

Robotic haptic adjective perception based on coupled sparse coding

Pengwen XIONG^{1*}, Kongfei HE¹, Aiguo SONG^{2*} & Peter X. LIU³

¹School of Advanced Manufacturing, Nanchang University, Nanchang 330031, China;

²School of Instrument Science and Engineering, Southeast University, Nanjing 210096, China;

³Department of Systems and Computer Engineering, Carleton University, Ottawa ON K1S 5B6, Canada

Received 10 March 2021/Revised 25 November 2021/Accepted 14 May 2022/Published online 5 January 2023

Citation Xiong P W, He K F, Song A G, et al. Robotic haptic adjective perception based on coupled sparse coding. *Sci China Inf Sci*, 2023, 66(2): 129201, https://doi.org/10.1007/s11432-021-3512-6

Objects, textures, and materials can be identified by extracting haptic interaction information [1,2]. For haptic adjective understanding, Gao et al. [3] adopted deep models as a unified way to learn information from vision and haptics modalities. Nevertheless, this approach depends on a transfer learning method designed for object classification. Unlike feature learning methods, Chu et al. [4] designed a hand-crafted feature, which is hand-designed and usually has a clear physical meaning, to distinguish haptic adjectives. In general, most scholars attempt to solve the abstract and tortuous haptic problems with an entirely single-feature extraction method: a hand-crafted feature or a learning feature. In this study, we attempt to fuse measurements to avoid these defects based on a new sparse coding model. Sparse coding [5] is the representation of the back neurons by the strong activation of a relatively small group of preceding neurons and is widely used in unconventional unsupervised learning methods. Sparse coding is a remarkable tool for multimeasurement fusion that aims to find a complete sparse representation of the input data. Sparse coding tries to solve the primary problem:

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (1)$$

where $\mathbf{y} \in \mathbb{R}^{m \times 1}$ is the original data, $\mathbf{D} \in \mathbb{R}^{m \times n}$ is the dictionary matrix, $\mathbf{x} \in \mathbb{R}^{n \times 1}$ is the sparse coefficient vector, λ is a hyperparameter, $m < n$, and $\|\cdot\|_1$ is the sum of absolute values of elements in the vector \mathbf{x} .

We divide haptic data such as in Figure 1 into scalar and electrode array signals. Formally, the matrix $\mathbf{S}_G \in \mathbb{R}^{d_S \times N}$ constituted by vectors comprising the features of scalar signals is derived as

$$\mathbf{S}_G = [\mathbf{S}_{1,1}, \mathbf{S}_{2,1}, \dots, \mathbf{S}_{N_1,1}, \dots, \mathbf{S}_{N_2,2}, \dots, \mathbf{S}_{N_o,o}],$$

where G denotes the different haptic adjectives and N denotes the number of samples in the training set. In the same way, the matrix $\mathbf{E}_G \in \mathbb{R}^{d_E \times N}$ constituted by vectors consisting of the features of electrode array signals is defined

as

$$\mathbf{E}_G = [\mathbf{E}_{1,1}, \mathbf{E}_{2,1}, \dots, \mathbf{E}_{N_1,1}, \dots, \mathbf{E}_{N_2,2}, \dots, \mathbf{E}_{N_o,o}].$$

The training set obtains $N = \sum N_i$ samples, where o indicates the last object and N_i indicates the number of samples of the i -th object. During the testing phase, all measurements of the test sample are simultaneously entered into the constructed classifier to obtain a common prediction label that lies in {positive, negative}.

Generally, for multisensor information fusion, a sparse coding method is used to attempt to solve the following optimization problem:

$$\min_{\mathbf{D}_S, \mathbf{D}_E, \mathbf{X}_S, \mathbf{X}_E} \Phi_R(\mathbf{D}_S, \mathbf{D}_E, \mathbf{X}_S, \mathbf{X}_E) + \Phi_P(\mathbf{X}_S, \mathbf{X}_E), \quad (2)$$

where $\Phi_R(\mathbf{D}_S, \mathbf{D}_E, \mathbf{X}_S, \mathbf{X}_E) = \|\mathbf{S}_G - \mathbf{D}_S \mathbf{X}_S\|_F^2 + \|\mathbf{E}_G - \mathbf{D}_E \mathbf{X}_E\|_F^2$, $\Phi_P(\mathbf{X}_S, \mathbf{X}_E)$ is the penalty term, which contains a traditional sparse term and other penalty terms in many studies. $\mathbf{D}_S \in \mathbb{R}^{d_S \times K}$ is the dictionary matrix for scalar signals, and $\mathbf{D}_E \in \mathbb{R}^{d_E \times K}$ is the dictionary matrix for electrode array signals. \mathbf{X}_S is the matrix comprising sparse coefficient vectors of scalar signals in training samples, and \mathbf{X}_E is the matrix comprising sparse coefficient vectors of electrode array signals in training samples.

