

Dynamical system's modelling for tic activity recognition

Jules Gottraux

EPFL

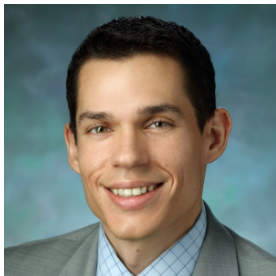
November 10th, 2020

Tic disorders:

- Conditions that induces motor and/or vocal spasms
- Pretty common, especially for young people
- Lowers the quality of life significantly in general
- Behavioral therapy as first-line treatment

→ Our goal is to facilitate this therapy by automating tic detection

Idea from Joseph F. McGuire and Joey Ka-Yee Essoe. Psychologists specialized in neuropsychiatric conditions such as tic disorder in the childhood.



Contribute by creating the dataset for tic detection. Creation of the dataset began in september 2020.



Setup for dataset

Our baseline is an activity recognition framework on pre-segmented dataset using linear dynamical systems. We'll build our methods from there.



Screenshot of tasks in JIGSAWS dataset

Fit a linear dynamical system to the video $\mathbf{Y} \in \mathbb{R}^{N \times P}$, $\mathbf{Y} = (y_1, \dots, y_N)^T$:

$$y_t = \mathbf{C}x_t$$

$$x_{t+1} = \mathbf{A}x_t$$

The projection is obtained by Principal Component Analysis (PCA) and is fixed. The matrix \mathbf{A} is then obtained via:

$$\mathbf{X}_- \mathbf{A} = \mathbf{X}_+$$

$$\Rightarrow \mathbf{A} = \mathbf{X}_-^\dagger \mathbf{X}_+$$

\mathbf{X} are the encoded frames and $\mathbf{X}_- = (x_1, \dots, x_{N-1})^T$, $\mathbf{X}_+ = (x_2, \dots, x_N)^T$

Our work consists in improving and adapting this baseline from activity recognition to tic detection and is composed of 4 main sections:

- Extension of the baseline using non-linear dimensionality reduction
- Extension of the baseline using joint learning of dynamical system components
- Extension of the techniques to online detection
- Test on Hopkins' dataset

All experiments are done on 256×256 videos converted to gray scale.

For a video with N frames: $\mathbf{Y} \in \mathbb{R}^{N \times P}$. We seek to find a mapping $\Phi_E : \mathbb{R}^P \rightarrow \mathbb{R}^R$, $R \ll P$ such that the frames are *well represented* in the latent space.

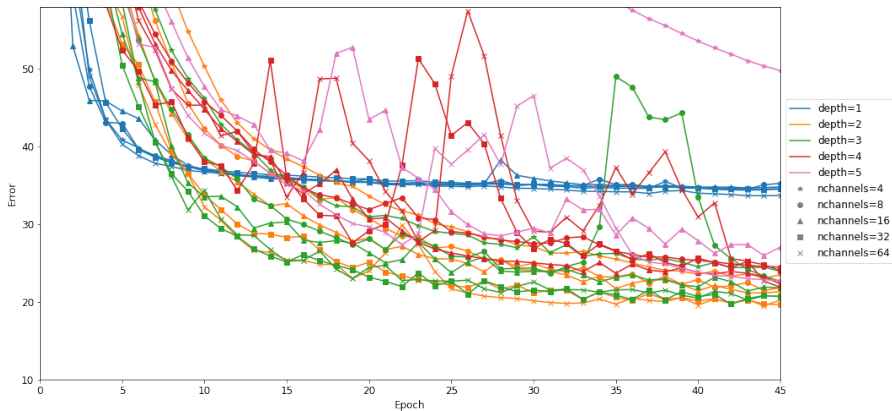
We minimize the reconstruction error on a video from Hopkins:

$$\mathcal{L}_{rec} = \frac{1}{N} \sum_{t=1}^N \|y_t - \Phi_E^{-1}(\Phi_E(y_t))\|_2^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2$$

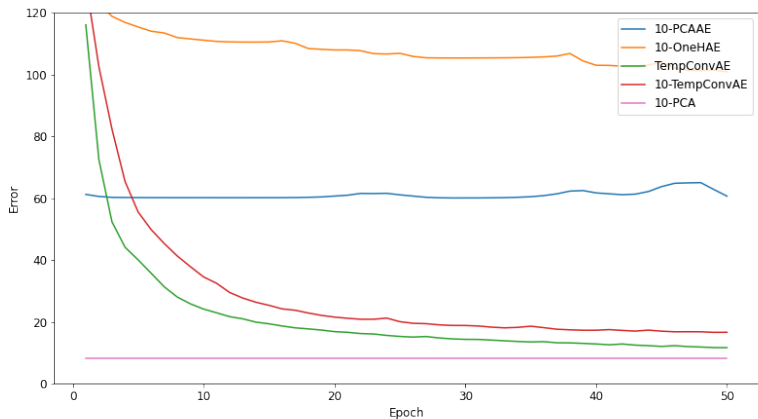
Where $\hat{\mathbf{Y}}$ is the reconstruction of all frames. And Φ_E^{-1} is simply Φ_D our mapping for the decoding.

