
Background Modeling Using Adaptive Pixelwise Kernel Variances in a Hybrid Feature Space

Abstract

Recent work on background subtraction has shown developments on two major fronts. In one, there has been increasing sophistication of probabilistic models, from mixtures of Gaussians at each pixel, to kernel density estimates at each pixel, and more recently to joint domain-range density estimates that incorporate spatial information. Another line of work has shown the benefits of increasingly complex feature representations, including the use of texture information, local binary patterns, and recently scale-invariant local ternary patterns. In this work, we use joint domain-range based estimates for background and foreground scores and show that dynamically choosing kernel variances in our kernel estimates at each individual pixel can significantly improve results. We give a heuristic method for selectively applying the adaptive kernel calculations which is nearly as accurate as the full procedure but runs much faster. We combine these modeling improvements with recently developed complex features and show significant improvements on a standard backgrounding benchmark.

1 Introduction

Background modeling is often an important step in detecting moving objects in video sequences. A common approach to background modeling is to define and learn a background distribution over feature values at each pixel location and then classify each image pixel as belonging to the background process or not. The distributions at each pixel may be modeled in a parametric manner using a mixture of Gaussians or using non-parametric kernel density estimation. More recently, models that allow a pixel's spatial neighbors to influence its distribution have been developed by joint domain-range density estimation. These models that allow spatial influence from neighboring pixels have been shown to perform better than earlier neighbor-independent models.

Also, the use of an explicit foreground model along with a background model can be useful. In a manner similar to theirs, we use a kernel estimate to obtain the background and foreground scores at each pixel location using data samples from a spatial neighborhood around that location from previous frames. The background score is computed as a kernel estimate depending on the distance in the joint domain-range space between the estimation point and the samples in the background model. A similar estimate is obtained for the foreground score. Each pixel is then assigned a (soft) label based on the ratio of the background and foreground scores.

The variance used in the estimation kernel reflects the spatial and appearance uncertainties in the scene. On applying our method to a data set with wide variations across the videos, we found that choosing suitable kernel variances during the estimation process is very important. With various experiments, we establish that the best kernel variance could vary for different videos and more importantly, even within a single video, different regions in the image should be treated with different variance values. For example, in a scene with a steady tree trunk and leaves that are waving in the wind, the trunk region can be explained with a small amount of spatial variance. The leaf regions may be better explained by a process with a large variance. Interestingly, when there is no wind, the leaf regions may also be explained with a low variance. The optimal variance hence changes for

each region in the video and also across time. This phenomenon is captured reasonably in MoG by use of different parameters for each pixel which adapt dynamically to the scene statistics, but the pixel-wise model does not allow a pixel’s neighbors to affect its distribution. address the phenomenon by updating the model with data samples from the most recent frame. We show that using location-specific variances in addition to updating the model greatly improves background modeling. Our approach with pixel-wise variances, which we call the variable kernel score (VKS) method results in significant improvement over uniform variance models and state of the art backgrounding systems.

The idea of using a pixel-wise variance for background modeling is not new. Although use a uniform variance, they discuss the use of variances that change as a function of the data samples or as a function of the point at which the estimation is made. Variance selection for KDE is a well studied problem with common solutions including mean integrated square error (MISE), asymptotic MISE (AMISE), and the leave-one-out-estimator based solutions. In the background subtraction context, there has been work on using a different covariance at each pixel. While require that the uncertainties in the feature values can be calculated in closed form, learn the covariances for each pixel from a training set of frames and keep the learned covariances fixed for the entire classification phase. We use a maximum-likelihood approach to select the best variance at each pixel location. For every frame of the video, at each pixel location, the best variance is picked from a set of variance values by maximizing the likelihood of the pixel’s observation under different variances. This makes our method a balloon estimator. By explicitly selecting the best variance from a range of variance values, we do not require the covariances to be calculable in closed-form and also allow for more flexibility at the classification stage.

Selecting the best of many kernel variances for each pixel means increased computation. One possible trade-off between accuracy and speed can be achieved by a caching scheme where the best kernel variances from the previous frame are used to calculate the scores for the current frame pixels. If the resulting classification is overwhelmingly in favor of either label, there is no need to perform a search for the best kernel variance for that pixel. The expensive variance selection procedure can be applied only to pixels where there is some contention between the two labels. We present a heuristic that achieves significant reduction in computation compared to our full implementation while maintaining the benefits of adaptive variance.

