# Nonparametric Statistics: Homework 7

蒋翌坤 20307100013

## Solution to Problem 1

From definition, we have $L_{ij} = l_j(x_i)$ and for any $x$, $\sum_{j=1}^n l_j(x) = 1$

$$
\begin{aligned}
Y_i - \hat{r}_{(-i)}(x_i) &= Y_i - \Big( \sum_{j=1}^n \frac{l_j(x_i)}{\sum_{k \neq i} l_k(x_i)} Y_j - \frac{l_i(x_i)}{\sum_{k \neq i} l_k(x_i)} Y_i \Big) \\
&= \Big( 1 + \frac{L_{ii}}{1 - L_{ii}} \Big) Y_i - \sum_{j=1}^n \frac{L_{ij}}{1 - L_{ii}} Y_j \\
&= \frac{Y_i - \hat{r}_n(x_i)}{1 - L_{ii}}
\end{aligned}
$$

Therefore, we have proved Theorem (5.34)

$$
\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{r}_{(-i)}(x_i))^2 = \frac{1}{n} \sum_{i=1}^n \Big( \frac{(Y_i - \hat{r}_n(x_i))}{(1 - L_{ii})} \Big)^2
$$

∎

# Solution to Problem 2

Several assumptions and methods used to solve the problem need pointing out.

- The chosen kernel $K(x)$ is the Gaussian kernel.

- Generalized cross validation is used to avoid the case where $L_{ii} = 1$.

- Optimal smoothing parameter is chosen by minimizing generalized cross validation score $\mathrm{GCV}(h)$ through analyzing $h$-$\mathrm{GCV}(h)$ plot (see Figure A.1 in the appendix) as well as using functions to find $\arg\min \mathrm{GCV}(h)$.

- In the spline estimator, we use each sample $x$ points as knots and optimize $\lambda$. If we encounter $x$ that has the same value, we take the average of $y$ corresponding to those $x$ and treat the sample points as a single point.

- Code for constructing estimator and producing figures can be found in the appendix.

A summary of four estimators is shown in Table 2.1.

| Estimator | Optimal smoothing parameter | Variance estimation | $c$ (for calculating CB) |
|---|---|---|---|
| Regressogram | $m = 30$ | 5.543 | 1.96 |
| Kernel | $h = 0.0418$ | 5.775 | 3.38 |
| Local Linear | $h = 0.0442$ | 5.769 | 3.39 |
| Spline | $\lambda = 2.176 \times 10^{-6}$ | 4.982 | 3.58 |

Table 2.1: Summary of four estimators

Confidence band and estimation of each estimator is shown in Figure 2.1.
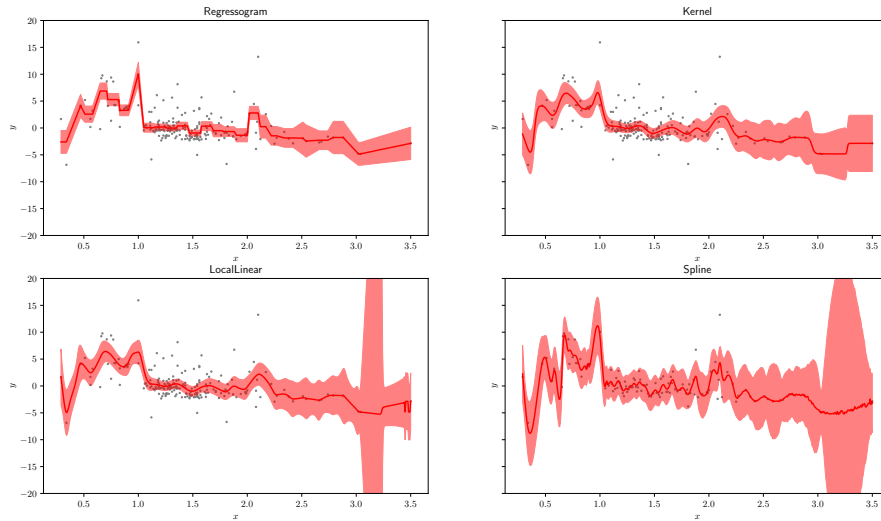


