Video Anomaly Detection by Estimating Likelihood of Representations

Yuqi Ouyang Victor Sanchez

Signal and Information Processing Lab, Department of Computer Science,

The University of Warwick

yuqi.ouyang@warwick.ac.uk v.f.sanchez-silva@warwick.ac.uk

Intro

Video anomaly detection refers to detecting abnormal activities or events in a scene. Since it is very challenging to collect and label examples of all possible types of abnormal events in a scene, video anomaly detection is usually solved as an unsupervised learning task, namely, training model with normal data for the objective of detecting outliers.

We propose a deep probabilistic model (named GMM-DAE) to transfer this task into a density estimation problem where latent manifolds are generated by a deep denoising autoencoder and clustered by expectation maximization. Evaluations on several benchmark datasets show the strengths of our model, achieving outstanding performance on challenging datasets.

Training to Reconstruct

$$L\left(\boldsymbol{X}, \hat{\boldsymbol{X}}\right) = \frac{1}{n} \sum_{i=1}^{n} \|x_i - \hat{x}_i\|_2^2 + \beta \|\boldsymbol{W}\|_2^2$$

X: n uncorrupted images.

 \widehat{X} : *n* output images.

 $\|W\|_2^2$: ℓ_2 -regularization on the weights by a factor β .

Experiments

Model	UCSD Ped2	CUHK Avenue	ShanghaiTech	
MPPCA [14]	69.3	-	-	Table 1: Erame lovel
MPPCA+SFA [25]	61.3	-	-	Table 1. Frame-level
MDT [25]	82.9	-	-	AUROC curve (%)
Unmasking [13]	82.2	80.6	-	comparison with
AMDN [42]	90.8	-	-	
FRCN action [10]	92.2	89.8	-	other baseline
Conv-AE [9]	90.0	70.2	60.9	models, on three
STAE [43]	91.2	77.1	-	henchmark datasets
GANs [31]	93.5	-	-	
FFP+MC [18]	95.4	85.1	72.8	(Higher is better).
LSA [1]	95.4	-	72.5	Current SOTA
MLAD [41]	99.21	71.54	-	porformanco is
sRNN-AE [22]	92.21	83.48	69.63	performance is
MemAE [8]	94.1	83.3	71.2	indicated in bold text.
UnetGAN [28]	96.2	86.9	-	(Please find all
MLEP-FP [19]	-	89.2	73.4	
MemAE2020 [30]	97.0	88.5	70.5	references in the
SDOR [29]	83.2	-	-	manuscript)
SAGC [26]	-	-	76.1	1,
GMM-DAE	96.5	89.3	81.2	

Method	AUROC
DAE+OP	90.1
DAE+OP+DP	93.9
DAE+OP+OL	94.2
DAE+OP+OL+DP+DL (GMM-DAE)	96.5
AE+OP+OL+DP+DL	95.8

Table 2: Frame-level AUROC curve (%) on the UCSD Ped2 dataset using different components of the proposed GMM-DAE model. O*: Frame patch; D*: Dynamic frame patch; *P: PSNR value; *L: Likelihood value.

Model: Overall



Figure 1: The proposed GMM-DAE model. An object detector is applied to generate patches. A dynamic image is computed by approximate rank pooling for motion info collection. Two denoising autoencoders (DAE) are trained to reconstruct data. Lowdimensional manifolds are clustered with two Gaussian Mixture Models (GMM). The final anomaly score is computed by fusing the reconstruction errors and latent likelihoods.

Training to Estimate Density



Performance Analysis







$$\hat{\phi}_{j} = \sum_{i=1}^{n} \frac{\gamma_{ij}}{n}$$

$$\hat{\phi}_{j} = \sum_{i=1}^{n} \frac{\gamma_{ij}}{n}$$

$$\hat{\mu}_{j} = \frac{\sum_{i=1}^{n} \gamma_{ij} z_{i}}{\sum_{i=1}^{n} \gamma_{ij}}$$

$$\hat{\Sigma}_{j} = \frac{\sum_{i=1}^{n} \gamma_{ij} (z_{i} - \hat{\mu}_{j}) (z_{i} - \hat{\mu}_{j})^{T}}{\sum_{i=1}^{n} \gamma_{ij}}$$

$$\hat{\nabla}_{j} = \frac{\sum_{i=1}^{n} \gamma_{ij} (z_{i} - \hat{\mu}_{j}) (z_{i} - \hat{\mu}_{j})^{T}}{\sum_{i=1}^{n} \gamma_{ij}}$$

$$M$$

$$M$$

$$K$$

$$(\hat{\phi}_{1}, \hat{\mu}_{1}, \hat{\Sigma}_{1}), N(\hat{\phi}_{2}, \hat{\mu}_{2}, \hat{\Sigma}_{2}), \cdots, N(\hat{\phi}_{k}, \hat{\mu}_{k}, \hat{\Sigma}_{k})$$

1750

Figure 3: Performance analysis on the UCSD Ped2 dataset. Higher anomaly score indicates more abnormal events happening. Blue areas indicate ground truth abnormal frames. Red boxes indicate correct detections. Activities from left to right in the first row: cycling, vehicle moving, skateboard riding & cycling.

Model: Denoise to Reconstruct



Figure 2: Architecture of the DAE in the GMM-DAE model. CONV: convolutional layer. BN: batch normalization.

Anomaly Inference

$$\begin{split} PSNR(x,\hat{x}) &= 10 \log_{10} \frac{max(x)}{MSE(x,\hat{x})} \\ P(z) &= \log \sum_{j=1}^{k} \frac{\hat{\phi}_{j} \cdot \exp - \frac{1}{2}(z-\hat{\mu}_{j})^{T} \hat{\Sigma}_{j}^{-1}(z-\hat{\mu}_{j})}{\sqrt{|2\pi\hat{\Sigma}_{j}|}} \\ A(x^{t}) &= - \left[\lambda_{1}P(z_{x^{t}}) + \lambda_{2} \cdot PSNR(x^{t},\hat{x}^{t}) + \lambda_{3}P(z_{d^{t}}) + \lambda_{4} \cdot PSNR(d^{t},\hat{d}^{t})\right] \\ A(I^{t}) &= max\{A(x_{1}^{t}), A(x_{2}^{t}), ..., A(x_{n}^{t})\} \\ PSNR(x,\hat{x}): \text{ Reconstruction accuracy.} \\ P(z): \text{ Likelihood of latent representation } z. \end{split}$$

 $A(x^t)$: Anomaly score of frame patch x^t .

 $A(I^t)$: Anomaly score of frame I^t .

Latent Space Visualization



Figure 4: Distribution of the generated manifolds based on the normalized anomaly scores using PSNR values (left: OI+PSNR) and latent likelihood values (right: OI+LL) on the UCSD Ped2 test videos. Each color represents a range of anomaly scores.