Development and improvement of the probabilistic models is one of the two main themes in background modeling research in recent years. The other theme is the development of complex features like local binary and ternary patterns that are more robust than color features for the task of background modeling. Scale-invariant local ternary patterns (SILTP) are recently developed features that have been shown to be very robust to lighting changes and shadows in the scene. By combining color features with SILTP features in our adaptive variance kernel model, we bring together the best ideas from both themes in the field and achieve state of the art results on a benchmark data set.

The main contributions of this paper are:

1. A practical scheme for pixel-wise variance selection for background modeling.
2. A heuristic for selectively updating variances to improve speed further.
3. Incorporation of complex SILTP features into the joint domain-range kernel framework to achieve state of the art results.

The paper is organized as follows. Section 2 discusses our background and foreground models. Dynamic adaptation of kernel variances is discussed in Section 3. Results and comparisons are in Section 4. An efficient algorithm is discussed in Section 5. We end with a discussion in Section 6.

2 Background and foreground models

In a video captured by a static camera, the pixel values are influenced by the background phenomenon, and new or existing foreground objects. We refer to any phenomenon that can affect image pixel values as a process. Like , we model the background and foreground processes using data samples from previous frames. The scores for the background and foreground processes at each pixel location are calculated using contributions from the data samples in each model. One major difference between and our model is that we allow “soft labeling”, i.e. the data samples contribute probabilistically to the background score depending on the samples’ probability of belonging to the background.

Let a pixel sample $a = [ax, ay, ar, ag, ab]$, where (ax, ay) are the location of the pixel and (ar, ag, ab) are the red, green, and blue values of the pixel. In each frame of the video, we compute background and foreground scores using pixel samples from the previous frames. The background model consists of the samples $B = b_i : i [1 : n_B]$ and foreground samples are $F = f_i : i [1 : n_F]$, with n_B and n_F being the number of background and foreground samples respectively, and b_i and f_i being pixel samples obtained from previous frames in the video. Under a KDE model, the likelihood of the sample under the background model is

$$P(a|bg; \sigma) = \frac{1}{n_B} \sum_{i=1}^{n_B} G(a - b_i; \sigma_B) \quad (1)$$

where $G(x; \sigma)$ is a multivariate Gaussian with zero mean and covariance B .

$$G(x; \sigma) = (2\pi)^{-\frac{D}{2}} |\sigma|^{-\frac{1}{2}} \exp(-\frac{1}{2} x^T \sigma^{-1} x), \quad (2)$$

where D is the dimensionality of the vector x .

In our model, we approximate the background score at sample a as

$$S_B(a; \sigma_B^d, \sigma_B^{rgb}) = \frac{1}{N_B} \sum_{i=1}^{N_B} G(a_{rgb} - b_{i_{rgb}}; \sigma_B^{rgb}) \times G(a_{xy} - b_{i_{xy}}; \sigma_B^d) \times P(bg|b_i) \quad (3)$$

N_B is the number of frames from which the background samples have been collected, B^d and B^{rgb} are two and three dimensional background covariance matrices in spatial and color dimensions respectively. A large spatial covariance allows neighboring pixels to contribute more to the score at a given pixel location. Color covariance allows for some color appearance changes at a given pixel location. Use of N_B in the denominator compensates for the different lengths of the background and foreground models.

The above equation basically sums the contribution from each background sample based on its distance in color space, weighted by its distance in spatial dimensions and the probability of the sample belonging to the background.

The use of $P(bg|b_i)$ in Equation 3 and normalization by the number of frames as opposed to the number of samples means that the score does not sum to 1 over all possible values of a . Thus, the score, although similar to the likelihood in Equation 1, is not a probability distribution.

A similar equation holds for the foreground score:

$$S_F(a; \sigma_F^d, \sigma_F^{rgb}) = \frac{1}{N_F} \sum_{i=1}^{N_F} G(a_{rgb} - f_{i_{rgb}}; \sigma_F^{rgb}) \times G(a_{xy} - f_{i_{xy}}; \sigma_F^d) \times P(fg|f_i) \quad (4)$$

N_F is the number of frames from which the foreground samples have been collected, F^d and F^{rgb} are the covariances associated with the foreground process.

