

Information Theory (UAI/500)



Lecture Notebook

Lecturer: Ivo Bukovsky

Something about Data and Systems that Generate them

- Data generating systems:
 - Stochastic (random process),
 - Deterministic
 - Chaotic (deterministic chaos, logistic eq. example)
 - Real (...omg)
- Anomaly detection: probabilistic or machine learning based ?

✓ Note on Deterministic Chaos

It can be good to aware of what deterministic chaos is (i.e., an unpredictable behavior of data generated by a deterministic governing law), just be aware it exists, ...)

Example of a very complex dynamical systems, or a very complex governing law in machine learning are recurrent neural networks (the ones with some kind of feedback, tapped delays,...). However do not panic, chaotic development in them can be very rare...

Example of Logistic Equation

The behavior of data (here scalar variable y) is driven by a deterministic governing law as follows

$$y(k) = a \cdot y(k-1) \cdot (1 - y(k-1)),$$

(1)

where k is discrete index of time, y is generated data (e.g., it has been demonstrated that extended to 3-D problem it works as an insect population [2]), a is called a bifurcation parameter (it changes qualitative behavior of generated data).

- What happens if we change a ?
- What happens if we keep $a = \text{const.}$ and change the initial condition $y(k=0)$ to $y(0) = y(0) + \text{"a butterfly"}$ where the "butterfly" is a very small number?

```

1 %matplotlib inline
2 %%matplotlib notebook
3 from numpy import *
4 from matplotlib.pyplot import *
5
6
7 a=4 # biffurcation parameter
8 N=10000
9 y=zeros(N)
10 y[0]=.1 # init. cond.
11 for k in range(1,N):
12     y[k]=a*y[k-1]*(1-y[k-1])
13
14 z=zeros(N)
15 z[0]=0.1000001 # init. cond.
16 for k in range(1,N):
17     z[k]=a*z[k-1]*(1-z[k-1])
18
19 figure(figsize=(10,4))
20 subplot(311)
21 plot(y, '.')
22 subplot(312)
23 plot(range(N-200,N),y[-200:], '-'),xlabel('k')
24 subplot(313)
25 plot(range(N-200,N),z[-200:], '-'),xlabel('k')
26 show()
27

```

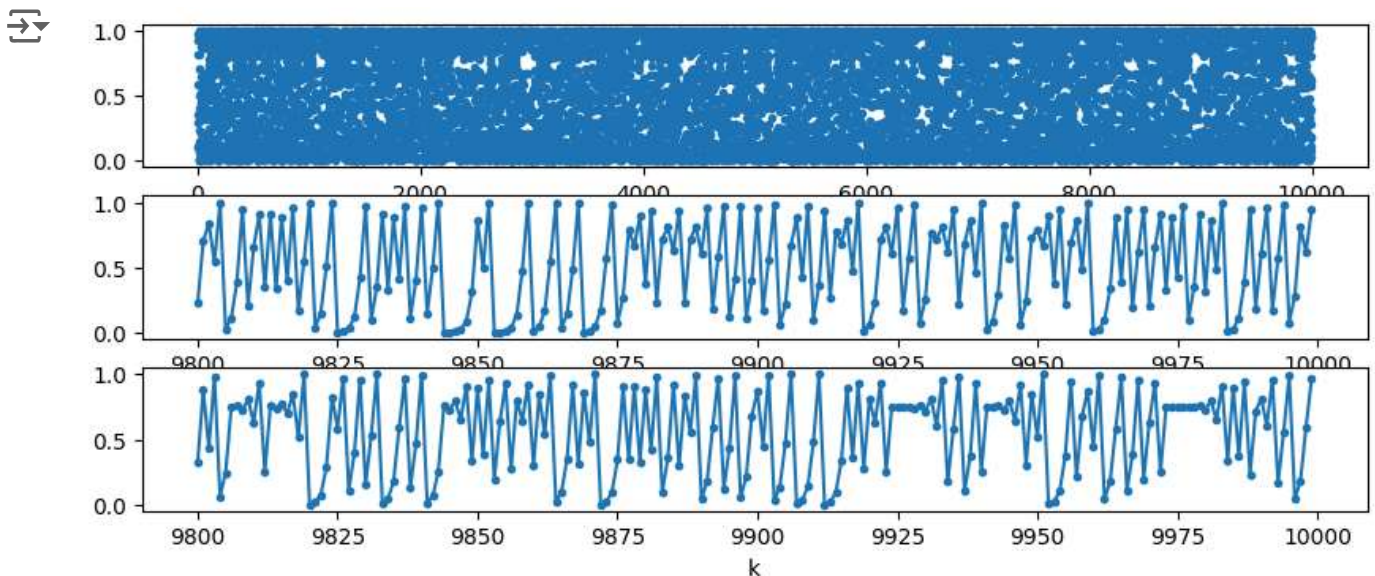
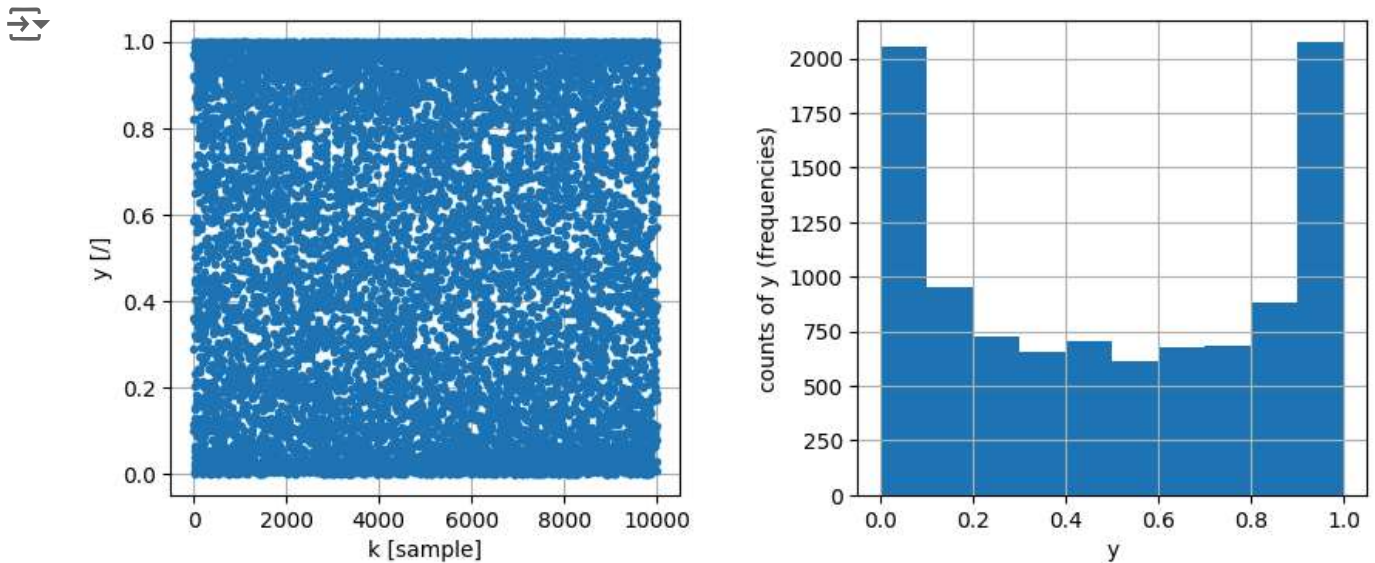