Figure 2.1: Confidence band and estimation of each estimator

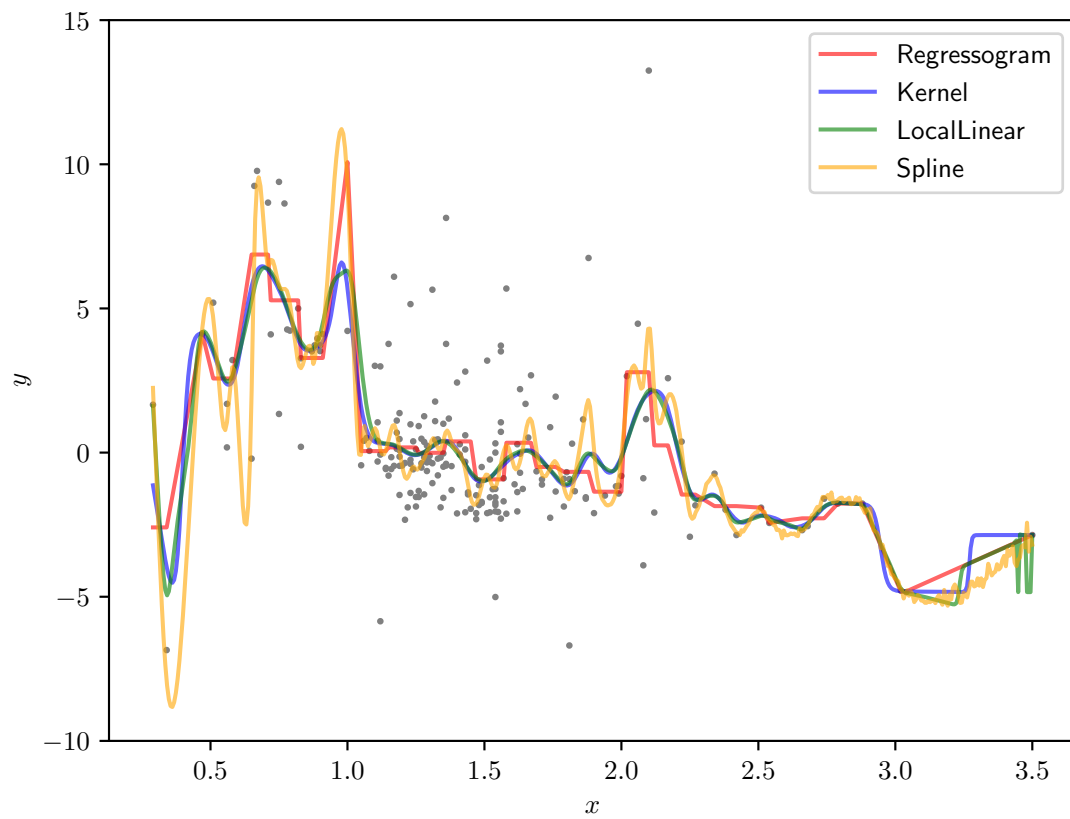A visual comparison of four estimators is shown in Figure 2.2.



