

Nonparametric Statistics: Homework 6

蒋翌坤 20307100013

Solution to Problem 1

From the data, we have $n = 1425$. Using cross-validation method, we have,

$$\hat{J}(h) = \int \left(\hat{f}_n(x) \right)^2 dx - \frac{2}{n} \sum_{i=1}^n \hat{f}_{(-i)}(X_i)$$

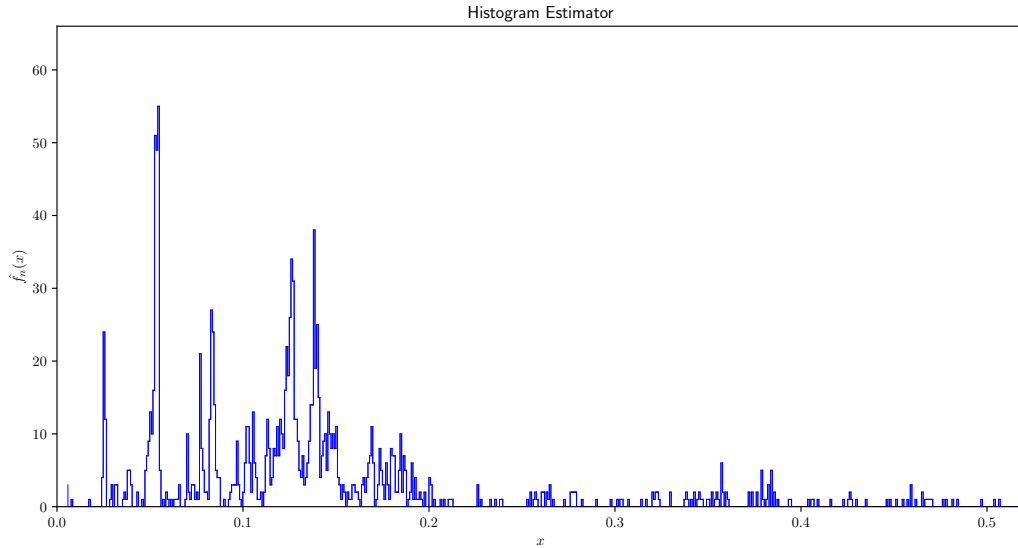
Histogram estimator

Here, $\hat{f}_n(x)$ is like the histogram of observations with $m = \frac{1}{h}$ bins. We also have,

$$\hat{J}(h) = \frac{2}{h(n-1)} - \frac{n+1}{h(n-1)} \sum_{j=1}^m \hat{p}_j^2$$

We can not deduce the optimal h from the formula, but, we can enumerate each m . In this way, we can find the minimum $\hat{J}(h)$ and the corresponding m . In this case, $\hat{J}(h)^* = -5.99$ and $m^* = \frac{1}{h^*} = 597$

We can therefore plot the graph of $\hat{f}_n(x)$,



Kernal estimator

Here, We use the Gaussian kernel: $K(x) = \frac{1}{\sqrt{2\pi}} \exp\{-\frac{x^2}{2}\}$.