Fig. 1: The detail of last 200 samples differs due to a very small distinction in initial conditions ($y(k=0) = 0.1, z(k=0) = y(k=0) + 0.000001$, i.e. this is the butterfly effect where 0.000001 is the "butterfly" here. Also, notice, that the flat data behavior around $k = 9925$ or

$k = 9975$ may seem to be an anomaly in sense of probability, but not in sense of the governing law that is the same for all k . (and may be learned from data ...if done properly...)

```
1 figure(figsize=(10,4))
2 subplots_adjust(wspace=.35)
3 subplot(121)
4 xlabel("k [sample]"); ylabel("y [/]")
5 plot(y, '.'); grid()
6 subplot(122)
7 hist(y,10)
8 grid( ); xlabel("y"); ylabel("counts of y (frequencies)")
9 show()
```



```
1 savetxt("y.txt",y)
```

✓ Example L-7

Let's play a little bit with theoretical data generated by chaotic system (1) and practice some basic statistical techniques.

Generate or use those 10 000 samples by system (1) for random initial condition $0 < y(0) < 1$ with $a = 4$ (as done above).

Now, let's forget about how you generated time series y and consider y be a continuous random variable observed at every sample time $k = 0, 1, 2, \dots, 9\,999$.

Sliding window exercise:

1. The time series is long $N = 10\,000$ data samples. Calculate and draw mean trajectory and variance trajectory for a sliding window of length 5, 10, 30, 100, 1000 samples with 1-sample

sliding. What is the shortest window where mean and variance are stationary? (stationary in a loose sense, show and discuss your result and your opinion)[0.6 Points]

1-D probability:

2. Estimate the probability mass function via histogram $p(y) \forall y$ and discuss your selection of number of bins (reason for your choice, how does it affect $p(y)$? [0.4 Points]
 3. Calculate and plot the probability of each individual sample point $y(k) \forall k$, i.e., horizontal axis is time index $k = 1, 2, \dots, N$, vertical axis shows $p(y(k))$ [0.6 Points]
-

2-D probability:

4. Estimate the 2-D probability mass function $p(y, y_{previous}) \forall y$, show the 2-D plot and briefly discuss your solution (you can use Numpy or some existing libraries (reference the library and briefly explain the principles that the library applies). [0.8 Points]
5. Estimate the 2-D conditional probability function $p(y|y_{previous}) \forall y$, show the 2-D plot and briefly discuss your solution (you can use Numpy or some existing libraries (reference the library and briefly explain the principles that the library applies):
 - a) Estimate $p(y|y_{previous})$ directly from the sequence of data (time series) [0.8 Points]
 - a) Estimate $p(y|y_{previous})$ by the Bayes rule using $p(y, y_{previous})$ and $p(y_{previous})$ [1 Points]
6. Calculate and plot the probability of each individual data point $\mathbf{x}(k) = [y(k), y(k-1)]$ where horizontal axis is time index k , and the vertical axis shows $p([y(k), y(k-1)])$ [0.6 Points]