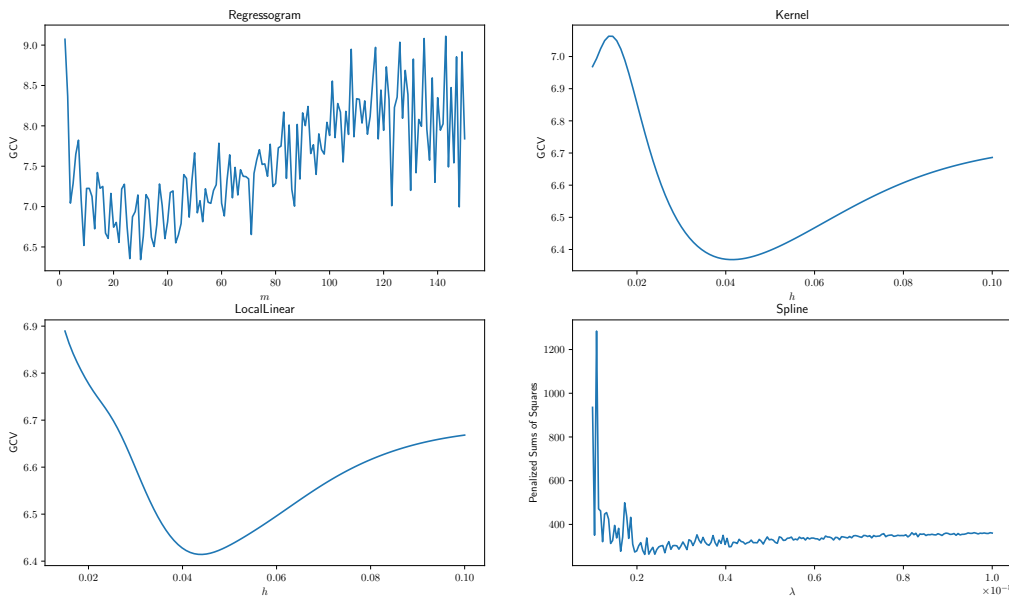Figure 2.2: Comparison among four estimators

# Appendix

**Figures**



Figure A.1: GCV or Penalized sum of squares of four different estimators

**Code**

Some explainations about the code.

- I use python to implement the code instead of R.

- Estimators are calculated using basic functions. Note that derivative and integral are computed using modules `numdifftools` and `scipy.integrate`.

- Estimators are constructed using class for better readability and convenience.

- In the spline estimator, we use nartual cubic spline basis instead of B-spline basis for simpler code implementation. Nartual cubic spline basis is roughly defined as the following:

$$f(x) = \sum_{k=1}^{n} \beta_k N_k(x), \quad N_1(x) = 1, \quad N_2(x) = x, \quad N_{k+2}(x) = d_k(x) - d_{n-1}(x)$$

where

$$d_k(x) = \frac{(x - \xi_k)_+^3 - (x - \xi_n)_+^3}{\xi_n - \xi_k} \quad \xi_k \text{ are knots}$$

- Note that some of the codes are not effective. Running all of the code may take as long as one hour.

