

STAT 510: Homework 04

David Dalpiaz

Due: Monday, February 21, 11:59 PM

Exercise 1 (Make It So)

Let $X_1, X_2, \dots, X_n \sim \text{Uniform}(0, \theta)$. Consider the estimator

$$\hat{\theta} = \max\{X_1, X_2, \dots, X_n\}.$$

Find the bias, variance, and MSE of this estimator. Assuming the estimator is biased, create a new estimator which is a simple function of $\hat{\theta}$ that is unbiased.

Solution

First, we find the CDF of $Y = \hat{\theta}$. (We use Y for ease of notation.)

$$P[\hat{\theta} < y] = (P[X < y])^n = \left(\frac{y}{\theta}\right)^n, \quad 0 < y < \theta$$

Thus, the PDF of $Y = \hat{\theta}$ is given by

$$f_{\hat{\theta}}(y) = \frac{n}{\theta} \left(\frac{y}{\theta}\right)^{n-1}, \quad 0 < y < \theta.$$

Next, we find the mean and variance

$$\mathbb{E}[\hat{\theta}] = \int_0^\theta y \cdot \frac{n}{\theta} \left(\frac{y}{\theta}\right)^{n-1} dx = \frac{n}{n+1} \cdot \theta$$

$$\mathbb{E}[(\hat{\theta})^2] = \int_0^\theta y^2 \cdot \frac{n}{\theta} \left(\frac{y}{\theta}\right)^{n-1} dx = \frac{n}{n+1} \cdot \theta^2$$

$$\mathbb{V}[\hat{\theta}] = \mathbb{E}[(\hat{\theta})^2] - (\mathbb{E}[\hat{\theta}])^2 = \frac{n\theta^2}{(n+1)^2 \cdot (n+2)}$$

We then obtain the requested quantities.

$$\text{bias}[\hat{\theta}] = \mathbb{E}[\hat{\theta}] - \theta = \boxed{-\frac{\theta}{n+1}}$$

$$\mathbb{V}[\hat{\theta}] = \boxed{\frac{n\theta^2}{(n+1)^2 \cdot (n+2)}}$$

$$\text{MSE} [\hat{\theta}] = \text{bias}^2 [\hat{\theta}] + \mathbb{V} [\hat{\theta}] = \boxed{\frac{2\theta^2}{(n+1) \cdot (n+2)}}$$

To obtain an unbiased estimator we define $\tilde{\theta} = \frac{n+1}{n}\hat{\theta}$. We then have

$$\text{bias} [\tilde{\theta}] = \frac{n+1}{n} \cdot \frac{n}{n+1} \cdot \theta = \theta.$$

Thus $\tilde{\theta} = \frac{n+1}{n}\hat{\theta}$ is an unbiased estimator for θ .

Exercise 2 (More Data, Less Problems)

Let $X_1, X_2, \dots, X_n \sim \text{Uniform}(0, \theta)$. Consider the estimator

$$\hat{\theta} = 2 \cdot \bar{X}_n$$

Find the bias, variance, and MSE of this estimator. Is this estimator consistent? Justify.

Solution

$$\mathbb{E} [\hat{\theta}] = 2 \cdot \mathbb{E} [\bar{X}_n] = 2 \cdot \frac{\theta - 0}{2} = \theta$$

$$\text{bias} [\hat{\theta}] = \boxed{0}$$

$$\mathbb{V} [\hat{\theta}] = 4 \cdot \mathbb{V} [\bar{X}_n] = 4 \cdot \frac{(\theta - 0)^2}{12n} = \boxed{\frac{\theta^2}{3n}}$$

$$\text{MSE} [\hat{\theta}] = \text{bias}^2 [\hat{\theta}] + \mathbb{V} [\hat{\theta}] = \boxed{\frac{\theta^2}{3n}}$$

Lastly, because we have

$$\text{MSE} [\hat{\theta}] \rightarrow 0$$

as $n \rightarrow \infty$, we see that $\hat{\theta}$ is a consistent estimator of θ .