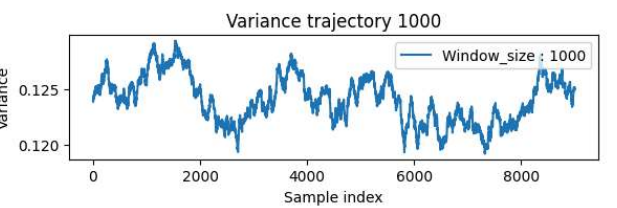
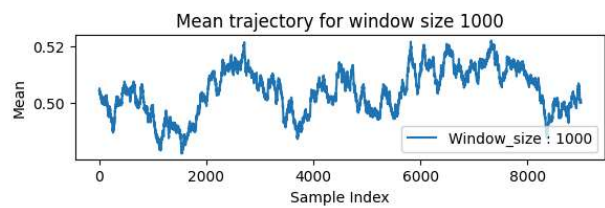
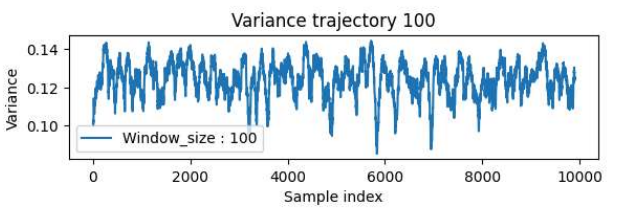
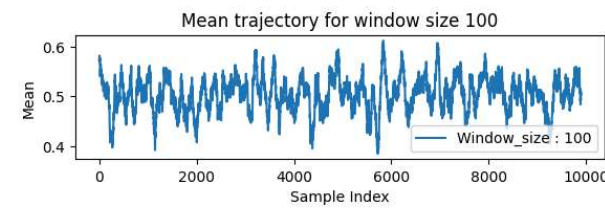
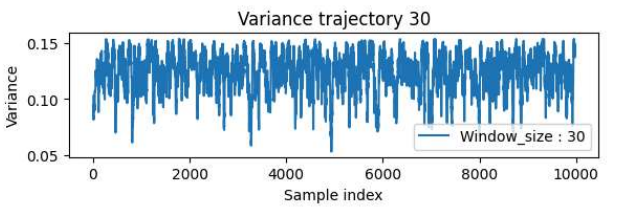
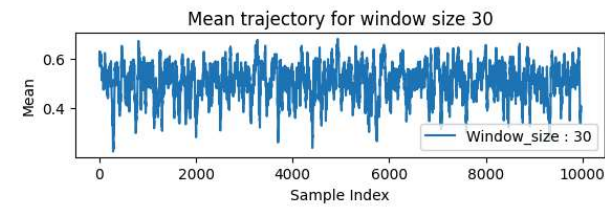
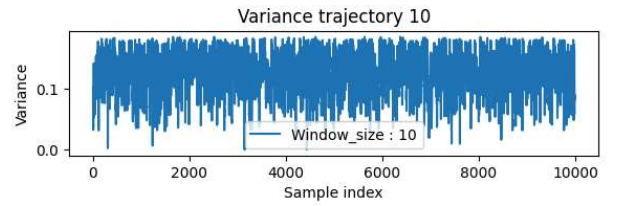
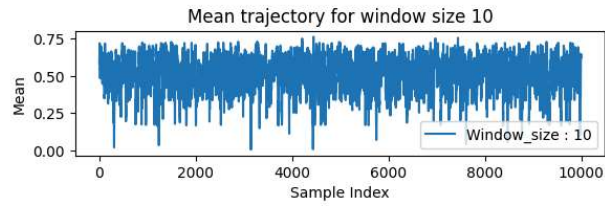
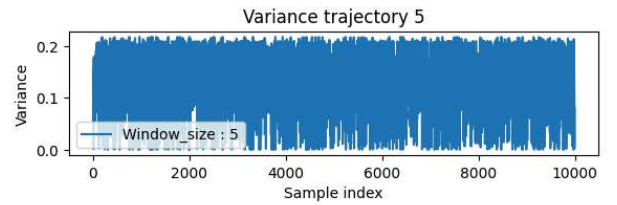
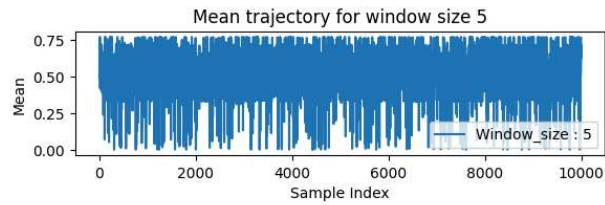


```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 """Sliding window excersize:
6 1) The time serie is long N=10 000 data samples. Calculate and draw mean trajectory
7 What is the shortest window where mean and variance are stationary?"""
8
9 a=4 # biffurcation parameter
10 N=10000
11
12 initial_condition = np.random.rand() # Random initial condition 0 < y(0) < 1
13 y = np.zeros(N)
14 y[0] = float(initial_condition)
15
16 for k in range(1,N):
17     y[k] = a * y[k-1] * (1-y[k-1]) # logistic map function
18
19 window_sizes = [5, 10, 30, 100, 1000] # window sizes
20
21 # Function to calculate mean and variance for time series data
22 def window_sliding(data, window_size):
```

```

23  series = pd.Series(data)
24  mean_rolling = series.rolling(window = window_size).mean().dropna().values
25  var_rolling = series.rolling(window = window_size).var().dropna().values
26  return mean_rolling, var_rolling
27
28 result = {}
29 for i in window_sizes:
30     mean, variance = window_sliding(y, i)
31     result[i] = (mean, variance)
32
33
34 # plotting the function
35 plt.figure(figsize=(15, 12))
36
37 for i, window_size in enumerate(window_sizes):
38     means, variances = result[window_size]
39
40     plt.subplot(5,2,2*i+1)
41     plt.plot(means, label = f'Window_size : {window_size}')
42     plt.title(f'Mean trajectory for window size {window_size}')
43     plt.xlabel('Sample Index')
44     plt.ylabel('Mean')
45     plt.legend()
46
47     plt.subplot(5,2,2*i+2)
48     plt.plot(variances, label = f'Window_size : {window_size}')
49     plt.title(f'Variance trajectory {window_size}')
50     plt.xlabel('Sample index')
51     plt.ylabel('Variance')
52     plt.legend()
53
54 plt.subplots_adjust(wspace=0.3, hspace=0.7)
55 plt.show()

