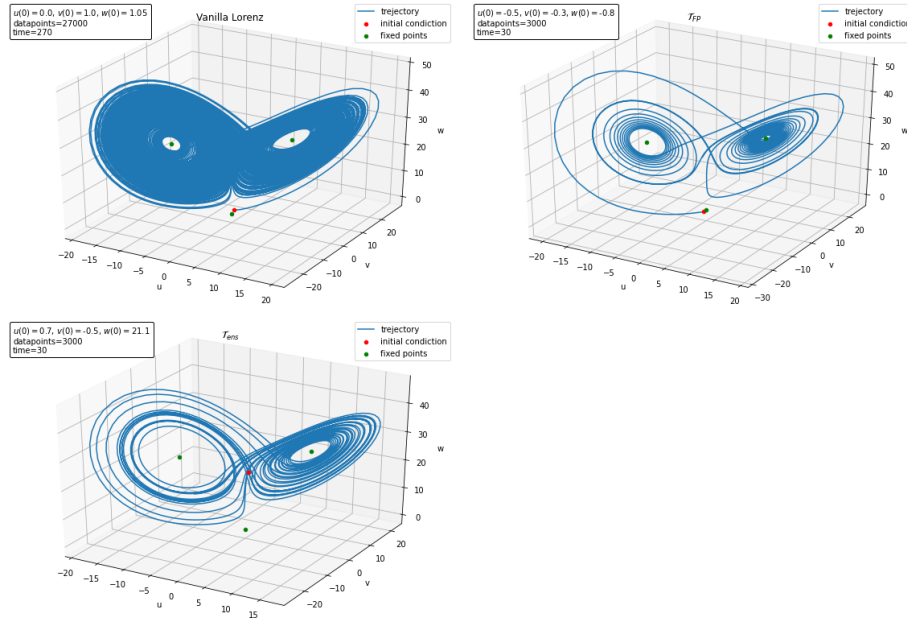


Entropy and Information in Lorenz Trajectories

December 2021

1 Entropy in Lorenz Trajectories

The following text is going to show, using different functions, why different trajectories in Lorenz Model seem to contain more information than others. Namely the performance of a LSTM trained with a dataset of trajectories with initial condition close to fixed points are better than others.



The idea behind this phenomena that we observe is that close to the fixed points the dynamics is far from the chaotic dynamics. We would like to measure the information that the different trajectories contain.

1.1 Shannon Entropy

The first function that we can use is the Shannon Entropy [1948]. This function is well known in Information Theory. The main idea of this function is that

it can measure the amount of information needed to reproduce the input (that is an event sampled from a distribution of probability). Namely, the minimum number of bits that you can use to store an events from a distribution.

Given X a random variable, with possible outcomes x_1, x_2, \dots, x_N , with probabilities $p(x_1), p(x_2), \dots, p(x_N)$, we can define the Shannon Entropy:

$$I_{sh} = I_{sh}(p(\mathbf{x}))$$

$$I_{sh} = - \sum_{i=1}^N p(x_i) \log_2(p(x_i)) \quad (1)$$

1.1.1 Technical Issues

Lorenz's systems trajectory is a timeseries with 'float64' numbers.

A classic Shannon Entropy algorithm has to compute the relative frequency of each number in the array to give an approximated value of $p(x_1)$. A known issue ^[1] when you have to use float64 values is that it is almost impossible to find two identical values in a timeseries.

Indeed the function "torch.load()" ^[2] creates a vector of 'float64' arrays with a precision of 16 decimals. That's why to compute entropy is useful to use another definition of equality.

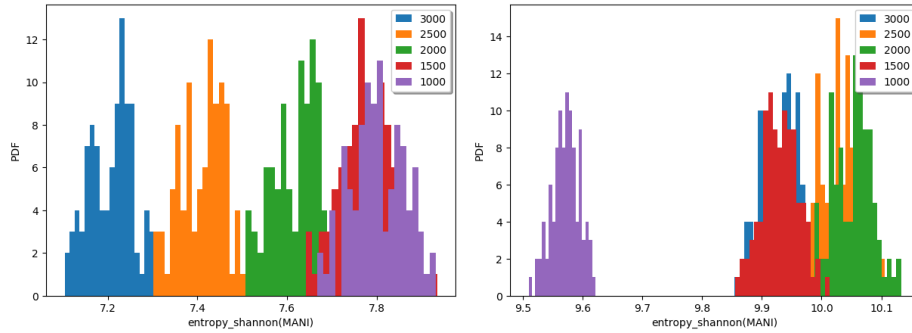
$$x(t_a) = a_0, a_1 a_2 \dots a_{16}$$

$$x(t_b) = b_0, b_1 b_2 \dots b_{16} \quad (2)$$

We say that $x(t_a) = x(t_b)$ with precision p if :

