

# LSTAT2150

## Non-parametric statistics

### Smoothing methods

Project number 7

December 7, 2020

LAMY Lionel  
1294-1700

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Quality assessment</b>	<b>2</b>
2.1	Bias . . . . .	2
2.2	Variance . . . . .	2
2.3	Mean squared error . . . . .	2
2.4	Mean integrated squared error . . . . .	2
2.5	Mean sum of squared error . . . . .	2
<b>3</b>	<b>Regression estimators</b>	<b>3</b>
3.1	Nadaya-Watson (nonparametric) . . . . .	3
3.2	Polynomial (parametric) . . . . .	5
<b>4</b>	<b>Investigation</b>	<b>6</b>
<b>5</b>	<b>Conclusion</b>	<b>8</b>
<b>A</b>	<b>Appendix</b>	<b>9</b>
A.1	Figures . . . . .	9
A.2	Code . . . . .	10

# 1 Introduction

Let  $\{(X_i, Y_i)\}_{i=1}^n$ ,  $n = 100$ , follow the regression model :

$$Y_i = m(X_i) + 0.5 \varepsilon_i, \quad i = 1, \dots, n$$

where  $X_i$  i.i.d  $\sim \text{Uniform}[0,1]$ ,  $\{\varepsilon_i\}$  i.i.d  $\sim \mathcal{N}(0, 1)$  and  $m(x) = (\sin(2 \pi x^3))^3$

The objective of this project is to compare the performance between a parametric polynomial regression estimator and a nonparametric Nadayara-Watson regression estimator<sup>1</sup>. To do so, we will first try to determine the optimal polynomial order and bandwidth by varying over a grid of possible values and selecting those values that minimize the MISE (or the MSSE) with respect to the true function  $m(x)$ . Then, we will inquire about the bias and the variance. Finally, we will investigate the MSE at locations we feel are relevant in order to conclude from the performance of our regressions.

## 2 Quality assessment

As announced in the introduction above, we will use estimators to judge and compare the quality of our models. As such, it seems important to us to define them beforehand.

### 2.1 Bias

The bias measures the systematic deviation of the estimator from the true density. If the estimator tends to zero when the number of observations tends towards infinity, then it is considered asymptotically unbiased.

$$\text{Bias}[\hat{m}(x)] = \text{E}[\hat{m}(x)] - m(x)$$

### 2.2 Variance

The variance represents the mean of the squares of the deviations from the true mean and is used to quantify the dispersion of the estimates.

$$\text{Var}[\hat{m}(x)] = \text{E}[(\hat{m}(x) - \text{E}[\hat{m}(x)])^2]$$

### 2.3 Mean squared error

The MSE is the average squared difference between the estimated values and the actual value which means that it characterizes the accuracy of an estimator.

$$\text{MSE}[\hat{m}(x)] = \text{E}[(\hat{m}(x) - m(x))^2]$$

### 2.4 Mean integrated squared error

We will mainly try to use this measurement to define the quality and accuracy of our models. It is the integration of the MSE on all possible values of  $x$  (in our case 0 to 1 as  $X_i$  follows an Uniform[0,1]). It can also be computed as the mean of integrated squared errors (ISE).

$$\begin{aligned} \text{MISE}[\hat{m}(\cdot)] &= \int \text{MSE}[\hat{m}(x)] dx \\ &= \int \text{E}[(\hat{m}(x) - m(x))^2] dx = \text{E} \int (\hat{m}(x) - m(x))^2 dx \end{aligned}$$

### 2.5 Mean sum of squared error

As an integral can be approximated by  $\frac{1}{n} \sum_{i=1}^n$ , we have a simpler equivalence of the MISE which saves us from having to solve an integral, which can sometimes be tricky.

$$\text{MSSE}[\hat{m}(\cdot)] = \frac{1}{n} \sum_{i=1}^n \text{MSE}[\hat{m}(x_i)]$$

---

<sup>1</sup>Called respectively 'LM' and 'NW' in this paper

### 3 Regression estimators

#### 3.1 Nadayara-Watson (nonparametric)

##### 3.1.1 Description

This estimator, published in 1964 by Nadaraya and Watson, proposes to consider the function  $m(x)$  as a locally weighted average of  $Y_i$  given observations  $X_i$  and try to estimate it using a kernel as a weighting function. The motivation behind it comes from :

$$m(x) = E[Y|X = x] = \int y \frac{f(x, y)}{f(x)} dy$$

Where the numerator  $f(x, y)$  and denominator  $f(x)$  can be estimated via a kernel density estimator :

$$f(x, y) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right) Y_i$$

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

Combining all in one, we therefore have the NW estimator, with  $h$  as the bandwidth of the kernel  $K$ :

$$\hat{m}_{NW}(x) = \frac{\sum_{i=1}^n K\left(\frac{X_i - x}{h}\right) Y_i}{\sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)}$$

The Kernel used in this work is the Gaussian, defined as  $K_G = \exp(-\frac{1}{2}u^2)(\sqrt{2\pi})^{-1}$  or simply **dnorm**( $u$ ) in R. Notice that the equations above show that the bandwidth is what affects the estimator the most, so the choice of a Kernel isn't that much of a matter.

##### 3.1.2 Optimal bandwidth selection

In order to find the best bandwidth, we code a R function that compute  $M$  regressions for each value of a sequence of possible bandwidths and returns the one where the average of the  $M$  minimize the MISE/MSSE. First, we tried bandwidths ranging from 0.01 to 0.50 in increments of 0.01 and with  $M=500$ . We obtained  $h = 0.04$  as result. For the sake of accuracy, we started the experiment again with possible values ranging from 0.03 to 0.05 per 0.001 increase. We thus found that 0.036 seems to be the best bandwidth.<sup>2</sup>

On the graph on the left, we drew the estimate of a single NW regression, while on the right, we drew the mean of  $M=500$  regressions. We notice that by bootstrapping, the regression fits slightly better to the true function but that the largest deviations still are at changes in convexity (second derivative).