Exercise 3 (A Little Bit of Bias Goes a Long Way)

Let Y have a binomial distribution with parameters n and p . Consider two estimators for p :

$$\hat{p}_1 = \frac{Y}{n}$$

and

$$\hat{p}_2 = \frac{Y+1}{n+2}$$

For what values of p does \hat{p}_2 achieve a lower mean square error than \hat{p}_1 ?

Solution:

Recall that for a binomial random variable we have

$$\begin{aligned} E[Y] &= np \\ \text{Var}[Y] &= np(1-p) \end{aligned}$$

We first calculate the bias, variance, and mean squared error of \hat{p}_1 .

$$\begin{aligned} \text{Bias}[\hat{p}_1] &= E[\hat{p}_1] - p = p - p = 0 \\ \text{Var}[\hat{p}_1] &= \frac{1}{n^2} \text{Var}[Y] = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n} \\ \text{MSE}[\hat{p}_1] &= \frac{p(1-p)}{n} \end{aligned}$$

Then we first calculate the bias, variance, and mean squared error of \hat{p}_2 .

$$\begin{aligned} E[\hat{p}_2] &= E\left[\frac{Y+1}{n+2}\right] = \frac{np+1}{n+2} \\ \text{Bias}[\hat{p}_2] &= E[\hat{p}_2] - p = \frac{np+1}{n+2} - p \\ \text{Var}[\hat{p}_2] &= \text{Var}\left[\frac{Y+1}{n+2}\right] = \frac{1}{(n+2)^2} \text{Var}[Y] = \frac{np(1-p)}{(n+2)^2} \\ \text{MSE}[\hat{p}_2] &= (\text{Bias}[\hat{p}_2])^2 + \text{Var}[\hat{p}_2] = \left(\frac{np+1}{n+2} - p\right)^2 + \frac{np(1-p)}{(n+2)^2} \end{aligned}$$

Finally, we solve for values of μ where $\text{MSE}[\hat{p}_2] < \text{MSE}[\hat{p}_1]$. We have,

$$\begin{aligned} \text{MSE}[\hat{p}_2] &< \text{MSE}[\hat{p}_1] \\ \left(\frac{np+1}{n+2} - p\right)^2 + \frac{np(1-p)}{(n+2)^2} &< \frac{p(1-p)}{n} \end{aligned}$$

Now we decide to be lazy and use [WolframAlpha](#). Thus, $\text{MSE}[\hat{p}_2] < \text{MSE}[\hat{p}_1]$ when

$$\boxed{\frac{1}{2} \left(1 - \sqrt{\frac{n+1}{2n+1}}\right) < p < \frac{1}{2} \left(1 + \sqrt{\frac{n+1}{2n+1}}\right)}$$

Note that this interval is symmetric about 0.5, which makes sense because \hat{p}_2 is biased towards 0.5

Exercise 4 (Minimizing MSE)

Suppose that $E[\hat{\theta}_1] = E[\hat{\theta}_2] = \theta$, $\text{Var}[\hat{\theta}_1] = \sigma_1^2$, $\text{Var}[\hat{\theta}_2] = \sigma_2^2$, and $\text{Cov}[\hat{\theta}_1, \hat{\theta}_2] = \sigma_{12}$. Consider the estimator

$$\hat{\theta}_3 = a\hat{\theta}_1 + (1-a)\hat{\theta}_2.$$

First, show this estimator is unbiased for all values of a . Then, what value should be chosen for the constant a in order to minimize the variance and thus mean squared error of $\hat{\theta}_3$ as an estimator of θ ?

Solution

We first verify that $\hat{\theta}_3$ is unbiased.

$$\mathbb{E}[\hat{\theta}_3] = \mathbb{E}[a\hat{\theta}_1 + (1-a)\hat{\theta}_2] = a\mathbb{E}[\hat{\theta}_1] + (1-a)\mathbb{E}[\hat{\theta}_2] = a\theta + (1-a)\theta = \theta$$

Next, we calculate the variance.

$$\begin{aligned}\text{Var}[\hat{\theta}_3] &= \text{Var}[a\hat{\theta}_1 + (1-a)\hat{\theta}_2] \\ &= a^2\text{Var}[\hat{\theta}_1] + (1-a)^2\text{Var}[\hat{\theta}_2] + 2(1-a)\text{Cov}[\hat{\theta}_1, \hat{\theta}_2] \\ &= a^2\sigma_1^2 + (1-a)^2\sigma_2^2 + 2a(1-a)\sigma_{12}\end{aligned}$$