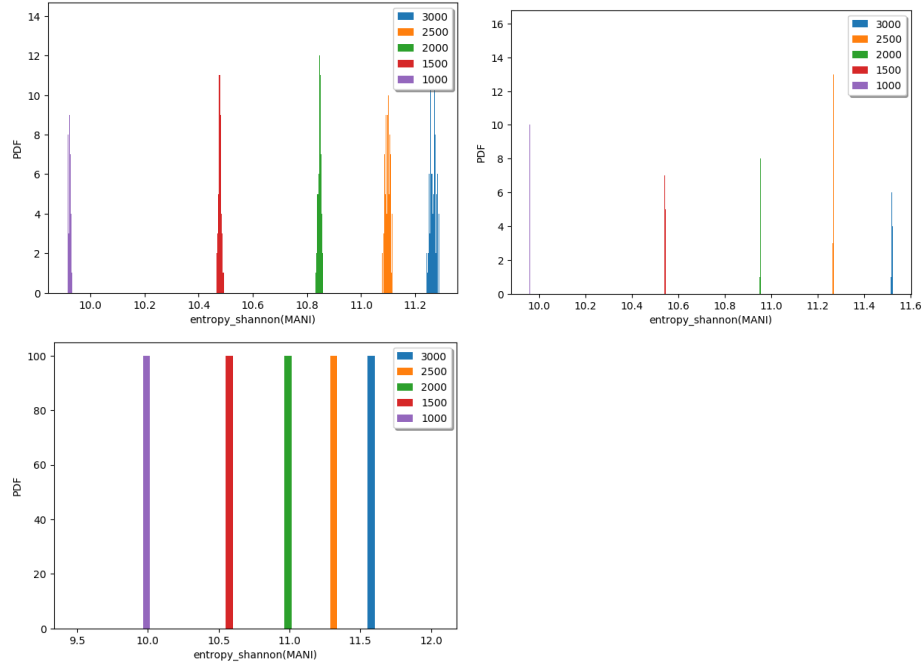
$$a_0 = b_0, a_1 = b_1, \dots, a_p = b_p \quad (3)$$

Following this idea we can make different plots of Shannon Entropy of the trajectories close to fixed points. With $p = 1, 2, 3, 4, 16$



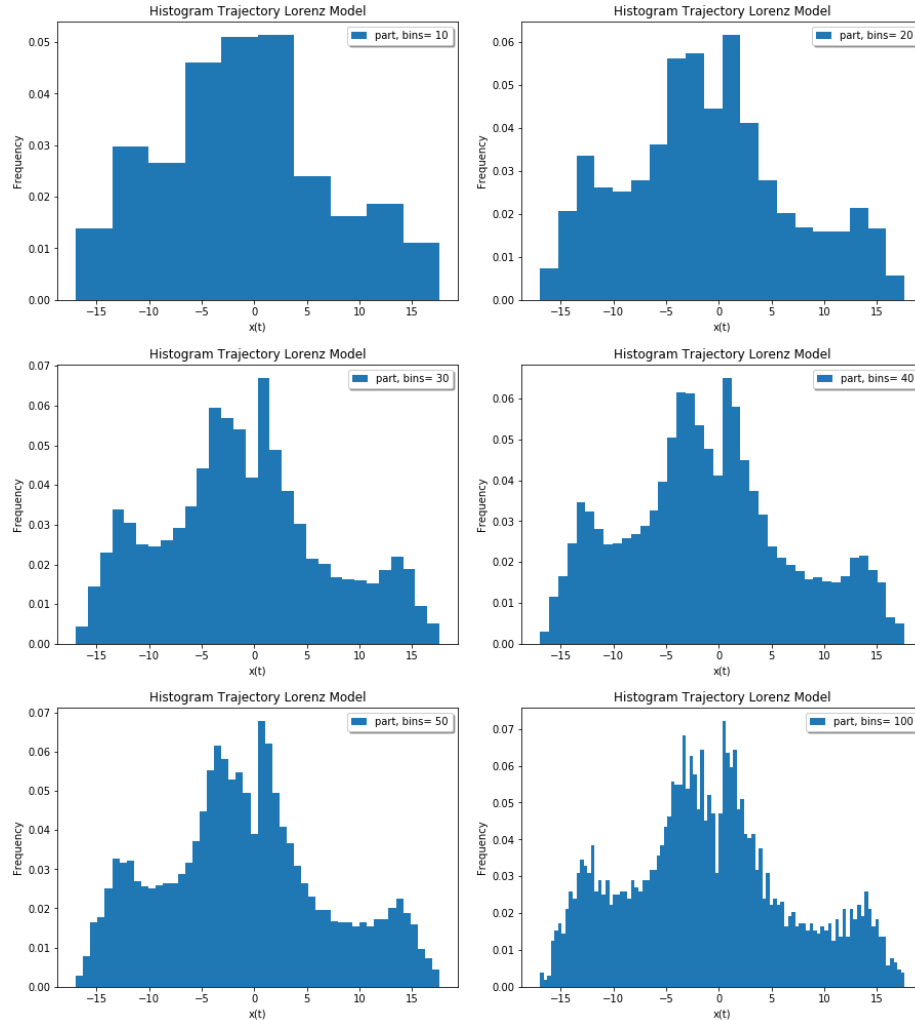
¹<https://stackoverflow.com/questions/18106975/compare-float-and-float64-in-python/18107037>