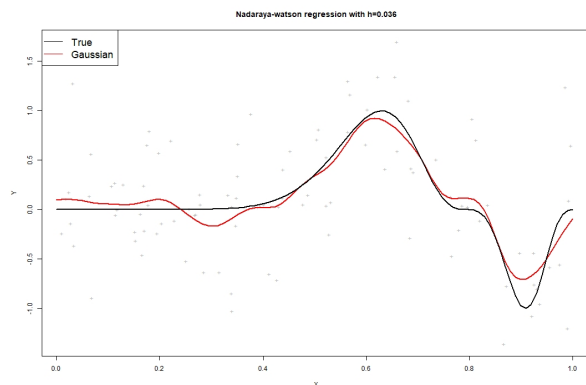


Figure 1: Fit of the NW estimator with the optimal  $h$  and  $n=100$

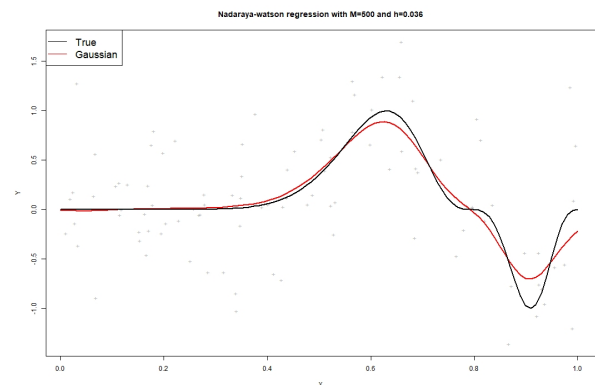


Figure 2: Fit of the NW estimator with the optimal  $h$ ,  $M=500$  and  $n=100$

<sup>2</sup>An animated example of the bandwidth evolution is [available here](#)

We also wanted to know what would be the best  $h$  with a smaller sample size, and to verify in the process that the MSSE is a good approximation of the MISE. This appears to be confirmed in the table below. The left part represents the value of the quality estimators while the right part corresponds to the associated  $h$  values. We observe that as the sample size increases, the estimators and bandwidths appear to decrease.

n	Minimized		Bandwidth	
	MISE	MSSE	MISE	MSSE
25	0.6493	0.6411	0.145	0.145
50	0.1073	0.1122	0.127	0.124
100	0.0314	0.0317	0.036	0.036

Table 1: MISE/MSSE and the optimal bandwidth by  $n$

### 3.1.3 Bias and variance

We will now focus on the bias and variance of the regression with the optimal bandwidth ( $h = 0.036$ ) for  $n=100$ . For this, we have kept the  $M$  predictions made in the above section. Thus, we can easily obtain the bias by subtracting the true value from the mean of the regressions for each of the points. The same is done for the variance, where we regressed once again and subtracted the mean of the regressions, all squared. See equations in sections 2.1 and 2.2 for a better understanding.

Let's first analyze the bias. Graphically, we notice that the bias is very important on the right side, corresponding indeed to the worst-estimated points in the Figure 2. We were expecting this result because theoretically, we know that the bias of a Kernel estimator  $\approx \frac{h^2}{2} m''(x) \int u^2 K(u) du$  and so is proportional to the second derivative of the density. Moreover, we are pleased to state that the bias only takes 0.2989 as maximum value, with a mean value of -0.0023 and median of 0.0098 and a absolute sum of 5.8867.

An additional plot with both regression and bias is available in appendix (Figure 14).

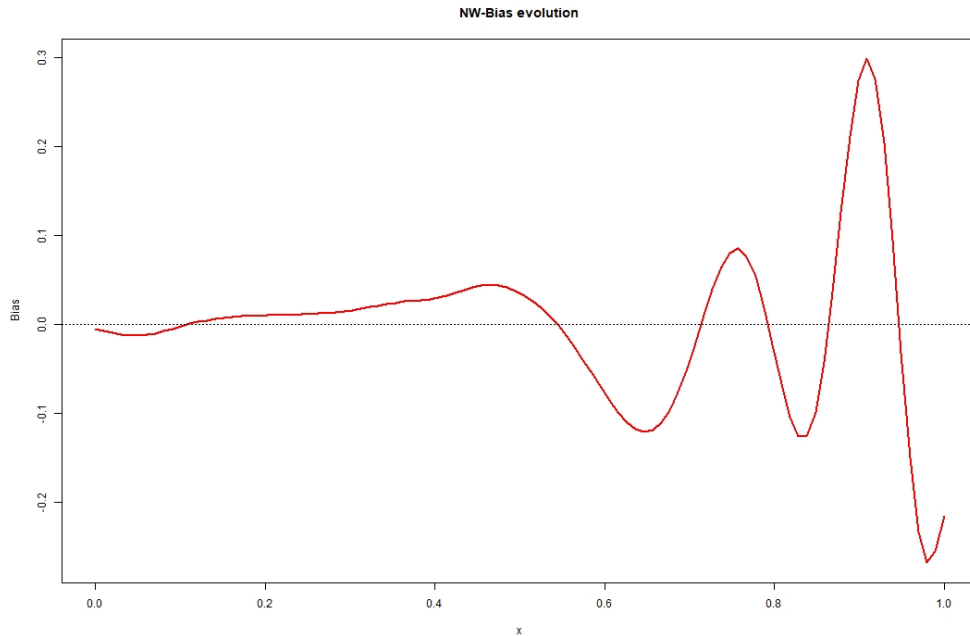


Figure 3: Bias evolution of the NW best regression

By generating a new regression on the same sample, we obtained a variance of 0.0038. In addition, by experimenting with different samples, we obtained a maximum variance of about 0.0450, which we believe to be a very good value.

## 3.2 Polynomial (parametric)

### 3.2.1 Optimal polynomial degree

We will use the following R command : `lm(Y ~ poly(X, degree))`. The latter allows us to avoid using a long formula with multiple powers, thanks to the `poly` function that returns a matrix of orthogonal polynomials. Our goal right now is to find out what is the best degree to use to fit the real  $m(x)$  function.

Unlike non-parametric regression, we were unable to find a way to calculate the MISE for linear regression and therefore sought to minimize the MSSE. Although this does not affect the results as we saw earlier, we wanted to make the point. So, we used the same strategy as for the Nadayara-Watson and generated  $M=500$  polynomial regressions, for each degree ranging from 1 to 20 and retained the one with the smallest MSSE.

