# **Prediction Under Uncertainty for Autonomous Driving**

#### **Anonymous Author(s)**

Affiliation Address email

#### Abstract

We consider the problem of prediction in a highly uncertain environment.

#### 2 1 Introduction

5

6

8

- 3 The contributions of this work are the following:
  - We introduce a new open-source environment to test methods for autonomous driving, based on a large dataset of real-world driving data.
    - We introduce a novel method for prediction under uncertainty which is simple to train, does not make assumptions about a prior distribution over latent variables or require sampling at training time, and is able to generate diverse predictions for hundreds of timesteps into the future.

## 10 2 Prediction Model

Our stochastic prediction model can be viewed as a conditional autoencoder paired with a non-parametric sampling procedure. The architecture consists of three neural networks: an encoder  $f_1$ , a decoder  $f_2$  and a latent variable network  $\phi$ . For each sample, the update equations are given by:

$$z_i = \phi(x_i, y_i)$$
  
$$\tilde{y}_i = f_2(f_1(x_i), z_i)$$

and all networks are trained by gradient descent to optimize the following objective:

$$\mathcal{L} = \sum_{i} \|y_i - \tilde{y}_i\|_2^2 \tag{1}$$

$$= \sum_{i} \|y_i - f(x_i, \phi(x_i, y_i))\|_2^2$$
 (2)

Note in particular that no sampling or reparamaterization is done at training time. After training, we extract all vectors  $z_i$  from the training set and use these as inputs to new inputs.

# 3 Related Work

- In recent years, several works have explored prediction of complex time series such as video [1].
- These typically train models to predict future frames with the goal of learning good representations

#### Algorithm 1 My algorithm

```
1: Input: Time series \{s_1...s_T\}.
 2: Train latent model:
 3: while not converged do
 4:
            z_t = f_{\phi}(s_{1:t}, s_{t+1})
           \tilde{s}_{t+1} = f_{\theta}(s_{1:t}, z_t)
\ell(\theta, \phi) = \|s_{t+1} - \tilde{s}_{t+1}\|_2^2
\theta \leftarrow \theta - \eta \nabla \theta
 5:
 6:
 7:
 8:
            \phi \leftarrow \phi - \eta \nabla \phi
 9: procedure EstimateLatentManifold
            V \leftarrow \{\}
10:
            E \leftarrow \{\}
11:
            for t = 1 : T do
12:
                 z_t = f_{\phi}(s_{1:t}, s_{t+1})

V[t] = z_t
13:
14:
            for t = 1 : T do
15:
                  E[t] \leftarrow \text{list of } k \text{ nearest neighbors of } V[t] = z_t \text{ in } V
16:
            return G = (V, E)
17: procedure Generate(s_0, s_1)
18:
            Initialize z_1 = f_{\phi}(s_0, s_1)
            for t = 1 : T \operatorname{do}^{q}
19:
20:
                  \tilde{s}_{t+1} = f_{\theta}(s_{1:t}, z_t)
```

- which disentangle factor of variation and can be used for unsupervised learning [2?, 3] or learn action-conditional forward models which can be used for planning [4, 5, 6?]. Several works have included latent variables as a means to model the uncertainty, using the framework of Variational Autoencoders [7, 8]. In contrast to VAEs, our model does not place any priors on the latent variable distribution, which removes the need for an additional loss term enforcing consistency between the prior and posterior. Additionally, our approach does not require any sampling at training time which reduces the variance in the gradients.
- Our method is closely related to Gated and Relational Autoencoders [9, 10], which were used to learn transformations between pairs of images in an unsupervised manner. The general architecture and loss is similar in both our works, however their focus was on representation learning for static images while ours is on video generation for use in planning.
- Video Prediction [1], stochastic: using the framework of Variational Autoencoders [11].
- Mixture Density Networks [12].
  - VQ-VAE, GLO, Gated Autoencoders

#### 4 Dataset

33

36

## 5 **Experiments**

 $\bullet$  Histogram of distances between consecutive z vs random pairs.

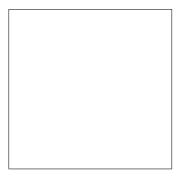


Figure 1: Sample figure caption.

