

STOCHASTIC LATENT RESIDUAL VIDEO PREDICTION

LOCUST @ Deezer

Tuesday 8th October, 2019

Jean-Yves Franceschi*, Edouard Delasalles*,
Mickael Chen, Sylvain Lamprier, Patrick Gallinari

Video Prediction

Video Prediction

Given a few conditioning frames, learning the distribution of future frames.

Challenges

- Account for uncertainty in the future
- Long-term prediction
- Generating realistic images

Applications

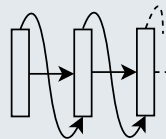
- Reinforcement Learning [Gre+19]
 - Robotics [Bab+18]
-) Challenges the ability of a model to capture visual and dynamic representations of the world

Related Work

Image Auto-Regressive Approaches [Bab+18] [DF18] [Lee+18]

Uses previously generated frames to predict next one

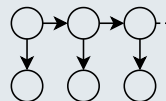
- (+) Easy to learn, based on powerful LSTMs, work with high-definition images
- () High computational cost, prediction tied to image generation



State-Space Models [Fra+16] [Kar+17] [Fra+17] [KSS17]

Prediction dynamics is fully latent. Latent states are independently decoded into frames.

- (+) Fast latent prediction, not all frames have to be decoded, interpretable latent space
- () Harder to train, only proposed on low dimensional data



Method

Our Approach

Latent Residual Dynamic Model

- Deterministic latent state \mathbf{y}_t
- Stochastic auxiliary variable \mathbf{z}_t
- Residual update of \mathbf{y}_t
 - ! Integration of ODEs in neural network architectures [Che+18]
 - ! Near-continuous dynamics
 - ! Generation at arbitrary frame rates
- End-to-end training with variational inference (VAE)

Content Variable

- Store “static” information (background, object shapes, etc...)
- Inferred from randomly sampled frame ! temporal invariance
- Skip-connections between encoder and decoder

Latent Residual Temporal Model

- \mathbf{x}_t : t -th frame
- \mathbf{y}_t : latent space corresponding to \mathbf{x}_t
- \mathbf{z}_t : random variable encapsulating stochasticity at time t
- f : residual function
- g : decoder

$$\begin{array}{ll}
 \mathbb{R} \mathbf{y}_1 & N(\mathbf{0}; I) & \text{(initial condition)} \\
 \mathbb{R} \mathbf{z}_{t+1} & N(\mathbf{y}_t; (\mathbf{y}_t)I) & \text{(random variable prediction)} \\
 \mathbb{R} \mathbf{y}_{t+1} & = \mathbf{y}_t + f(\mathbf{y}_t; \mathbf{z}_{t+1}) & \text{(latent state prediction)} \\
 \mathbb{R} \mathbf{x}_t & N(g(\mathbf{y}_t); I) & \text{(decoding)}
 \end{array}$$

Latent Residual Temporal Model

- \mathbf{x}_t : t -th frame
- \mathbf{y}_t : latent space corresponding to \mathbf{x}_t
- \mathbf{z}_t : random variable encapsulating stochasticity at time t
- f : residual function
- g : decoder

$$\begin{array}{ll}
 \mathbb{R} \mathbf{y}_1 \sim N(\mathbf{0}; I) & \text{(initial condition)} \\
 \mathbb{R} \mathbf{z}_{t+1} \sim N(\mathbf{0}; I) & \text{(random variable prediction)} \\
 \mathbb{R} \mathbf{y}_{t+1} = \mathbf{y}_t + f(\mathbf{y}_t; \mathbf{z}_{t+1}) & \text{(latent state prediction)} \\
 \mathbb{R} \mathbf{x}_t \sim N(g(\mathbf{y}_t); I) & \text{(decoding)}
 \end{array}$$

Extension Leveraging the Euler Analogy

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \Delta t f(\mathbf{y}_t; \mathbf{z}_{t+1})$$