MSSE	Bias	Degree
0.0597	9.0993	7
0.0784	5.4018	19

Table 2: MSSE, bias and degree of polynomial regression

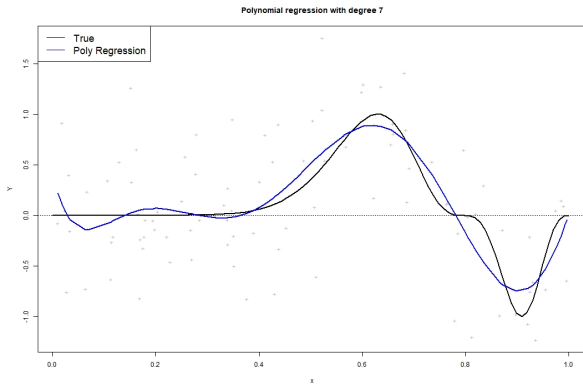


Figure 4: Fit with degree 7

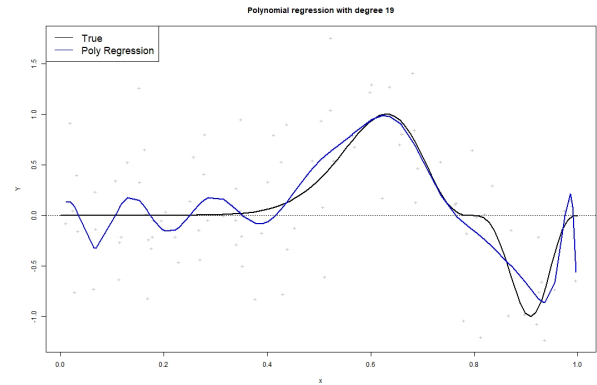


Figure 5: Fit with degree 19

However, by trying to minimize the bias rather than the MSSE, we noticed that the degree returned is higher and as you will see below in Figure 7, the mean prediction of  $M=500$  simulations is a bit more accurate with degree 19 than degree 7 (Figure 6) despite a lower MSSE. This might be explained by the fact that with such a high degree, the regression is highly dependent on the sample (the distribution of  $X$ ) and is therefore highly variable (overfit), leading to an increase of the MSSE. We observe this effect in figures 6 and 7, where the dotted lines are the minimum and maximum values taken by the  $M$  predictions (like a confidence interval). With the lower degree, we don't have the best results but each individual prediction is more stable and less far from the real function.

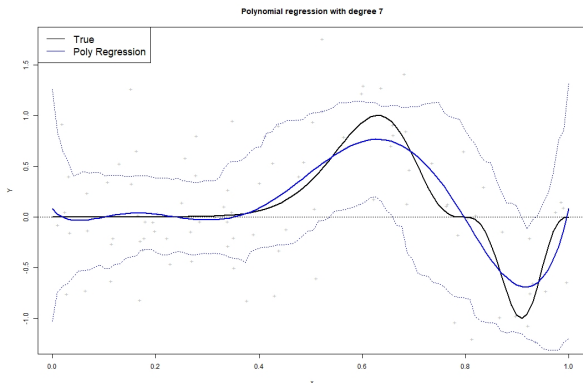


Figure 6: Mean fit with degree 7,  $M=500$

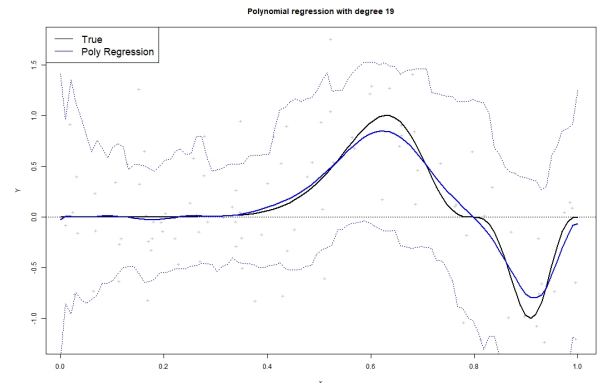


Figure 7: Mean fit with degree 19,  $M=500$

By precaution, and in order to follow the instructions given in the introduction, we have decided to keep `degree = 7` for our next analyses.

### 3.2.2 Bias and variance

We observe the same limitations and weaknesses as with the Nadayara-Watson estimator yet with a slightly larger amplitude. The maximum bias is 0.3124, while the mean and the median are respectively 0.0004 and 0.0019. The sum, however is 9.1032, which corresponds to an increase of 155% compared to the bias obtained with NW which was 5.8867. The variance is 0.0184.

An additional plot with both regression and bias is also available in appendix (Figure 15)

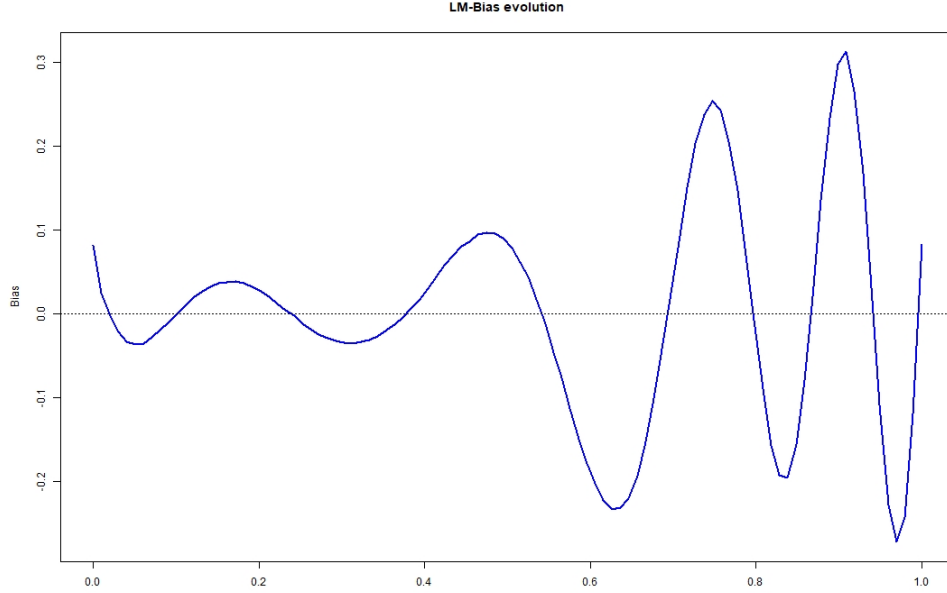


Figure 8: Bias evolution of the LM best regression

## 4 Investigation

In this section, we are going to investigate numerically with 5000 Monte-Carlo simulations the MSE of some points. The points that we find interesting to explore are  $[0, 0.2, 0.45, 0.63, 0.78, 0.91, 1]$  and are represented on the Figure 9 below by a red dashed line. We have chosen the two points at the extremities, the point between the maximum and the minimum, as well as these last two, and an additional point near the center.

