

HW3

Teryl Schmidt | tschmidt6@wisc.edu | 9072604920

9/27/2018

Problem 1

A chemical supply company ships a certain solvent in 10-gallon drums. Let X represent the number of drums ordered by a randomly chosen customer. Assume X has the following probability mass function (pmf). The mean and variance of X is : $\mu_X = 2.3$ and $\sigma^2_X = 1.81$:

X	1	2	3	4	5
$p(X=x)$	0.4	0.2	0.2	0.1	0.1

- a. Find $P(X \leq 2)$ and describe what it means in the context of the problem.

$$P(X \leq 2) = P(X = 1) + P(X = 2) = 0.4 + 0.2$$

$$P(X \leq 2) = 0.6$$

- b. Let Y be the number of gallons ordered, so $Y = 10X$. Find the probability mass function (pmf) of Y .

define new table (second table)

X	1	2	3	4	5
$Y = 10x$	$10 \cdot 1$	$10 \cdot 2$	$10 \cdot 3$	$10 \cdot 4$	$10 \cdot 5$
$p(X=x)$	0.4	0.2	0.2	0.1	0.1

$Y = 10x$	10	20	30	40	50
$p(X=x)$	0.4	0.2	0.2	0.1	0.1

- c. Find the mean number of gallons ordered μ_Y using the pmf from part (b) and a second time using the expectation of linear combination formula shown in class.

$$\begin{aligned}\mu_Y &= E(Y) \\ &= E(10X) \\ &= 10E(X) \\ &= 10 \cdot \mu_X \\ &= 10 \cdot 2.3 \\ &= 23\end{aligned}$$

- d. Find the standard deviation of the number of gallons ordered σ_Y using the pmf from part (b) and a second time using the standard deviation of linear combination formulas shown in class.

$$\begin{aligned}\sigma_Y &= \sqrt{\sigma^2_Y} \\ &= \sqrt{181} \\ &= 13.454\end{aligned}$$

- e. What possible values of Y are within two standard deviations of the mean value (that is, in the interval from $(\mu - 2\sigma, \mu + 2\sigma)$)? What is the probability that the observed value of Y is within two standard deviations of the mean value?

$$\begin{aligned}&(23 - 13.454, 23 + 13.454) \\ &(9.546, 36.454) \\ &P(X = 1) + P(X = 2) + P(X = 3) \\ &= 0.8\end{aligned}$$

Problem 2

Consider a large population which has true mean μ and true standard deviation σ . We take a sample of size 3 from this population, thinking of the sample as the RVs X_1, X_2, X_3 where X_i can be considered iid. We are interested in estimating μ .

- a. Consider the estimator $\hat{\mu}_1 = X_1 + X_2 - X_3$. Is this estimator unbiased? Calculate the bias.

$$\begin{aligned}E(\hat{\mu}) &= E(X_1 + X_2 - X_3) \\ &= E(X_1) + E(X_2) - E(X_3) \\ &= \mu + \mu - \mu \\ &= \mu\end{aligned}$$

Therefore $\hat{\mu}_1 = X_1 + X_2 - X_3$ is an unbiased estimator for μ

- b. Find the variance of $\hat{\mu}_1$.

$$\begin{aligned}
V(\hat{\mu}) &= V(X_1 + X_2 - X_3) \\
&= V(X_1) + V(X_2) - V(X_3) \\
&= \sigma^2 + \sigma^2 + \sigma^2 \\
&= 3\sigma^2 \\
&= \text{The variance of } \hat{\mu}_1 \text{ is } 3\sigma^2
\end{aligned}$$

c. When estimating μ , explain why someone would prefer the estimator $\bar{X} = (X_1 + X_2 + X_3) / 3$ over $\hat{\mu}_1$

$$\begin{aligned}
\hat{\mu}_2 &= (X_1 + X_2 + X_3) / 3 \\
E(\bar{X}) &= \mu \\
E(\hat{\mu}) &= E(X_1 + X_2 + X_3) / 3 \\
&= \frac{1}{3} [E(X_1) + E(X_2) + E(X_3)] \\
&= \frac{1}{3} [\mu + \mu + \mu] \\
&= \frac{3\mu}{3} \\
&= \mu
\end{aligned}$$