Visual Representation

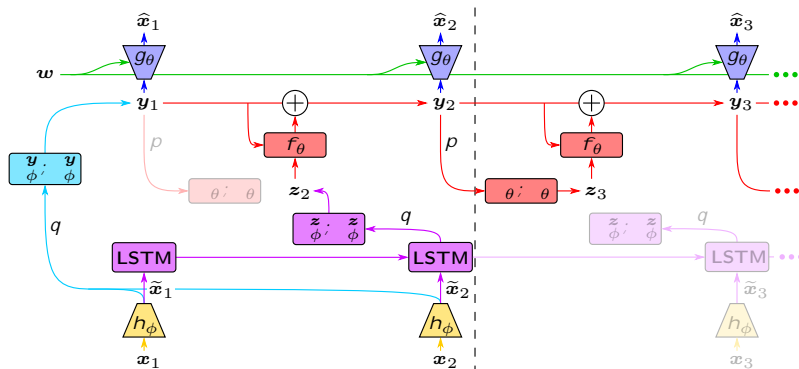
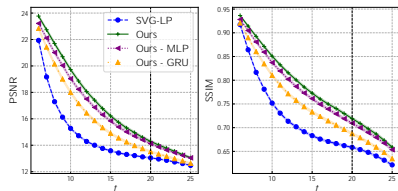


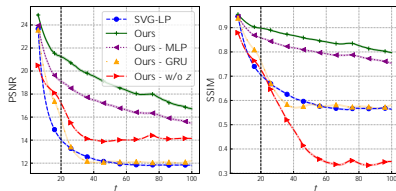
Figure 1: Model () and inference (, LSTM) architecture on a test sequence.

Experiments

Stochastic Moving MNIST



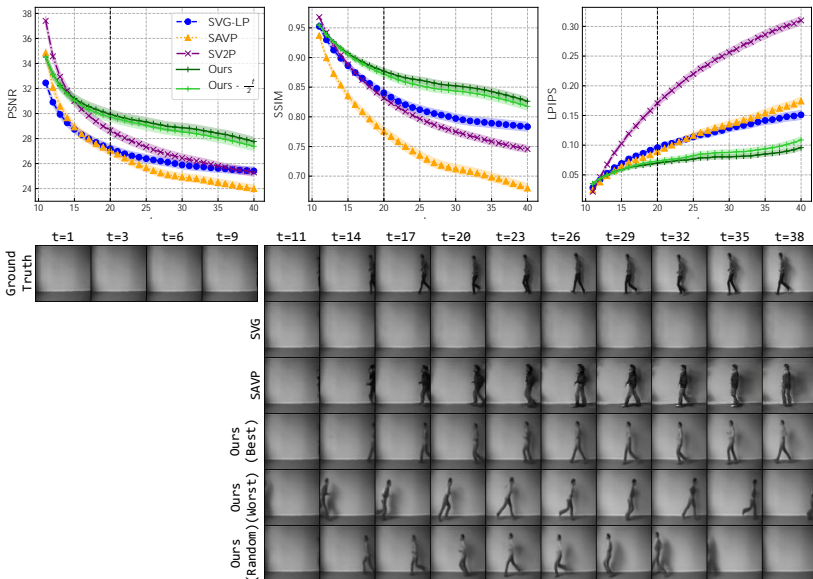
(a) Stochastic Moving MMNIST.