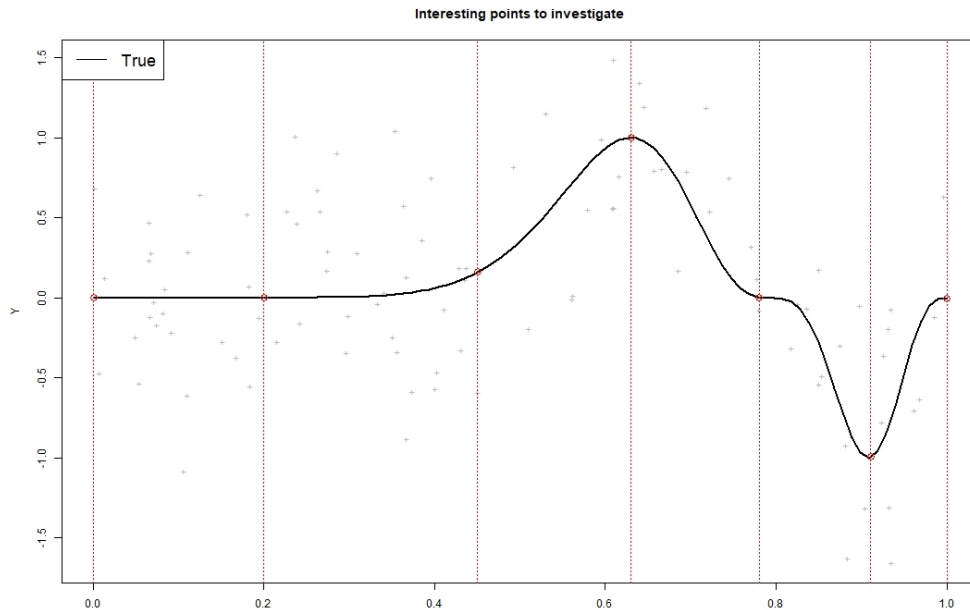


Figure 9: Interesting points marked by a red line and a circle

By calculating the MSE for each point of interest for  $n = 25, 50, 100$  and  $1000$  we notice several things: i) the mean MSE of all points decreased drastically between  $n=25$  and  $n=1000$  for both polynomial and Nadaraya-Watson regression. ii) the two points on the right are always with the highest mean squared error. iii) the Nadaraya-Watson estimator seems to have systematically less error, with the exception of point 0.2.

	Nadaraya-Watson						Polynomial				
Points	n=25	n=50	n=100	n=1000	Mean		n=25	n=50	n=100	n=1000	Mean
<b>0</b>	0.1441	0.0889	0.0424	0.0040	0.0699		0.1941	0.1643	0.1191	0.0270	0.1261
<b>0.2</b>	0.0909	0.0440	0.0212	0.0020	0.0395		0.0713	0.0338	0.0179	0.0032	0.0316
<b>0.45</b>	0.0992	0.0450	0.0220	0.0026	0.0422		0.1261	0.0769	0.0513	0.0100	0.0662
<b>0.63</b>	0.1120	0.0572	0.0308	0.0111	0.0528		0.2339	0.1483	0.0819	0.0327	0.1392
<b>0.78</b>	0.0996	0.0458	0.0235	0.0034	0.0431		0.3828	0.2292	0.1406	0.0224	0.1938
<b>0.91</b>	0.2459	0.1562	0.1146	0.0864	0.1508		0.7049	0.2811	0.1394	0.0641	0.2974
<b>1</b>	0.3630	0.1951	0.1033	0.0454	0.1767		0.4177	0.2481	0.1435	0.0931	0.2256
<b>Mean</b>	0.1650	0.0903	0.0511	0.0221	0.0821		0.3130	0.1688	0.0991	0.0362	0.1543

Table 3: MSE evolution and comparison at specifics points

Graphically, we see more easily the very fast decrease in error at the beginning and we also notice that an asymptotic trend is emerging. The mean squared error progressively decreases. To remedy this and get even closer to zero, the optimal bandwidth and degree values would have to be recalculated with a sample size larger than  $n=100$ .

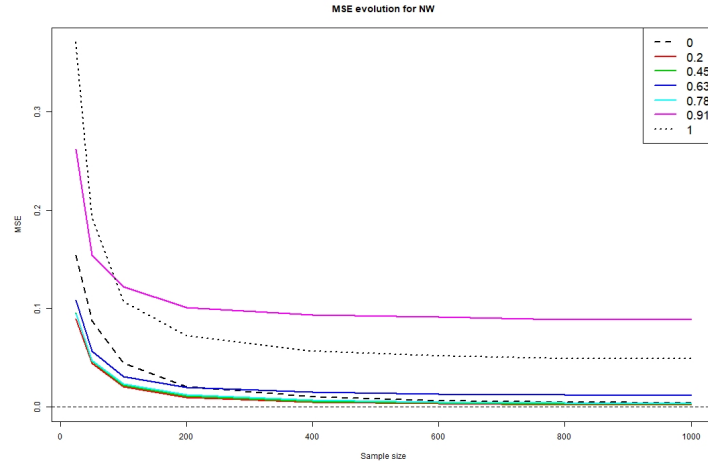


Figure 10: NW-MSE evolution

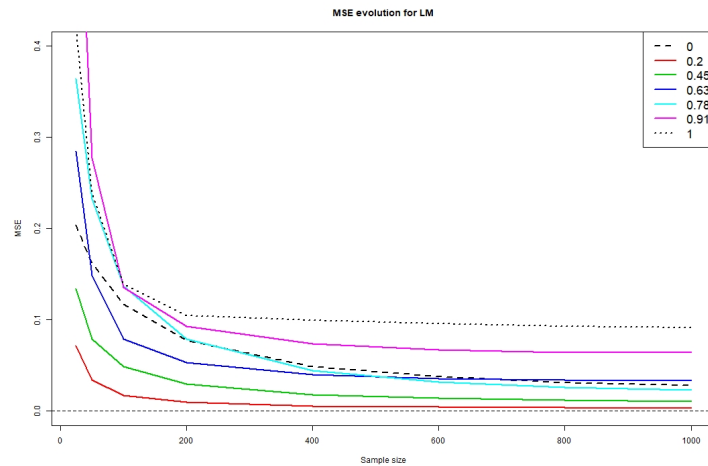


Figure 11: LM-MSE evolution

## 5 Conclusion

To conclude this work and this comparison, we will end with two final graphs with the two estimators. In our case the **Nadaya-Watson estimator did slightly better to estimate the true function  $m(x)$** . Nevertheless, we believe that since the true function  $m(x)$  is relatively simple, with a higher  $n$  and a higher degree, the parametric regression could be better. In addition, a perhaps more interesting approach to explore would be the use of spline, via a GAM model for example. However, we consider ourselves satisfied and leave this possibility open for future work.