unbiased

$$\begin{aligned}
V(\bar{X}) &= V(X_1 + X_2 - X_3) / 3 \\
&= V(X_1) + V(X_2) - V(X_3) \\
&= \sigma^2 + \sigma^2 + \sigma^2 / 3 \\
&= 3\sigma^2 / 3 \\
&= \sigma^2 \text{ smaller variance}
\end{aligned}$$

Therefore $\hat{\mu}_2 = (X_1 + X_2 + X_3) / 3$ is an unbiased estimator for μ . And we choose \bar{X} over $\hat{\mu}_1$ because \bar{X} has a smaller variance $(3\sigma^2 / 3) = \sigma^2$, as oppose to $3\sigma^2$

d. Now consider the estimator $\hat{\mu}_2 = (X_1 + X_2 + X_3) / 2$. Is this estimator unbiased? Calculate the bias.

$$\begin{aligned}
\hat{\mu}_2 &= (X_1 + X_2 + X_3) / 2 \\
E(\bar{X}) &= \mu \\
E(\hat{\mu}) &= E(X_1 + X_2 + X_3) / 2 \\
&= \frac{1}{2} [E(X_1) + E(X_2) + E(X_3)] \\
&= \frac{1}{2} [\mu + \mu + \mu] \\
&= \frac{3\mu}{2} \neq \mu
\end{aligned}$$

Therefore $\hat{\mu}_2 = (X_1 + X_2 + X_3) / 2$ is a biased estimator for μ

e. Compute the MSE for $\hat{\mu}_2$.

$$\begin{aligned}
\hat{\mu}_2 &= (X_1 + X_2 + X_3) / 2 \\
\text{Var}(\hat{\mu}) &= \text{Var}(X_1 + X_2 + X_3) / 2 \\
&= \frac{1}{4} [\text{Var}(X_1) + \text{Var}(X_2) + \text{Var}(X_3)] \\
&= \frac{1}{4} [\sigma^2 + \sigma^2 + \sigma^2] \\
&= \frac{3\sigma^2}{4} \\
\text{MSE}(\hat{\mu}) &= \text{Var}(\hat{\mu}) + [\text{Bias}(\hat{\mu})]^2 \\
&= \text{MSE}(\hat{\mu}) = \frac{3\sigma^2}{4} + \left(\frac{\mu}{2}\right)^2 \\
&= \text{MSE}(\hat{\mu}) = \frac{3\sigma^2}{4} + \left(\frac{\mu^2}{4}\right) \\
&= \text{MSE}(\hat{\mu}) = \frac{3\sigma^2 + \mu^2}{4}
\end{aligned}$$

Problem 3

Let F be an RV that represents the operating temperature in Fahrenheit of one instance of a manufacturing process, and assume $F \sim N(100, \text{Var}(F) = 5^2)$. Let C be an RV that represents the same process, but measured in Celsius.

Fahrenheit can be converted to Celsius using $C = \frac{5}{9} (F - 32)$. Using normal table in Canvas, solve for the following (You can also check your answers using R):

a. Find the probability that one randomly selected instance of the process will have operating temperature greater than 98.6 Fahrenheit.

$$\begin{aligned}
&P(X > 98.6) \\
&P(X - 100 / 5 > 98.6 - 100 / 5) \\
&1 - P(Z \leq -0.28) \\
&1 - 0.3897 \\
&= 0.6103
\end{aligned}$$

b. Find the distribution of C . (Hint: $C \sim N(\mu, \sigma^2)$)

distribution of will be normal with the mean

$$\begin{aligned}\mu_c &= E\left(\frac{5}{9}(X - 32)\right) \\ &= \frac{5}{9}E(E(X) - 32) \\ &= \frac{5}{9}E(100 - 32) \\ &= 340 / 9 = 37.77778\end{aligned}$$

and variance

$$\begin{aligned}\sigma_c^2 &= (5/9)^2[\text{Var}(X) + 0] \\ &= (5/9)^2[5^2 + 0] \\ &= 7.716\end{aligned}$$

$$C \sim N(37.77, 7.716)$$

- c. Find the probability that one randomly selected instance of the process will have operating temperature below 32 Celsius.