- Code can also be found in https://thisiskunmeng.github.io/nonparametric/hw7.html.

```python
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   from scipy.stats import norm
5   import scipy.integrate as integrate
6   from scipy.optimize import fsolve
7   import numdifftools as nd
8   from tqdm import tqdm
9
10
11  class Estimate:
12      """
13      Base Class for nonparametric regression estimator
14      """
15
16      def __init__(self, y: pd.Series, x: pd.Series):
17          self.y = y
18          self.x = x
19          self.n = self.x.shape[0]
20          self.parameter = None
21          self.optimal_parameter = None
22          self.effective_kernel = None   # L
23          self.v = None  # trace of L
24          self.v_hat = None  # trace of L^T L
25          self.kappa0 = None  # use to calculate confidence interval coefficient
26          self.c = None  # confidence interval coefficient
27          self.method = None  # The method used to estimate the parameter
28
29      def l_x(self, i) -> np.ndarray:
30          pass
31
32      def get_effective_kernel(self):
33          self.effective_kernel = np.array([self.l_x(i) for i in self.x])
34          self.v = np.trace(self.effective_kernel)
35          self.v_hat = np.trace(self.effective_kernel.T @ self.effective_kernel)
36
37      def set_optimal_parameter(self, para):
38          self.optimal_parameter = para
39          self.set_parameter(para)
40          self.get_effective_kernel()
41
42      def set_parameter(self, para):
43          self.parameter = para
44          self.get_effective_kernel()
45
46      def estimate(self, x) -> float | np.ndarray:
47          """
48          :return: .. math:: \hat{r}_n(x)
49          """
50          if isinstance(x, np.ndarray):
51              if x.ndim == 0:
52                  return np.sum(self.y.values * self.l_x(x))
53              return np.array([self.estimate(i) for i in x])
54          return np.sum(self.y.values * self.l_x(x))
55
56      def plot_cv_score(self, low, high, num=100):
```

```
57          h = np.linspace(low, high, num)
58          score_f = self.cross_validation()
59          y = [score_f(i) for i in tqdm(h, desc=self.method + " cross validation score")]
60          plt.plot(h, y)
61          return np.array([h, y])
62
63      def get_kappa0(self):
64          self.kappa0 = integrate.quad(self._get_tx_norm, self.x.min(), self.x.max())[0]
65
66      def get_c(self, alpha=0.05):
67          self.get_kappa0()
68          self.c = fsolve(self._solve_c, np.array(1.96))[0]
69
70      def _solve_c(self, x):
71          return 2 * (1 - norm.cdf(x)) + self.kappa0 / np.pi * np.exp(-x ** 2 / 2) - 0.05
72
73      def _get_tx_norm(self, x):
74          return np.linalg.norm(self._get_t_derivative_x(x))
75
76      def _get_tx(self, x):
77          l_ix = self.l_x(x)
78          t_ix = l_ix / np.linalg.norm(l_ix)
79          return t_ix
80
81      def _get_t_derivative_x(self, x):
82          d = nd.Derivative(self._get_tx)
83          return d(x)
84
85      def cross_validation(self):
86          def score(h):
87              self.set_parameter(h)
88              r = np.sum(
89                  ((self.y.values - self.estimate(self.x.values)) / (1 - self.v / self.n)) ** 2
90              ) / self.n
91              self.set_optimal_parameter(h)
92              return r
93
94          return score
95
96      def variance_estimate(self) -> float:
97          sigma = np.sum((self.y - self.estimate(self.x.values)) ** 2) / (self.n - 2 * self.v + self.v_hat)
98          return sigma
99
100     def confidence_band(self):
101         self.get_c()