²Used to analyze the entropy of the trajectories



We can see that in the first two plots the histograms shows to have a higher standard deviation than the others, this happens because the probability to have two identical number decrease with precision. The last plot shows this problem, every number in the timeseries is unique: that's why we have histograms with a very small standard deviation, namely the entropy is the same for each trajectory.

1.1.2 Histograms



1.1.3 Infomatic implementation of Shannon Entropy

To evaluate the entropy we use this library ^[3].

```
1 def shannon_entropy(time_series):
2     """Return the Shannon Entropy of the sample data.
3     Args:
4         time_series: Vector or string of the sample data
5     Returns:
6         The Shannon Entropy as float value
7     """
8
```

³<https://github.com/nikdon/pyEntropy/blob/master/pyentrp/entropy.py>

```

9      # Check if string
10     if not isinstance(time_series, str):
11         time_series = list(time_series)
12
13     # Create a frequency data
14     data_set = list(set(time_series))
15     freq_list = []
16     for entry in data_set:
17         counter = 0.
18         for i in time_series:
19             if i == entry:
20                 counter += 1
21         freq_list.append(float(counter) / len(time_series))
22
23     # Shannon entropy
24     ent = 0.0
25     for freq in freq_list:
26         ent += freq * np.log2(freq)
27     ent = -ent
28     return ent

```

1.2 Theory svd entropy

Single Value Decomposition is a factorization of a matrix. This method has been used in different applications. In the context of Deep Learning it is used in Principal Component Analysis [1936] that is useful in the dimensionality reduction. The key point behind this is to find the useful features that you need in the training process, namely the features that explain the variance of your signal. We can find online the documentation of the function svd-entropy [4].

$$S_{svd} = S(\mathbf{x}(t), order, delay) \quad (4)$$

We denote with \mathbf{x} the dataset that we want to study.

$$\begin{aligned} \mathbf{x}(t_0) &= x_0, \\ \mathbf{x}(t) &= (x_1, x_2, x_3, \dots, x_N) \end{aligned} \quad (5)$$

This function creates a matrix Y with the dataset.

$$\mathbf{y}(i) = (x_i, x_{i+delay}, \dots, x_{i+(order-1)delay}) \quad (6)$$

$$Y = \begin{pmatrix} \mathbf{y}(1) \\ \mathbf{y}(2) \\ \vdots \\ \mathbf{y}(N - (ord - 1)del) \end{pmatrix}$$

⁴https://raphaelvallat.com/entropy/build/html/generated/entropy.svd_entropy.html

If we use $order = 3$ and $delay = 1$ we obtain the following matrix⁵:

$$Y = \begin{pmatrix} x_1 & x_2 & x_3 \\ x_2 & x_3 & x_4 \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & x_{N-1} & x_N \end{pmatrix}$$

The next step is to write this matrix using the SVD decomposition.

$$Y = U\Sigma V^* \quad (7)$$

Where U is an unitary matrix, V is also an unitary matrix, and Σ is a diagonal matrix. Now we denotes the eigenvalues of the matrix Σ with σ_i . Now we can compute the average eigenvalues:

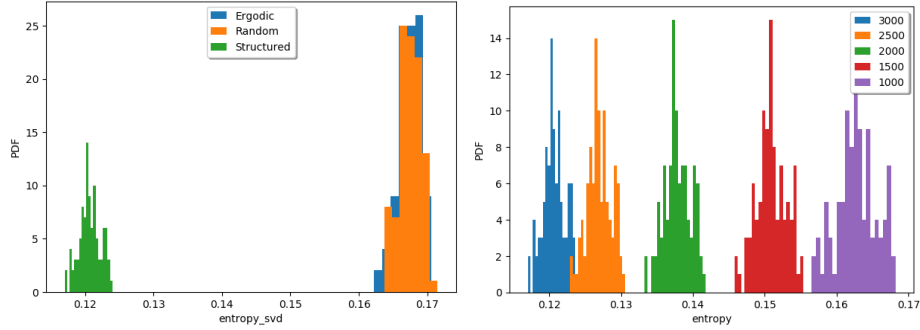
$$\bar{\sigma}_k = \frac{\sigma_k}{\sum_j^M \sigma_j} \quad (8)$$

Where M is the number of eigenvalues.

After this we can compute the SVD Entropy:

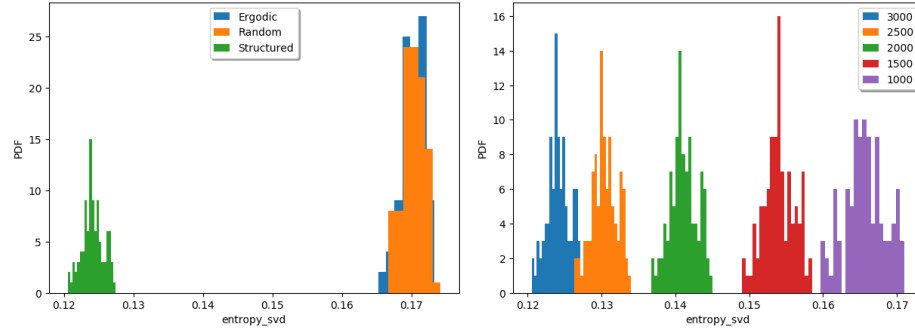
$$S_{svd} = - \sum_k^M \bar{\sigma}_k \log_2(\bar{\sigma}_k) \quad (9)$$

1.2.1 SVD with order 3

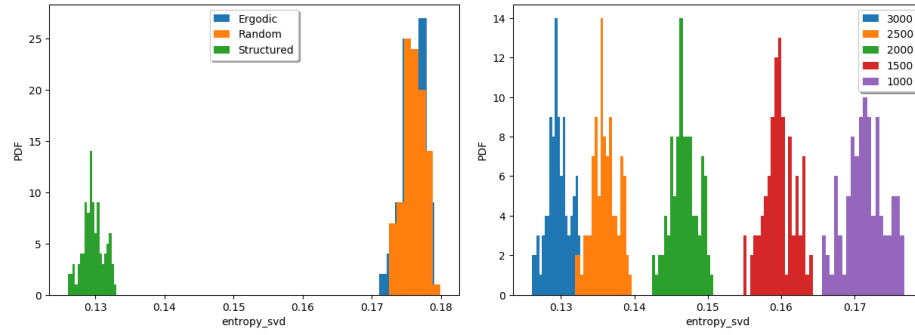


⁵To choose the delay and the order we can use the mutual information [1986] and to find the embedding dimension we can use the False Nearest Neighbours algorithm 2015, even if the qualitative behaviour of the results does not change even if we use different values of order and delay