Thus we want to find a to minimize the function

$$f(a) = a^2\sigma_1^2 + (1-a)^2\sigma_2^2 + 2a(1-a)\sigma_{12}$$

Taking the derivative, we have

$$f'(a) = 2a\sigma_1^2 - 2\sigma_2^2 + 2a\sigma_2^2 + 2\sigma_{12} - 4a\sigma_{12}$$

Setting this equal to 0 and solving for a gives

$$a = \frac{\sigma_2^2 - \sigma_{12}}{\sigma_1^2 + \sigma_2^2 - 2\sigma_{12}}$$

when $\sigma_1^2 + \sigma_2^2 - 2\sigma_{12} \neq 0$.

Finally, we check that this gives a minimum by evaluating the second derivative.

$$f''(a) = 2\sigma_1^2 + 2\sigma_2^2 - 4\sigma_{12}$$

To verify that a is a minimum, we need to show that

$$2\sigma_1^2 + 2\sigma_2^2 - 4\sigma_{12} > 0$$

which is equivalent to

$$\sigma_1^2 + \sigma_2^2 - 2\sigma_{12} > 0$$

First, recall that

$$\rho = \frac{\sigma_{12}}{\sigma_1\sigma_2}$$

Rearranged, we have

$$\rho\sigma_1\sigma_2 = \sigma_{12}$$

Since $-1 \leq \rho \leq 1$ we can state

$$\sigma_1\sigma_2 \geq \sigma_{12}$$

Now starting from the always true fact that $(\sigma_1 - \sigma_2)^2 \geq 0$, we have

$$\begin{aligned} (\sigma_1 - \sigma_2)^2 &\geq 0 \\ \sigma_1^2 + \sigma_2^2 - 2\sigma_1\sigma_2 &> 0 \\ \sigma_1^2 + \sigma_2^2 &> 2\sigma_1\sigma_2 \\ \sigma_1^2 + \sigma_2^2 &> 2\sigma_1\sigma_2 \geq 2\sigma_{12} \\ \sigma_1^2 + \sigma_2^2 &> 2\sigma_{12} \\ \sigma_1^2 + \sigma_2^2 - 2\sigma_{12} &> 0 \end{aligned}$$

Thus, we conclude that

$$a = \frac{\sigma_2^2 - \sigma_{12}}{\sigma_1^2 + \sigma_2^2 - 2\sigma_{12}}$$

gives a minimum.

Exercise 5 (Dependence in the Empirical Distribution)

Let x and y be two distinct points. Find

$$\text{Cov}(\hat{F}_n(x), \hat{F}_n(y)).$$

Solution

Without loss of generality, we will assume that $x < y$. We first evaluate the quantity

$$\text{Cov}(I(X_i < x), I(X_j < y)).$$

When $i \neq j$, this quantity is 0 since we assume that we are in the usual setup where X_1, \dots, X_n are independent.

When $i = j$, we have

$$\begin{aligned} \text{Cov}(I(X_i < x), I(X_j < y)) &= \text{Cov}(I(X_i < x), I(X_i < y)) \\ &= \mathbb{E}[I(X_i < x) \cdot I(X_i < y)] - \mathbb{E}[I(X_i < x)] \cdot \mathbb{E}[I(X_i < y)] \\ &= \mathbb{E}[I(X_i < x)] - \mathbb{E}[I(X_i < x)] \cdot \mathbb{E}[I(X_i < y)] \\ &= F(x) - F(x)F(y) \end{aligned}$$

Now, backing up, we have

$$\begin{aligned}
\text{Cov}(\hat{F}_n(x), \hat{F}_n(y)) &= \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n I(X_i < x), \frac{1}{n} \sum_{i=1}^n I(X_i < y)\right) \\
&= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{Cov}(I(X_i < x), I(X_j < y)) \\
&= \frac{n}{n^2} [F(x) - F(x)F(y)] \\
&= \frac{1}{n} [F(x) - F(x)F(y)] \\
&= \boxed{\frac{1}{n} [F(x) \cdot (1 - F(y))]}
\end{aligned}$$