102         sigma = np.sqrt(self.variance_estimate())
103
104         def band(x):
105             s = self.c * np.sqrt(sigma) * np.linalg.norm(self.l_x(x))
106             return [self.estimate(x) - s, self.estimate(x) + s, self.estimate(x)]
107
108         return band
109
110     def draw_confidence_band(self, flag=True, ax=None):
111         if ax is None:
112             ax = plt
113         band = self.confidence_band()
114         if flag:
115             x = np.linspace(self.x.min(), self.x.max(), 500)
```

```
116            else:
117                x = self.x.sort_values()
118            y = np.array([band(i) for i in tqdm(x, desc=self.method + " confidence band")])
119            ax.plot(self.x, self.y, "o", color="grey", ms=2)
120            ax.plot(x, y[:, 0], color="red", label="lower")
121            ax.plot(x, y[:, 1], color="green", label="upper")
122            ax.plot(x, y[:, 2], color="blue", label="estimate")
123            ax.legend()
124            return np.array([x, y[:, 0], y[:, 1]])
125
126        @staticmethod
127        def _k(x):
128            """
129            Gaussian kernel
130            """
131            return np.exp(-x ** 2 / 2) / np.sqrt(2 * np.pi)
132
133
134  class Regressogram(Estimate):
135        def __init__(self, y: pd.Series, x: pd.Series):
136            super().__init__(y, x)
137            self.range = self.x.max() - self.x.min()
138            self.method = "Regressogram"
139
140        def get_effective_kernel(self):
141            if self.parameter is None:
142                raise ValueError("bins is not set")
143
144            def i_th_effective_kernel(i):
145                x_v = self.x.values
146                ek = np.zeros(self.n)
147                w = np.arange(self.x.min(), self.x.max(), self.range / self.parameter)
148                w = np.append(w, self.x.max())
149                for j in range(len(w) - 1):
150                    if w[j] <= self.x[i] <= w[j + 1]:
151                        if j == len(w) - 2:
152                            ek[(x_v >= w[j]) & (x_v <= w[j + 1])] = 1 / len(self.x[(self.x >= w[j]) & (self.x <= w[j + 1])])
153                        else:
154                            ek[(x_v >= w[j]) & (x_v < w[j + 1])] = 1 / len(self.x[(self.x >= w[j]) & (self.x < w[j + 1])])
155                return ek
156
157            self.effective_kernel = np.array([i_th_effective_kernel(i + 1) for i in range(self.n)]).T
158            self.v = np.trace(self.effective_kernel)
159            self.v_hat = np.trace(self.effective_kernel.T @ self.effective_kernel)
160
161        def estimate(self, x):
162            if isinstance(x, np.ndarray):
163                return np.array([self.estimate(i) for i in x])
164            w = np.arange(self.x.min(), self.x.max(), self.range / self.parameter)
165            w = np.append(w, self.x.max())
166            if x < w[1]:
167                return self.y[self.x <= w[1]].mean()
168            if x > w[-2]:
169                return self.y[self.x >= w[-2]].mean()
170            for j in range(len(w) - 1):
171                if w[j] <= x < w[j + 1]:
172                    return self.y[(self.x >= w[j]) & (self.x < w[j + 1])].mean()
173
174        def plot_cv_score(self, low=2, high=150, num=149):
```

```
175            h = np.linspace(low, high, num)
176            score_f = self.cross_validation()
177            y = [score_f(i) for i in h]
178            plt.plot(h, y)
179            return np.array([h, y])
180
181        def confidence_band(self):
182            self.get_c()
183            sigma = np.sqrt(self.variance_estimate())
184
185            def band(x):
