

Probabilistic Kernels for the Classification of Auto-Regressive Visual Processes

Antoni B. Chan and Nuno Vasconcelos
Department of Electrical and Computer Engineering
University of California, San Diego
abchan@ucsd.edu, nuno@ece.ucsd.edu

Abstract

We present a framework for the classification of visual processes that are best modeled with spatio-temporal autoregressive models. The new framework combines the modeling power of a family of models known as dynamic textures and the generalization guarantees, for classification, of the support vector machine classifier. This combination is achieved by the derivation of a new probabilistic kernel based on the Kullback-Leibler divergence (KL) between Gauss-Markov processes. In particular, we derive the KL-kernel for dynamic textures in both 1) the image space, which describes both the motion and appearance components of the spatio-temporal process, and 2) the hidden state space, which describes the temporal component alone. Together, the two kernels cover a large variety of video classification problems, including the cases where classes can differ in both appearance and motion and the cases where appearance is similar for all classes and only motion is discriminant. Experimental evaluation on two databases shows that the new classifier achieves superior performance over existing solutions.

1. Introduction

Figure 1 presents a sample from a large collection of visual processes that have proven remarkably challenging for traditional motion representations, based on modeling of the individual trajectory of pixels [1, 2], particles [3], or objects in a scene. Since most of the information required for the perception of these processes is contained in the interaction between the many motions that compose them, they require a holistic representation of the associated motion field capable of capturing its variability without the need for segmentation or tracking of individual components. Throughout the years, some representations appeared particularly promising in this respect, e.g. the representation of the motion field as a collection of layers [4]. However, only recently some real success has been demonstrated through the modeling of these processes as dynamic textures, i.e. realizations of an auto-regressive stochastic process with both a spatial and temporal component [5, 6]. Like many other re-

cent advances in vision, the success of these methods derives from the adoption of representations based on generative probabilistic models that can be learned from collections of training examples.

In the context of classification, detection, and recognition problems, the probabilistic representation has various properties that are known to be assets for perception [7], e.g. existence of principled inference formalisms that allow the fusion of diverse sources of information, the ability to incorporate domain-knowledge in the form of prior beliefs, etc. There are, nevertheless, core aspects in which it also has strong shortcomings. In particular, while it can lead to optimal classifiers by simple application of Bayesian decision theory, these classifiers have weak generalization guarantees, and can be quite sensitive to the dimensionality of the underlying feature space, or prone to over-fitting when the models have large numbers of parameters. This is a source of particular concern for the problems of Figure 1 since spatio-temporal autoregressive modeling tends to require high dimensional feature and state spaces.

An alternative classification framework [8], which delivers large-margin classifiers of much better generalization ability (e.g. the now popular support vector machine), does exist but has strong limitations of its own. For the classification of spatio-temporal data-streams, the most restrictive among these is a requirement for the representation of those data-streams as points in Euclidean space. These points are then mapped into a high-dimensional feature space by a kernel function that transforms Euclidean distances in domain space into distances defined over a manifold embedded in range space. The Euclidean representation is particularly troublesome for spatio-temporal processes, where different instances of a process may have different temporal extents (e.g. two similar video streams with different numbers of frames), or be subject to simple transformations that are clearly irrelevant for perception and classification (e.g. a change of sampling rate), but can map the same data-stream into very different points of Euclidean space.

Recent developments in the area of probabilistic kernels have shown significant promise to overcome these limitations. Probabilistic kernels are kernels that act on pairs of generative probabilistic models, enabling simultaneous sup-



Figure 1: Examples of visual processes that are challenging for traditional spatio-temporal representations: fire, smoke, the flow of a river stream, or the motion of an ensemble of objects, e.g. a flock of birds, a bee colony, a school of fish, the traffic on a highway, or the flow of a crowd.

port for complex statistical inference, which is characteristic of probabilistic representations, and good generalization guarantees, which are characteristic of large-margin learning. Although the feasibility of applying these kernels to vision problems has been demonstrated on relatively simple recognition tasks, e.g. the recognition of objects presented against a neutral background [9], we believe that this is unsatisfactory in two ways. First, the greatest potential for impact of probabilistic kernels is in the solution of classification problems where 1) simple application of Bayesian decision theory is likely to fail, e.g. problems involving large state-space models and high-dimensional features, and 2) the inappropriateness of the Euclidean representation makes the traditional large-margin solutions infeasible. Second, many of the recognition problems for which there are currently no good solutions in the vision literature, e.g. those involving the processes of Figure 1, are exactly of this type.