Exercise 6 (Empirical Distribution Properties)

For any fixed value of x , show each of the following.

$$\mathbb{E}[\hat{F}_n(x)] = F(x)$$

$$\mathbb{V}[\hat{F}_n(x)] = \frac{F(x) \cdot (1 - F(x))}{n}$$

$$\text{MSE}[\hat{F}_n(x)] = \frac{F(x) \cdot (1 - F(x))}{n} \rightarrow 0$$

$$\hat{F}_n(x) \xrightarrow{P} F(x)$$

Solution

$$\begin{aligned}
\mathbb{E}[\hat{F}_n(x)] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n I(X_i < x)\right] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[I(X_i < x)] \\
&= \frac{1}{n} \cdot n \cdot F(x) \\
&= F(x)
\end{aligned}$$

$$\begin{aligned}
\mathbb{V}[\hat{F}_n(x)] &= \mathbb{V}\left[\frac{1}{n} \sum_{i=1}^n I(X_i < x)\right] \\
&= \frac{1}{n^2} \sum_{i=1}^n \mathbb{V}[I(X_i < x)] \\
&= \frac{1}{n^2} \sum_{i=1}^n \left(\mathbb{E}[I(X_i < x)^2] - (\mathbb{E}[I(X_i < x)])^2 \right) \\
&= \frac{1}{n^2} \sum_{i=1}^n \left(\mathbb{E}[I(X_i < x)] - (\mathbb{E}[I(X_i < x)])^2 \right) \\
&= \frac{1}{n^2} \sum_{i=1}^n (F(x) - F(x)^2) \\
&= \frac{F(x) \cdot (1 - F(x))}{n}
\end{aligned}$$

Since we see that $\hat{F}_n(x)$ is unbiased, we have

$$\text{MSE}[\hat{F}_n(x)] = \frac{F(x) \cdot (1 - F(x))}{n} \rightarrow 0.$$

Lastly, we see that

$$P\left(|\hat{F}_n(x) - F(x)| > \epsilon\right) \leq \frac{\mathbb{V}[\hat{F}_n(x)]}{\epsilon^2} = \frac{F(x) \cdot (1 - F(x))}{n\epsilon^2} \rightarrow 0.$$

Thus we have

$$\hat{F}_n(x) \xrightarrow{P} F(x)$$

Exercise 7 (Limiting Distribution of Empirical Distribution)

Let $X_1, X_2, \dots, X_n \sim F$. Given the empirical distribution function $\hat{F}_n(x)$ and a fixed point x , use the central limit theorem to find the limiting distribution of $\sqrt{n}(\hat{F}_n(x) - F(x))$.

Solution

Using the CLT, and the results of the previous exercise, we obtain

$$\sqrt{n}(\hat{F}_n(x) - F(x)) \xrightarrow{D} N(0, F(x) \cdot (1 - F(x)))$$

Exercise 8 (Using Statistical Functionals)

Let $X_1, X_2, \dots, X_n \sim F$ and let $\hat{F}_n(x)$ be the empirical distribution function. Let fixed numbers $a < b$ and define

$$\theta = T(F) = F(b) - F(a).$$

Find the estimated standard deviation of

$$\hat{\theta} = T\left(\hat{F}_n(x)\right) = \hat{F}_n(b) - \hat{F}_n(a).$$

Solution

$$\begin{aligned}\mathbb{V}\left[\hat{\theta}\right] &= \mathbb{V}\left[\hat{F}_n(b) - \hat{F}_n(a)\right] \\ &= \mathbb{V}\left[\hat{F}_n(b)\right] + \mathbb{V}\left[\hat{F}_n(a)\right] - 2 \cdot \text{Cov}\left[\hat{F}_n(b), \hat{F}_n(a)\right] \\ &= \frac{F(b) \cdot (1 - F(b))}{n} + \frac{F(a) \cdot (1 - F(a))}{n} - 2 \cdot \frac{F(a) \cdot (1 - F(b))}{n} \\ \text{s.d.}\left[\hat{\theta}\right] &= \sqrt{\frac{\hat{F}(b) \cdot (1 - \hat{F}(b))}{n} + \frac{\hat{F}(a) \cdot (1 - \hat{F}(a))}{n} - 2 \cdot \frac{\hat{F}(a) \cdot (1 - \hat{F}(b))}{n}}\end{aligned}$$