Table 1: Sample table title

	Part	
Name	Description	Size $(\mu m)$
Dendrite Axon Soma	Input terminal Output terminal Cell body	$\begin{array}{c} \sim \! 100 \\ \sim \! 10 \\ \text{up to } 10^6 \end{array}$

- 37 5.1 Figures
- 38 **5.2 Tables**

## **39 6** Final instructions

# 40 7 Appendix

- Our model can also be viewed through the lens of Variational Autoencoders, as a type of conditional
- 42 VAE with a zero-variance posterior network and a uniform categorical prior. In this case, the KL term
- reduces to a constant which can be ignored during training; details can be found in the Appendix.
- 44 Fixed Prior

$$p_{\psi}(z_j|x_i, y_i) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$$
 (3)

$$p_{\phi}(z_j|x_i) = \frac{1}{M} \tag{4}$$

(5)

The KL divergence is therefore:

$$KL(p_{\phi}(z|x_i)||p_{\psi}(z|x_i, y_i)) = \sum_{j} p_{\psi}(z_j|x_i, y_i) \log \frac{p_{\psi}(z_j|x_i, y_i)}{p_{\phi}(z_j|x_i)}$$
(6)

$$=1\cdot\log\frac{1}{\frac{1}{M}}\tag{7}$$

$$= \log M \tag{8}$$

# 46 Learned Prior

$$p_{\psi}(z_j|x_i, y_i) = \begin{cases} 1 & \text{if } i = j\\ 0 & \text{else} \end{cases}$$
 (9)

$$p_{\psi}(z_{j}|x_{i}, y_{i}) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$$

$$p_{\phi}(z_{j}|x_{i}) = \frac{h(z_{j}; x_{j}, \theta)}{\sum_{k=1}^{M} h(z_{k}; x_{j}, \theta)}$$
(10)

The KL divergence is therefore:

$$KL(p_{\phi}(z|x_i)||p_{\psi}(z|x_i, y_i)) = \sum_{j} p_{\psi}(z_j|x_i, y_i) \log \frac{p_{\psi}(z_j|x_i, y_i)}{p_{\phi}(z_j|x_i)}$$
(11)

$$= 1 \cdot \log \frac{1}{\frac{h(z_{j}; x_{j}, \theta)}{\sum_{k=1}^{M} h(z_{k}; x_{j}, \theta)}}$$

$$= -\log \frac{h(z_{j}; x_{j}, \theta)}{\sum_{k=1}^{M} h(z_{k}; x_{j}, \theta)}$$
(13)

$$= -\log \frac{h(z_j; x_j, \theta)}{\sum_{k=1}^{M} h(z_k; x_j, \theta)}$$
 (13)

$$= -\log h(z_j; x_j, \theta) + \log \sum_{k=1}^{M} h(z_k; x_j, \theta)$$
 (14)

The last term is a constant and can be ignored.

#### Acknowledgments 49

Use unnumbered third level headings for the acknowledgments. All acknowledgments go at the end 50 of the paper. Do not include acknowledgments in the anonymized submission, only in the final paper. 51

#### References 52

- [1] Michaël Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond 53 mean square error. CoRR, abs/1511.05440, 2015. 54
- [2] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. Unsupervised learning of video 55 representations using 1stms. CoRR, abs/1502.04681, 2015. 56
- [3] Emily Denton and Vighnesh Birodkar. Unsupervised learning of disentangled representations 57 from video. CoRR, abs/1705.10915, 2017. 58
- [4] Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard L. Lewis, and Satinder P. Singh. Action-59 conditional video prediction using deep networks in atari games. CoRR, abs/1507.08750, 60 2015. 61
- [5] Chelsea Finn, Ian J. Goodfellow, and Sergey Levine. Unsupervised learning for physical 62 interaction through video prediction. CoRR, abs/1605.07157, 2016. 63
- [6] Pulkit Agrawal, Ashvin Nair, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Learning to 64 65 poke by poking: Experiential learning of intuitive physics. CoRR, abs/1606.07419, 2016.
- [7] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. 66 Stochastic variational video prediction. CoRR, abs/1710.11252, 2017. 67
- [8] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. CoRR, 68 abs/1802.07687, 2018. 69
- [9] Roland Memisevic and Geoffrey Hinton. Unsupervised learning of image transformations, 07 70 71
- [10] R. Memisevic. Learning to relate images. IEEE Transactions on Pattern Analysis and Machine 72 Intelligence, 35(8):1829–1846, Aug 2013. 73

- [11] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. 12 2013.
- 75 [12] Christopher Bishop. Mixture density networks. Technical report, January 1994.

# 76 References

- 77 References follow the acknowledgments. Use unnumbered first-level heading for the references. Any
- choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font
- size to small (9 point) when listing the references. Remember that you can use more than eight
- pages as long as the additional pages contain *only* cited references.