```



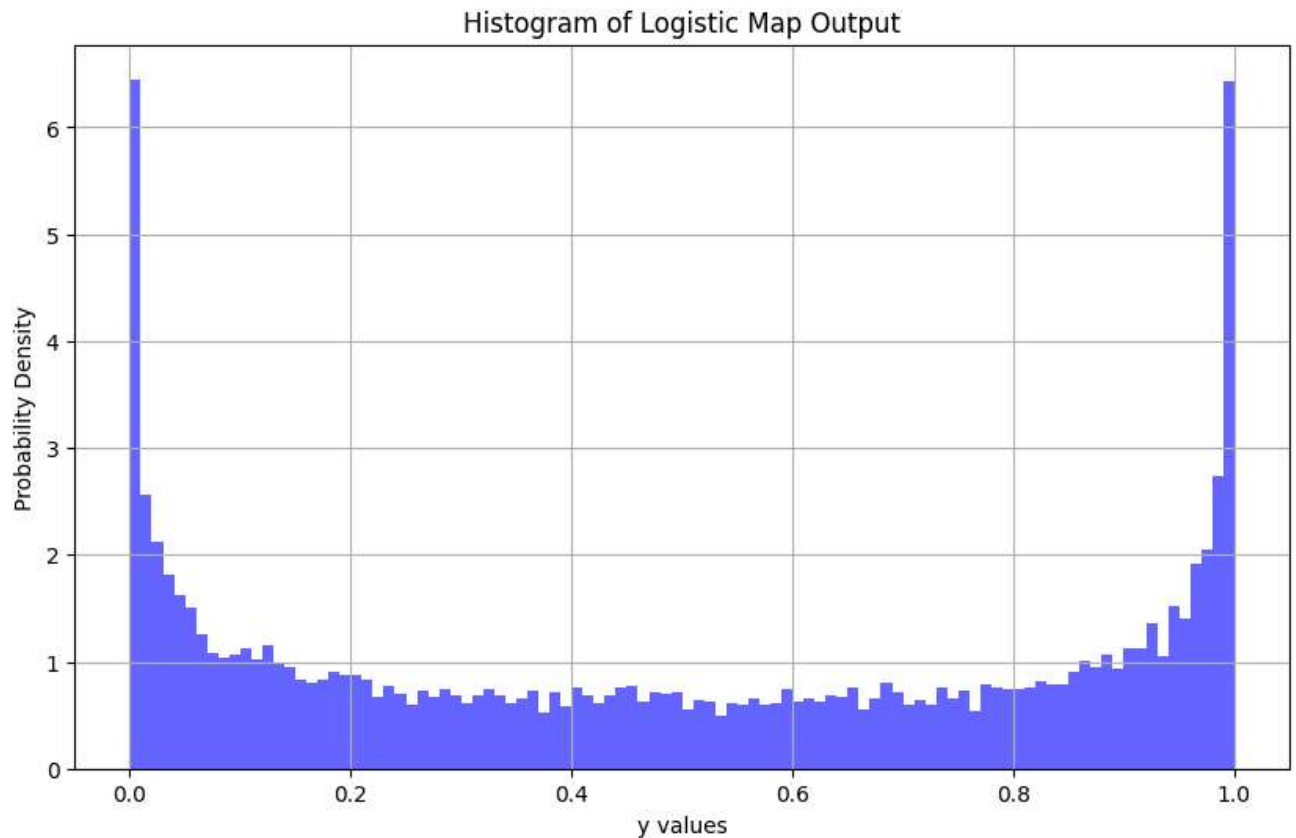
Window size 5 and 10 - The trajectories exhibit high fluctuations due to the small window size. The time series shows chaotic behavior and these small windows are not enough to smooth out the noise. The variance fluctuates wildly, indicating non-stationarity.

Window size 30 - The fluctuations in both mean and variance start to decrease, but there is still some noticeable variation. This suggests some initial stabilization, but the system has not fully reached stationarity.

Window size 100 - Both the mean and variance appear to be more stable compared to smaller window sizes. Although there are still minor variations, the system is much closer to stationarity.

Window size 1000 - Here, the mean and variance show very little change over time. This large window size effectively smooths the data, suggesting that the system has become stationary.

```
1 """2) Estimate the probability mass function via histogram  $p(y) \forall y$  and discuss your
2 """
3
4 num_bins = int(np.sqrt(N))
5
6 plt.figure(figsize=(10, 6))
7 plt.hist(y, bins=num_bins, density=True, alpha=0.6, color='b')
8
9 # Adding titles and labels
10 plt.title('Histogram of Logistic Map Output')
11 plt.xlabel('y values')
12 plt.ylabel('Probability Density')
13 plt.grid()
14 plt.show()
```



I have used the square root of the sample size as the appropriate choice for selecting the bins as it balances the detail and smoothness in the distribution.

If the number of bins is too low, the histogram will be overly simplistic, failing to capture the nuances of the distribution. Important features such as peaks and valleys may be lost, making it hard to discern the actual shape of the PMF.

Conversely, using too many bins can result in a noisy histogram that fluctuates widely. This is particularly problematic with limited data points, as random variations may appear as significant features of the distribution which can mislead interpretations.