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right)$$

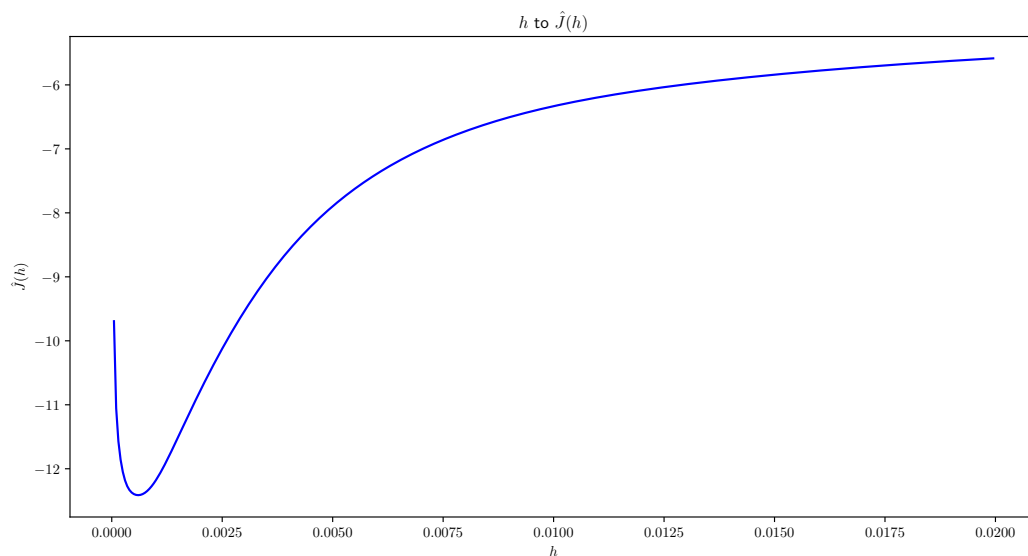
Then,

$$\begin{aligned} \hat{J}(h) &= \frac{1}{hn^2} \sum_i \sum_j K^*\left(\frac{X_i - X_j}{h}\right) + \frac{2}{nh} K(0) + O\left(\frac{1}{n^2}\right) \\ &= \frac{1}{hn^2} \sum_i \sum_j \left[K^{(2)}\left(\frac{X_i - X_j}{h}\right) - 2K\left(\frac{X_i - X_j}{h}\right) \right] + \frac{2}{nh\sqrt{2\pi}} + O\left(\frac{1}{n^2}\right) \end{aligned}$$

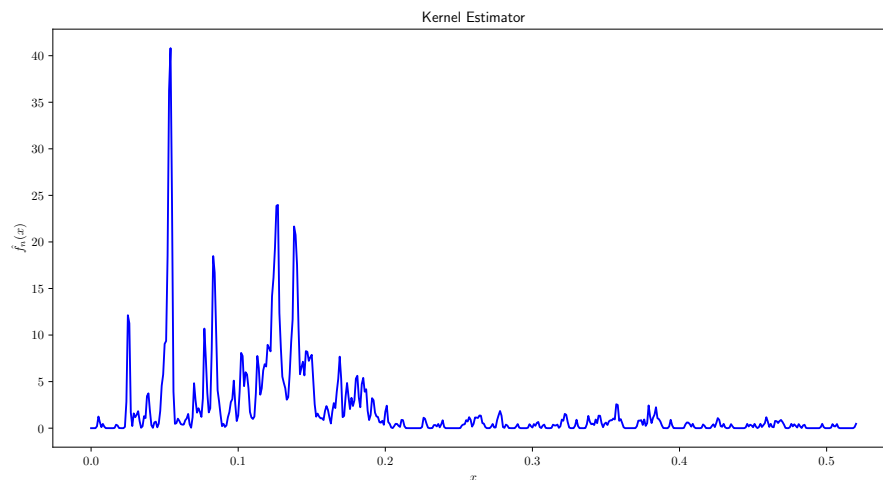
where,

$$K^{(2)}(x) = \frac{1}{2\sqrt{2\pi}} \exp\left\{-\frac{x^2}{8}\right\}$$

We can therefore plot the graph of $\hat{J}(h)$ with respect to h ,



From the graph, we can find the minimum $\hat{J}(h)$ and the corresponding m . In this case, $\hat{J}(h)^* \approx -12.41$ and $h^* \approx 0.0006$. We can therefore plot the graph of $\hat{f}_n(x)$ using $h^* = 0.006$,



