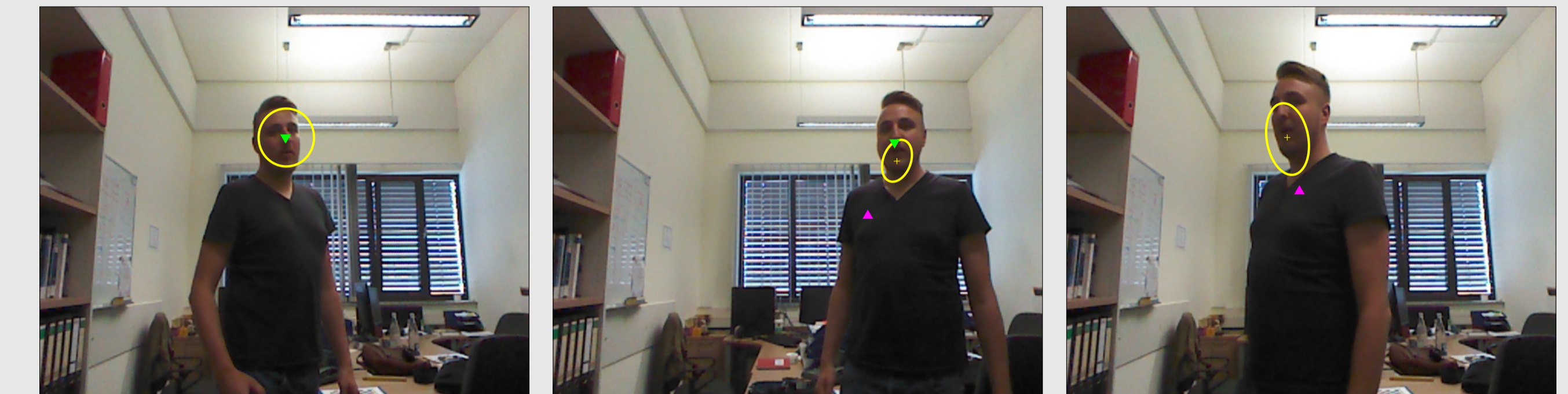


Figure 1: Exemplary speaker localization and tracking situations using oracle dynamic stream weights. Video observations are shown as green and audio observations as magenta triangles. The estimated state mean and covariance matrix are depicted in yellow.

Evaluation

Speaker localization experiments were conducted using 150 audiovisual recordings of the humanoid robot Nao in a reverberant room with $T_{60} \approx 450$ ms. DPD-MUSIC [2] was used for acoustic localization and visual speaker locations were estimated with the Viola-Jones algorithm [3]. The evaluation metric is localization RMSE. Systematic errors in visual localization were generated by adding zero-mean Gaussian noise at different standard deviations.

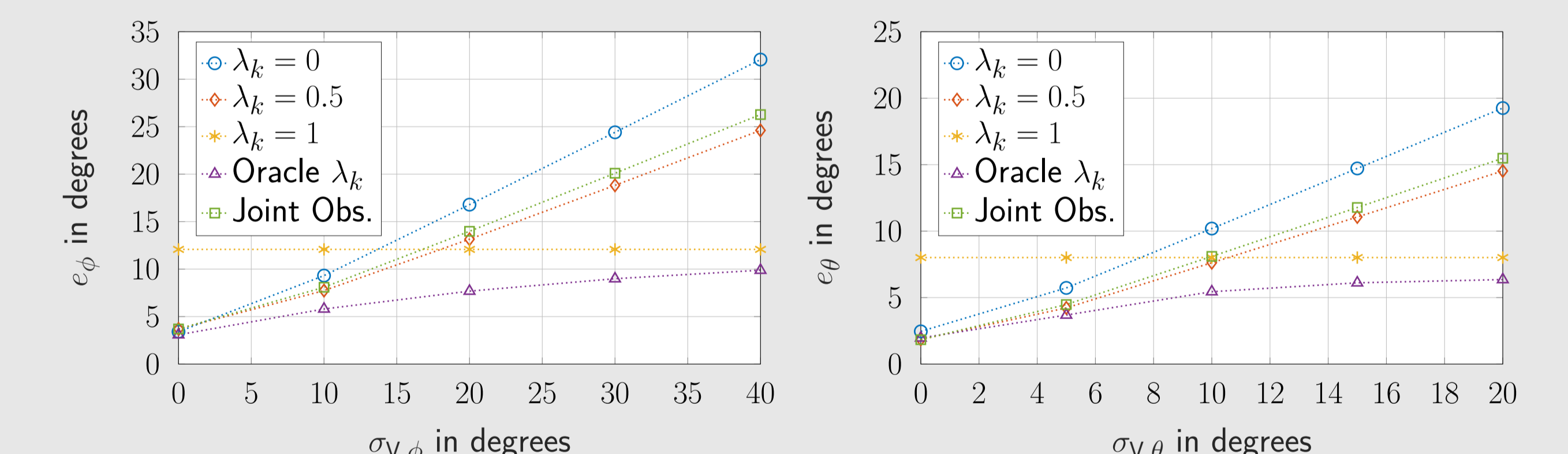


Figure 2: Localization errors obtained with different dynamic stream weight settings at varying levels of disturbance for the visual modality. Azimuth errors are shown in the left plot, elevation errors are depicted in the right plot. "Joint Obs." shows the baseline performance obtained with a conventional Kalman filter.

References

- [1] A. H. Abdelaziz, S. Zeiler, and D. Kolossa, "A new EM estimation of dynamic stream weights for coupled-HMM-based audio-visual ASR," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014.
- [2] O. Nadiri and B. Rafaely, "Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2014.
- [3] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proc. of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2001.