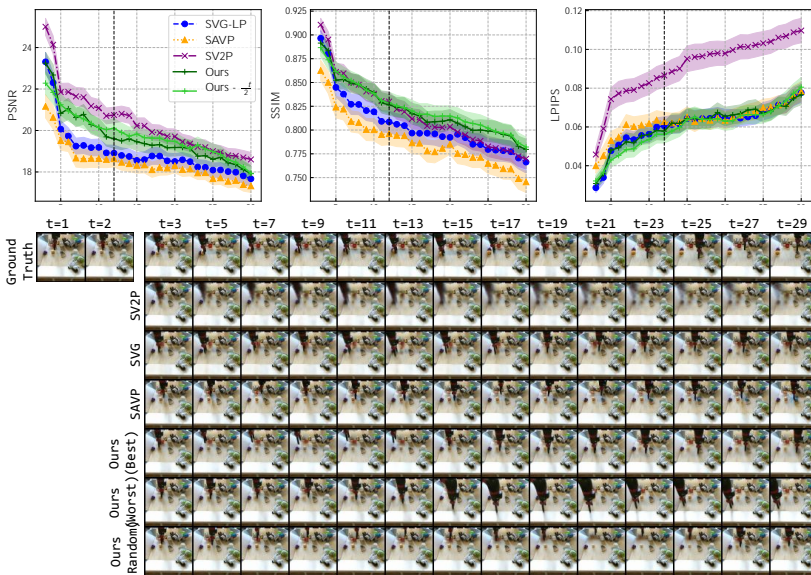
(b) Deterministic Moving MMNIST.



KTH Action Dataset








BAIR Robot Pushing Dataset







References

References I

-  Babaeizadeh, Mohammad et al. (2018). "Stochastic Variational Video Prediction". In: *International Conference on Learning Representations*.
-  Chen, Tian Qi et al. (2018). "Neural Ordinary Differential Equations". In: *Advances in Neural Information Processing Systems 31*. Ed. by Samy Bengio et al. Curran Associates, Inc., pp. 6571–6583.
-  Denton, Emily and Rob Fergus (July 2018). "Stochastic Video Generation with a Learned Prior". In: *Proceedings of the 35th International Conference on Machine Learning*. Ed. by Jennifer Dy and Andreas Krause. Vol. 80. Proceedings of Machine Learning Research. Stockholm, Sweden: PMLR, pp. 1174–1183.
-  Fraccaro, Marco et al. (2016). "Sequential Neural Models with Stochastic Layers". In: *Advances in Neural Information Processing Systems 29*. Ed. by Daniel D. Lee et al. Curran Associates, Inc., pp. 2199–2207.
-  Fraccaro, Marco et al. (2017). "A Disentangled Recognition and Nonlinear Dynamics Model for Unsupervised Learning". In: *Advances in Neural Information Processing Systems 30*. Ed. by Isabelle Guyon et al. Curran Associates, Inc., pp. 3601–3610.

References II

-  Gregor, Karol et al. (2019). "Temporal Difference Variational Auto-Encoder". In: *International Conference on Learning Representations*.
-  Karl, Maximilian et al. (2017). "Deep Variational Bayes Filters: Unsupervised Learning of State Space Models from Raw Data". In: *International Conference on Learning Representations*.
-  Krishnan, Rahul G., Uri Shalit, and David Sontag (2017). "Structured Inference Networks for Nonlinear State Space Models". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 31, pp. 2101–2109.
-  Lee, Alex X. et al. (2018). "Stochastic Adversarial Video Prediction". In: *arXiv preprint arXiv:1804.01523*.

Appendix

Prior

$$\begin{aligned}
 \mathbb{P} \mathbf{y}_1 & \sim N(\mathbf{0}; I) && \text{(initial condition)} \\
 \mathbb{P} \mathbf{z}_{t+1} & \sim N(\mathbf{y}_t; (\mathbf{y}_t)^\top I) && \text{(random variable prediction)} \\
 \mathbb{P} \mathbf{y}_{t+1} & = \mathbf{y}_t + f(\mathbf{y}_t; \mathbf{z}_{t+1}) && \text{(latent state prediction)} \\
 \mathbb{P} \mathbf{x}_t & \sim N(g(\mathbf{y}_t); I) && \text{(decoding)}
 \end{aligned}$$

The above model induces the following factorization:

$$\begin{aligned}
 & p(\mathbf{x}_{1:T}; \mathbf{z}_{2:T}; \mathbf{y}_{1:T} | \mathbf{w}) \\
 & = p(\mathbf{y}_1) \prod_{t=1}^{T-1} p(\mathbf{z}_{t+1} | \mathbf{y}_t) p(\mathbf{y}_{t+1} | \mathbf{y}_t; \mathbf{z}_{t+1}) \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{y}_t; \mathbf{w});
 \end{aligned}$$

with $p(\mathbf{y}_{t+1} | \mathbf{y}_t; \mathbf{z}_{t+1})$ being a Dirac.