186                w = np.arange(self.x.min(), self.x.max(), self.range / self.parameter)
187                w = np.append(w, self.x.max())
188                l_x = np.zeros(self.n)
189                l_x_norm = 0
190                if x < w[1]:
191                    l_x[self.x <= w[1]] = 1 / len(self.x[self.x <= w[1]])
192                    l_x_norm = np.linalg.norm(l_x)
193                if x > w[-2]:
194                    l_x[self.x >= w[-2]] = 1 / len(self.x[self.x >= w[-2]])
195                    l_x_norm = np.linalg.norm(l_x)
196                for j in range(len(w) - 1):
197                    if w[j] <= x < w[j + 1]:
198                        l_x[(self.x >= w[j]) & (self.x < w[j + 1])] = 1 / len(
199                            self.x[(self.x >= w[j]) & (self.x < w[j + 1])])
200                        l_x_norm = np.linalg.norm(l_x)
201                s = self.c * np.sqrt(sigma) * l_x_norm
202                return [self.estimate(x) - s, self.estimate(x) + s, self.estimate(x)]
203
204            return band
205
206        def get_kappa0(self):
207            self.kappa0 = 0
208
209        def get_c(self, alpha=0.05):
210            self.c = norm.ppf(1 - alpha / 2)
211
212
213    class Kernel(Estimate):
214        def __init__(self, y: pd.Series, x: pd.Series):
215            super().__init__(y, x)
216            self.method = "Kernel"
217
218        def l_x(self, x) -> np.ndarray:
219            return self._k((x - self.x) / self.parameter) / self._k((x - self.x) / self.parameter).sum()
220
221
222    class LocalLinear(Estimate):
223        def __init__(self, y: pd.Series, x: pd.Series):
224            super().__init__(y, x)
225            self.method = "Local Linear"
226
227        def l_x(self, x):
228            b_x = self._b_x(x)
229            s = np.sum(b_x)
230            l_ix = b_x / s
231            return l_ix
232
233        def _s_nj(self, x, j: {1, 2}):
```

```python
234            return np.sum(self._k((self.x - x) / self.parameter) * (self.x - x) ** j)
235
236    def _b_x(self, x):
237        return self._k((self.x - x) / self.parameter) * (self._s_nj(x, 2) - (self.x - x) * self._s_nj(x, 1))
238
239
240 class Spline(Estimate):
241    def __init__(self, y: pd.Series, x: pd.Series):
242        super().__init__(y, x)
243        self.data = None
244        self.process_same_value()
245        self.method = "Spline"
246        self._G()
247        self._Omega()
248        self.inv_np = None
249        self.beta = None
250
251    def process_same_value(self):
252        self.data = pd.DataFrame({"x": self.x, "y": self.y})
253        self.x = pd.Series(self.data["x"].value_counts().index.sort_values().values)
254        self.y = pd.Series([self.data[self.data["x"] == i]["y"].mean() for i in self.x])
255        self.n = len(self.x)
256
257    def set_parameter(self, para):
258        self.parameter = para
259        self._get_inv()
260        self.get_effective_kernel()
261
262    def l_x(self, x) -> np.ndarray:
263        return self._g_x(x)[0:self.n].T @ self.inv_np @ self.G.T
264
265    def _get_inv(self):
266        A = self.G.T @ self.G + self.parameter * self.Omega
267        self.inv_np = np.linalg.inv(A)
268        self.beta = self.inv_np @ self.G.T @ self.y.values.T
269
270    def _g_x(self, x):
271        g = np.zeros(self.n + 2)
272        d_K_1 = self._g_x_i(x, self.x.values[-2])
273        g[0] = 1
274        g[1] = x
275        g[2:] = np.array(
276            [
277                self._g_x_i(x, i) - d_K_1 if i < x else 0 for i in self.x
278            ])
279        return g
280
281    def _g_x_i(self, x, i):
282        if x < i:
283            temp = 0
