

Background on Martin distance

Jules Gottraux

November 2020

Classification using subspace angles

As the first classification algorithm, we will use a nearest neighbor classifier between the transition matrix of the models fitted from the training set and the model of the test video. This algorithm will, as distances, use the squared Martin and Frobenius distance on the subspace angles $\theta_1, \theta_2, \dots, \theta_{2n}$. Let models \mathcal{M}_1 and \mathcal{M}_2 be $\mathcal{M}_i = (A_i, C_i), i \in \{1, 2\}$, then if $\|A_i\|_2 < 1$ the subspace angles can directly be computed from the eigenvalues of matrix P^{ij} :

$$P^{ij} = P_{ii}^{-1} P_{ij} P_{jj}^{-1} P_{ji} \quad (1)$$

with P_{ij} the solution to:

$$P_{ij} = A_i^T P_{ij} A_j + C_i^T C_j \iff (A_i^T)^{-1} P_{ij} + P_{ij} (-A_j) = (A_i^T)^{-1} C_i^T C_j$$

Which is a Sylvester equation, i.e. $A_s X + X B_s = C_s$ of parameters:

$$\begin{aligned} A_s &= (A_i^T)^{-1} \\ B_s &= (-A_j) \\ C_s &= (A_i^T)^{-1} C_i^T C_j \end{aligned}$$

Once we have the matrices P^{ij} and P^{ji} with their corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ we define:

$$\theta_k = \begin{cases} \cos^{-1}(\sqrt{\lambda_i^{ij}}), & \text{if } k \in \{1, \dots, n\} \\ \cos^{-1}(\sqrt{\lambda_i^{ji}}), & \text{if } k \in \{n+1, \dots, 2n\} \end{cases}$$

For the model definition $\mathcal{M}_i = (A_i, C_i), i \in 1, 2$, if model i does not uses a linear mapping but rather a neural network, we set the corresponding $C_i = I_{R \times P}$. Then from the subspace angles, the distances can be computed as follows:

$$\begin{aligned} d_M^2(\mathcal{M}_1, \mathcal{M}_2) &= -\log \prod_{i=1}^{2n} \cos(\theta_i)^2 \\ d_F^2(\mathcal{M}_1, \mathcal{M}_2) &= 2 \sum_{i=1}^{2n} \sin(\theta_i)^2 \end{aligned}$$

We can also use a kernel SVM algorithm to classify the models from those distances. If we have S samples points (corresponding of models \mathcal{M}_s fitted on video fragments) the associated kernel K_{RBF} can be defined as:

$$K_{RBF,ij} = e^{-\gamma d_X^2(\mathcal{M}, \mathcal{M})}, \quad \gamma > 0, \quad X \in M, F$$

Enforcing small intra-cluster distance

When learning the transition matrices A_i , we want to cluster matrices from the same activity $g \in G$ together. Meaning that if we have two videos i and j from activity g_i and g_j we want $d_X^2(\mathcal{M}, \mathcal{M})$ to be small if $g_i = g_j$ and high if $g_i \neq g_j$. We can enforce that by minimizing not only the prediction error but also this intra-cluster distance. So once we have obtained the points $\mathcal{M} = (A_i, C_i)$ for all video fragment we can fix the projection C_i and continue to learn A_i by minimizing the following loss. For each activity $g \in G$, with $i, j \in \{k : \text{video } k \text{ has activity } g\}$:

$$\mathcal{L} = \alpha \sum_i \|X_i^+ - A_i X_i^-\|_F^2 + (1 - \alpha) \sum_{i \neq j} d_X^2(\mathcal{M}_i, \mathcal{M}_j) \quad (2)$$

Then at test time for a given video fragment, we fit the projection and the matrix using loss (??). Then, two approaches are possible, we stop here and classify as before using k-nearest neighbors or SVM. Or we can fix the projection C_t and continue learning the matrix A_t . To continue the learning of A_t , we will have $|G|$ models: $\mathcal{M}_t^g = (A_t^g, C_t)$. For each activity g we have training samples: $\mathcal{M}_i = (A_i, C_i), i \in \mathcal{I} = \{i : \text{video } i \text{ is in activity } g\}$ and we minimize:

$$\mathcal{L}_g = \alpha \|X_t^+ - A_t^g X_t^-\|_F^2 + \frac{(1 - \alpha)}{|\mathcal{I}|} \sum_i d_X^2(\mathcal{M}_i, \mathcal{M}_t^g) \quad (3)$$

Then, the out of those $|G|$ models, we predict the gesture of the model that gives the "best" approximation.

Derivative of the Martin distance

Let $\mathcal{M} = (\mathbf{C}, \mathbf{A})$ and $\mathcal{M}_i = (\mathbf{C}_i, \mathbf{A}_i)$ be the parameters of two dynamical systems. We want to find the transition matrix \mathbf{A} of model \mathbf{A} that minimizes the Martin distance:

$$\min_{\mathbf{A}} d_M^2(\mathcal{M}, \mathcal{M}_i). \quad (4)$$

A natural strategy to minimize over \mathbf{A} would be to use a gradient descent method. For that purpose, we need the derivative of the Martin distance w.r.t. matrix \mathbf{A} . Assuming that \mathbf{C} has orthonormal columns (*i.e.*, $\mathbf{C}^T \mathbf{C} = \mathbf{I}$), it can be shown that derivative w.r.t. the jk -th entry of \mathbf{A} is given by:

$$\frac{\partial d_M^2(\cdot)}{\partial A_{jk}} = \text{Tr} \left((\mathbf{X}^{-1})^T \frac{\partial \mathbf{X}}{\partial A_{jk}} \right) - 2 \text{Tr} \left((\mathbf{X}_i^\dagger)^T \frac{\partial \mathbf{X}_i}{\partial A_{jk}} \right), \quad (5)$$

where \cdot^\dagger denote pseudo-inverse and where

$$\mathbf{X} = \mathbf{A}^T \mathbf{X} \mathbf{A} + \mathbf{I} \quad (6)$$

$$\mathbf{X}_i = \mathbf{A}_i^T \mathbf{X}_i \mathbf{A}_i + \mathbf{I} \quad (7)$$

$$\frac{\partial \mathbf{X}}{\partial A_{jk}} = \mathbf{A}^T \frac{\partial \mathbf{X}}{\partial A_{jk}} \mathbf{A} + \mathbf{J}_{jk}^T \mathbf{X} \mathbf{A} + \mathbf{A}^T \mathbf{X} \mathbf{J}_{jk} \quad (8)$$

$$\frac{\partial \mathbf{X}_i}{\partial A_{jk}} = \mathbf{A}^T \frac{\partial \mathbf{X}_i}{\partial A_{jk}} \mathbf{A}_i + \mathbf{J}_{jk}^T \mathbf{X}_i \mathbf{A}_i, \quad (9)$$

with \mathbf{J}_{jk} being the single-entry matrix (*i.e.*, all zeros except 1 at the jk -th entry). Note that the above relations correspond to the solutions of discrete Lyapunov or Sylvester equations.

Note: for computation purposes, we can rewrite the derivative of martin:

$$\frac{\partial d_M^2(\cdot)}{\partial A_{jk}} = \sum_{j,k} (X^{-1} \circ \frac{\partial \mathbf{X}}{\partial A_{jk}}) - 2 \sum_{j,k} (X_i^\dagger \circ \frac{\partial \mathbf{X}_i}{\partial A_{jk}}) \quad (10)$$