Exercise 9 (More Coverage)

Let $X_1, X_2, \dots, X_n \sim \text{Bernoulli}(p)$. Set $n = 100$ and $\alpha = 0.05$. Consider two confidence intervals for p . For both, define

$$\hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

First, consider the interval from the previous homework that we justified via Hoeffding's inequality.

$$C_n^H = \left(\hat{p}_n - \sqrt{\frac{1}{2n} \log\left(\frac{2}{\alpha}\right)}, \hat{p}_n + \sqrt{\frac{1}{2n} \log\left(\frac{2}{\alpha}\right)} \right)$$

Second, consider the “normal” interval,

$$C_n^N = \left(\hat{p}_n - z_{\alpha/2} \sqrt{\frac{\hat{p}_n(1 - \hat{p}_n)}{n}}, \hat{p}_n + z_{\alpha/2} \sqrt{\frac{\hat{p}_n(1 - \hat{p}_n)}{n}} \right).$$

Use simulation to check these intervals' coverage and expected length. Report your results using appropriate plots. Consider as many values of p as you can, but at minimum use

$$p \in (0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9).$$

Comment on the validity of these intervals and the interval lengths.

Solution

Note that this is a subset of a larger and classic analysis performed by [Brown, Cai, and DasGupta](#).

```
# sequence of values of p to consider
p = seq(from = 0.0, to = 1, by = 0.01)

# function to perform single simulation for a given p
sim_cover_and_length = function(p) {

  data = rbinom(n = 100, size = 1, prob = p)
  phat = mean(data)
```



```

margin_hoef = sqrt((1 / (2 * 100)) * log(2 / 0.05))
margin_norm = qnorm(1 - 0.025) * sqrt(phat * (1 - phat)/(100))
interval_hoef = phat + c(-1, 1) * margin_hoef
interval_norm = phat + c(-1, 1) * margin_norm

c(in_hoef = interval_hoef[1] < p & p < interval_hoef[2],
  in_norm = interval_norm[1] < p & p < interval_norm[2],
  len_hoef = diff(interval_hoef),
  len_norm = diff(interval_norm))
}

# wrapper around single simulation to repeat process
repeat_sim = function(p) {
  rowMeans(replicate(n = 20000, sim_cover_and_length(p = p)))
}

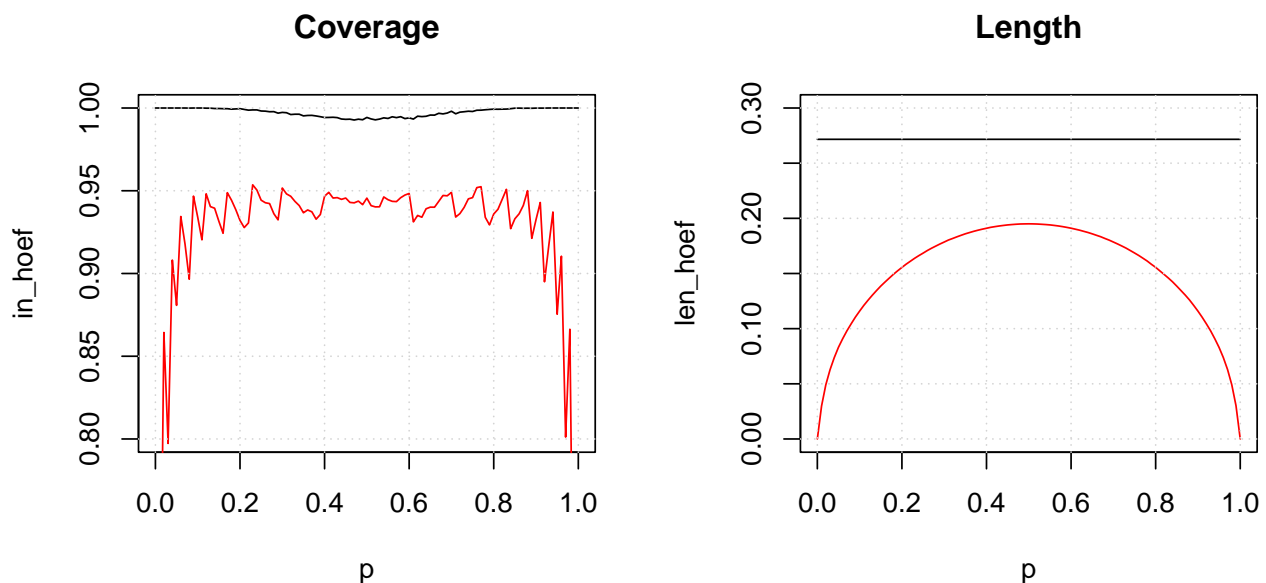
# calculate and store results for all values of p
res = cbind(p, as.data.frame(t(sapply(p, repeat_sim))))

# plot results
par(mfrow = c(1, 2))

# coverage
plot(in_hoef ~ p, data = res, type = "l", ylim = c(0.80, 1), main = "Coverage")
lines(in_norm ~ p, data = res, type = "l", col = "red")
grid()

# length
plot(len_hoef ~ p, data = res, type = "l", ylim = c(0, 0.30), main = "Length")
lines(len_norm ~ p, data = res, type = "l", col = "red")
grid()