284        elif x < self.x.values[-1]:
285            temp = ((x - i) ** 3) / (self.x.values[-1] - i)
286        elif x >= self.x.values[-1]:
287            temp = ((x - i) ** 3 - (x - self.x.values[-1]) ** 3) / (self.x.values[-1] - i)
288        return temp
289
290    def d2_g_x(self, x):
291        g = np.zeros(self.n + 2)
292        d2d_K_1 = self._d2_g_x_i(x, self.x.values[-2])
```

```
293            g[0] = 0
294            g[1] = 0
295            g[2:] = np.array(
296                [
297                    self._d2_g_x_i(x, i) - d2d_K_1 if i < x else 0 for i in self.x
298                ])
299            return g[0:self.n]
300
301        def _d2_g_x_i(self, x, i):
302            if x < i:
303                temp = 0
304            elif x < self.x.values[-1]:
305                temp = (6 * (x - i)) / (self.x.values[-1] - i)
306            elif x >= self.x.values[-1]:
307                temp = (6 * ((x - i) - (x - self.x.values[-1]))) / (self.x.values[-1] - i)
308            return temp
309
310        def _G(self):
311            self.G = np.array([self._g_x(i) for i in self.x])[0:self.n:, 0:self.n]
312
313        def _Omega(self):
314            self.Omega = np.array([self._omega_i(i) for i in tqdm(range(self.n), desc=self.method + " Omega")])
315
316        def _omega_i(self, i):
317            return integrate.quad_vec(lambda x: self.d2_g_x(x) * self.d2_g_x(x)[i], self.x.min(), self.x.max())[0]
318
319        def cross_validation(self):
320            def score(h):
321                self.set_parameter(h)
322                mse = np.sum(
323                    (self.y.values - self.estimate(self.x.values)) ** 2)
324                penalty = self.parameter * \
325                        integrate.quad(lambda x: np.sum(self.beta * self.d2_g_x(x)), self.x.min(), self.x.max())[0]
326                return mse + penalty
327            return score
```

```
1    df = pd.read_csv("https://www.stat.cmu.edu/~larry/all-of-nonpar/=data/glass.dat", sep="\s+", header=0)
2    target = df["RI"]
3    x_al = df["Al"]
4
5    reg = Regressogram(target, x_al)
6    bin_cv = reg.plot_cv_score()
7    plt.savefig("./fig/Regressogram_cv_score.pdf")
8    plt.clf()
9
10   reg.set_optimal_parameter(int(bin_cv[0, np.argmin(bin_cv[1])]))
11
12   bin_cb = reg.draw_confidence_band(flag=False)
13   plt.savefig("./fig/Regressogram_confidence_band.pdf")
14   plt.clf()
15
16   print(f"optimal smoothing parameter: {np.argmin(bin_cv[1]) + 2}")
17   print(f"Regressogram variance estimation: {reg.variance_estimate()}")
18   print(f"Regressogram c: {reg.c}")
19
20   ker = Kernel(target, x_al)
21   ker_cv = ker.plot_cv_score(0.01, 0.1)
```

```python
22    plt.savefig("./fig/Kernel_cv_score.pdf")
23    plt.clf()
24
25    ker.set_optimal_parameter(ker_cv[0, np.argmin(ker_cv[1])])
26
27    ker_cb = ker.draw_confidence_band()
28    plt.savefig("./fig/Kernel_confidence_band.pdf")
29    plt.clf()
30
31    print(f"optimal smoothing parameter: {ker_cv[0, np.argmin(ker_cv[1])]}")
32    print(f"Kernel variance estimation: {ker.variance_estimate()}")
33    print(f"Kernel c: {ker.c}")
34
35    print(f"Kernel c: {ker.c}")
36
37    ll = LocalLinear(target, x_al)
38    ll_cv = ll.plot_cv_score(0.015, 0.1)
39    plt.savefig("./fig/LocalLinear_cv_score.pdf")
40    plt.clf()
41
42    ll.set_optimal_parameter(ll_cv[0, np.argmin(ll_cv[1])])
43
44    ll_cb = ll.draw_confidence_band()
45    plt.savefig("./fig/LocalLinear_confidence_band.pdf")
46    plt.clf()
47
48    print(f"optimal smoothing parameter: {ll_cv[0, np.argmin(ll_cv[1])]}")
49    print(f"LocalLinear variance estimation: {ll.variance_estimate()}")