Normal reference rule

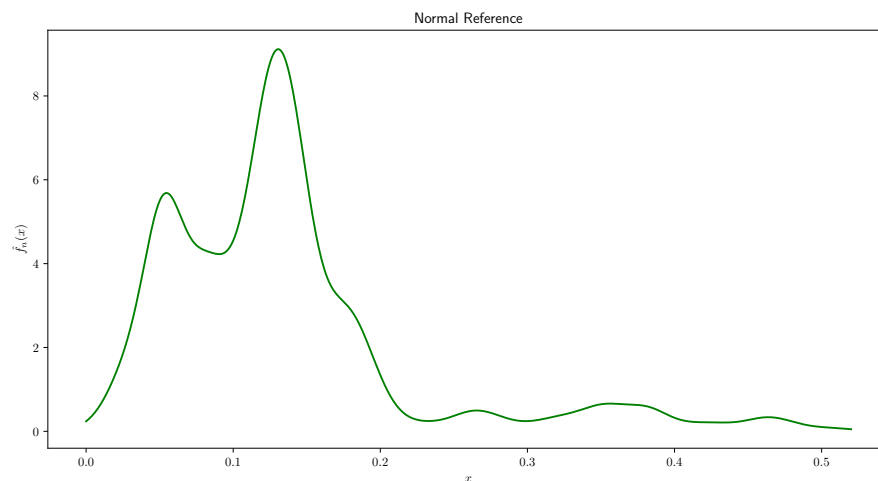
We can also apply Normal reference rule to pick a bandwidth for the kernel,

$$\hat{\sigma} = \min\left\{s, \frac{\text{Sample Interquantile range}}{1.34}\right\} = \min\{0.089, 0.053\} = 0.053$$

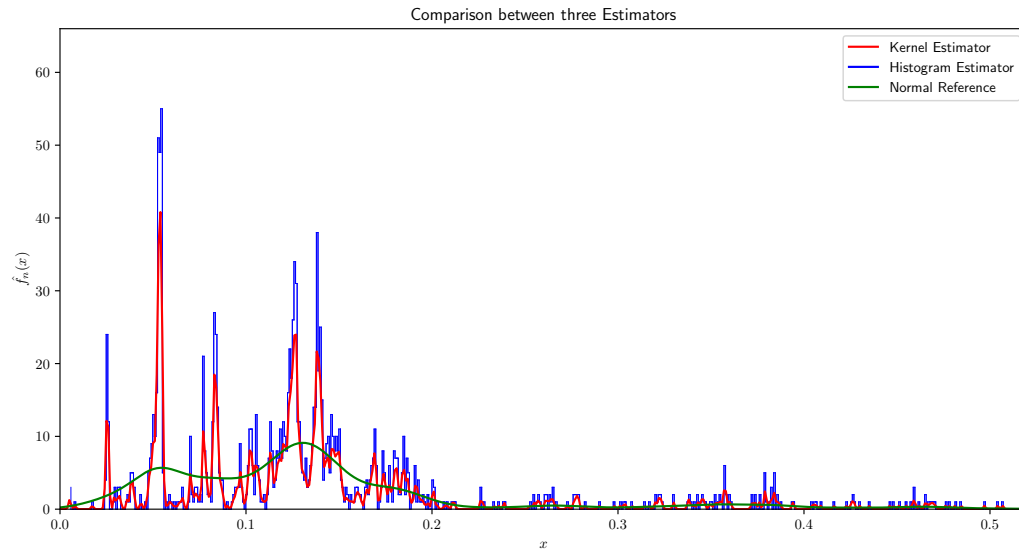
$$h_n = \frac{1.06\hat{\sigma}}{n^{1/5}} = 0.0131$$

Here, h_n is very different from the optimal h^* we get from enumeration.

