

## CSE 544, Spring 2020, Probability and Statistics for Data Science

### Assignment 3: Non Parametric Inference

Due: 3/4, in class

(7 questions, 75 points total)

I/We understand and agree to the following:

- (a) Academic dishonesty will result in an 'F' grade and referral to the Academic Judiciary.
- (b) Late submission, beyond the 'due' date/time, will result in a score of 0 on this assignment.

(write down the name of all collaborating students on the line below)

---

1. (5 pts)

$$\begin{aligned}MSE &= E[(\hat{\theta} - \theta)^2] = E[\hat{\theta}^2 + \theta^2 - 2\theta\hat{\theta}] \\&= E[\hat{\theta}^2] + E[\theta^2] - 2\theta E[\hat{\theta}] \quad - (1)\end{aligned}$$

$$Var(\hat{\theta}) = E[\hat{\theta}^2] - E^2[\hat{\theta}] \quad - (2)$$

$$Bias(\hat{\theta}) = E[\hat{\theta}] - \theta$$

$$Bias^2(\hat{\theta}) = E^2[\hat{\theta}] + \theta^2 - 2\theta E[\hat{\theta}] \quad - (3)$$

From (3) + (2), we get

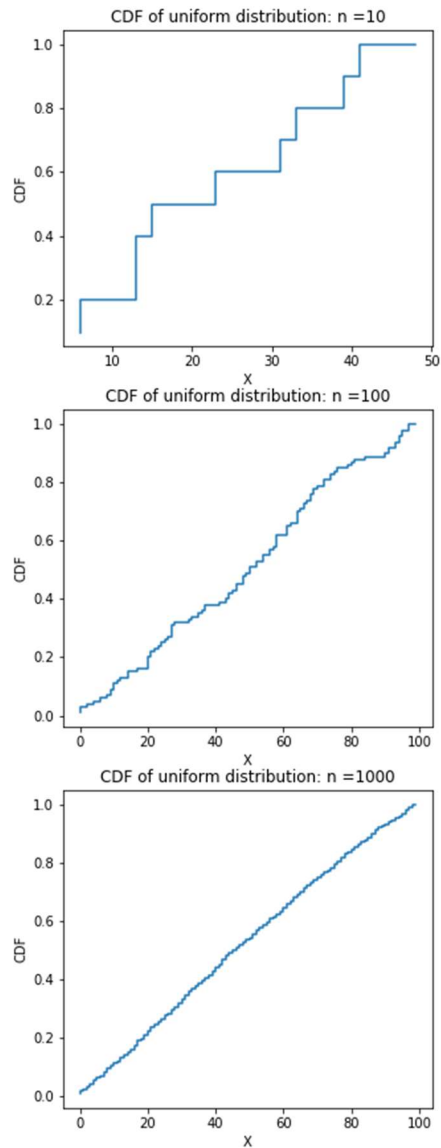
$$Var(\hat{\theta}) + Bias^2(\hat{\theta}) = E[\hat{\theta}^2] + \theta^2 - 2\theta E[\hat{\theta}] \quad - (4)$$

From (1) and (4), we get

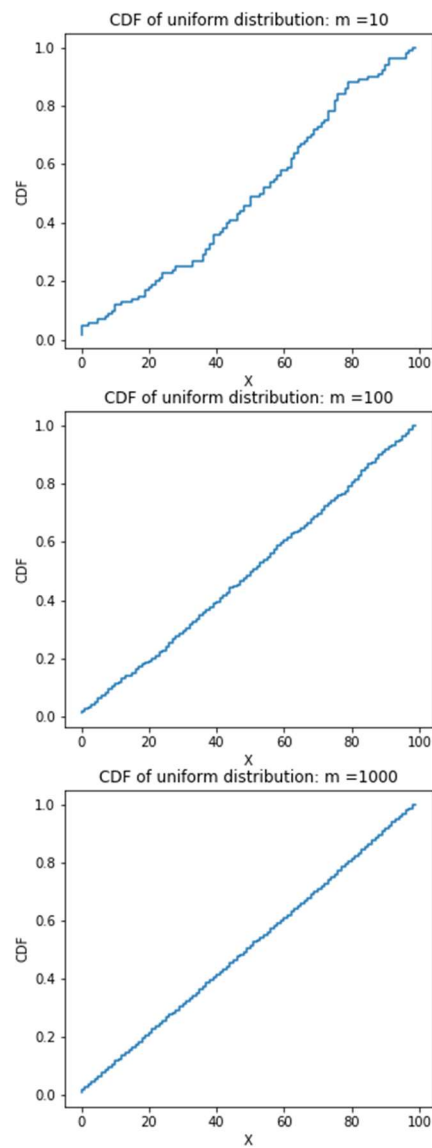
$$MSE = Var(\hat{\theta}) + Bias^2(\hat{\theta})$$