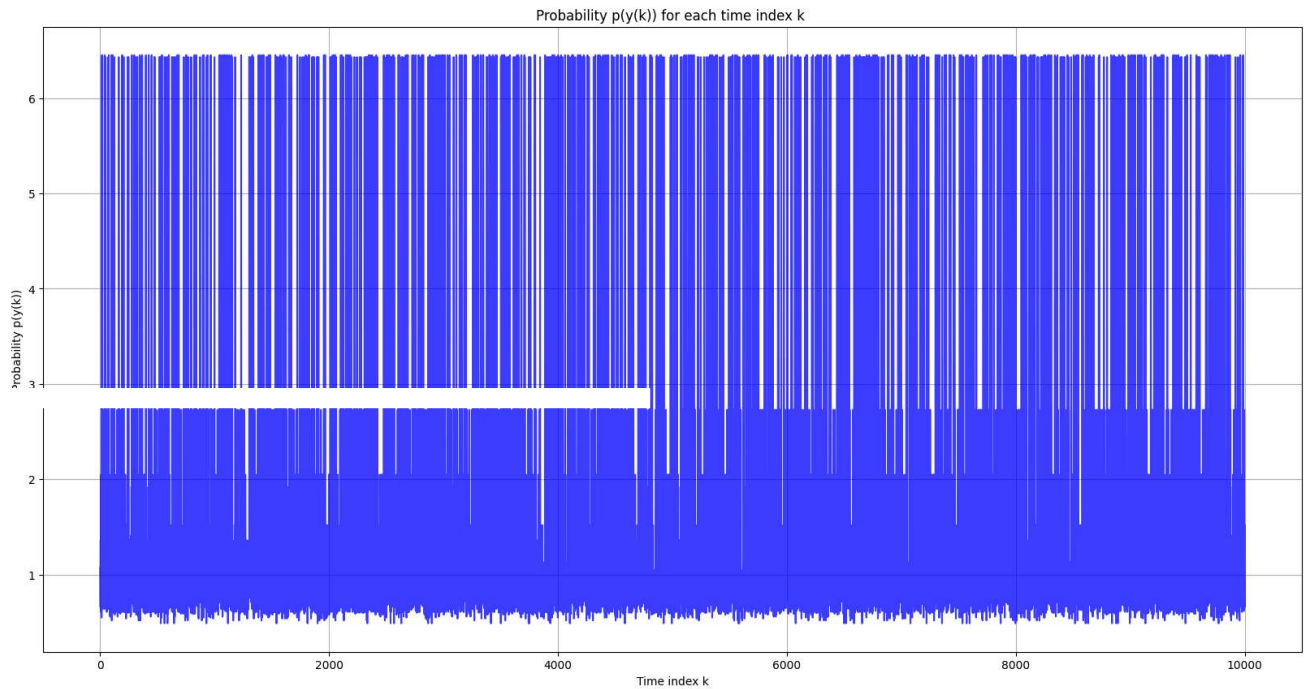
```
1 """ 3)Calculate and plot the probability of each individual sample point  $y(k) \forall k$ , i.
2 """
3
4 num_bins = int(np.sqrt(N)) # Number of bins
5 hist, bin_edges = np.histogram(y, bins=num_bins, density=True)
6
7 bin_indices = np.digitize(y, bin_edges[:-1]) - 1 # Bins are indexed from 0
8 bin_indices[bin_indices >= num_bins] = num_bins - 1 # Correct out-of-bound indices
9
10 # Get the probability for each  $y(k)$  based on the corresponding bin
11 prob_yk = hist[bin_indices]
```



```

12
13 # Step 3: Plot  $p(y(k))$  for each  $k$  (time index)
14 plt.figure(figsize=(20, 10))
15 plt.plot(range(1, N+1), prob_yk, color='blue', alpha=0.75)
16 plt.title('Probability  $p(y(k))$  for each time index  $k$  ')
17 plt.xlabel('Time index  $k$ ')
18 plt.ylabel('Probability  $p(y(k))$ ')
19 plt.grid(True)
20
21 plt.show()