Both of these points are addressed in this work, which makes contributions at two levels. On one hand, we introduce a procedure for the design of large-margin classifiers for spatio-temporal autoregressive processes. This includes the derivation of a discriminant distance function (the Kullback-Leibler divergence) for this class of processes and its application to the design of probabilistic kernels. On the other, we demonstrate the practical feasibility of large margin classification for vision problems involving complex spatio-temporal visual stimuli, such as the classification of dynamic textures or the classification of patterns of highway traffic flow under variable environmental conditions. The new large-margin solution is shown to perform well above the state of the art and to produce quite promising results for difficult problems, such as monitoring highway congestion.

2. Modeling motion flow

Various representations of a video sequence as a spatio-temporal texture have been proposed in the vision literature over the last decade. Earlier efforts were aimed at the

extraction of features that capture both the spatial appearance of a texture and the associated motion flow field. For example, in [10], temporal textures are represented by the first and second order statistics of the normal flow of the video. These types of strictly feature-based representation can be useful for recognition but do not provide a probabilistic model that could be used for kernel design.

More recently, various authors proposed to model a temporal texture as a generative process, resulting in representations that can be used for both synthesis and recognition. One example is the multi-resolution analysis tree method of [11], which represents a temporal texture as the hierarchical multi-scale transform associated with a 3D wavelet. The conditional probability distributions of the wavelet coefficients in the tree are estimated from a collection of training examples and the texture is synthesized by sampling from this model. Another possibility is the spatio-temporal autoregressive (STAR) model of [5], which models the interaction of pixels within a local neighborhood over both space and time. By relying on spatio-temporally localized image features these representations are incapable of abstracting the video into a pair of holistic appearance and motion components.

This problem is addressed by the dynamic texture model of [6], an auto-regressive random process (specifically, a linear dynamical system) that includes a hidden state variable, in addition to the observation variable that determines the appearance component. The motion flow of the video is captured by a dynamic generative model, from which the hidden state vector is drawn. The observation vector is then drawn from a second generative model, conditioned on the state variable. Both the hidden state vector and the observation vector are representative of the entire image, enabling a holistic characterization of the motion for the entire sequence. For this reason, we adopt the dynamic texture model in the remainder of this work.



Figure 2: Example of the dynamic texture model: (top-left) frames from a traffic sequence; (bottom-left) the first three principal components; (right) the state space trajectory of the corresponding coefficients.

2.1. The dynamic texture model

The dynamic texture model [6] is defined by

$$x_{t+1} = Ax_t + Bv_t \quad (1)$$

$$y_t = Cx_t + w_t \quad (2)$$

where $x_t \in \mathbb{R}^n$ is a sequence of n -dimensional hidden (state) random variables, $y_t \in \mathbb{R}^m$ a m -dimensional sequence of observed (video frame) random variables, $v_t \sim_{iid} \mathcal{N}(0, I_{n_v})$ a n_v -dimensional driving process (typically $n \ll m$ and $n_v \leq n$), and $w_t \sim_{iid} \mathcal{N}(0, R)$ an observation noise process. The model is parameterized by $\Theta = (A, B, C, R, x_0)$, where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times n_v}$, $C \in \mathbb{R}^{m \times n}$, $R \in \mathbb{R}^{m \times m}$, and x_0 is a known initial state. Note that $Bv_t \sim \mathcal{N}(0, Q)$ where $Q = BB^T$. The covariance of the observation noise is assumed to be i.i.d, i.e. $R = \sigma^2 I_m$.

The sequence $\{y_t\}$ encodes the appearance component of the video (video frames), and the motion component is encoded into the state sequence $\{x_t\}$. The hidden state is modeled as a first-order Gauss-Markov process, where the state at time $t + 1$, x_{t+1} , is determined by the transition matrix A , the state at time t , x_t , and the driving process v_t . The image at time t , y_t , is a linear combination of the principal components of the entire video sequence, stored in the columns of C , with each component weighted by the corresponding coefficient in the state vector x_t . Figure 2 shows an example of a traffic sequence, its first three principal components, and the corresponding state space coefficients.