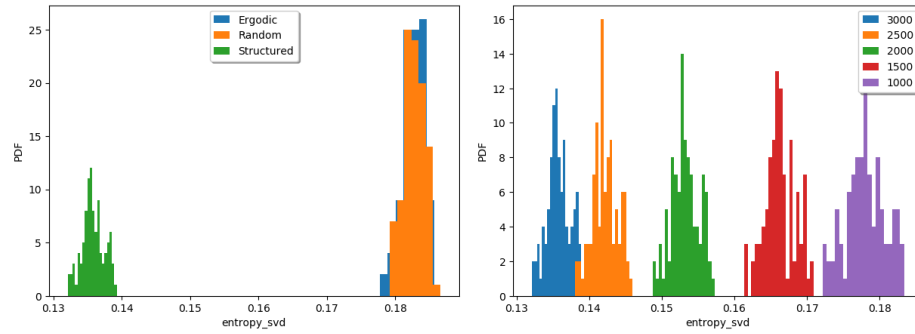
1.2.2 SVD with order 4



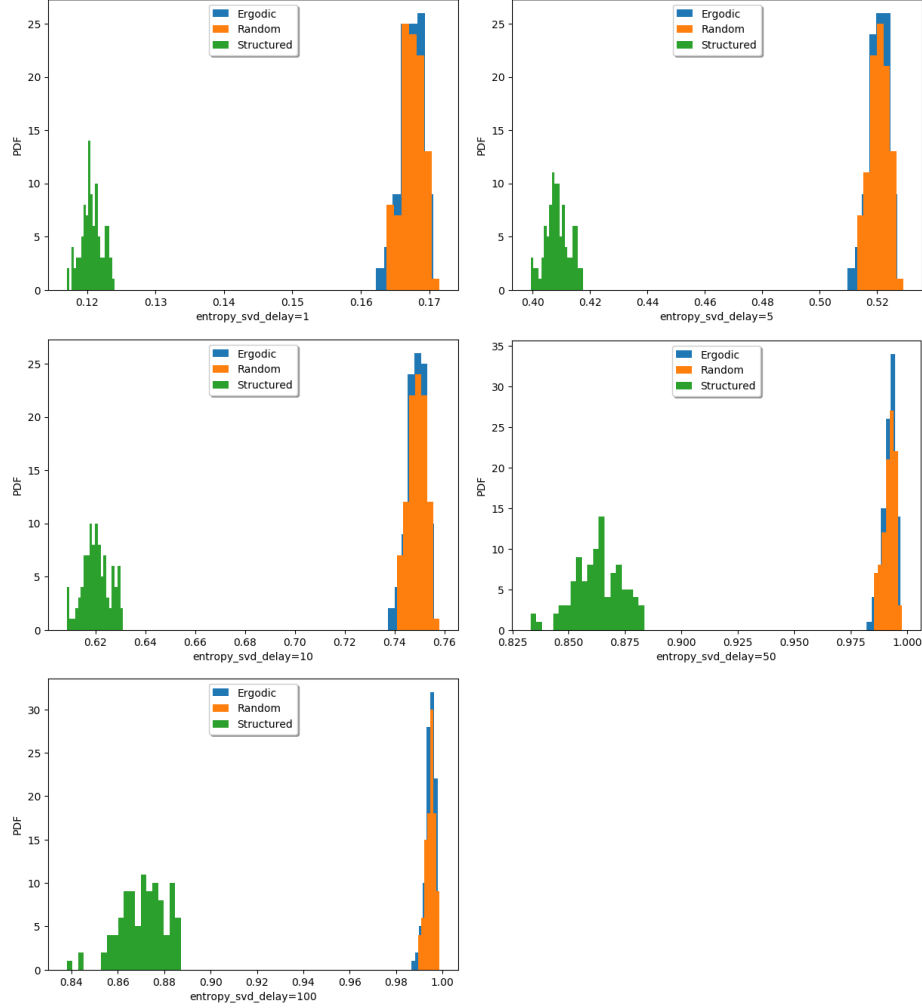
1.2.3 SVD with order 5



1.2.4 SVD with order 6



1.2.5 SVD time delay



1.3 Fisher Information

Fisher information is used together with svd decomposition.[2006] We can define the Fisher Information in the following way^[6]:

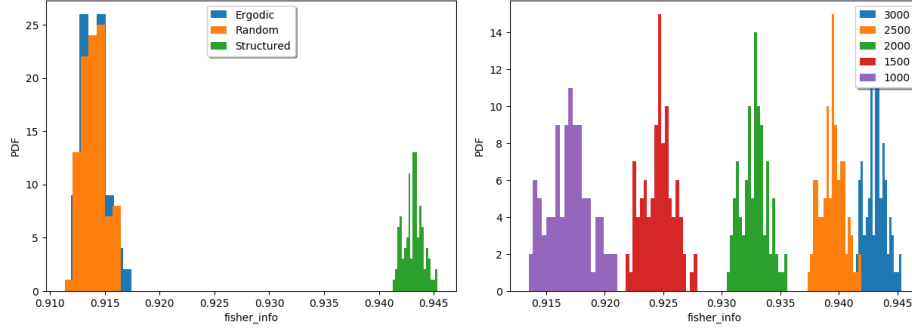
$$I_F = \sum_k^{M-1} \frac{(\bar{\sigma}_{k+1} - \bar{\sigma}_k)^2}{\bar{\sigma}_{k-1}} \quad (10)$$

The quantities are exactly the same that we have defined in the previous section. The difference between this function and the svd entropy is that FI behaves