```



```

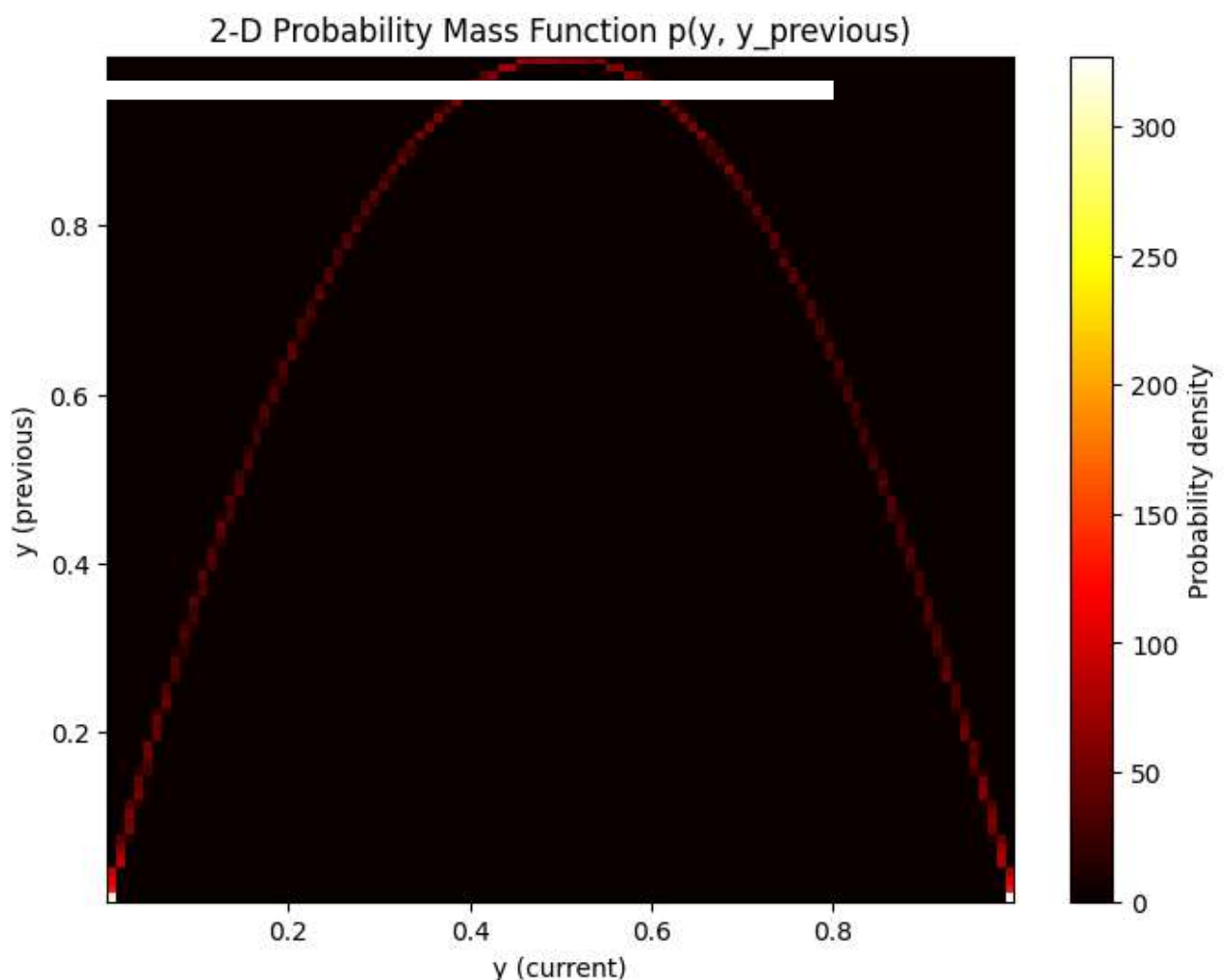
1 """4) Estimate the 2-D probability mass function  $p(y, y_{\text{previous}})$   $\forall y$  , show the 2-D p
2 """
3

```

```

4 y_previous = y[:-1] # y(k-1)
5 y_current = y[1:]    # y(k)
6
7 # Step 2: Compute the 2-D histogram to estimate joint distribution p(y, y_previous)
8 num_bins = int(np.sqrt(N)) # Set the number of bins
9 hist, xedges, yedges = np.histogram2d(y_previous, y_current, bins=num_bins, density=True)
10
11 # Step 3: Plot the 2-D PMF as a heatmap
12 plt.figure(figsize=(8, 6))
13 plt.imshow(hist.T, origin='lower', extent=[xedges[0], xedges[-1], yedges[0], yedges[-1]])
14 plt.colorbar(label='Probability density')
15 plt.title('2-D Probability Mass Function p(y, y_previous)')
16 plt.xlabel('y (current)')
17 plt.ylabel('y (previous)')
18 plt.grid(False)
19 plt.show()