2. (3+2+3+2+3+4 pts)

2. b



2. d

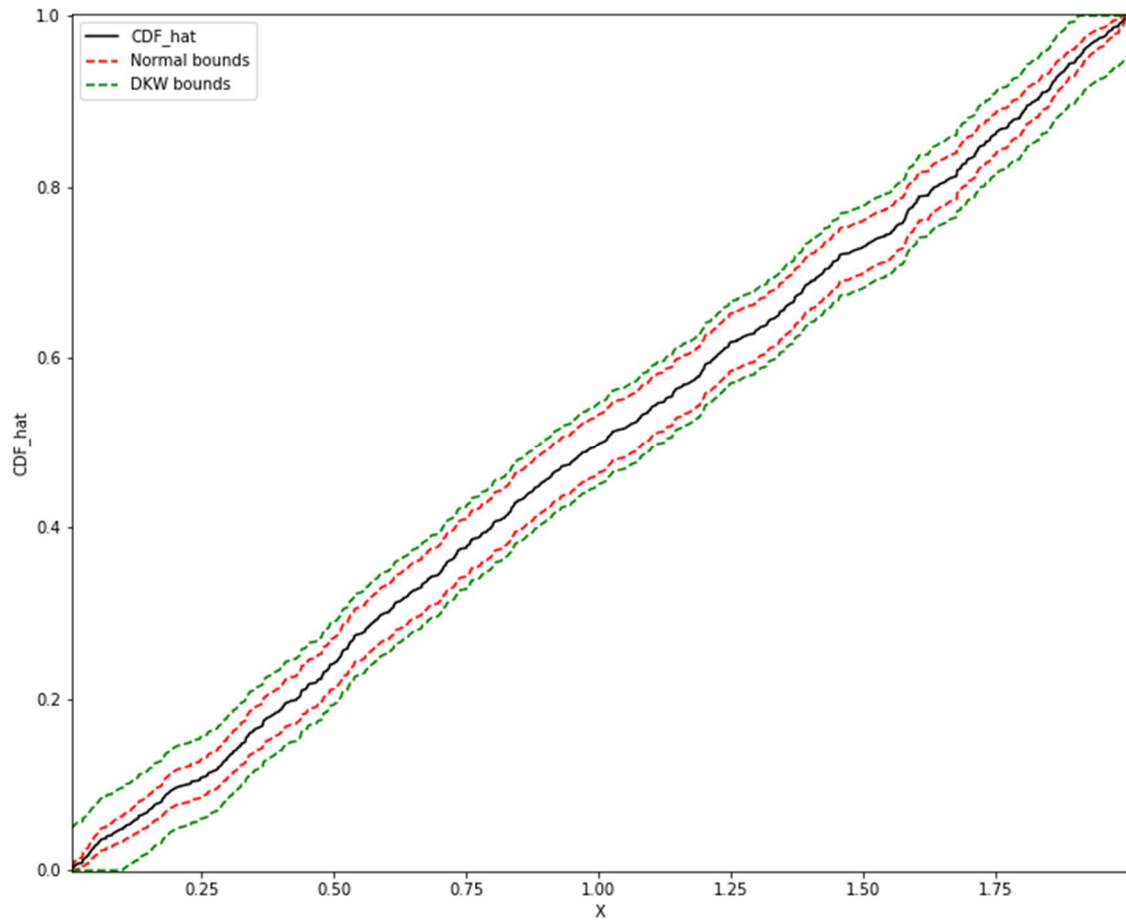


#### OBSERVATIONS

2.b: As value of  $n$  (sample size) increases the CDF estimate becomes smoother and the estimated CDF approaches the true CDF.