⁶for reference see here:<https://www.mdpi.com/1099-4300/23/11/1424>

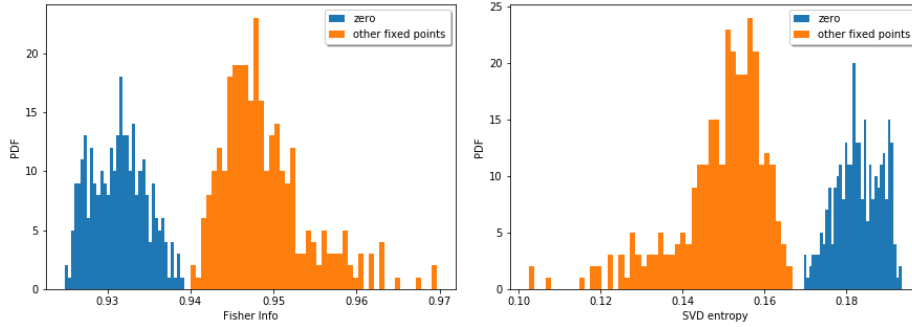
contrary to the entropy. We can see this in the plots.[2021]

1.3.1 FISHER INFORMATION

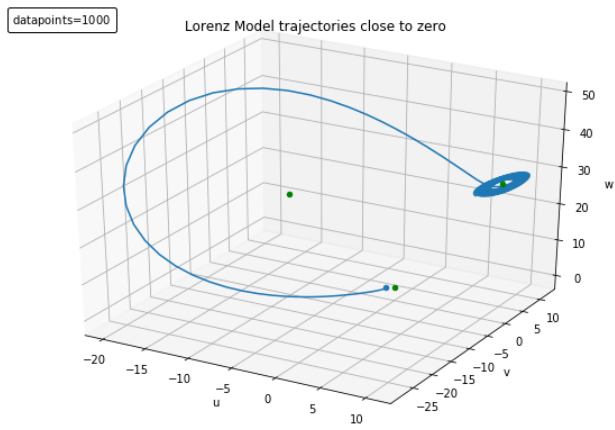
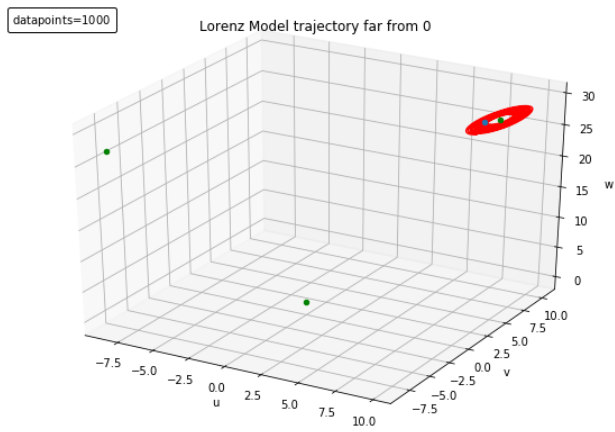


1.4 Is zero the most informative point ?

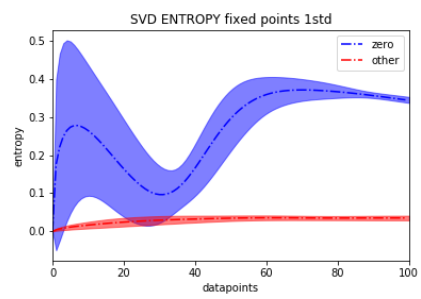
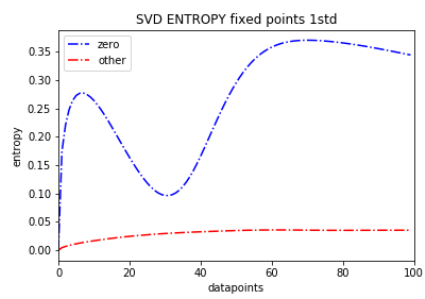
We can try to measure the amount of information that we obtain if we use a trajectory close to zero to the train our model. By following the same ideas that helped us in the previous section we are going to see the results of this computation. It has been used the same amount of data to obtain these two plots.

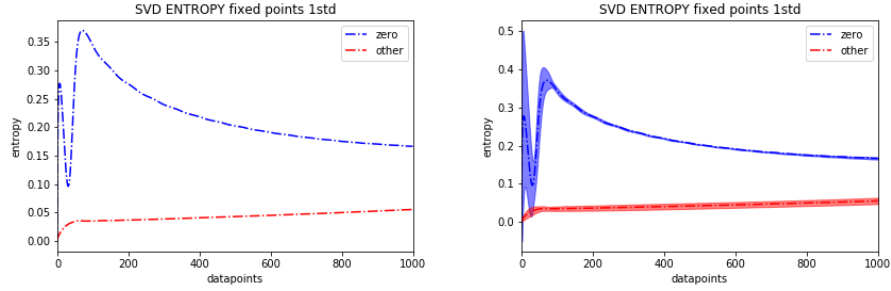


Here we can see also two different trajectories in the attractor.