We can therefore plot the graph of $\hat{f}_n(x)$ using $h_n = 0.0131$,



We can combine three graphs together, and get,



Appendix: Code

Python code and output used to solve the problem can also be found at:

<https://thisiskunmeng.github.io/nonparametric/hw6.html>

The following provides Python code I use to complete this homework:

```

1 import numpy as np
2 import pandas as pd
3 import requests
4 from io import StringIO
5 import matplotlib.pyplot as plt
6 data = requests.get("https://www.stat.cmu.edu/~larry/all-of-nonpar/=data/galaxy.dat", verify=False)
7 df = pd.read_csv(StringIO(data.text), sep="\s+", header=None, names=["ra", "declination", "redshift"])
8 redshift = df["redshift"]
9 print(len(redshift))

```

1425

Histogram estimator

```

1 def get_jh_hist(m, n, dat):
2     yj = np.histogram(dat, bins=m)[0]
3     pj = yj / n
4     first = 2 * m / (n - 1)
5     second = (n + 1) * m / (n - 1)
6     third = np.sum(pj ** 2)
7     return first - second * third
8
9 def get_jh_hist_all(dat):
10     jh = []
11     for i in range(len(dat)):
12         jh.append(get_jh_hist(i + 1, len(dat), dat))
13     return jh
14
15 jh_all = get_jh_hist_all(redshift)
16 print(np.min(jh_all))
17 m_star = np.argmin(jh_all) + 1
18 print(m_star)

```

-5.989204025439116

597

```

1 plt.rcParams["text.usetex"] = True
2 fig, ax = plt.subplots(figsize=(12, 6))
3 height = np.histogram(redshift, bins=m_star)[0] / len(redshift)
4 x = (np.histogram(redshift, bins=m_star)[1])

```

```

5
6 for i in range(1, m_star-1):
7     ax.add_line(plt.Line2D([x[i], x[i+1]], [height[i], height[i]], color='blue', linewidth=0.5))
8     ax.add_line(plt.Line2D([x[i], x[i]], [height[i - 1], height[i]], color='blue', linewidth=0.5))
9     ax.add_line(plt.Line2D([x[i+1], x[i+1]], [height[i], height[i + 1]], color='blue', linewidth=0.5))
10 ax.set_xlim(0, np.max(redshift))
11 ax.set_ylim(0, 1.2 * np.max(height))
12 ax.set_xlabel(r"$x$")
13 ax.set_ylabel(r"$\hat{f}_n(x)$")
14 ax.set_title(r"Histogram Estimator")
15 plt.savefig("./fig/histogram_estimator.pdf")

```

Kernal estimator

```

1 def kernel_function(xx):
2     return 1 / np.sqrt(2 * np.pi) * np.exp(-xx ** 2 / 2)
3
4 def kernel_function_2(xx):
5     return 1 / (2 * np.sqrt(2 * np.pi)) * np.exp(-xx ** 2 / 8)
6
7 def hat_f(xx, hx, dat):
8     return 1 / len(dat) * np.sum(kernel_function((xx - dat) / hx) / hx)
9
10 datx, daty = np.meshgrid(redshift, redshift)
11 datx_x = datx - daty
12
13 def jh_kernel(hx, dat):
14     n = len(dat)
15     the_sum = np.sum(kernel_function_2(datx_x / hx) - 2 * kernel_function(datx_x / hx))
16     the_jh = 1 / (hx * n ** 2) * the_sum + 2 / (hx * n * np.sqrt(2 * np.pi))
17     return the_jh
18
19 step = 0.00005
20 hn = np.arange(step, 0.02, step)
21 jh = []
22 for i in hn:
23     jh.append(jh_kernel(i, redshift))
24 jh = np.array(jh)
25
26 print(np.min(jh))
27 print(np.argmax(jh))
28 print((np.argmax(jh) + 1) * step)

```

```

-12.412517973369557
11
0.0006000000000000001

```

```

1 fig, ax = plt.subplots(figsize=(12, 6))
2 ax.plot(hn, jh, color="blue")
3 ax.set_xlabel(r"$h$")
4 ax.set_ylabel(r"$\hat{J}(h)$")

```

```

5 ax.set_title(r"$h$ to $\hat{J}(h)$")
6 plt.savefig("./fig/h_to_jh.pdf")