2.2. Parameter estimation

Given an image sequence (y_1, \dots, y_N) , it is possible to learn the parameters of the dynamic texture which best models the image observations. While asymptotically optimal solutions, in the maximum likelihood sense, exist (e.g. N4SID [12]) the high dimensionality of the observed image space makes such solutions infeasible for dynamic texture models. A suboptimal (but tractable) alternative

[6] is to learn the spatial and temporal parameters separately. If $Y_1^N = [y_1, \dots, y_N] \in \mathbb{R}^{m \times N}$ is the matrix of observed video frames, its singular value decomposition (SVD) $Y_1^N = U\Sigma V^T$ is a natural decomposition into 1) principal components (columns of U) and 2) corresponding state vectors (columns of ΣV^T). It is therefore natural to rely on estimates of the form

$$\hat{C} = U \quad \hat{X}_1^N = \Sigma V^T \quad (3)$$

where $\hat{X}_1^N = [\hat{x}_1, \dots, \hat{x}_N]$ is a matrix of state estimates for each frame. Given these state estimates, the transition matrix is computed using the least-squares estimate of the linear dependence of the state between consecutive time steps (assuming the state random variables have zero mean),

$$\hat{A} = \hat{X}_2^N (\hat{X}_1^{N-1})^\dagger \quad (4)$$

where $M^\dagger = M^T(MM^T)^{-1}$ is the pseudo-inverse of M . Finally, the estimate of the covariance of the driving process is

$$\hat{Q} = \frac{1}{N-1} \sum_{i=1}^{N-1} \hat{v}_i \hat{v}_i^T \quad (5)$$

where $\hat{v}_t = \hat{x}_{t+1} - \hat{A}\hat{x}_t$.

3. Support vector machines and probabilistic kernels

A support vector machine (SVM) [8] is a discriminative classifier that constructs a maximum-margin hyperplane between two classes using a set of training examples $\{x_1, \dots, x_N\} \in \mathcal{X}$. The SVM provides strong generalization guarantees for learning and usually leads to improved performance, outside the training set, when compared to classical methods based on Bayesian decision theory [13]. The training examples that are most difficult to classify are

referred to as support vectors, and determine the separating hyperplane.

The SVM can be augmented by using the “kernel” trick, which maps the training examples into a high-dimensional non-linear feature space. This feature space transformation is defined by the kernel function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. One interpretation of the kernel function is that $K(x_i, x_j)$ measures the similarity between the two points x_i and x_j in the space \mathcal{X} . A popular example is the Gaussian kernel, defined as $K_g(x_i, x_j) = \exp(-\gamma\|x_i - x_j\|^2)$.

If the training examples are represented as probabilistic models (e.g. dynamic textures), the kernel becomes a measure of similarity between probability distributions. A probabilistic kernel is thus defined as a mapping $K : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$, where \mathcal{P} is the space of probability distributions. One such kernel is the Kullback-Leibler kernel [14], defined as

$$K_{KL}(p, q) = e^{-\gamma(D(p\|q) + D(q\|p))} \quad (6)$$

where $D(p\|q)$ is the Kullback-Leibler (KL) divergence between the probability distributions $p(x)$ and $q(x)$ [16]

$$D(p\|q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx. \quad (7)$$

The KL divergence is a natural distance measure between two probability distributions, and the KL kernel in probability space is analogous to the Gaussian kernel in Euclidean space. The KL kernel has been shown to achieve very good results in the domains of object [9] and speech [14] recognition.

4. Probabilistic kernels for dynamic textures

In this section we introduce a probabilistic kernel for visual processes that can be modeled as dynamic textures.

4.1. General considerations

The dynamic texture model provides a probability distribution of the texture in both image and state space. This allows the derivation of two kernels that can ground the classification in either the appearance or the flow of the dynamic texture. Grounding the classification on image space tends to favor iconic pixel matches and leads to best performance when the goal is to differentiate between dynamic textures of different visual appearance (e.g. a flock of birds from a school of fish in Figure 1). Under this approach, two sequences of distinct textures subject to similar motion can be correctly identified. It is, however, not clear that the dynamic texture model is of great significance in this context: a simple appearance classifier based, for example, on the principal component decomposition of the sequences may

achieve good results. Ideally, the kernel based on the dynamic texture model should achieve performance at least as good as that of a static kernel, when this is the case.

An alternative classification scenario is that where the different classes have similar appearance and all the discriminant information is contained in the motion flow. For example, problems such as determining the level of traffic on a highway, or detecting outliers and unusual events (e.g. cars speeding or committing other traffic violations). Since for these cases the iconic pixel matching inherent to existing static kernels is clearly inappropriate, these are the cases where dynamic kernels have the greatest potential for improvement over the state of the art.

In summary, depending on the specific classification problem, it may be advisable to ground the classification on either the state space or the image space components of the dynamic texture model. In the remainder of this section we derive the KL kernel for these two representations.