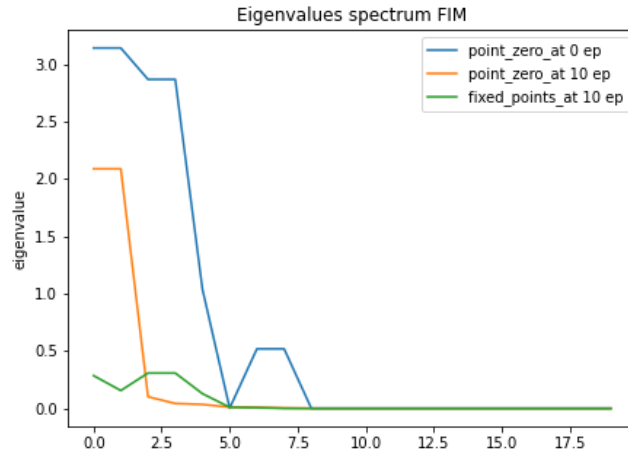
These two trajectories have different svd entropy growth.





1.4.1 FIM and information of zero

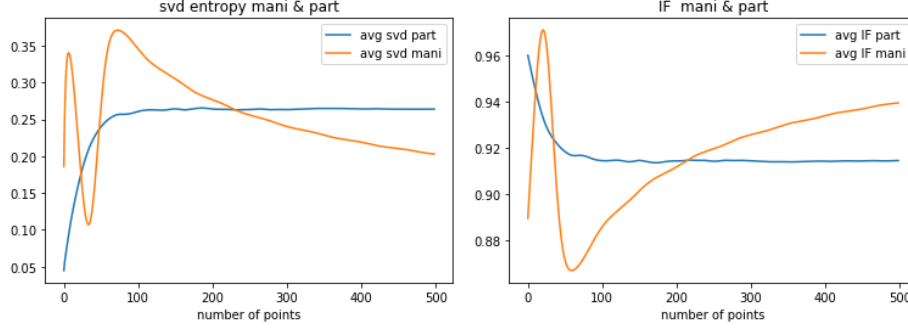
To answer to the question of this section we can use the FIM approach to see how this matrix changes if we use as training set the trajectories close to zero, and the others. We can find these results in the following plots.



As we can see the dataset of trajectories closes to zero seems to contain less information than the others. Indeed we see that after 10 epochs the green curve is close to zero.

1.5 Growth of Information and Entropy

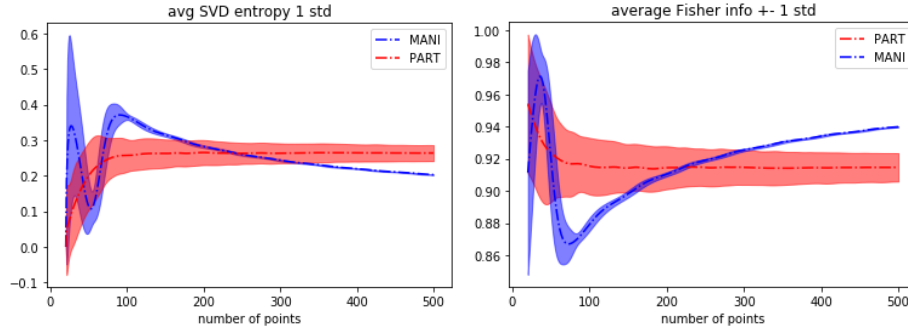
We can try to see the growth of information when we increase the number of trajectory's points. We are going to see the behaviour of Fisher Information and SVD entropy in two type of dataset: part (random sampling of Lorenz Orbit Trajectories) and mani (wich is a dataset of trajectories close to fixed points).



As we can see and as we can expect the plots behave in opposite ways. If we look at this by a statistical point of view, it could be useful to understand the average behaviour of the IF and SVD entropy, and also to understand the fluctuation that we have. We can try to compare the different shapes of the average amount of information (or entropy) of different trajectories when we increase the number of points.

It is interesting to notice that around 200 number of points we can see that the amount of FI is bigger in the manifold data than in the part. We have the same behaviour with the svd entropy, but as we expect it is in the opposite direction.

It is important also to study the fluctuation of information that we have. We can also notice by looking at these plots that the part trajectories seems to have a more stable variation of information than the manifold. In particular we see that at the very beginning of the curve, namely around 200 points the svd entropy of the manifold dataset increases a lot.