$$\begin{aligned}P(X < 32) \\ &= P(32 - 37.77778 / \sqrt{7.716}) \\ &= P(Z < -2.0799) \\ &= 0.0188\end{aligned}$$

- d. Above what temperature (in Celsius) is the top 10% of operating temperatures?

$$\begin{aligned}P(X > a) &= 0.90 \\ P(x - 37.77 / \sqrt{7.716} > a - 37.77 / \sqrt{7.716}) &= 0.1 \\ (a - 37.77 / \sqrt{7.716}) &= 0.90 \\ &= 1.28 \text{ (table)} \\ a &= 37.77 + 1.28 + \sqrt{7.716} \\ a &= 41.34\end{aligned}$$

- e. Find the probability in a sample of 6 instances, more than 4 instances have operating temperature above 32 Celsius. (Assuming observations in the sample are independent)

$$\begin{aligned}P(X > 32) \\ &= 1 - P(32 - 37.77778 / \sqrt{7.716}) \\ &= 1 - P(Z > -2.0799) \\ &= 1 - 0.0188 \\ &= 0.9812\end{aligned}$$

$$\begin{aligned}P(X > 32) &= 0.9812 \\ B &= \# \text{ of instances samples} > 32 \\ X &\sim \text{Bin}(6, 0.9812) \\ P(X > 4) &= P(X = 5) + P(X = 6) \\ &= \binom{6}{5} * (0.98)^5 * (1 - 0.98) + \binom{6}{6} * (0.98)^6 \\ &= 0.9943\end{aligned}$$

- f. Find the distribution of \bar{X} for $n=6$, then find the probability that the average operating temperature in a sample of 6 instances is above 32 Celsius.

$$\begin{aligned}\sigma^2 / n \\ &= 7.716 / 6 = 1.29 \\ \bar{X} &\sim N(37.77, 1.29) \\ P(\bar{X} > 32) \\ &= P(\bar{X} - 37.77 / \sqrt{1.29} > 32 - 37.77 / \sqrt{1.29}) \\ &= 1 - P(Z > -5.09) \\ &= 0.99999 \text{ or } 1\end{aligned}$$

Problem 4

Retail stores experience their heaviest returns on December 26th and December 27th each year. The distribution for the Number of Items Returned (X) for Hildale Macy's on those days last year is given in the table below. It has mean: $\mu = 2.61$ and variance $\sigma^2 = 1.80$. Assume the probability distribution also holds for this year. Number of Items Probability

1		0.25
2		0.28
3		0.20
4		0.17
5		0.08
6		0.02

- a. In this year, a random sample of size 45 returns is selected for review. Describe the sampling distribution of the sample mean (shape, center, and spread).

$$\mu_{\bar{X}} = \mu = 2.61$$

$$\sigma_{\bar{X}} = \sigma / \sqrt{n} = \sqrt{(1.80 / 45)} = 0.20$$

The sample follows a normal distribution

$$X \sim N(2.61, 0.20)$$

- b. What is the probability that the sample mean will be greater than 2.9 items?

$$P(\bar{X} > 2.9) = P(Z > 2.9 - 2.61 / 0.20)$$

$$P(Z > 1.45) = 0.5 - 0.4265 = 0.0735$$

- c. Find an upper bound b such that the total number of items returned by 45 customers will be less than b with probability 0.95.