We compare three models, each one is a neural network with an autoencoder structure:

- PCAA (PCA autoencoder): autoencoder using a linear projection
- OneHAE (one hidden layer autoencoder): autoencoder using a one hidden layer neural network structure
- TempConvAE (temporal convolutional autoencoder): autoencoder using 3d (a.k.a. temporal) convolutional layers



Training error for all temporal convolutional networks, $R = 10$



Training error for all models, $R = 10$

Notes:

- Further analysis showed that enforcing the latent dimension with a linear mapping were hurting the model's power
- Better performance could be obtained using known video compression algorithm or network (e.g. U-Net)
- We'll stick to linear projections, as no improvement have been observed

Initial goals with this dataset:

- Measure the capability of the method to detect an activity on raw frames
- Assess whether learning **A** and **C** jointly helps the algorithm

Instead of fixing the projection and obtaining **A** as in the baseline, we initialize the models with the parameters from the baseline and minimize:

$$\mathcal{L}_{pred} = \sqrt{\frac{1}{N-1} \sum_{t=2}^N \|y_t - \mathbf{C}^{-1}(\mathbf{AC}(y_{t-1}))\|_2^2}$$

Evaluation:

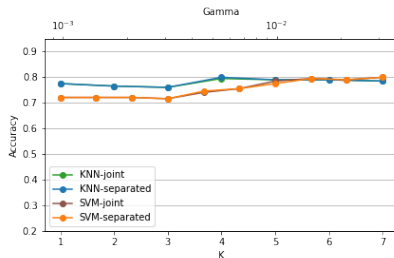
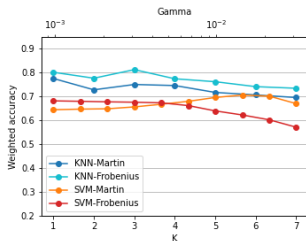
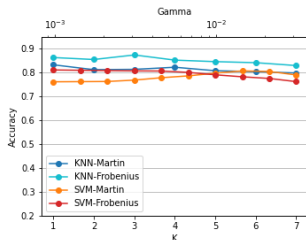
- We have fragments of videos, each with a given activity and fitted dynamical system
- Classification task, the features are the dynamical system's matrices and the labels are the corresponding activity

The classification uses metrics based on subspace angles between matrices of the dynamical systems. These metrics are the Frobenius and Martin distance.

These metrics tells us how much the dynamics of two different fragments differ.

On top of these metrics, we use K-Nearest Neighbors (KNN) or Support Vector Machine with radial basis function kernel (SVM) for the classification

Results for one task (suturing):



Notes:

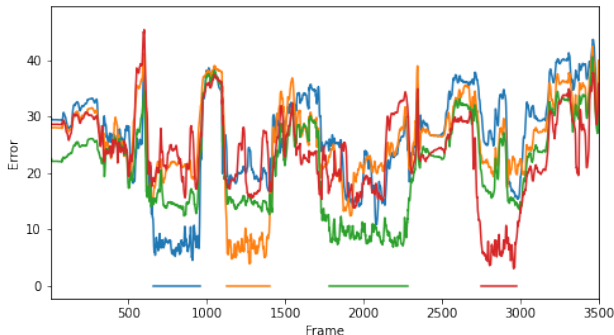
- Joint learning performance is indistinguishable from separated learning. Hence, separated learning is nearly optimal for this framework
- The model's performance seems enough for our detection task

We test two approach for online detection:

- Detection based on the reconstruction error of the prediction of models
- Detection based on the distance between models and a model fitted in a moving window fashion

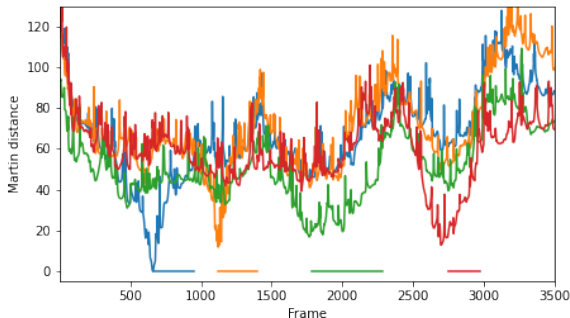
Experiments are done on a video picked from the JIGSAWS dataset.

- Pick a particular gesture
- Do the models of this gesture capture the dynamics of the others as well?
- If yes we could use this information to construct an online classifier



Reconstruction error of the prediction of models from same activity

- Pick a particular gesture
- Can we detect the occurrence of this gesture based on the distance between a moving window model and the fitted models?



Martin distance between models of the gesture and moving window model