2.d: As value of  $m$  (# of rows) increases the CDF estimate becomes smoother and approaches the true CDF even with small sample size because  $m$  list of  $n$  samples is equivalent to a sample of size  $n.m$ .

2.e, f



From the figure we can see that Normal bound is tighter than DKW bound.

a) Let  $\hat{\sigma}^2$  be plugin estimator of  $\sigma^2$  &  $\bar{X}_n$  be plugin estimator for mean  $\mu$ .

We know,  $E[X] = \sum_i x_i p(x_i)$ . Using plugin estimator for  $p(x_i)$

we get  $p(x_i) = 1/n$  where  $n = \text{sample size}$ .

$$\therefore E[X] = \frac{1}{n} \sum_i x_i = \bar{X}_n ; E[X^2] = \sum_i x_i^2 p(x_i) = \sum_i x_i^2 \hat{p}(x_i) \\ = \frac{1}{n} \sum_i x_i^2$$

$$\therefore \hat{\sigma}^2 = E[X^2] - (E[X])^2 = \frac{1}{n} \sum_i x_i^2 - \left(\frac{1}{n} \sum_i x_i\right)^2 = \frac{1}{n} \sum_i x_i^2 - \bar{X}_n^2 \quad \text{--- (1)}$$

$$\text{R.H.S.} = \frac{1}{n} \sum_i (x_i - \bar{X}_n)^2 = \frac{1}{n} \sum_i (x_i^2 - 2x_i \bar{X}_n + \bar{X}_n^2)$$

$$= \frac{1}{n} \sum_i x_i^2 - \frac{2\bar{X}_n}{n} \sum_i x_i + \frac{\bar{X}_n^2}{n} \sum_i 1$$

$$= \frac{1}{n} \sum_i x_i^2 - 2\bar{X}_n \cdot \bar{X}_n + \bar{X}_n^2 \cdot \frac{n}{n}$$

$$= \frac{1}{n} \sum_i x_i^2 - \bar{X}_n^2 \quad \text{--- (2)}$$

From (1) & (2)

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i (x_i - \bar{X}_n)^2$$

(b)

The bias of estimator  $\hat{\sigma}^2$  is

$$\begin{aligned} E[\hat{\sigma}^2] &= E\left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right] \\ &= E\left[\frac{1}{n} \sum_{i=1}^n ((x_i - \mu) - (\bar{x}_n - \mu))^2\right] \\ &= E\left[\frac{1}{n} \sum_i (x_i - \mu)^2 - \frac{2}{n} (\bar{x}_n - \mu) \sum_i (x_i - \mu) + (\bar{x}_n - \mu)^2\right] \\ &= E\left[\frac{1}{n} \sum_i (x_i - \mu)^2 - \frac{2}{n} (\bar{x}_n - \mu) \cdot n \cdot (\bar{x}_n - \mu) + (\bar{x}_n - \mu)^2\right] \\ &= E\left[\frac{1}{n} \sum_i (x_i - \mu)^2 - (\bar{x}_n - \mu)^2\right] \\ &= E\left[\frac{1}{n} \sum_i (x_i - \mu)^2\right] - E[(\bar{x}_n - \mu)^2] \\ &= \sigma^2 - E[(\bar{x}_n - \mu)^2] \\ &= \sigma^2 - \text{Var}(\bar{x}_n) = \sigma^2 - \text{Var}\left(\frac{1}{n} \sum_i x_i\right) \\ &= \sigma^2 - \frac{1}{n^2} \sum_i \text{Var}(x_i) = \sigma^2 - \frac{1}{n} \sigma^2 \end{aligned}$$

Thus bias is  $E[\hat{\sigma}^2] - \sigma^2 = -\frac{1}{n} \sigma^2$

(C) Let  $\hat{\sigma}^2$ ,  $\hat{\mu}$  be plugin estimator for  $\sigma^2$  &  $\mu$ .