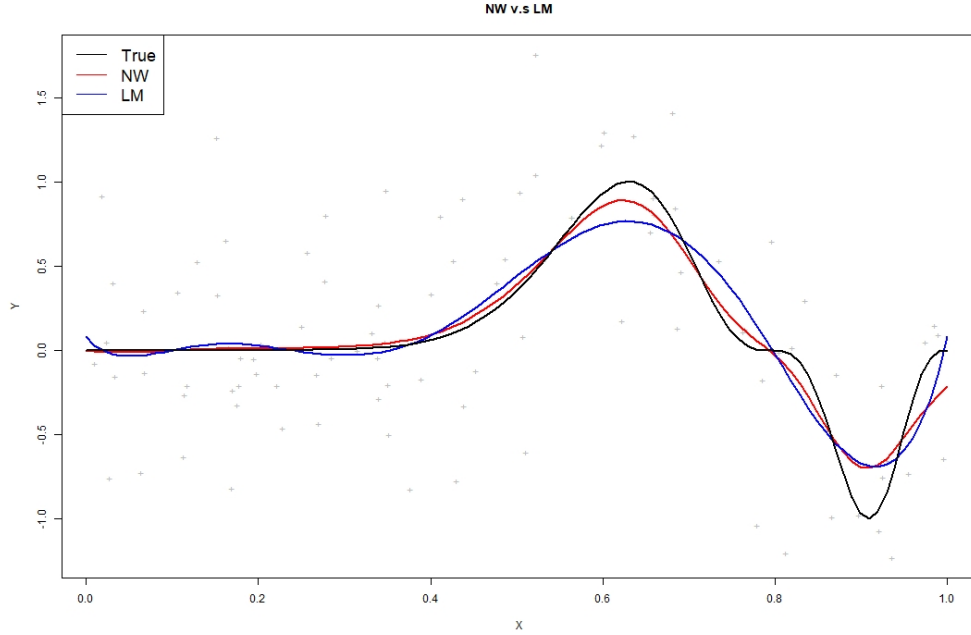


Figure 12: Comparison of the two estimators

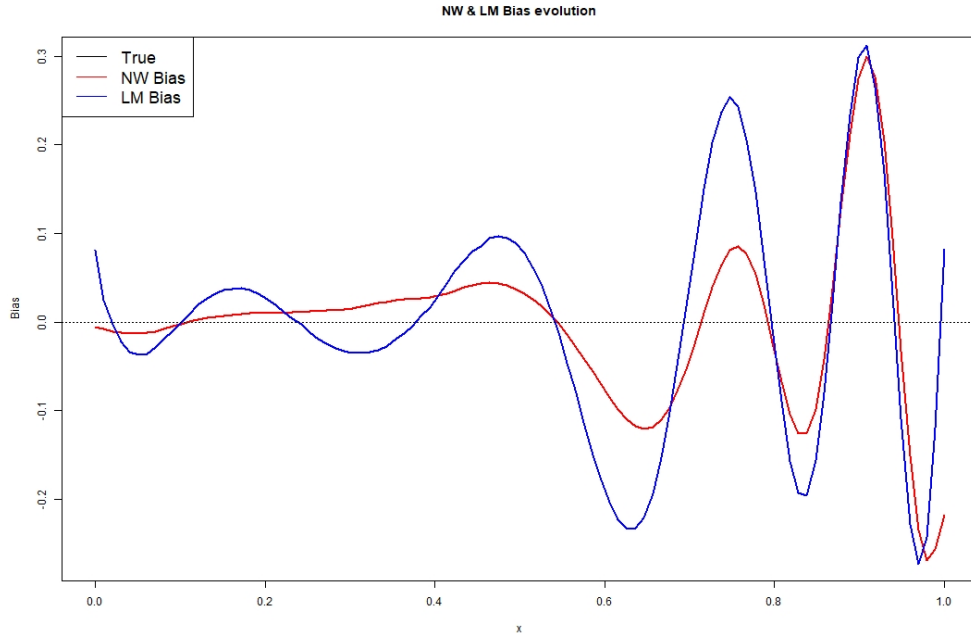


Figure 13: NW and LM bias together



# A Appendix

**Info:** All the content (including full-sized figures) and the code can be found on github via <https://github.com/lamylio/LSTAT2150-Project>