4.2. Probability distributions

We start by obtaining the probability distributions of the Gauss-Markov process [15] that models the state of the dynamic texture. The conditional probability of state x_t given state x_{t-1} follows from (1),

$$\begin{aligned} p(x_t|x_{t-1}) &= G(x_t, Ax_{t-1}, Q) & (8) \\ &= \frac{1}{\sqrt{(2\pi)^n |Q|}} e^{-\frac{1}{2}\|x_t - Ax_{t-1}\|_Q^2} & (9) \end{aligned}$$

where $\|x\|_Q^2 = x^T Q^{-1} x$. Recursively substituting into (1),

$$x_t = A^t x_0 + \sum_{i=1}^t A^{t-i} B v_i \quad (10)$$

where the initial state x_0 is known. Since x_t is the sum of $t - 1$ Gaussian random variables, it is also Gaussian

$$p(x_t) = G(x_t, \mu_t, S_t) \quad (11)$$

with mean and covariance given by

$$\mu_t = A^t x_0 = A \mu_{t-1} \quad (12)$$

$$S_t = \sum_{i=0}^{t-1} A^i Q (A^i)^T = A S_{t-1} A^T + Q \quad (13)$$

Let $x_1^\tau = (x_1, \dots, x_\tau)$ be a sequence of τ state vectors. The probability of a state sequence is also Gaussian, and can be expressed using conditional probabilities as

$$p(x_1^\tau) = p(x_1) \prod_{i=2}^{\tau} p(x_i|x_{i-1}) \quad (14)$$

$$= G(x_1^\tau, \mu, \Sigma) \quad (15)$$

where $\mu = [\mu_1^T \ \cdots \ \mu_\tau^T]^T$ and the covariance is

$$\Sigma = \begin{bmatrix} S_1 & (AS_1)^T & \cdots & (A^{\tau-1}S_1)^T \\ AS_1 & S_2 & \cdots & (A^{\tau-2}S_2)^T \\ \vdots & \vdots & \ddots & \vdots \\ A^{\tau-1}S_1 & A^{\tau-2}S_2 & \cdots & S_\tau \end{bmatrix} \quad (16)$$

The image sequence y_1^τ is a linear transformation of the state sequence, and is thus given by

$$p(y_1^\tau) = G(y_1^\tau, \gamma, \Phi) \quad (17)$$

where $\gamma = \mathbf{C}\mu$ and $\Phi = \mathbf{C}\Sigma\mathbf{C}^T + \mathbf{R}$, and \mathbf{C} and \mathbf{R} are block diagonal matrices formed from C and R respectively.

4.3. Projection between state spaces

The KL divergence between state spaces cannot be computed directly because each dynamic texture uses a different PCA space. Instead, one state space must be projected into the other by applying a sequence of two transformations: 1) from the original state space into image space, and 2) from image space into the target state space. If the original state space is that of x_1 and the target that of x_2 , this is the transformation $\hat{x}_1 = Fx_1$ with $F = C_2^T C_1$. From (1),

$$x_{t+1} = A_1 x_t + B_1 v_t \quad (18)$$

$$F x_{t+1} = F A_1 F^{-1} F x_t + F B_1 v_t \quad (19)$$

$$\hat{x}_{t+1} = \hat{A}_1 \hat{x}_t + \hat{B}_1 v_t \quad (20)$$

and the transformation of a Gauss-Markov process with parameters (A_1, B_1, x_{01}) is a Gauss-Markov process with parameters $\hat{A}_1 = (C_2^T C_1) A_1 (C_2^T C_1)^{-1}$, $\hat{B}_1 = (C_2^T C_1) B_1$, and $\hat{x}_{01} = (C_2^T C_1) x_{01}$. The KL divergence between state spaces can now be computed with this transformed state model.

4.4. KL divergence between state spaces

The KL divergence rate between two random processes with distributions, $p(X)$ and $q(X)$ over $X = (x_1, x_2, \dots)$, is defined as

$$D(p(X) \| q(X)) = \lim_{t \rightarrow \infty} \frac{1}{\tau} D(p(x_1^\tau) \| q(x_1^\tau)). \quad (21)$$

If $p(x_1^\tau)$ and $q(x_1^\tau)$ are the state probability distributions of two dynamic textures parameterized by (A_1, Q_1, x_{01}) and (A_2, Q_2, x_{02}) , the KL divergence on the RHS of (21) is (see Appendix for derivation),

$$\begin{aligned} \frac{1}{\tau} D(p(x_1^\tau) \| q(x_1^\tau)) &= \frac{1}{2} \left[\log \frac{|Q_2|}{|Q_1|} \right. \\ &+ \text{tr}(Q_2^{-1} Q_1) - n + \frac{1}{\tau} \|A_1 x_{01} - A_2 x_{02}\|_{Q_2}^2 \\ &\left. + \frac{1}{\tau} \sum_{i=2}^{\tau} \text{tr}(\bar{A}^T Q_2^{-1} \bar{A} (S_{i-1} + \mu_{i-1} \mu_{i-1}^T)) \right] \end{aligned} \quad (22)$$