```



- The above plots show the Hoeffding inspired intervals on black, and the normal interval in red.
- First note that the Hoeffding interval consistently over-covers, but is always valid. It's length is constant, which is not surprising, as the margin does not depend on p .

- The normal interval's coverage is at times close to 0.95, but especially as p moves closer to 0 or 1, this interval significantly under-covers. The length of these intervals are longest at $p = 0.5$, which should come as not surprise given that the margin is a function of p , in particular $p \cdot (1 - p)$ which has exactly this shape.

Exercise 10 (Empirical Distribution Confidence Bands)

The following code simulates data from three different distributions.

```
set.seed(42)
data_1 = rexp(n = 100)
data_2 = rnorm(n = 25)
data_3 = rt(n = 500, df = 3)
```

For each, plot the empirical distribution with 95% confidence bands. For each, overlay the true cumulative distribution function. Do not use R's `ecdf()` function or anything similar. You may use R's `stepfun()` function.

Solution

```
plot_ecdf = function(data, alpha, main) {

  x = sort(data)
  n = length(data)
  f = c(0, 1:n / n)
  e = sqrt(1 / (2 * n) * log(2 / alpha))

  z = plot(stepfun(x, f), do.points = FALSE, lwd = 0,
           main = main,
           ylab = "F(x)",
           ylim = c(-0.001, 1.001))

  lower = stepfun(x, pmax(f - e, 0))
  upper = stepfun(x, pmin(f + e, 1))

  lines(lower, lwd = 4, do.points = FALSE, col = "dodgerblue")
  lines(upper, lwd = 4, do.points = FALSE, col = "dodgerblue")

  plot_rect = function(i) {
    rect(xleft = z$t[i], ybottom = lower(z$t[i]),
         xright = z$t[i + 1], ytop = upper(z$t[i]),
         col = "gray90", border = "gray90", lwd = 2)
  }

  lapply(1:(length(z$t) - 1), plot_rect)
  grid()
  lines(stepfun(x, f), do.points = FALSE)

}

add_legend = function() {
  legend("bottomright", legend = c("Empirical", "Confidence Bands", "True"),
        col = c("black", "dodgerblue", "darkorange"),
        lty = c(1, 1, 2), cex = 0.8)
}
```

```

set.seed(42)
data_1 = rexp(n = 100)
data_2 = rnorm(n = 25)
data_3 = rt(n = 500, df = 3)

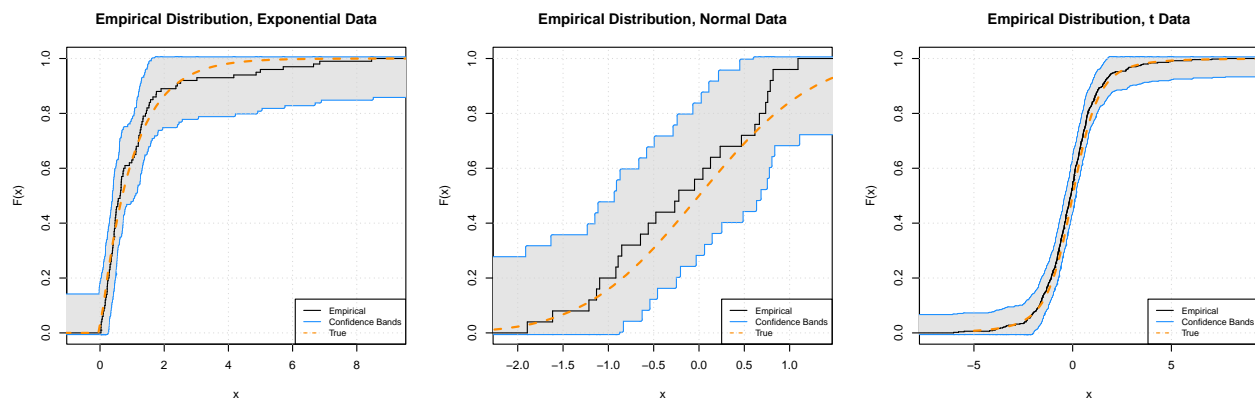
par(mfrow = c(1, 3))

plot_ecdf(data_1, alpha = 0.05,
           main = "Empirical Distribution, Exponential Data")
curve(pexp(x), col = "darkorange", add = TRUE,