## A.1 Figures

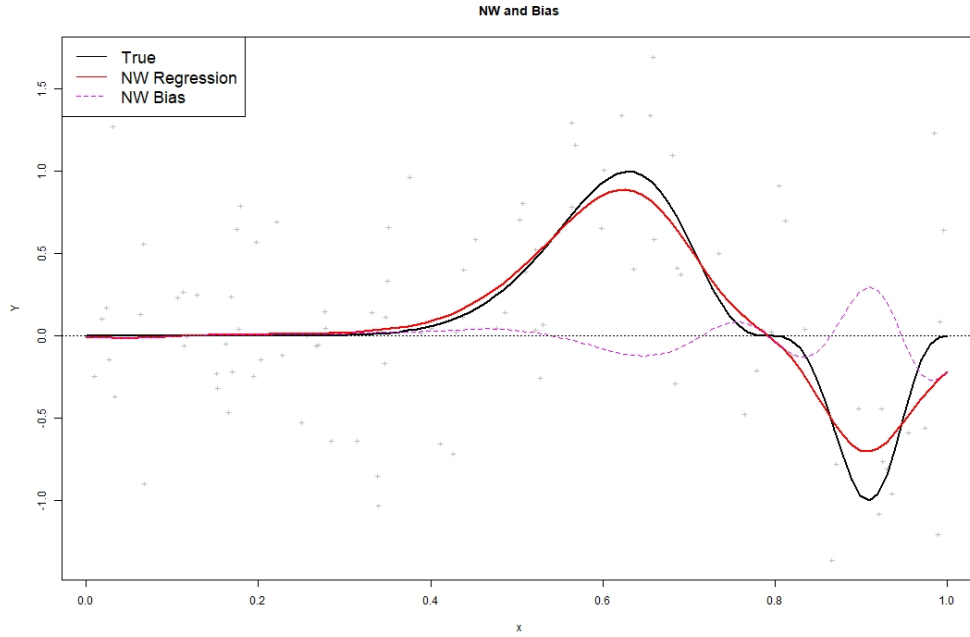


Figure 14: Bias with the NW best regression in same time

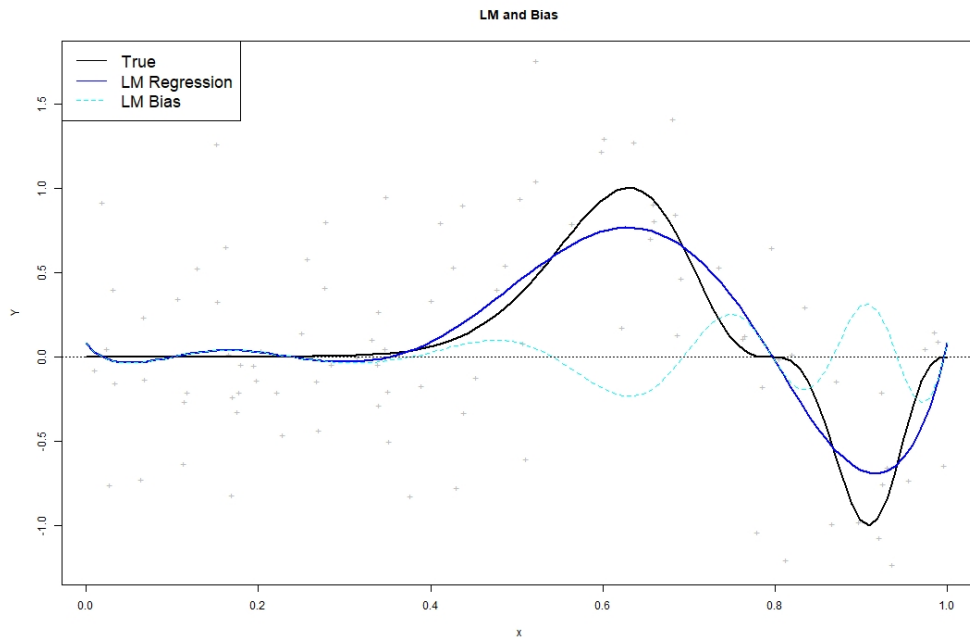


Figure 15: Bias with the LM best regression in same time

## A.2 Code

### A.2.1 Main

```
1 # =====
2 # By: LAMY Lionel
3 # Noma: 1294-17-00
4 # https://github.com/lamylio/LSTAT2150-Project
5 # =====
6
7
8 # =====
9 # Please see Setup
10 # =====
11
12 set.seed(122020)
13 source("./sources/setup.r")
14 source("./sources/drawing.r")
15 source("./sources/minimization.r")
16 source("./sources/investigate.R")
17
18 # Global X and Y
19 CP.X = sort(X(n))
20 CP.Y = Y(CP.X)
21 CP.x = seq(0, 1, length = n)
22
23 # Kernel
24 Knorm = function(u) dnorm(u) #Gaussian kernel
25
26 # NW Estimator
27 NW.regEst <- function(x, X, Y, h, K) sum(Y * K((x - X)/h))/sum(K((x - X)/h))
28
29
30
31 # =====
32 # Please see Minimize
33 # =====
34
35 # ---
36 # Warning: might take a while
37 # Use ni=100, M=100 by default but better results with higher values
38 # ---
39
40 # NW.minimizeMISE(h.test = seq(0.03, 0.05, 0.001), M=500)
41
42 NW.optimal = NW.minimizeMSSE(h.test = seq(0.03, 0.05, 0.001), M=500)
43 LM.optimal = LM.minimizeMSSE(M=500)
44
45 # Optimal bandwidth
46 NW.h = NW.optimal$h[1] # 0.036
47 # Optimal degree
48 LM.poly = LM.optimal$poly[1] # 7
49
50 # Models (for futur use)
51 NW.optimal.model = sapply(seq(0,1,length=nrow(NW.optimal)), function(x) NW.regEst(x, CP.X, CP.Y))
52 LM.optimal.model = lm(CP.Y ~ poly(CP.X, LM.poly))$fitted
53
54
55 # As X.optimal already contains the bias and the MSSE
56 # we only need to compute the variance, via MReg
57
58 NW.optimal$variance = round(
59   mean((NW.optimal.model - NW.optimal$MReg)^2)
60 ,4)
```

```

61
62 LM.optimal$variance = round(
63   mean((LM.optimal.model - LM.optimal$MReg)^2)
64 ,4)
65
66 # head(NW.optimal, 1)
67 # head(LM.optimal, 1)
68
69
70
71 # =====
72 # Please see Investigate
73 # =====
74
75 CP.to_check = c(0,0.2,0.45,0.63,0.78,0.91,1)
76 CP.investigation.samples = c(25,50,100,200,400,600,800,1000) # smoother lines
77
78 # Draw which points
79 plot(CP.X, CP.Y, pch="+", col = "grey", xlab = "x", ylab="Y")
80 lines(CP.x, m(CP.x), col = 1, lwd=2)
81 abline(v=CP.to_check, col=2, lty=3)
82 title("Interesting points to investigate")
83 legend("topleft",
84       legend = c("True"),
85       col = c(1),
86       lty = c(1),
87       cex = 1.5
88 )
89
90 # Compute all the MSE
91 NW.investigation.results = CP.investigation.toframe(NW.investigation)
92 LM.investigation.results = CP.investigation.toframe(LM.investigation)
93
94
95
96 # =====
97 # Please see Draw
98 # =====
99
100 # ---
101 # Draw the evolution of MSE
102 CP.drawMSE(NW.investigation.results, "MSE evolution for NW", save=T)
103 CP.drawMSE(LM.investigation.results, "MSE evolution for LM", save=T)
104
105 # ---
106 # Plot both LM and NW best regression
107 CP.drawComparison(NW.optimal$MReg, LM.optimal$MReg, 100, save=T)
108
109 # ==
110
111 # ---
112 # Plot the bias alone
113 jpeg("./plots/bias/NW-Alone.jpg", quality = 100, width = 1080, height = 720)
114 plot(CP.x, NW.optimal$Bias, col=2, lwd=2, lty=1, type="l", xlab = "x", ylab="Bias", main="NW-B")
115 abline(h=0, col=1, lty=3)
116 dev.off()
117
118 jpeg("./plots/bias/LM-Alone.jpg", quality = 100, width = 1080, height = 720)
119 plot(CP.x, LM.optimal$Bias, col=4, lwd=2, lty=1, type="l", xlab = "x", ylab="Bias", main="LM-B")
120 abline(h=0, col=1, lty=3)
121 dev.off()
122
123 # Both in the same time