However, for the foreground process, to account for emergence of new colors in the scene, we mix in a constant contribution independent of the estimation point's and data samples' color values. We assume that each data sample in a pixel's spatial neighborhood contributes a constant value u to the foreground score. The constant contribution $U_F(a)$ is given by

$$U_F(a; \sigma_F^d) = \sum_{i=1}^{N_F} u \times G(a_{xy} - f_{i_{xy}}; \sigma_F^d) \quad (5)$$

We get a modified foreground score by including the constant contribution:

$$\hat{S}_F(a; \sigma_F^d, \sigma_F^{rgb}) = \alpha_F \times U_F(a; \sigma_F^d) + (1 - \alpha_F) \times S_F(a; \sigma_F^d, \sigma_F^{rgb}). \quad (6)$$

F is a parameter that represents the amount of mixing between the constant contribution and the color dependent foreground score. u is set to 106 and α_F is set to 0.5 for our experiments.

To classify a particular sample as background or foreground, we can use a Bayes-like formula:

$$P(bg|a) = \frac{S_B(a; \sigma_B^d, \sigma_B^{rgb})}{S_B(a; \sigma_B^d, \sigma_B^{rgb}) + \hat{S}_F(a; \sigma_F^d, \sigma_F^{rgb})} \quad (7)$$

$$P(fg|a) = 1 - P(bg|a). \quad (8)$$

Adding the constant factor U to the foreground score (and hence to the denominator of the Bayes-like equation) has the interesting property that when either one of the foreground or background scores is significantly larger than U , U has little effect on the classification. However, if both the background and foreground scores are less than U , then Equation 7 will return a low value as $P(bg|a)$. Hence, an observation that has very low background and foreground scores will be classified as foreground. This is desirable because if a pixel observation is not well explained by either model, it is natural to assume that the pixel is a result of a new object in the scene and is hence foreground. In terms of likelihoods, adding the constant factor to the foreground likelihood is akin to mixing it with a uniform distribution.

2.1 Model initialization and update

To initialize the models, it is assumed that the first few frames (typically 50) are all background pixels. The background model is populated using pixel samples from these frames. In order to improve efficiency, we sample 5 frames at equal time intervals from these 50 frames. The foreground model is initialized to have no samples. The modified foreground score (Equation 6) enables colors that are not well explained by the background model to be classified as foreground, thus bootstrapping the foreground model. Once the pixel at location (ax, ay) from a new frame is classified using Equation 7, the background and foreground models at the location (ax, ay) can then be updated with the new sample a . Background and foreground samples at location (ax, ay) from the oldest frame in the models are replaced by a . Samples from the previous 5 frames are maintained in memory as the foreground model samples. The label probabilities of the background/foreground from Equation 7 are also saved along with the sample values for subsequent use in the Equations 3 and 4.

One consequence of the update procedure described above is that when a large foreground object occludes a background pixel at (ax, ay) for more than 50 frames, all the background samples in the spatial neighborhood of (ax, ay) are replaced by these foreground samples that have very low $P(bglbi)$ values. This causes the pixel at (ax, ay) to be misclassified as foreground even when the occluding foreground object has moved away (because the background score will be extremely low due to the influence of $P(bglbi)$ in Equation 3). To avoid this problem, we replace the background sample from location (ax, ay) in the oldest frame in the background model with the new sample a from the current frame only if $P(bg|a)$ estimated from Equation 7 is greater than 0.5.

In our chosen evaluation data set, there are several videos with moving objects in the first 50 frames. The assumption that all these pixels are background is not severely limiting even in these videos. The model update procedure allows us to recover from any errors that are caused by the presence of foreground objects in the initialization frames.

2.2 Using MRF to clean the classification

Similar to , we use a Markov random field (MRF) defined over the posterior label probabilities of the 4-neighbors of each pixel and perform the min-cut procedure to post-process the labels. The interaction factor between the nodes was set to 1 for all our experiments.