```

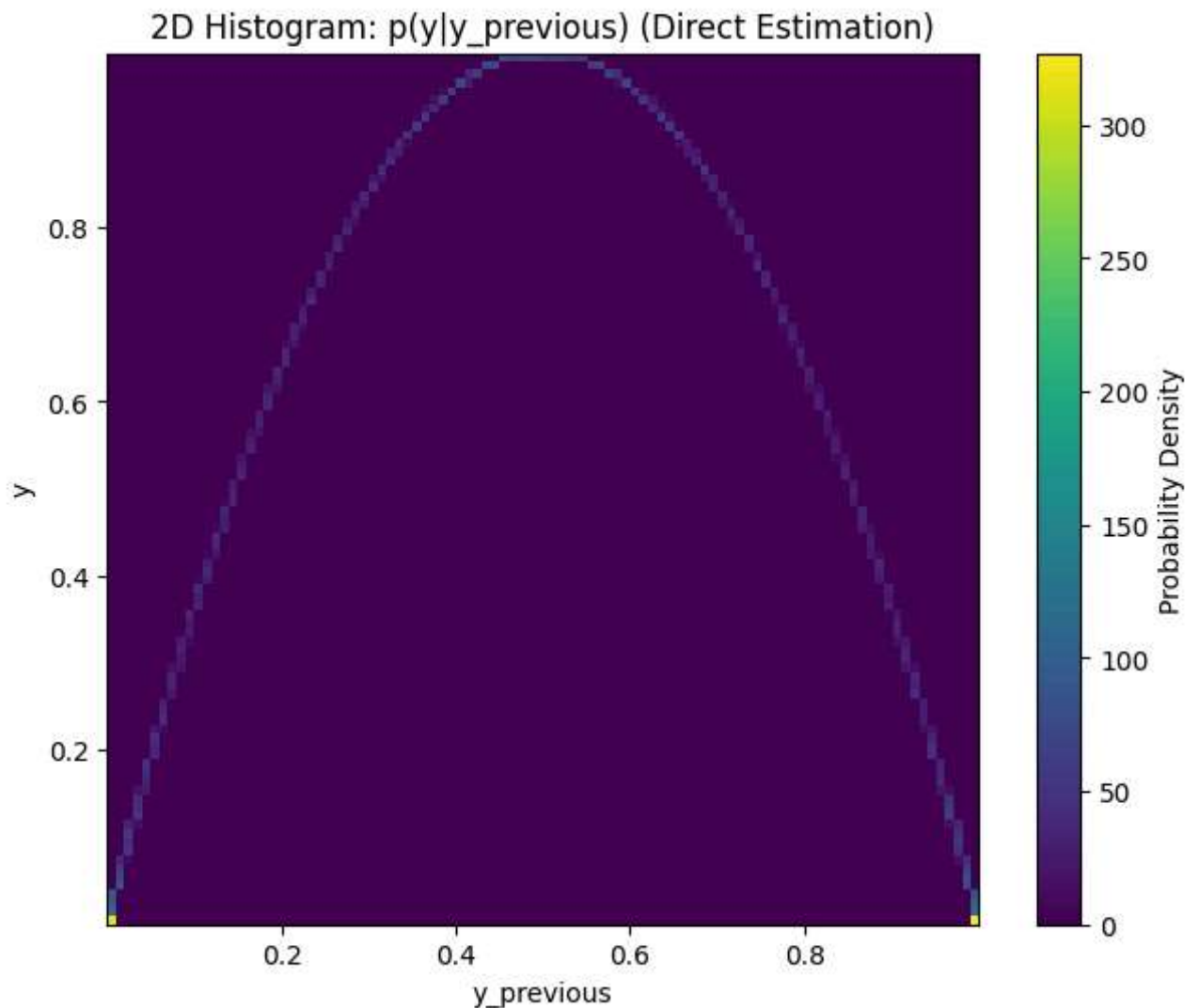


The plot exhibits a structured pattern, particularly with prominent regions in the lower-middle of the graph. This behavior is a result of the chaotic nature of the logistic map at $a=4$, where the system oscillates between different regions of the state space.

The relationship between y and y_{previous} shows significant density in certain areas, reflecting where the system spends more time.

We used `numpy.histogram2d` to compute the joint histogram of y and $y_{previous}$. By normalizing the counts into probabilities, we obtain a 2D PMF, which is visualized with a heatmap using `matplotlib`. This approach allows us to estimate the density of points in the phase space, capturing the dynamics of the system.

```
1 """ 5) Estimate the 2-D conditional probability function  $p(y|y_{previous}) \forall y$ , show t
2 a) Estimate  $p(y|y_{previous})$  directly from the sequence of data (time series)
3
4 """
5 y_previous = y[:-1] #  $y(k-1)$ 
6 y_current = y[1:]   #  $y(k)$ 
7
8 num_bins = int(np.sqrt(N)) # Set the number of bins for histogram
9
10 # 2D histogram (joint distribution of  $y$  and  $y_{previous}$ )
11 hist, xedges, yedges = np.histogram2d(y_previous, y_current, bins=num_bins, density=True)
12
13 # Plotting the 2D histogram
14 plt.figure(figsize=(8, 6))
15 plt.imshow(hist.T, origin='lower', cmap='viridis', extent=[xedges[0], xedges[-1], yedges[0], yedges[-1]])
16 plt.colorbar(label='Probability Density')
17 plt.title('2D Histogram:  $p(y|y_{previous})$  (Direct Estimation)')
18 plt.xlabel('y_previous')
19 plt.ylabel('y')
20 plt.show()
```



This method simply counts the occurrences of each pair (y_{previous}, y) and normalizes the counts to represent the probability density.

In the direct estimate, the conditional probabilities are approximated but may not sum to 1 for each y_{previous} . This could lead to inaccuracies in interpreting the probability distribution.