```

```

124 jpeg("./plots/bias/NW-LM-Alone.jpg", quality = 100, width = 1080, height = 720)
125 plot(CP.x, NW.optimal$Bias, col=2, lwd=2, lty=1,type="l", xlab = "x", ylab="Bias", main="NW &
126 lines(CP.x, LM.optimal$Bias, col=4, lwd=2, lty=1)
127 abline(h=0, col=1, lty=3)
128 legend("topleft",
129     legend = c("True", "NW Bias", "LM Bias"),
130     col = c(1, 2, 4),
131     lty = c(1, 1, 1),
132     cex = 1.5
133 )
134 dev.off()
135
136 # ---
137 # Plot bias and regression
138 jpeg("./plots/bias/NW-Both.jpg", quality = 100, width = 1080, height = 720)
139 plot(CP.X, CP.Y, pch="+", col = "grey", xlab = "x", ylab="Y", main="NW and Bias")
140 lines(CP.x, m(CP.x), col = 1, lwd=2)
141 lines(CP.x, NW.optimal$MReg, col=2, lty=1, lwd=2)
142 lines(CP.x, NW.optimal$Bias, col=6, lty=2)
143 abline(h=0, col=1, lty=3)
144 legend("topleft",
145     legend = c("True", "NW Regression", "NW Bias"),
146     col = c(1, 2, 6),
147     lty = c(1, 1, 2),
148     cex = 1.5
149 )
150 dev.off()
151
152 jpeg("./plots/bias/LM-Both.jpg", quality = 100, width = 1080, height = 720)
153 plot(CP.X, CP.Y, pch="+", col = "grey", xlab = "x", ylab="Y", main="LM and Bias")
154 lines(CP.x, m(CP.x), col = 1, lwd=2)
155 lines(CP.x, LM.optimal$MReg, col=4, lty=1, lwd=2)
156 lines(CP.x, LM.optimal$Bias, col=5, lty=2)
157 abline(h=0, col=1, lty=3)
158 legend("topleft",
159     legend = c("True", "LM Regression", "LM Bias"),
160     col = c(1, 4, 5),
161     lty = c(1, 1, 2),
162     cex = 1.5
163 )
164 dev.off()
165
166
167 # =====
168 # END
169 # =====

```

## A.2.2 Setup

```

1 # =====
2 # Given in the project statement
3 # =====
4
5 n=100
6
7 X = function(n=100) runif(n)
8 m = function(x) (sin(2*pi*x^3))^3
9 Y = function(X=runif(n)) m(X) + rnorm(length(X), 0, 0.5) # transformed into a function

```

### A.2.3 Minimize

```

1 library(progress)
2
3 # ---
4 # Only return the minimized MISE and the bandwidth associated
5 NW.minimizeMISE = function(Kernel = Knorm, h.test = seq(0.01, 0.15, 0.001), M = 100, ni = 100)
6   set.seed(122020)
7   best = Inf
8   best_h = 0
9
10  pb = progress_bar$new(
11    format = "[NW MISE] Best: :bmise (:bh) | Current: :cmise (:ch) [:bar] :current/:total |:p
12    total = length(h.test)
13  )
14  Sys.sleep(0.1)
15  pb$tick(0)
16
17  for (h in h.test){
18
19    ISE = matrix(NA, M, 1)
20
21    for (i in 1:M){
22      boot.X = runif(ni)
23      boot.Y = Y(boot.X)
24      boot.ISE = integrate(
25        function(xi) (sapply(xi, function(x) NW.regEst(x, boot.X, boot.Y, h, Kernel)) - m(xi)
26        lower=0, upper=1)$value
27      ISE[i,] = boot.ISE
28    }
29
30    MISE = round(mean(ISE), 6)
31
32    if (MISE < best){
33      best = MISE
34      best_h = h
35    }
36    pb$tick(tokens=list(bmise=best, bh=best_h, cmise=MISE, ch=h))
37  }
38  sprintf("[NW MISE]: %.6f with %.6f as bandwidth", best, best_h)
39  return (c(best, best_h))
40 }
41
42 # ---
43 # This one returns MSSE, bandwidth, Bias and also the mean of the regressions (for variance)
44 NW.minimizeMSSE = function(Kernel = Knorm, h.test = seq(0.01, 0.15, 0.001), M = 100, ni = 100)
45   set.seed(122020)
46   boot.x = seq(0, 1, length = ni)
47   res = data.frame(MSSE=Inf, h=NA, Bias=matrix(NA, ni, 1), MReg=matrix(NA, ni, 1))
48
49   pb = progress_bar$new(
50     format = "[NW MSSE] Best: :bmsse (:bh) | Current: :cmsse (:ch) [:bar] :current/:total |:p
51     total = length(h.test)
52   )
53   Sys.sleep(0.1)
54   pb$tick(0)
55   for (h in h.test){
56     SE = matrix(NA, M, ni)
57     Reg = matrix(NA, M, ni)
58
59     for (i in 1:M){
60       boot.X = runif(ni)
61       boot.Y = Y(boot.X)
62       boot.reg = sapply(boot.x, function(xi) NW.regEst(xi, boot.X, boot.Y, h, Kernel))