3 Pixel-wise adaptive kernel variance selection

Background and foreground kernels. use the same kernel parameters for background and foreground models. Given the different nature of the two processes, it is reasonable to use different kernel parameters. For instance, foreground objects typically move between 5 and 10 pixels per frame in the data set, whereas background pixels are either stationary or move very little. Hence, it is useful to have a larger spatial variance for the foreground model than for the background model.

Optimal kernel variance for all videos. In the results section, we show that for a data set with large variations like , a single value for kernel variance for all videos is not sufficient to capture the variability in all the videos.

Variable kernel variance for a single video. As explained in the introduction, different parts of the scene may have different statistics and hence need different kernel variance values. For example, in Figure 1a to 1d, having a high spatial dimension kernel variance helps in accurate classification of

the water surface pixels, but doing so causes some pixels on the person’s leg to become part of the background. Ideally, we would have different kernel variances for the water surface pixels and the rest of the pixels. Similarly in the second video (Figure 1e to 1h), having a high kernel variance allows accurate classification of some of the fountain pixels as background at the cost of misclassifying many foreground pixels. The figure also shows that while the medium kernel variance may be the best choice for the first video, the low kernel variance may be best for the second video.

Optimal kernel variance for classification. Having different variances for the background and foreground models reflects the differences between the expected uncertainty in the two processes. However, having different variances for the two processes could cause erroneous classification of pixels. Figure 2 shows a 1-dimensional example where using a very wide kernel (high variance) or very narrow kernel for the background process causes misclassification. Assuming that the red point (square) is a background sample and the blue point (triangle) is a foreground sample, having a very low variance kernel (dashed red line) or a very high variance (solid red line) for the background process makes the background likelihood of the center point ‘x’ lower than the foreground likelihood. Thus, it is important to pick the optimal kernel variance for each process during classification.

In order to address all four issues discussed above, we propose the use of location-specific variances. For each location in the image, a range of kernel variances is tried and the variance which results in the highest score is chosen for the background and the foreground models separately.

The background score with location-dependent variances is

$$S_B(a; \sigma_{B_d, x, y}, \sigma_{B_{rgb}, x, y}) = \frac{1}{N_B} \sum_{i=1}^{N_B} G(a_{rgb} - b_{i_{rgb}}; \sigma_{B_{rgb}, x, y}) \times G(a_{xy} - b_{i_{xy}}; \sigma_{B_d, x, y}) \times P(bg|b_i) \quad (9)$$

where B_d, x, y and B_{rgb}, x, y represent the location-specific spatial and color dimension variances at location (x, y) .

For each pixel location (ax, ay) , the optimal variance for the background process is selected by maximizing the score of the background label at sample a under different variance values:

$$\{\sigma_{B_d, ax, ay}^*, \sigma_{B_{rgb}, ax, ay}^*\} = \underset{\sigma_{B_d, ax, ay}, \sigma_{B_{rgb}, ax, ay}}{\operatorname{argmax}} S_B(a; \sigma_{B_d, ax, ay}, \sigma_{B_{rgb}, ax, ay}). \quad (10)$$

Here, B_d and B_{rgb} represent the set of spatial and color dimension variances from which to choose the optimal variance.

A similar procedure may be followed for the foreground score. However, in practice, it was found that the variance selection procedure yielded large improvements when applied to the background model and little improvement in the foreground model. Hence, our final implementation uses an adaptive kernel variance procedure for the background model and a fixed kernel variance for the foreground model.

4 Results

For comparisons, we use the data set which consists of 9 videos taken using a static camera in various environments. The data set offers various challenges including dynamic background like trees and waves, gradual and sudden illumination changes, and the presence of multiple moving objects. Ground truth for 20 frames in each video is provided with the data set. The F-measure is used to measure accuracy.

The effect of choosing various kernel widths for the background and foreground models is shown in Table 1. The table shows the F-measure for each of the videos in the data set for various choices of the kernel variances. The first 5 columns correspond to using a constant variance for each process at all pixel locations in the video. Having identical kernel variances for the background and foreground models (columns 1, 2) is not as effective as having different variances (all other columns). Comparing columns 2 and 3 shows that using a larger spatial variance for the foreground model than for the background model is beneficial. Changing the spatial variance from 3 (column 3) to 1 (column 4) helps the overall accuracy in one video (Fountain). Using a selection procedure where the best kernel variance is chosen from a set of values gives the best results for most videos (column 6) and frames.