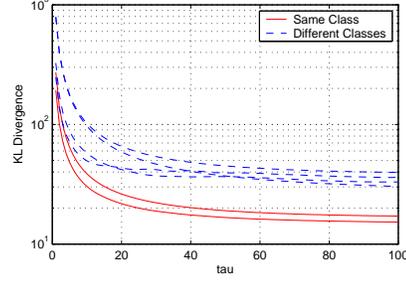


Figure 3: Example of the convergence of the KL divergence rate versus τ . The clustering of generative models in the same and different classes can also be seen.

where $\bar{A} = A_1 - A_2$, and S_{i-1} and μ_{i-1} are the covariance and mean associated with the state x_{i-1} of the first dynamic texture, as calculated in (12) and (13). In practice, the KL divergence rate can be estimated by setting τ to some large number. Figure 3 shows a graph of some examples of the KL between dynamic textures in state space. The KL rate converges as τ increases, and the clustering between generative models in the same and different classes can be seen.

4.5. KL divergence in image space

The KL divergence rate between two image sequence distributions, $p(Y)$ and $q(Y)$ over $Y = (y_1, y_2, \dots)$ is

$$D(p(Y) \| q(Y)) = \lim_{t \rightarrow \infty} \frac{1}{\tau} D(p(y_1^\tau) \| q(y_1^\tau)) \quad (23)$$

The image probabilities $p(y_1^\tau)$ and $q(y_1^\tau)$ are both distributed as Gaussians with means γ_1 and γ_2 and covariance Φ_1 and Φ_2 respectively, thus the KL divergence between them is,

$$\begin{aligned} D(p(y_1^\tau) \| q(y_1^\tau)) &= \frac{1}{2} \left[\log \frac{|\Phi_2|}{|\Phi_1|} \right. \\ &\left. + \text{tr}(\Phi_2^{-1} \Phi_1) + \|\gamma_1 - \gamma_2\|_{\Phi_2}^2 - m\tau \right] \end{aligned} \quad (24)$$

Direct computation of the KL divergence between image sequences is intractable since the covariance matrices are $m\tau \times m\tau$, where m is the number of pixels in a frame. Using several matrix identities, it is possible to rewrite the terms of the image KL into a recursive form that is computationally efficient and only requires storing $n\tau \times n\tau$ matrices (recall $n \ll m$). For brevity, we omit the details here and refer the reader to a companion tech report [17].

5. Experimental evaluation

We evaluate the performance of motion flow recognition using the KL kernel on two video databases. The first database contains many visually distinct classes. The second database, based on traffic video, contains visually similar classes, but with varying temporal characteristics.

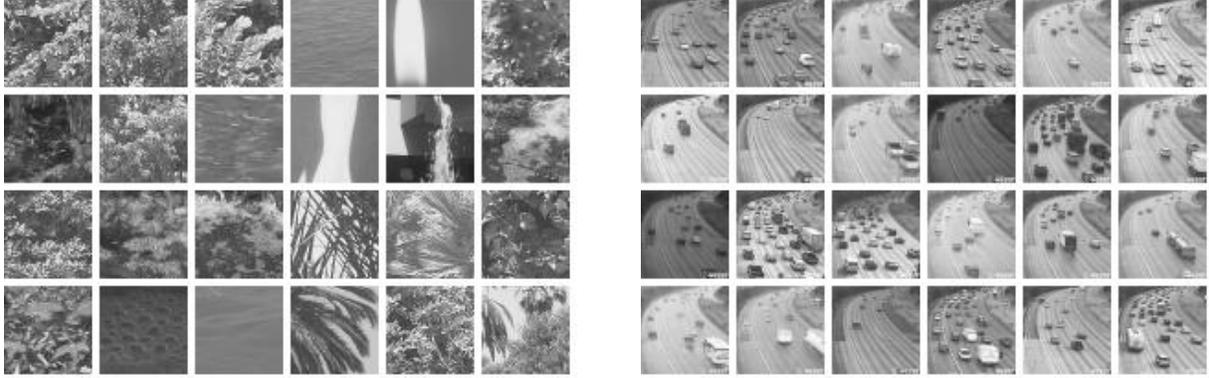


Figure 4: Examples from the databases used for evaluation: (left) the dynamic texture database; (right) the traffic video database.