50    print(f"LocalLinear c: {ll.c}")
51
52    spline = Spline(target, x_al)
53    spline_cv = spline.plot_cv_score(1e-06, 1e-05, 200)
54    plt.savefig("./fig/Spline_cv_score.pdf")
55    plt.clf()
56
57    spline.set_optimal_parameter(spline_cv[0, np.argmin(spline_cv[1])])
58
59    spline_cb = spline.draw_confidence_band()
60    plt.savefig("./fig/Spline_confidence_band.pdf")
61    plt.clf()
62
63    print(f"optimal smoothing parameter: {spline_cv[0, np.argmin(spline_cv[1])]}")
64    print(f"Spline variance estimation: {spline.variance_estimate()}")
65    print(f"Spline c: {spline.c}")
```

```
optimal smoothing parameter: 30
Regressogram variance estimation: 5.5430334733388245
Regressogram c: 1.959963984540054

optimal smoothing parameter: 0.0441919191919192
LocalLinear variance estimation: 5.768970862437912
LocalLinear c: 3.3862804226633734

optimal smoothing parameter: 2.175879396984925e-06
Spline variance estimation: 4.981818961988229
Spline c: 3.578015201230727
```

```
1   fig, ax = plt.subplots(2, 2, figsize=(16, 9))
2   ax[0, 0].set_title('Regressogram')
3   ax[0, 1].set_title('Kernel')
4   ax[1, 0].set_title('LocalLinear')
5   ax[1, 1].set_title('Spline')
6   ax[0, 0].set_xlabel(r'$m$')
7   ax[0, 1].set_xlabel(r'$h$')
8   ax[1, 0].set_xlabel(r'$h$')
9   ax[1, 1].set_xlabel(r'$\lambda$')
10  ax[0, 0].set_ylabel('GCV')
11  ax[0, 1].set_ylabel('GCV')
12  ax[1, 0].set_ylabel('GCV')
13  ax[1, 1].set_ylabel('Penalized Sums of Squares')
14  ax[0, 0].plot(bin_cv[0], bin_cv[1])
15  ax[0, 1].plot(ker_cv[0], ker_cv[1])
16  ax[1, 0].plot(ll_cv[0], ll_cv[1])
17  ax[1, 1].plot(spline_cv[0], spline_cv[1])
18  plt.savefig('./fig/cv_score.pdf')
```

```
1   fig, ax = plt.subplots(2, 2, figsize=(16, 9), sharey='all')
2   ax[0, 0].set_ylim(-20, 20)
3   ax[0, 0].set_title('Regressogram')
4   ax[0, 1].set_title('Kernel')
5   ax[1, 0].set_title('LocalLinear')
6   ax[1, 1].set_title('Spline')
7   ax[0, 0].set_xlabel(r'$x$')
8   ax[0, 1].set_xlabel(r'$x$')
9   ax[1, 0].set_xlabel(r'$x$')
10  ax[1, 1].set_xlabel(r'$x$')
11  ax[0, 0].set_ylabel(r'$y$')
12  ax[0, 1].set_ylabel(r'$y$')
13  ax[1, 0].set_ylabel(r'$y$')
14  ax[1, 1].set_ylabel(r'$y$')
15  ax[0, 0].scatter(reg.x, reg.y, s=2, color='grey')
16  ax[0, 1].scatter(ker.x, ker.y, s=2, color='grey')
17  ax[1, 0].scatter(ll.x, ll.y, s=2, color='grey')
18  ax[1, 1].scatter(spline.x, spline.y, s=2, color='grey')
19  ax[0, 0].fill_between(bin_cb[0], bin_cb[1], bin_cb[2], alpha=.5, color='red')
20  ax[0, 0].plot(bin_cb[0], reg.estimate(bin_cb[0]), color='red')
21  ax[0, 1].fill_between(ker_cb[0], ker_cb[1], ker_cb[2], alpha=.5, color='red')
22  ax[0, 1].plot(ker_cb[0], ker.estimate(ker_cb[0]), color='red')
23  ax[1, 0].fill_between(ll_cb[0], ll_cb[1], ll_cb[2], alpha=.5, color='red')
24  ax[1, 0].plot(ll_cb[0], ll.estimate(ll_cb[0]), color='red')
25  ax[1, 1].fill_between(spline_cb[0], spline_cb[1], spline_cb[2], alpha=.5, color='red')
26  ax[1, 1].plot(spline_cb[0], spline.estimate(spline_cb[0]), color='red')
27  plt.savefig('./fig/confidence_band.pdf')
```

```
1   fig, ax = plt.subplots()
2   ax.set_ylim(-10, 15)
3   ax.scatter(x_al, target, s=2, color='grey')
4   ax.plot(bin_cb[0], reg.estimate(bin_cb[0]), color='red', alpha=0.6, label='Regressogram')
5   ax.plot(ker_cb[0], ker.estimate(ker_cb[0]), color='blue', alpha=0.6, label='Kernel')
6   ax.plot(ll_cb[0], ll.estimate(ll_cb[0]), color='green', alpha=0.6, label='LocalLinear')
7   ax.plot(spline_cb[0], spline.estimate(spline_cb[0]), color='orange', alpha=0.6, label='Spline')
```

```
8    ax.set_xlabel(r'$x$')
9    ax.set_ylabel(r'$y$')
10   ax.legend()
11   plt.savefig('./fig/compare.pdf')
```