Comparison of our selection procedure to a baseline method of using a standard algorithm for variance selection in KDE (AMISE criterion) shows that the standard algorithm is not as accurate as our

method (column 7). Our choice for the variance values for spatial dimension reflects no motion ($B_d = 1/4$) and very little motion ($B_d = 3/4$) for the background, and moderate amount of motion ($F_d = 12/4$) for the foreground. For the color dimension, the choice is between little variation ($B_{rgb} = 5/4$), moderate variation ($B_{rgb} = 15/4$), and high variation ($B_{rgb} = 45/4$) for the background, and moderate variation ($F_{rgb} = 15/4$) for the foreground. These choices are based on our intuition about the processes involved. For videos that differ significantly from the videos we use, it is possible that the baseline AMISE method would perform better.

We would like to point out that ideally the variance value sets should be learned automatically from a separate training data set. In absence of suitable training data for these videos in particular and for background subtraction research in general, we resort to manually choosing these values. This also appears to be the common practice among researchers in this area.

Benchmark comparisons are provided for selected existing methods - MOG, the complex foreground model (ACMMM03), and SILTP. To evaluate our results, the posterior probability of the background label is thresholded at a value of 0.5 to get the foreground pixels. Following the same procedure as , any foreground 4-connected components smaller than a size threshold of 15 pixels are ignored.

Figure 3 shows qualitative results for the same frames that were reported by . We present results for our kernel method with uniform variances and adaptive variances with RGB features (Uniform-rgb and VKS-rgb respectively), and adaptive variances with a hybrid feature space of LAB color and SILTP features (VKS-lab+siltp). Except for the Lobby video, the VKS results are better than other methods. The Lobby video is an instance where there is a sudden change in illumination in the scene (turning a light switch on and off). Due to use of an explicit foreground model, our kernel methods misclassify most of the pixels as foreground and take a long time to recover from this error. A possible solution for this case is presented later. Compared to the uniform variance kernel estimates, we see that VKS-rgb has fewer false positive foreground pixels.

Quantitative results in Table 3 compare the F-measure scores for our method against MoG, ACMMM03, and SILTP results as reported by . The table shows that methods that share spatial information (uniform kernel and VKS) with RGB features give significantly better results than methods that use RGB features without spatial sharing. Comparing the variable kernel method to a uniform kernel method in the same feature space (RGB), we see a significant improvement in performance for most videos. Scale-invariant local ternary pattern (SILTP) is a recent texture feature that is robust to soft shadows and lighting changes. We believe SILTP represents the state of the art in background modeling and hence compare our results to this method. Scale-invariant local states is a slight variation in the representation of the SILTP feature. For comparison, we use SILTP results from because in human judgement was used to vary a size threshold parameter for each video. We believe results from the latter fall under a different category of human-assisted backgrounding and hence do not compare to our method where no video-specific hand-tuning of parameters was done. Table 3 shows that SILTP is very robust to lighting changes and works well across the entire data set. Blue entries in Table 3 correspond to videos where our method performs better than SILTP. VKS with RGB features (VKS-rgb) performs well in videos that have few shadows and lighting changes. Use of color features that are more robust to illumination change, like LAB features in place of RGB helps in successful classification of the shadow regions as background. Texture features are robust to lighting changes but not effective on large texture-less objects. Color features are effective on large objects, but not very robust to varying illumination. By combining texture features with LAB color features, we expect to benefit from the strengths of both feature spaces. Such a combination has proved useful in earlier work. Augmenting the LAB features with SILTP features (computed at 3 resolutions) in the VKS framework (VKS-lab+siltp) results in an improvement in 7 out of 9 videos (last column). The variance values used in our implementation are given in Table 2.