We use two projection matrices that can perfectly project raw features to a higher dimensional space:

$$\mathbf{S}_G \rightarrow \mathbf{P}_S^T \mathbf{S}_G, \quad \mathbf{E}_G \rightarrow \mathbf{P}_E^T \mathbf{E}_G,$$

where $\mathbf{P}_S \in \mathbb{R}^{d_S \times d}$ and $\mathbf{P}_E \in \mathbb{R}^{d_E \times d}$ are the projection matrices, $\mathbf{P}_S^T \mathbf{P}_S = \mathbf{P}_E^T \mathbf{P}_E = \mathbf{I}_d$, and d is the subspace dimension. With the development of convex optimization methods, this function can be easily solved by off-the-shelf solvers. The reconstruction error term can be reconstructed with projection matrices as

$$\Phi_R(\mathbf{D}, \mathbf{X}_S, \mathbf{X}_E) = \left\| \mathbf{P}_S^T \mathbf{S}_G - \mathbf{D} \mathbf{X}_S \right\|_F^2 + \left\| \mathbf{P}_E^T \mathbf{E}_G - \mathbf{D} \mathbf{X}_E \right\|_F^2, \quad (3)$$

* Corresponding author (email: steven.xpw@ncu.edu.cn, a.g.song@seu.edu.cn)

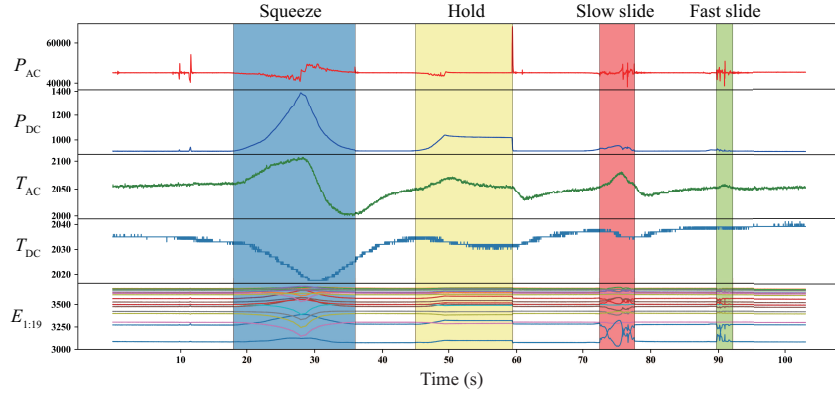


Figure 1 (Color online) Haptic signals were recorded from a PR2 robot trial using two BioTac sensors installed in its left hand. A single record contains pressure (P_{AC} , P_{DC}), temperature (T_{AC} , T_{DC}), and spatially distributed impedance-measuring electrodes ($E_{1:19}$), and each integrated process contains four exploratory procedures, including squeeze, hold, slow slide, and fast slide.

where $\mathbf{D} = [\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_K] \in \mathbb{R}^{d \times K}$ is a dictionary matrix with $d < K$. The proposed fusion model can be defined as the following optimization problem:

$$\min_{\mathbf{D}, \mathbf{X}_S, \mathbf{X}_E, \mathbf{P}_S, \mathbf{P}_E} \Phi_R(\mathbf{D}, \mathbf{X}_S, \mathbf{X}_E) + \Phi_P(\mathbf{X}_S, \mathbf{X}_E). \quad (4)$$

$\|\mathbf{d}_k\|_2 \leq 1$ is a typical constraint that can prevent obtaining an oversized solution, and $\mathbf{P}_S^T \mathbf{P}_S = \mathbf{P}_E^T \mathbf{P}_E = \mathbf{I}_d$.

For our haptic adjective recognition model, the constraint of the proposed fusion model can be defined as

$$\Phi_P(\mathbf{X}_S, \mathbf{X}_E) = \partial \|\mathbf{X}_S, \mathbf{X}_E\|_{1,1} + \Phi_L(\mathbf{X}_S, \mathbf{X}_E), \quad (5)$$

where $\|\mathbf{X}_S, \mathbf{X}_E\|_{1,1}$ is the sum of absolute values of elements in $[\mathbf{X}_S, \mathbf{X}_E]$. The label of training samples is defined as

$$q_{*,i} = \begin{cases} +\delta, & \text{when object } i \text{ is labeled as positive,} \\ -\delta, & \text{when object } i \text{ is labeled as negative,} \end{cases} \quad (6)$$

where δ is a positive number usually set as one, and * refers to any sample in object i . Labels of training samples can be combined as

$$\mathbf{Q} = [q_{1,1}, q_{2,1}, \dots, q_{N_1,1}, \dots, q_{N_2,2}, \dots, q_{N_o,o}]. \quad (7)$$