```

```

63     Reg[i,] = boot.reg
64     SE[i,] = (boot.reg - m(boot.x))^2
65 }
66
67 MReg = colMeans(Reg)
68 Bias = MReg - m(boot.x)
69 MSE = colMeans(SE)
70 MSSE = round(mean(MSE),6)
71
72 if (MSSE < mean(res$MSSE)){
73     res$MSSE = MSSE
74     res$h = h
75     res$Bias = Bias
76     res$MReg = MReg
77 }
78 pb$tick(tokens=list(bmsse=res$MSSE, bh=res$h, cmsse=MSSE, ch=h))
79 }
80 return (res)
81 }
82
83 # ---
84 # Haven't found a way to compute MISE with LM (yet)
85 # But as MSSE ~ MISE, this is fine.
86 # This one returns MSSE, bandwidth, Bias, the mean of the M regressions
87 # And also the min and max values taken
88 # ---
89 LM.minimizeMSSE = function(poly.test = seq(1, 20, 1), M = 200){
90     set.seed(122020)
91     res = data.frame(MSSE=Inf, poly=NA, MSE=matrix(Inf, n, 1), Bias=matrix(Inf, n, 1), MReg=matrix(Inf, n, 1),
92                     low = matrix(NA, n, 1), up = matrix(NA, n, 1))
93
94     pb = progress_bar$new(
95         format = "[LM MSSE] Best: :bmsse (:bp) | Current: :cmsse (:cp) | Bias: :bias [:bar] :curr",
96         total = length(poly.test)
97     )
98     pb$tick(0)
99
100    for (p in poly.test){
101        Reg = matrix(NA, M, n)
102        SE = matrix(NA, M, n)
103
104        for (i in 1:M){
105            boot.X = sort(X(n))
106            boot.Y = Y(boot.X)
107            boot.reg = lm(boot.Y ~ poly(boot.X, degree = p, raw=T))
108            boot.pred = predict(boot.reg, newdata=poly(CP.x, degree=p), type="response")
109            Reg[i,] = boot.pred
110            SE[i,] = (boot.pred - m(CP.x))**2
111        }
112
113        # Add "confidence interval"
114        # (this is simply the min / max values taken not real CI)
115        CI = matrix(NA, n, 2)
116
117        for (i in 1:n){
118            CI[i,1] = min(Reg[,i])
119            CI[i,2] = max(Reg[,i])
120        }
121
122        MReg = colMeans(Reg)
123        Bias = MReg - m(CP.x)
124
125        MSE = colMeans(SE)

```

```

126     MSSE = round(mean(MSE),6)
127
128     if (mean(MSE) < mean(res$MSE)){
129         res$MSE = MSE
130         res$MSSE = MSSE
131         res$poly = p
132         res$Bias = Bias
133         res$MReg = MReg
134         res$low = CI[,1]
135         res$up = CI[,2]
136     }
137
138     pb$tick(tokens=list(bias=sum(abs(Bias)), bmsse=res$MSSE, bp=res$poly, cmsse=MSSE, cp=p))
139 }
140 return (res)
141 }

```

#### A.2.4 Investigate

```

1 # ---
2 # Compute the MSE by M Monte-carlo simulations
3
4 NW.investigation = function(M=5000, n=100){
5     SE = matrix(NA, M, length(CP.to_check))
6
7     for (i in 1:M){
8         boot.X = runif(n)
9         boot.Y = Y(boot.X)
10        SE[i,] = (sapply(CP.to_check,function(xi) NW.regEst(xi, boot.X, boot.Y, NW.h, Knorm))- m(
11    })
12
13    MSE = colMeans(SE)
14    return(MSE)
15 }
16
17 LM.investigation = function(M=5000, n=100){
18     SE = matrix(NA, M, length(CP.to_check))
19
20     for (i in 1:M){
21         boot.X = sort(runif(n))
22         boot.Y = Y(boot.X)
23         boot.reg = lm(boot.Y ~ poly(boot.X, degree = LM.poly, raw=ifelse(n<50, T, F)))
24         boot.pred = fitted(boot.reg)[c(1,CP.to_check*n)]
25         SE[i,] = (boot.pred - m(CP.to_check))^2
26     }
27
28     MSE = colMeans(SE)
29     return(MSE)
30 }
31
32 # ---
33 # Create the dataframe for each sample size
34
35 CP.investigation.toframe = function(investigation_method){
36     set.seed(122020)
37     res = NULL
38     pb = progress_bar$new(
39         format = "[Investigation] Sample: :sample | [:bar] :current/:total |:percent | :elapsed",
40         total = length(CP.investigation.samples)
41     )
42     pb$tick(0)
43     Sys.sleep(0.1)
44     for (s in CP.investigation.samples){

```

```

45     pb$tick(tokens=list(sample=s))
46     res = cbind(res, investigation_method(n=s))
47 }
48 return (res)
49 }

```

### A.2.5 Draw

```

1 # ---
2 # Draw the comparison of LM and NW
3
4 CP.drawComparison = function(BestNW, BestLM, n=100, save=F){
5   set.seed(122020)
6   if(save) {jpeg("./plots/comparison.jpg", quality = 100, width = 1080, height = 720)}
7   plot(CP.X, CP.Y, pch="+", col = "grey", xlab="X", ylab="Y")
8   lines(CP.x, BestNW, col = 2, lwd=ifelse(save, 2, 1))
9   lines(CP.x, BestLM, col=4, lwd=ifelse(save, 2, 1), pch=19)
10  lines(CP.x, m(CP.x), col = 1, lwd=ifelse(save, 2, 1))
11  legend("topleft",
12        legend = c("True", "NW", "LM"),
13        col = c(1, 2, 4),
14        lty = c(1, 1, 1),
15        #pch = c(NA, 19, 19),
16        cex = ifelse(save, 1.5, 0.8)
17  )
18  title("NW v.s LM")
19  if(save) {dev.off()}
20 }
21
22 # ---
23 # Draw the evolution of the MSE
24
25 CP.drawMSE = function(investigation_frame, title, save=F){
26   set.seed(122020)
27   middle = nrow(investigation_frame)-1
28   if(save) {jpeg(paste0("./plots/investigation/",title,".jpg"), quality = 100, width = 1080,
29   plot(CP.investigation.samples,
30       head(investigation_frame, 1), col = 1, lwd = ifelse(save, 2, 1), lty=2,
31       type="l", xlab="Sample size", ylab="MSE", ylim = c(-0.0001, min(max(investigation_fram
32   )
33
34   for (i in 2:middle){
35     lines(CP.investigation.samples, investigation_frame[i,], col = i, lwd = ifelse(save, 2, 1
36   }
37
38   lines(CP.investigation.samples, tail(investigation_frame, 1), col = 1, lwd = ifelse(save, 2
39   abline(h = 0, lty = 2)
40   legend("topright",
41         legend = CP.to_check,
42         col = c(1, 2:middle, 1),
43         lty = c(2, rep(1,middle-1), 3),
44         lwd = 2,
45         cex = ifelse(save, 1.5, 0.8)
46   )
47   title(title)
48   if(save) {dev.off()}
49 }

```