From part A we know  $\hat{\sigma}^2 = \frac{1}{n} \sum_i (X_i - \bar{X}_n)^2$  — (1)

$$\& \hat{\mu} = \bar{X}_n$$

$$\begin{aligned} E[(X - \mu)^4] &= \sum_i (X_i - \mu)^4 \cdot P(X_i) = \sum_i (X_i - \mu)^4 \cdot \hat{P}(X_i) \\ &= \frac{1}{n} \sum_i (X_i - \hat{\mu})^4 \quad \{ \text{Plugin estimator} \} \quad \text{--- (2)} \end{aligned}$$

From (1)

$$\hat{\sigma}^4 = \frac{1}{n^2} \left( \sum_i (X_i - \bar{X}_n)^2 \right)^2 \quad \text{--- (3)}$$

From (2) & (3) & by definition of Kurt[X], we have:

$$\begin{aligned} \text{Kurt}[X] &= \frac{\frac{1}{n} \sum_i (X_i - \bar{X}_n)^4}{\frac{1}{n^2} \left( \sum_i (X_i - \bar{X}_n)^2 \right)^2} \\ &= \frac{n \sum_i (X_i - \bar{X}_n)^4}{\left( \sum_i (X_i - \bar{X}_n)^2 \right)^2} \end{aligned}$$

4.a

$$\begin{aligned}\hat{F}(\alpha) &= \frac{\sum_{i=1}^n I(X_i < \alpha)}{n} \\ E[\hat{F}(\alpha)] &= E\left[\frac{\sum_{i=1}^n I(X_i < \alpha)}{n}\right] \\ E[\hat{F}(\alpha)] &= \frac{\sum_{i=1}^n E[I(X_i < \alpha)]}{n} \quad \text{By L.O.E} \\ E[\hat{F}(\alpha)] &= \frac{n \cdot E[I(X_i < \alpha)]}{n} \quad \because X_i \text{ s are iid} \\ E[\hat{F}(\alpha)] &= E[I(X_i < \alpha)] \\ E[\hat{F}(\alpha)] &= \Pr(X_i < \alpha) \\ E[\hat{F}(\alpha)] &= F(\alpha) \quad - (1)\end{aligned}$$

4.b

$$\text{Bias}(\hat{F}(\alpha)) = E[\hat{F}(\alpha)] - F(\alpha) \quad - (2)$$

From (1) and (2), we get  $\text{Bias}(\hat{F}(\alpha)) = 0$

4.c

$$\begin{aligned}\text{Se}(\hat{F}(\alpha)) &= \sqrt{\text{Var}(\hat{F}(\alpha))} \\ \text{Var}(\hat{F}(\alpha)) &= \text{Var}\left(\frac{\sum_{i=1}^n I(X_i < \alpha)}{n}\right) \\ \text{Var}(\hat{F}(\alpha)) &= \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n I(X_i < \alpha)\right) \\ \text{Var}(\hat{F}(\alpha)) &= \frac{n}{n^2} \text{Var}(I(X_i < \alpha)) \quad \because X_i \text{ s are iid} \\ \text{Var}(\hat{F}(\alpha)) &= \frac{1}{n} \text{Var}(I(X_i < \alpha)) \\ \text{Var}(\hat{F}(\alpha)) &= \frac{(1 - \Pr(X_i < \alpha)) \cdot \Pr(X_i < \alpha)}{n} \\ \text{Var}(\hat{F}(\alpha)) &= \frac{F(\alpha)(1 - F(\alpha))}{n} \\ \text{Se}(\hat{F}(\alpha)) &= \sqrt{\frac{F(\alpha)(1 - F(\alpha))}{n}}\end{aligned}$$

4.d

As  $n \rightarrow \infty$   $\text{Bias}(\hat{F}(\alpha)) = 0$  and  $\text{Se}(\hat{F}(\alpha)) \rightarrow 0$ ,  $\therefore \hat{F}$  is consistent estimator of  $F$ .

5. (3+4+3+3)

a. Let  $x_i$  be the start point of the bin  $B_i$  which is of size  $b$ .

Let  $p_i$  be the probability that a point  $x$  lies in  $B_i$ , i.e.  $p_i = \Pr(x \in (x_i, x_i + b))$

$$\therefore p_i = \frac{1}{n} \sum_{j=1}^n I(x_j \in (x_i, x_i + b)) \quad (1)$$