Inspired by [6], we linearly transform all coefficient vectors to a common label for pairing, which can be defined as

$$\begin{aligned} +\delta &= \mathbf{W}_{S,T}^T \mathbf{x}_{S,T} = \mathbf{W}_E^T \mathbf{x}_{E,T}, \\ -\delta &= \mathbf{W}_{S,F}^T \mathbf{x}_{S,F} = \mathbf{W}_E^T \mathbf{x}_{E,F}, \end{aligned} \quad (8)$$

where $\mathbf{x}_{S,T}$ is the coefficient vector of scalar signals of any positively labeled sample, and $\mathbf{x}_{S,F}$ is the coefficient vector of scalar signals of any negatively labeled sample. Similarly, $\mathbf{x}_{E,T}$ is the coefficient vector of electrode array signals of any positively labeled sample, and $\mathbf{x}_{E,F}$ is the coefficient vector of electrode array signals of any negatively labeled sample. According to the above discussion, the entire optimization problem can be constructed as

$$\begin{aligned} \min_{\mathbf{D}, \mathbf{X}_S, \mathbf{X}_E, \mathbf{P}_S, \mathbf{P}_E} & \Phi_R(\mathbf{D}, \mathbf{X}_S, \mathbf{X}_E) + \Phi_P(\mathbf{X}_S, \mathbf{X}_E), \\ \Phi_R(\mathbf{D}, \mathbf{X}_S, \mathbf{X}_E) &= \left\| \mathbf{P}_S^T \mathbf{S}_G - \mathbf{D} \mathbf{X}_S \right\|_F^2 \\ &+ \left\| \mathbf{P}_E^T \mathbf{E}_G - \mathbf{D} \mathbf{X}_E \right\|_F^2, \\ \Phi_P(\mathbf{X}_S, \mathbf{X}_E) &= \partial \|\mathbf{X}_S, \mathbf{X}_E\|_{1,1} \end{aligned}$$

$$+ \beta \left\| \begin{bmatrix} \mathbf{Q} \\ \mathbf{Q} \end{bmatrix} - \begin{bmatrix} \mathbf{W}_S^T \mathbf{X}_S \\ \mathbf{W}_E^T \mathbf{X}_E \end{bmatrix} \right\|_F^2 + \chi \|\mathbf{W}_S, \mathbf{W}_E\|_F^2, \quad (9)$$

where $\|\mathbf{W}_S, \mathbf{W}_E\|_F^2$ can prevent the overfitting problem. β and χ are weight parameters. The above approximation solution can be obtained by iterating a process that immobilizes different variables according to a fixed sequence. With the optimal solution, \mathbf{D}^* , \mathbf{W}_S^* , \mathbf{W}_E^* , \mathbf{P}_S^* , and \mathbf{P}_E^* are obtained. With the feature \mathbf{s}_G of the scalar signal and the feature \mathbf{e}_G of the electrode array signal of a testing sample, the label l of the testing sample can be decided by

$$l = \begin{cases} \text{positive,} & \mathbf{W}_S^{*T} \mathbf{x}_S^* + \mathbf{W}_E^{*T} \mathbf{x}_E^* \geq 0, \\ \text{negative,} & \mathbf{W}_S^{*T} \mathbf{x}_S^* + \mathbf{W}_E^{*T} \mathbf{x}_E^* < 0. \end{cases} \quad (10)$$

Conclusion. Traditional methods cannot fuse multimodal haptic measurements well; to solve this problem, we attempt to translate multimodal sparse codes obtained using a unified dictionary into a shared label. The proposed model not only preserves the original multimodal information but also transforms the raw data into a shared feature space.

Acknowledgements This work was partially supported by National Natural Science Foundation of China (Grant Nos. 62163024, 61903175, 61663027) and Academic and Technical Leaders Foundation of Major Disciplines of Jiangxi Province (Grant No. 20204BCJ23006).

References

- Liu H, Yu Y, Sun F, et al. Visual-tactile fusion for object recognition. *IEEE Trans Automat Sci Eng*, 2017, 14: 996–1008
- Strese M, Schuwerk C, Iepure A, et al. Multimodal feature-based surface material classification. *IEEE Trans Haptics*, 2017, 10: 226–239
- Gao Y, Hendricks L A, Kuchenbecker K J, et al. Deep learning for tactile understanding from visual and haptic data. In: *Proceedings of 2016 IEEE International Conference on Robotics and Automation (ICRA)*, 2016. 536–543
- Chu V, McMahon I, Riano L, et al. Robotic learning of haptic adjectives through physical interaction. *Robot Auton Syst*, 2015, 63: 279–292
- Huang D, Wang Y F. Coupled dictionary and feature space learning with applications to cross-domain image synthesis and recognition. In: *Proceedings of 2013 IEEE International Conference on Computer Vision*, 2013. 2496–2503
- Jiang Z, Lin Z, Davis L S. Learning a discriminative dictionary for sparse coding via label consistent K-SVD. In: *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2011. 1697–1704