5.1. Dynamic texture database

The dynamic texture database used in [18] contains 50 classes of various texture, including boiling water, fountains, fire, rippling water, waterfalls, and plants and flowers swaying in the wind. Each class contains four grayscale video sequences with 75 frames¹ of 160×110 pixels. Each sequence was clipped to a 48×48 window that contained the representative motion. Figure 4 (left) shows several examples of the video patches from the dynamic texture database. Since almost all of the classes are visually distinct, the appearance component of the model is likely to be as important for classification as the motion component.

5.2. Traffic video database

The traffic video database consists of 254 video sequences of highway traffic in Seattle, collected from a single stationary traffic camera over two days [19]. The database contains a variety of traffic patterns and weather conditions (e.g. overcast, raining, sunny, rain drops on the camera lens). Each video was recorded in color with a resolution of 320×240 pixels with between 42 to 52 frames at 10 fps. Each sequence was converted to grayscale, resized to 80×60 pixels, and then clipped to a 48×48 window over the area with the most total motion. Finally, for each video clip, the mean image was subtracted and the pixel intensities were normalized to have unit variance. This was done to reduce the impact of the different lighting conditions.

The database was labeled by hand with respect to the amount of traffic congestion in each sequence. In total there were 44 sequences of heavy traffic (slow or stop and go speeds), 45 of medium traffic (reduced speed), and 165 of light traffic (normal speed). Figure 4 (right) shows a representative set of clips from this database. All clips are very similar in that the views are obtained with a fixed camera

facing the same stretch of road, and the motion is always in the same direction and confined to the same area. Thus, an effective classifier for this problem must be able to distinguish between the different patterns of flow, i.e. the underlying temporal process.

5.3. Experiment setup

The parameters of the dynamic texture model were learned for each video clip using the method of Section 2.2. To ensure that the KL divergence converges, the transition matrix A was scaled so that the largest eigenvalues lie on the unit circle. In addition, the covariance of the driving process was regularized to prevent problems with singular matrices, i.e. we set $Q' = Q + I_n$. All classification results were averaged over four trials. In each trial the data set was split differently with 75% used for training and cross-validation, and 25% reserved for testing.

For the dynamic texture database, SVMs were trained using the KL kernel in image space ($\tau = 25$), and for the traffic video database, the SVMs were trained with the KL kernel in state space ($\tau = 250$). A one-versus-all scheme was used to learn the multi-class problem, and the C and γ parameters were selected using 3-fold cross-validation over the training set. The SVM training and testing was performed using the `libsvm` software package [20]. We also tested a nearest neighbor (NN) classifier using the image space and state space KL as distance measures. Finally, for comparison with the state-of-the-art, a nearest neighbor classifier was implemented using the Martin distance [21, 22] as suggested in [18]. The Martin distance is related to the principal angles between the subspaces of the extended observability matrices of two dynamic textures. For this experiment, the extended observability matrices were approximated with $\tau = 250$.

¹The four videos in each class originate from 2 videos with 150 frames each.

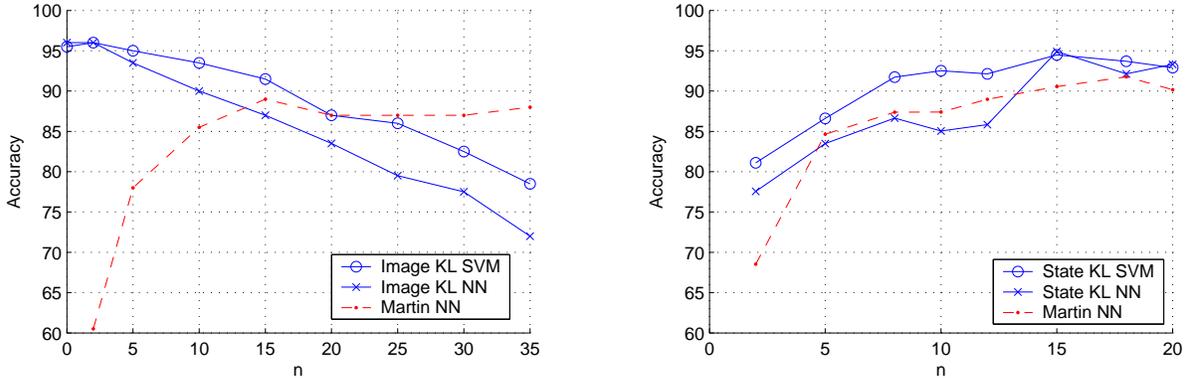


Figure 5: Evaluation of the KL kernel on two databases: (left) classification accuracy on the dynamic texture database using the SVM with the image KL kernel; (right) classification accuracy on the traffic video database using the SVM with the state KL kernel. In both plots, the accuracy of nearest neighbors classification using the appropriate KL distance and the Martin distance is also shown. The x-axis is the number of principal components (n) used in the dynamic texture model.