      lwd = 2, from = -5, to = 10, lty = 2)
add_legend()

plot_ecdf(data_2, alpha = 0.05,
           main = "Empirical Distribution, Normal Data")
curve(pnorm(x), col = "darkorange", add = TRUE,
      lwd = 2, from = -5, to = 10, lty = 2)
add_legend()

plot_ecdf(data_3, alpha = 0.05,
           main = "Empirical Distribution, t Data")
curve(pt(x, df = 3), col = "darkorange", add = TRUE,
      lwd = 2, from = -5, to = 10, lty = 2)
add_legend()

```



Exercise 11 (Estimating Functionals with the Empirical Distribution)

The following code simulates data from a [Weibull distribution](#).

```

set.seed(42)
some_data = rweibull(n = 250, shape = 2, scale = 3)

```

Use the empirical distribution function to create plug-in estimates of the following:

- Mean
- Variance
- Skewness
- Median

Compare these results to their true values given the data generating process defined above. Report your results as a table.

Solution

```
# helper function for plug-in estimator of variance
calc_var = function(x) {
  mean((x - mean(x)) ^ 2)
}

# helper function for plug-in estimator of skewness
calc_skew = function(x) {
  mean((x - mean(x)) ^ 3) / (sqrt(calc_var(x))) ^ 3
}

# calculate estimated values
mean_est = mean(some_data)
var_est = calc_var(some_data)
skew_est = calc_skew(some_data)
med_est = quantile(some_data, probs = 0.5, type = 1)

# calculate true values
mu = 3 * gamma(1 + 1 / 2)
s2 = sigma_2 = 3 ^ 2 * (gamma(1 + 2 / 2) - (gamma(1 + 1 / 2)) ^ 2)
skew_true = (gamma(1 + 3 / 2) * 3 ^ 3 - 3 * mu * s2 - mu ^ 3) / (sqrt(s2)) ^ 3
med_true = 3 * (log(2)) ^ (1 / 2)

# arrange results in data frame
res = data.frame(
  Mean = c(mean_est, mu),
  Variance = c(var_est, s2),
  Skewness = c(skew_est, skew_true),
  Median = c(med_est, med_true),
  row.names = c("Estimated", "Truth")
)

# print data frame as table
knitr::kable(res)
```

	Mean	Variance	Skewness	Median
Estimated	2.651403	2.189681	1.0098122	2.383722
Truth	2.658681	1.931417	0.6311107	2.497664

Note that, as we have requested the plug-in estimates, we expect the exact results above, especially for the variance. (Unless you specifically noted that you are substituting the sample variance.) Since it is trickier to deal with, we will accept the more readily available answer for the median.

```
# also acceptable
median(some_data)
```

```
## [1] 2.40646
```