b. Given  $x \in B_j = (x_j, x_j + b)$

From (1), histogram estimate of  $x$  is

$$\begin{aligned}\hat{h}(x) &= \frac{\hat{p}_j}{b} = \frac{1}{nb} \sum_{k=1}^n I(x_k \in (x_j, x_j + b)) \\ E[\hat{h}(x)] &= \frac{1}{nb} \sum_{k=1}^n E[I(x_k \in (x_j, x_j + b))] \quad \{\text{By L.O.E}\} \\ E[\hat{h}(x)] &= \frac{1}{nb} \sum_{k=1}^n \Pr(x_k \in (x_j, x_j + b)) \\ E[\hat{h}(x)] &= \frac{1}{nb} \sum_{k=1}^n \Pr(x_k \in (x_j, x_j + b)) \\ E[\hat{h}(x)] &= \frac{1}{nb} \sum_{k=1}^n [F(x_j + b) - F(x_j)] \\ E[\hat{h}(x)] &= \frac{[F(x_j + b) - F(x_j)]}{b}\end{aligned}$$

As  $b \rightarrow 0$ , we can say

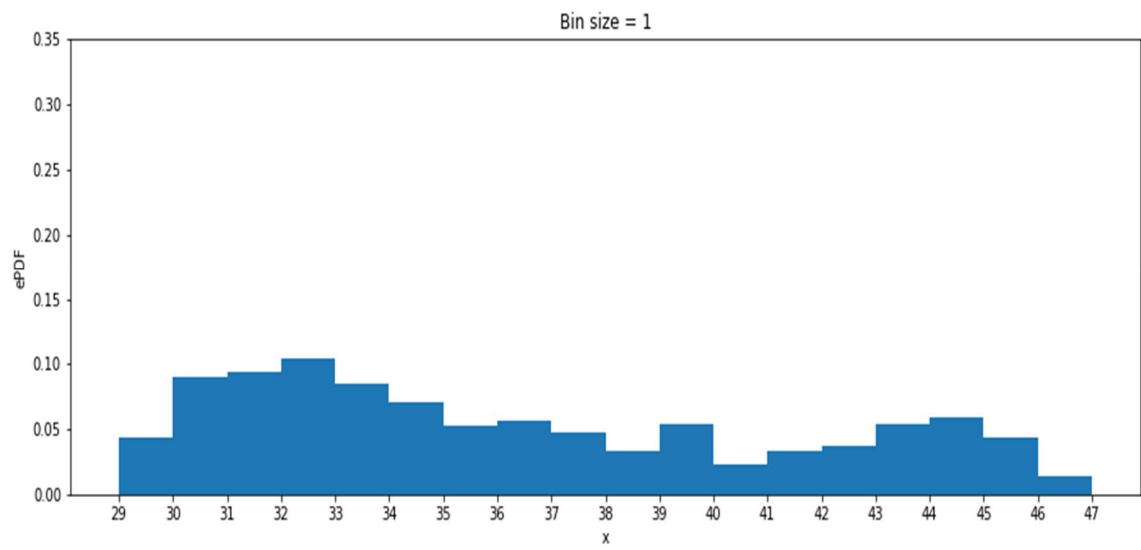
$$E[\hat{h}(x)] = \lim_{b \rightarrow 0} \frac{[F(x_j + b) - F(x_j)]}{b} = f(x_j)$$

As PDF is derivative of CDF and  $x \in (x_j, x_j + b)$ , we can say that  $f(x_j) \rightarrow f(x)$