5.4. Results

The results on the dynamic texture database, presented in Figure 5 (left), show that the image-based KL classifiers performs significantly better than the Martin distance classifier (an improvement of the best classification accuracy from 89% to 96%). Note that the accuracies of the image KL classifiers improve as the number of principal components n decreases. This was expected since the dynamic texture database contains many visually distinct classes for which the appearance components is more discriminant than the motion. In fact, in the degenerate case of $n = 0$, the video is modeled as a Gaussian whose mean is the mean image of the video sequence, and covariance is the deviation of the frames from this mean. Note how, by achieving top performance for a small number of components, the image-based KL classifiers virtually become static classifiers. In contrast, the Martin distance nearest neighbors classifier does rather poorly with a small number of components. Hence, although performance improves as n increases, it never reaches an accuracy comparable to that of the KL-based classifiers.

Figure 5 (right) presents the results obtained on the traffic video database. It can be seen from this figure that the two state KL classifiers outperform the Martin NN classifier on this database. Furthermore, all classifiers improve as the number of principal components increases, confirming the fact that a static classifier would do rather poorly on this database. Comparing the performance of the state KL classifiers versus the Martin NN counterpart it can be concluded that 1) the state KL-SVM combination is consistently better, and 2) the state KL-NN combination is better for $n \geq 15$ and also achieves a higher maximum accuracy.

Overall, the image and state KL classifiers outperform the Martin distance nearest neighbor method in classification tasks with both visually distinct video textures, and vi-

usually similar, but temporally distinct, video textures. The KL classifiers are also capable of spanning the gamut from static to highly-varying dynamic classifier and, therefore, provide a generic framework for the classification of a large variety of video streams. Comparing the performance of the two KL classifiers, it is clear that SVM-KL combination achieves better classification performance than NN-KL. In particular, the greater robustness of the SVM classifier to a poor selection of the number of components indicates that it has better generalization ability.

Finally, we tested the robustness of the dynamic texture and KL classification framework by using the trained classifiers to label a set of 12 sequential traffic videos that spanned an hour at night, and contained a traffic jam (Figure 6). The NN-KL and SVM-KL correctly labeled the 12 video sequences, including the event of the traffic jam (heavy traffic), and the events of reduced speed (medium traffic leading up to and immediately following the traffic jam). This is particularly interesting since the classifiers were trained with daytime images containing normally lit cars, yet they are able to correctly label nighttime images where the cars are represented as headlights and a pair of tail lights. These results provide evidence that the dynamic texture model is indeed extracting relevant motion information, and that the proposed classification framework is capable of using the motion model to discriminate between classes of motion. We are currently exploring how this framework could be used for tracking highway congestion, and detection of outlier events such as speeding cars and accidents.

Appendix

Given that $p(x)$ and $q(x)$ are distributions of Markov processes, (21) can be simplified using the chain rule of diver-

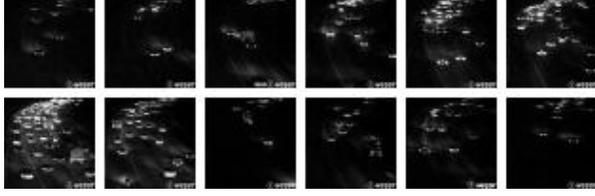


Figure 6: Twelve nighttime traffic videos with a traffic jam.

gence [23],

$$D(p(x_1^T) \| q(x_1^T)) = D(p(x_1) \| q(x_1)) + \sum_{i=2}^{\tau} D(p(x_i | x_{i-1}) \| q(x_i | x_{i-1}))$$

For two dynamic textures, the KL of the first state vector is

$$D(p(x_1) \| q(x_1)) = \frac{1}{2} \|A_1 x_{01} - A_2 x_{02}\|_{Q_2}^2 + \frac{1}{2} \log \frac{|Q_2|}{|Q_1|} + \frac{1}{2} \text{tr}(Q_2^{-1} Q_1) - \frac{n}{2}$$