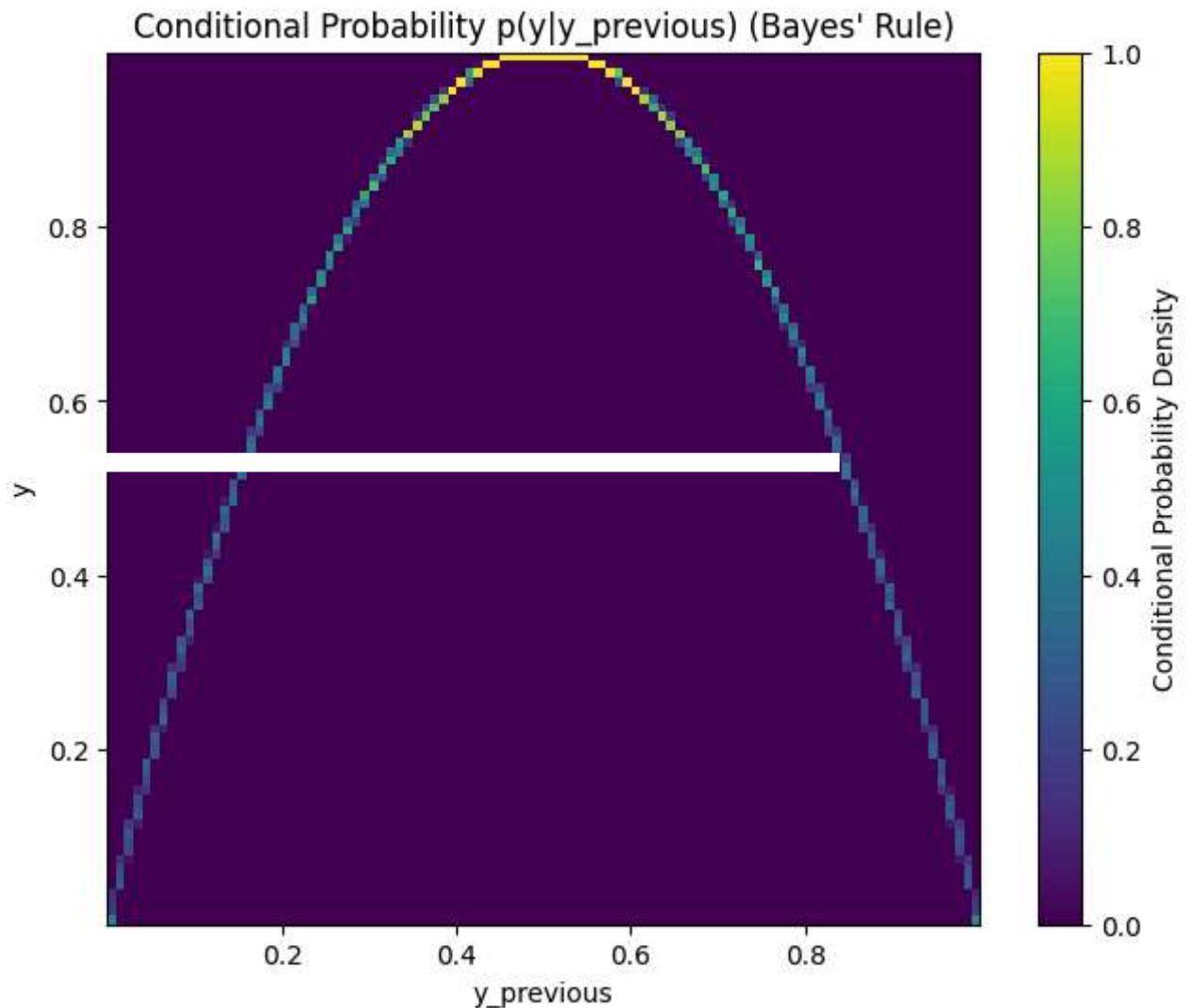
The direct estimate might reflect the raw frequencies more accurately, but it's also more susceptible to sampling noise.

```
1 """5 b) Estimate  $p(y|y_{\text{previous}})$  by the Bayes rule using  $p(y, y_{\text{previous}})$  and  $p(y_{\text{previous}})$ 
2 """
3
4 # Calculate the marginal distribution  $p(y_{\text{previous}})$ 
5 p_y_previous = np.sum(hist, axis=1) # Marginalize over y (sum across rows)
6
7 # Avoid division by zero by ensuring no zero values in p_y_previous
8 p_y_previous = np.where(p_y_previous == 0, 1e-10, p_y_previous)
9
10 # Use Bayes' rule to calculate  $p(y|y_{\text{previous}})$ 
11 p_y_given_y_previous = hist / p_y_previous[:, np.newaxis]
12
13 # Plot the conditional probability  $p(y|y_{\text{previous}})$  using Bayes' rule
```

```

14 plt.figure(figsize=(8, 6))
15 plt.imshow(p_y_given_y_previous.T, origin='lower', cmap='viridis', extent=[xedges[0],
16 plt.colorbar(label='Conditional Probability Density')
17 plt.title('Conditional Probability p(y|y_previous) (Bayes\' Rule)')
18 plt.xlabel('y_previous')
19 plt.ylabel('y')
20 plt.show()

```



This method ensures that for each value of y_{previous} , the sum of the conditional probabilities across all possible values of y is 1.

By explicitly normalizing the joint distribution using the marginal $p(y_{\text{previous}})$, this method gives a more accurate and consistent representation of the conditional probabilities.

The resulting graph is smoothed and better conforms to the probabilistic interpretation, removing noise and better highlighting the regions where the probabilities are highest.

```

1 """ 6) Calculate and plot the probability of each individual data point  $x(k)=[y(k),y$ 
2     where horizontal axis is time index  $k$  , and the vertical axis shows  $p([y(k),y$ 
3 """
4
5 y_previous = y[:-1] #  $y(k-1)$ 
6 y_current = y[1:]   #  $y(k)$ 

```

```

7
8 # Step 2: Compute the 2-D histogram to estimate joint distribution p(y, y_previous)
9 num_bins = int(np.sqrt(N)) # Set the number of bins
10 hist, xedges, yedges = np.histogram2d(y_current, y_previous, bins=num_bins, density=1)
11
12 # Step 3: Assign probabilities to each data point x(k) = [y(k), y(k-1)]
13 # Find the bin indices for each (y(k), y(k-1)) pair
14 bin_x_indices = np.digitize(y_current, xedges[:-1]) - 1 # for y(k)
15 bin_y_indices = np.digitize(y_previous, yedges[:-1]) - 1 # for y(k-1)
16
17 # Correct out-of-bound indices
18 bin_x_indices[bin_x_indices >= num_bins] = num_bins - 1
19 bin_y_indices[bin_y_indices >= num_bins] = num_bins - 1
20
21 # Get the probability for each x(k) = [y(k), y(k-1)] based on the corresponding bin
22 prob_xk = hist[bin_x_indices, bin_y_indices]
23
24 # Step 4: Plot p([y(k), y(k-1)]) for each k (time index)
25 plt.figure(figsize=(20, 10))
26 plt.plot(range(1, N), prob_xk, color='blue', alpha=0.75)
27 plt.title('Probability p([y(k), y(k-1)]) for each time index k')
28 plt.xlabel('Time index k')
29 plt.ylabel('Probability p([y(k), y(k-1)])')
30 plt.grid(True)
31
32 plt.show()

```