1.5.1 FISHER INFORMATION MATRIX

To study the differences between two type of trajectories we can analyze the difference of the amount of information that the LSTM extracts during the training process.

1.5.2 Theory of Information Matrix

We have a Dataset: $\mathcal{T} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ and we want to study the amount of information that a model (for instance LSTM) extracts during the training process from these data. We denote with $f(\hat{\mathbf{y}}; \theta)$ the density function of a random variable $\hat{\mathbf{Y}}$ parametrized by θ . We are interested in the shape of likelihood function. That is why we need the Fisher Matrix, that it is defined in this expression:

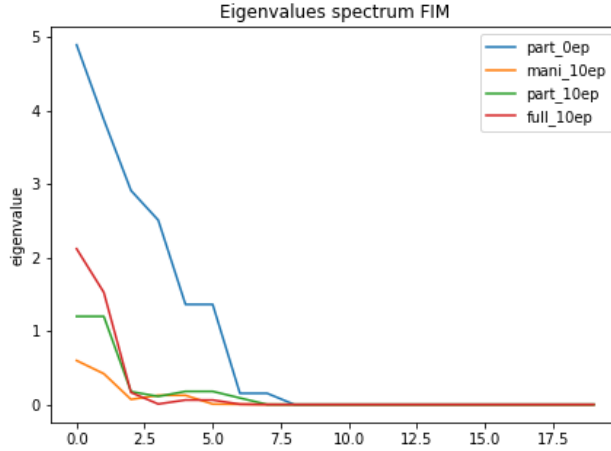
$$I(\theta) := \mathbb{E}_{\hat{\mathbf{Y}}} \left[(\nabla_{\theta} \log f(\hat{\mathbf{y}}; \theta))^2 \right] \quad (11)$$

This matrix, with its eigenvalues is able to describe the 'direction' that contains more information in the space of parameters. In practice, in a ML approach, we want to know the gradient of the loss function, $\mathcal{L}(\theta)$, with respect to model parameters.

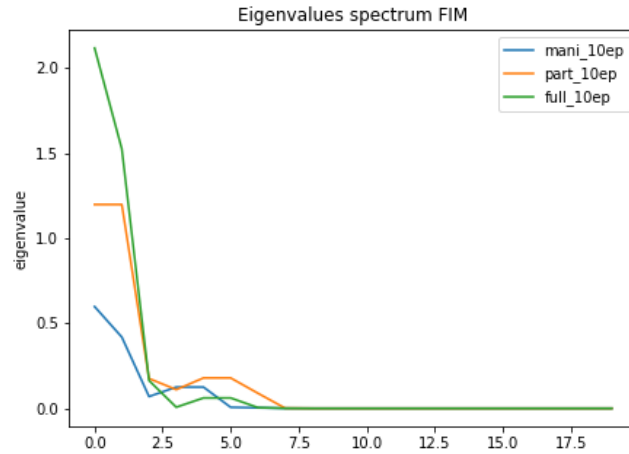
$$\nabla_{\theta} \log \hat{f}(\hat{\mathbf{y}}_i; \theta) \sim -\nabla_{\theta} \mathcal{L}_i + \frac{\sum_k \exp(-\mathcal{L}_k/\sigma^2) \nabla_{\theta} \mathcal{L}_k}{\sum_j \exp(-\mathcal{L}_j/\sigma^2)} \quad (12)$$

1.5.3 Fisher Matrix Plots

Then we can see how a RNN extracts information in different dataset type after 10 epochs. Using the same amount of training points we can see different plots.



We can see that in this plot the amount of information extracted from the dataset "mani", namely trajectories closes to fixed points, is greater than then the others dataset.



Here we have a zoom of the spectrum of the eigenvalues of the FIM of different trajectories after 10 epochs. As we can see the dataset that contains the bigger amount of information is "mani". Indeed the eigenvalues of the FIM are linked with the shape of the space in wich the "gradient descent" moves.

References

- Fraser A. M., Swinney H. L., 1986, Phys. Rev. A, 33, 1134
- Hotelling H., 1936, Biometrika, 28, 321
- Mayer A. L., Pawlowski C. W., Cabezas H., 2006, Ecological Modelling, 195, 72
- Raubitzek S., Neubauer T., 2021, Entropy, 23
- Shannon C. E., 1948, The Bell System Technical Journal, 27, 379
- Tang Y., Krakovská A., Mezeiová K., Budáčová H., 2015, Journal of Complex Systems, 2015, 932750