and the conditional KL term is

$$\begin{aligned} & D(p(x_i | x_{i-1}) \| q(x_i | x_{i-1})) \\ &= \int p(x_{i-1}) \int G(x_i, A_1 x_{i-1}, Q_1) \\ & \quad \cdot \log \frac{G(x_i, A_1 x_{i-1}, Q_1)}{G(x_i, A_2 x_{i-1}, Q_2)} dx_i dx_{i-1} \\ &= \int p(x_{i-1}) \frac{1}{2} \left[\|(A_1 - A_2)x_{i-1}\|_{Q_2}^2 + \log \frac{|Q_2|}{|Q_1|} \right. \\ & \quad \left. + \text{tr}(Q_2^{-1} Q_1) - n \right] dx_{i-1} \\ &= \frac{1}{2} \left[\text{tr}(\bar{A}^T Q_2^{-1} \bar{A} (S_{i-1} + \mu_{i-1} \mu_{i-1}^T)) + \log \frac{|Q_2|}{|Q_1|} \right. \\ & \quad \left. + \text{tr}(Q_2^{-1} Q_1) - n \right] \end{aligned}$$

where $\bar{A} = A_1 - A_2$, and in the last line we have used the property that if $p(x)$ has mean μ and covariance Σ ,

$$\int p(x) \|Ax\|_B^2 dx = \text{tr}(A^T B^{-1} A (\Sigma + \mu \mu^T))$$

Finally, (22) is obtained by summing the initial KL term and the conditional KL terms from above.

Acknowledgments

We wish to thank Gianfranco Doretto and Stefano Soatto for providing the dynamic texture database used in [18].

References

- [1] B. Horn and B. Schunk. Determining Optical Flow. *Artificial Intelligence*, Vol. 17, 1981.
- [2] B. Lucas and T. Kanade. An Iterative Image Registration Technique with an Application to Stereo Vision. In *Proc. DARPA Image Understanding Workshop*, 1981.
- [3] M. Isard and A. Blake. Condensation – Conditional Density Propagation for Visual Tracking. *International Journal of Computer Vision*, Vol. 29(1), pp. 5-28, 1998.
- [4] J. Wang and E. Adelson. Representing Moving Images with Layers. *IEEE Trans. on Image Processing*, Vol. 3, September 1994.
- [5] M. Szummer and R. Picard. Temporal texture modeling. In *IEEE Conference on Image Processing*, volume 3, pages 823–6, 1996.
- [6] G. Doretto, A. Chiuso, Y. N. Wu, and S. Soatto. Dynamic textures. *International Journal of Computer Vision*, (2):91–109, 2003.
- [7] D. Knull and W. Richards. *Perception as Bayesian Inference*. Cambridge Univ. Press, 1996.
- [8] V. Vapnik. *The Nature of Statistical Learning Theory*. Springer Verlag, 1995.
- [9] N. Vasconcelos, P. Ho, and P. Moreno. The Kullback-Leibler Kernel as a Framework for Discriminant and Localized Representations for Visual Recognition. In *Proc. European Conference on Computer Vision, Prague, Czech Republic*, 2004.
- [10] R. Polana and R. C. Nelson. Recognition of motion from temporal texture. In *IEEE Conference on Computer Vision and Pattern Recognition, Proceedings*, pages 129–34, 1992.
- [11] Z. Bar-Joseph, R. El-Yaniv, D. Lischinski, and M. Werman. Texture mixing and texture movie synthesis using statistical learning. *IEEE Transactions on Visualization and Computer Graphics*, 7(2):120–35, 2001.
- [12] P. Van Overschee and B. De Moor. N4sid: Subspace algorithms for the identification of combined deterministic-stochastic systems. *Automatica*, 30:75–93, 1994.
- [13] R. Duda, P. Hart, and D. Stork. *Pattern Classification*. John Wiley and Sons, 2001.
- [14] P. J. Moreno, P. Ho, and N. Vasconcelos. A kullback-leibler divergence based kernel for svm classification in multimedia applications. In *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec 2003.
- [15] S. M. Kay. *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice-Hall Signal Processing Series, 1993.
- [16] S. Kullback. *Information Theory and Statistics*. Dover Publications, New York, 1968.
- [17] A. B. Chan and N. Vasconcelos. Efficient Computation of the KL Divergence between Dynamic Textures. Technical Report SVCL-TR-2004-02, <http://www.svcl.ucsd.edu/>, November 2004.
- [18] P. Saisan, G. Doretto, Y. Wu, and S. Soatto. Dynamic texture recognition. In *IEEE Conference on Computer Vision and Pattern Recognition, Proceedings*, volume 2, pages 58–63, 2001.
- [19] <http://www.wsdot.wa.gov>
- [20] Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: a library for support vector machines*, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [21] R. J. Martin. A metric for arma processes. *IEEE Transactions on Signal Processing*, 48(4):1164–70, April 2000.
- [22] K. De Cock and B. De Moor. Subspace angles between linear stochastic models. In *IEEE Conference on Decision and Control, Proceedings*, pages 1561–6, December 2000.
- [23] T. Cover and J. Thomas. *Elements of Information Theory*. John Wiley, 1991.