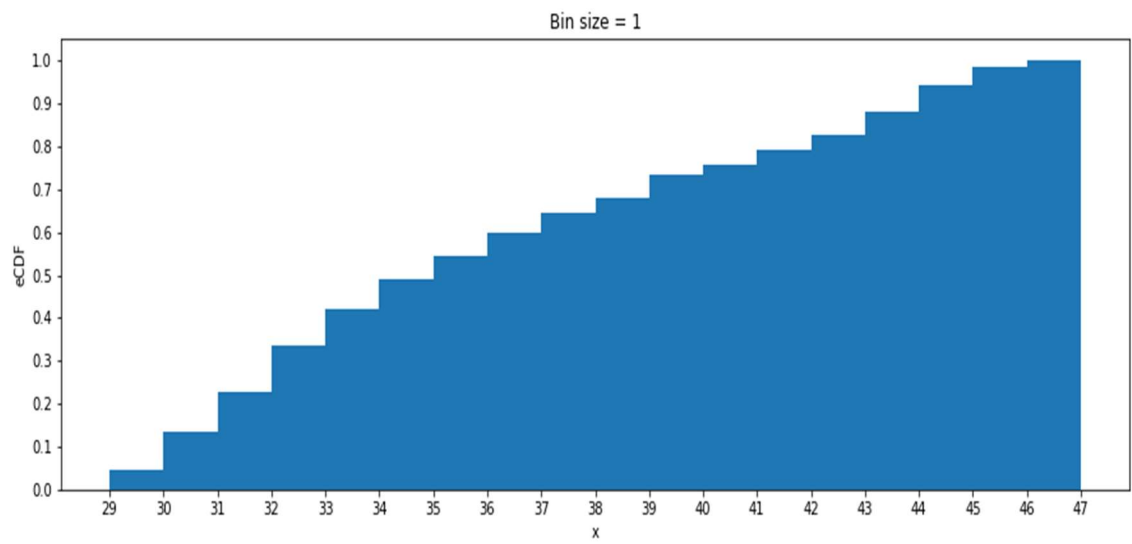
$$\therefore E[\hat{h}(x)] = f(x)$$



**c**



**d**



6. (5 pts)

$$Bias(\hat{\theta}) = E[\hat{\theta}] - \theta$$

$$\Rightarrow Bias(\hat{\theta}) = E\left[\frac{1}{n}\sum_i X_i\right] - \theta$$

$$\Rightarrow Bias(\hat{\theta}) = \frac{1}{n}\sum_i E[X_i] - \theta \quad [By L.O.E]$$

$$\Rightarrow Bias(\hat{\theta}) = \frac{1}{n}n\theta - \theta$$

$$\Rightarrow Bias(\hat{\theta}) = 0$$

$$Var(\hat{\theta}) = Var\left(\frac{1}{n}\sum_i X_i\right)$$

$$\Rightarrow Var(\hat{\theta}) = \frac{1}{n^2}\sum_i Var(X_i) \quad [\because X_i \text{ s are iid}]$$

$$\Rightarrow Var(\hat{\theta}) = \frac{1}{n^2}n\theta = \frac{\theta}{n}$$

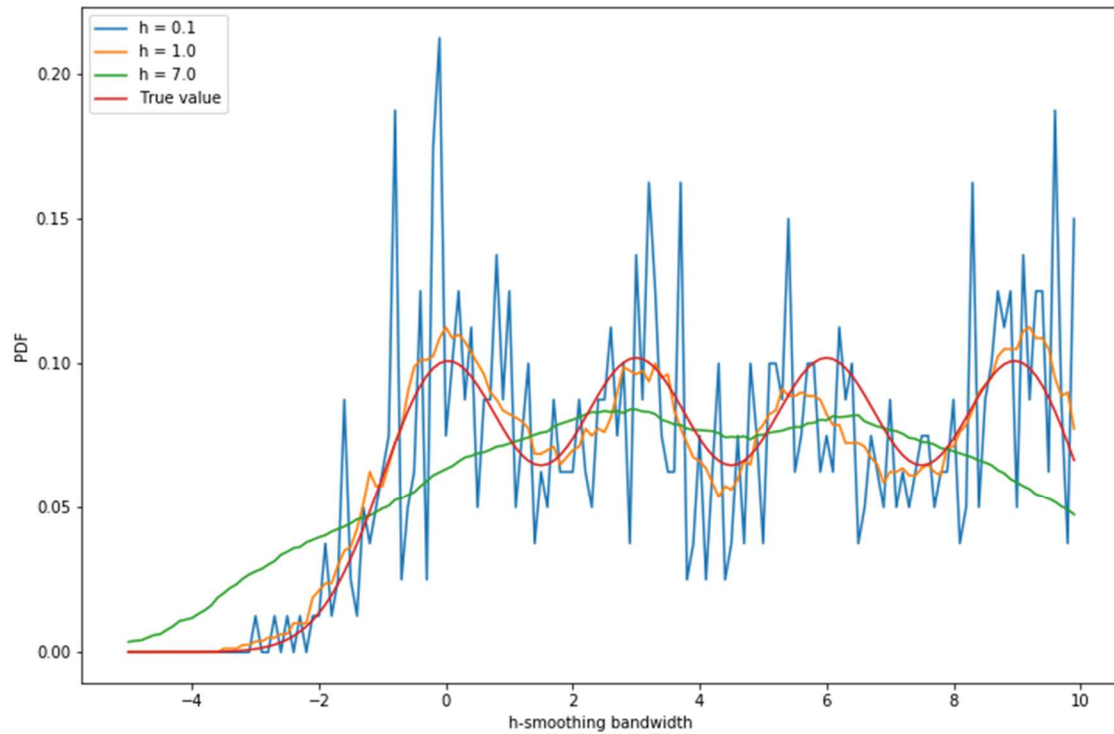
$$\Rightarrow Se(\hat{\theta}) = \sqrt{Var(\hat{\theta})} = \sqrt{\frac{\theta}{n}}$$