We also compare our results (VKS-lab+siltp) to the 5 videos that were submitted as supplementary material by . Figure 4 highlights some key frames that highlight the strengths and weaknesses of our system versus the SILTP results. The common problems with our algorithm are shadows being classified as foreground (row e) and initialization errors (row e shows a scene where the desk was occluded by people when the background model was initialized. Due to the explicit foreground model, VKS takes some time to recover from the erroneous initialization). A common drawback with SILTP is that large texture-less objects have “holes” in them (row a). Use of color features helps avoid these errors. The SILTP system also loses objects that stop moving (rows b, c, d, f). Due to the explicit modeling of the foreground, VKS is able to detect objects that stop moving.

The two videos in the data set where our algorithm performs worse than SILTP are the Escalator video (rows g, h) and the Lobby video (rows i, j). In the Escalator video, our algorithm fails at the escalator steps due to large variation in color in the region.

In the Lobby video, at the time of sudden illumination change, many pixels in the image get classified as foreground. Due to the foreground model, these pixels continue to be misclassified for a long duration (row j). The problem is more serious for RGB features (Figure 3 column 2). One method to address the situation is to observe the illumination change from one frame to the next. If more than half the pixels in the image change in illumination by a threshold value of TI or more, we throw away all the background samples at that instance and begin learning a new model from the subsequent 50 frames. This method allows us to address the poor performance in the Lobby video with resulting F-measure values of 86.77 for uniform-rgb, 78.46 for VKS-rgb, and 77.76 for VKS-lab+silt. TI of 10 and 2.5 were used for RGB and LAB spaces respectively. The illumination change procedure does not affect the performance of VKS on any other video in the data set.

5 Caching optimal kernel variances from previous frame

A major drawback with trying multiple variance values at each pixel to select the best variance is that the amount of computation per pixel increases significantly. In order to reduce the complexity the algorithm, we use a scheme where the current frame's optimal variance values for each pixel location for both the background and foreground processes is stored ($B_{cache} x,y$, $F_{cache} x,y$) for each location (x, y) in the image. When classifying pixels in the next frame, these cached variance values are first tried. If the resulting scores are very far apart, then it is very likely that the pixel has not changed its label from the previous frame. The expensive variance selection procedure is performed only at pixels where the resulting scores are close to each other. Algorithm 1 for efficient computation results in a reduction in computation in about 80

6 Discussion

By applying kernel estimate method to a large data set, we have established, as do , that the use of spatial information is extremely helpful. Some of the important issues pertaining to the choice of kernel parameters for data sets with wide variations have been addressed. Having a uniform kernel variance for the entire data set and for all pixels in the image results in a poor overall system. Dynamically adapting the variance for each pixel results in a significant increase in accuracy.

Using color features in the joint domain-range kernel estimation approach can complement complex background model features in settings where the latter are known to be inaccurate. Combining robust color features like LAB with texture features like SILTP in a VKS framework yields a highly accurate background classification system.

For future work, we believe our method could be explained more elegantly in a probabilistic framework where the scores are replaced by likelihoods and informative priors are used in the Bayes rule classification.

Column num	(1)	(2)	(3)	(4)	(5)	(6)	(7)
4*B d →	3	3	3	1	3	[1 3]	AMISE
4*B rgb→	15	45	45	45	15	[5 15 45]	AMISE
4*F d →	3	3	12	12	12	[12]	[12]
4*F rgb→	15	45	45	45	15	[15]	[15]
AirportHall	40.72	59.53	67.07	63.53	47.21	70.44	53.01
Bootstrap	49.01	57.90	63.04	58.39	51.49	71.25	63.38
Curtain	66.26	83.33	91.91	89.52	81.54	94.11	52.00
Escalator	20.92	30.24	34.69	28.58	22.65	48.61	32.02
Fountain	41.87	51.89	73.24	74.58	67.60	75.84	28.50
ShoppingMall	55.19	60.17	64.95	62.18	63.85	76.48	70.14
Lobby	22.18	23.81	25.79	25.69	25.06	18.00	36.77
Trees	30.14	58.41	73.53	47.03	67.80	82.09	64.30
WaterSurface	85.82	94.04	94.93	92.91	94.64	94.83	30.29
Average	45.79	57.70	65.46	60.27	52.98	70.18	47.82

Table 1: F-measure for different kernel variances. Using our selection procedure (Column 6) results in the highest accuracy.