```

```

1 h_star = (np.argmin(jh) + 1) * step
2 fig, ax = plt.subplots(figsize=(12, 6))
3 x = np.arange(0, np.max(redshift), 0.001)
4 y = []
5 for i in x:
6     y.append(hat_f(i, h_star, redshift))
7 ax.plot(x, y, color="blue")
8 ax.set_xlabel(r"$x$")
9 ax.set_ylabel(r"$\hat{f}_n(x)$")
10 ax.set_title(r"Kernel Estimator")
11 plt.savefig("./fig/kernel_estimator.pdf")

```

```

1 # Comparison between histogram estimator and kernel estimator
2 fig, ax = plt.subplots(figsize=(12, 6))
3 height = np.histogram(redshift, bins=m_star)[0]
4 x = (np.histogram(redshift, bins=m_star)[1])
5 x_k = np.arange(0, np.max(redshift), 0.001)
6 y_k = []
7 for i in x_k:
8     y_k.append(hat_f(i, h_star, redshift))
9 for i in range(1, m_star-1):
10     ax.add_line(plt.Line2D([x[i], x[i+1]], [height[i], height[i]], color='blue', linewidth=0.5))
11     ax.add_line(plt.Line2D([x[i], x[i]], [height[i - 1], height[i]], color='blue', linewidth=0.5))
12     ax.add_line(plt.Line2D([x[i+1], x[i+1]], [height[i], height[i + 1]], color='blue', linewidth=0.5))
13 ax.plot(x_k, y_k, color="red", label="Kernel Estimator")
14 ax.plot(0,0, color='blue', label="Histogram Estimator")
15 ax.set_xlim(0, np.max(redshift))
16 ax.set_ylim(0, 1.2 * np.max(height))
17 ax.legend()
18 ax.set_xlabel(r"$x$")
19 ax.set_ylabel(r"$\hat{f}_n(x)$")
20 ax.set_title(r"Comparison between Histogram Estimator and Kernel Estimator")
21 plt.savefig("./fig/histogram_and_kernel_estimator.pdf")

```

Normal reference rule

```

1 s = np.std(redshift)
2 iqr = np.quantile(redshift, 0.75) - np.quantile(redshift, 0.25)
3 h = 1.06 * np.min([s, iqr / 1.34]) * len(redshift) ** (-1 / 5)
4 print(s, iqr / 1.34, h)

```

```
0.0889070103989908 0.05297985074626865 0.0131417627598199
```

```

1 fig, ax = plt.subplots(figsize=(12, 6))
2 x_n = np.arange(0, np.max(redshift), 0.001)
3 y_n = []
4 for i in x_n:
5     y_n.append(hat_f(i, h_norm, redshift))
6 ax.plot(x_n, y_n, color="green")
7 ax.set_xlabel(r"$x$")
8 ax.set_ylabel(r"$\hat{f}_n(x)$")
9 ax.set_title(r"Normal Reference")
10 plt.savefig("./fig/normal_reference.pdf")

```

```

1 # Comparisons
2 fig, ax = plt.subplots(figsize=(12, 6))
3 for i in range(1, m_star-1):
4     ax.add_line(plt.Line2D([x[i], x[i+1]], [height[i], height[i]], color='blue', linewidth=0.5))
5     ax.add_line(plt.Line2D([x[i], x[i]], [height[i - 1], height[i]], color='blue', linewidth=0.5))
6     ax.add_line(plt.Line2D([x[i+1], x[i+1]], [height[i], height[i + 1]], color='blue', linewidth=0.5))
7 ax.plot(x_k, y_k, color="red", label="Kernel Estimator")
8 ax.plot(0,0, color='blue', label="Histogram Estimator")
9 ax.plot(x_n, y_n, color="green", label="Normal Reference")
10 ax.set_xlim(0, np.max(redshift))
11 ax.set_ylim(0, 1.2 * np.max(height))
12 ax.legend()
13 ax.set_xlabel(r"$x$")
14 ax.set_ylabel(r"$\hat{f}_n(x)$")
15 ax.set_title(r"Comparison between three Estimators")
16 plt.savefig("./fig/normal_reference_comp.pdf")

```