Posterior

We factorize the posterior as follows:

$$\begin{aligned}
 q_{Z;Y} &, q(\mathbf{z}_{2:T}; \mathbf{y}_{1:T} \mid \mathbf{x}_{1:T}; \mathbf{w}) \\
 &= q(\mathbf{y}_1 \mid \mathbf{x}_{1:k}) \prod_{t=2}^T q(\mathbf{z}_t \mid \mathbf{x}_{1:t}) \mathbf{y}_{t-1} + f(\mathbf{y}_{t-1}; \mathbf{z}_t) \mathbf{y}_t :
 \end{aligned}$$

- \mathbf{y}_1 is inferred from the first k frames
- \mathbf{z}_t is inferred from the past and present frames only (works best in practice)
- \mathbf{y}_t is deterministically computed from the previously inferred \mathbf{y}_{t-1} and \mathbf{z}_t

Evidence Lower Bound

$$\begin{aligned}
 \log p(\mathbf{x}_{1:T} | \mathbf{w}) &= L(\mathbf{x}_{1:T}; \mathbf{w}; \cdot; \cdot) \\
 &= E_{(\tilde{\mathbf{z}}_{2:T}; \tilde{\mathbf{y}}_{1:T})} q_{\mathbf{Z}; \mathbf{Y}} \sum_{t=1}^T \log p(\mathbf{x}_t | \mathbf{y}_t; \mathbf{w}) \\
 &\quad - D_{\text{KL}} q(\mathbf{y}_1 | \mathbf{x}_{1:k}) - p(\mathbf{y}_1) \\
 &\quad - E_{(\tilde{\mathbf{z}}_{2:T}; \tilde{\mathbf{y}}_{1:T})} q_{\mathbf{Z}; \mathbf{Y}} \sum_{t=2}^T D_{\text{KL}} q(\mathbf{z}_t | \mathbf{x}_{1:t}) - p(\mathbf{z}_t | \mathbf{y}_{t-1})
 \end{aligned}$$

Loss Function

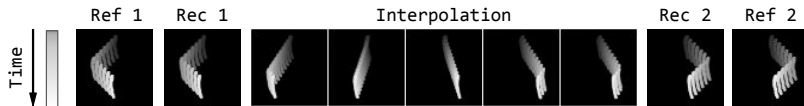
We add to the evidence lower bound an ℓ_2 regularization on the residues to stabilize the dynamics, giving the following loss function:

$$\arg \max_{\theta} \mathbb{E}_{\mathbf{x}_c^{(k)}} \sum_{\mathbf{x} \in \mathcal{X}} \frac{1}{4} L(\mathbf{x}_{1:T}; \mathbf{c}, \mathbf{x}_c^{(k)}) + \frac{\lambda}{2} \sum_{t=2}^T \mathbb{E}_{(\mathbf{z}_{2:T}; \mathbf{y}_{1:T})} \left[\sum_{\mathbf{z}; \mathbf{y}} \sum_{t=2}^T f(\mathbf{y}_{t-1}; \mathbf{z}_t) \right]$$

Latent Space Interpolation

Interpolation of y_1 on Moving MNIST

- Take 2 trajectories x^s and x^t
- Infer Latent initial conditions y_1^s and y_1^t
- Generate frame sequences from \tilde{y}_1 linearly interpolated between y_1^s and y_1^t .



Deep State Space Model

Without content variable, we can compare our model to Deep State Space Models.

Table 1: ELBO for DVBF [Kar+17], KVAE [Fra+17] and our model on the Pendulum dataset.

Score	DVBF	KVAE	Ours
ELBO	798.56	807.02	806.12