$$MSE(\hat{\theta}) = Bias^2(\hat{\theta}) + Var(\hat{\theta})$$

$$MSE(\hat{\theta}) = \frac{\theta}{n}$$

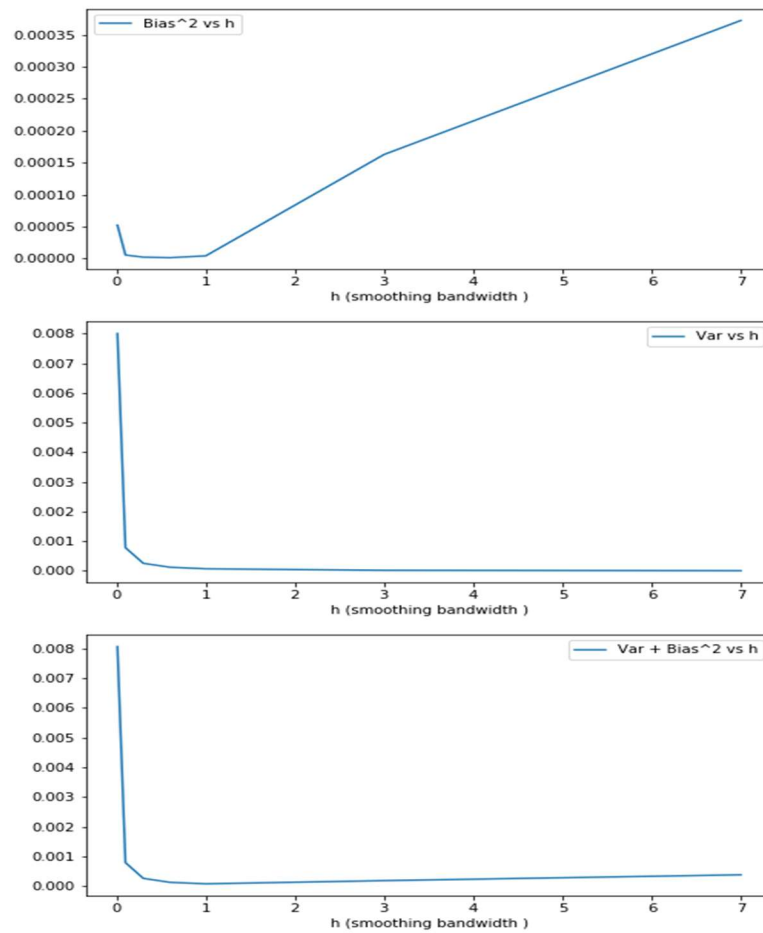
7. (8+ (5+2) pts)

a.



Observations:

1. Size of the smoothing bandwidth is inversely proportional to degree of smoothness of density estimate.
2. Smaller window leads to density estimate based on local configuration of data points.
3. In ideal case, when # of data points  $n \rightarrow \infty$ , window size  $h \rightarrow 0$  the density estimate asymptotically converges to the true distribution.



7.b Observation: Size of the smoothing window is directly proportional to Bias and inversely proportional to the variance of the density estimate.

Small window size is highly sensitive to noise whereas large window size does not take into account the fine-grained local information.

h	0.01	0.1	0.3	0.6	1.0	3.0	7.0
Bias <sup>2</sup>	5.23E-06	5.61E-06	2.27E-06	1.49E-06	4.27E-06	1.63E-04	3.72E-04
Variance	8.01E-03	7.83E-04	2.56E-04	1.23E-04	7.07E-05	1.89E-05	5.23E-06

Based on the given measure of Bias<sup>2</sup>+Variance, the **optimal value of h is 1.0**.