For 95% the Z value is 1.645

Therefore the upperbound b is:

$$b = \mu + Z\sigma = 2.61 + 1.645 * 0.20 = 2.939$$

$$(Z \leq x - 2.61 / 0.20) = 0.95$$

$$x = 2.939$$

$$b = 45 * 2.939 = 132$$

Problem 5

We will be exploring the difference between using the standard deviation formula: $s1 = \sqrt{(\bar{X} - \bar{X})^2 / (n - 1)}$

$s2 = \sqrt{(\bar{X} - \bar{X})^2 / \sqrt{n}}$ through a simulation.

In the code below, I have defined a population of values, named pop1. I have also written a function sample.sd to compute the sample standard deviation on a set of numbers passed in.

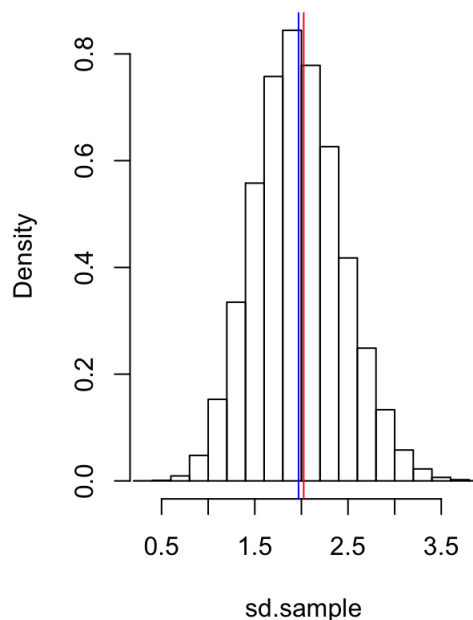
- Copy and paste the entire chunk of code from set.seed(1) (so that we're all using the same data) through par(mfrow=c(1,1)) (so we reset the graphics pane). Then copy and paste the three lines of code within the samp.sd function into the pop.sd function. Update the three lines as necessary so the pop.sd function will calculate the population standard deviation formula for a set of values.
- Run the entire chunk of code (from set.seed(1) through par(mfrow=c(1,1))). (i) What do you notice about the average of the standard deviations produced using the samp.sd function compared to the pop.sd function compared to the true population standard deviation? (ii) Why might we prefer to use the sample.sd formulation when we have a sample of data and are interested in estimating the population standard deviation? (You can compare the resulting histograms to help you answer the question.)
- The curve of the sample.sd is more like a bell curve.
- There is less variance in the sample.sd function.

```

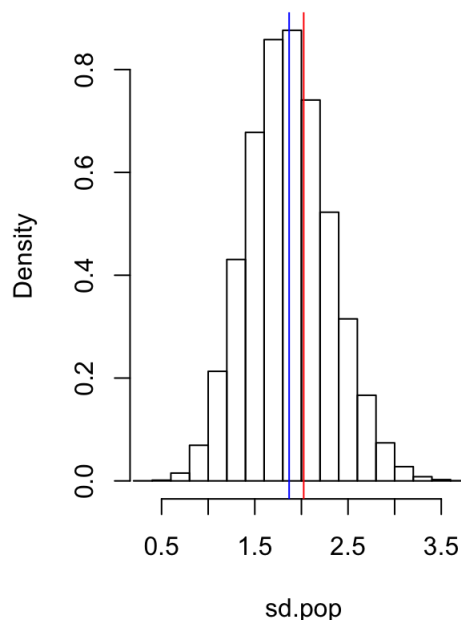
set.seed(1)
pop1<-rnorm(10000, 4, 2)
#Writing functions to calculate sample sd and population sd
samp.sd<-function(data){
  n<-length(data)
  sum.sq.devs<-sum((data-mean(data))^2)
  av.dev<-sqrt(sum.sq.devs/(n-1))
  return(av.dev)
}
pop.sd<-function(data){
  n<-length(data)
  sum.sq.devs<-sum((data-mean(data))^2)
  av.dev<-sqrt(sum.sq.devs/(n))
  return(av.dev)
}
#Simulation Section
#Building sampling distribution of pop standard deviation estimators
#estimator 1 is sd.sample and estimator 2 is sd.pop formulation
nsamples<-100000
sd.sample<-rep(0, nsamples)
sd.pop<-rep(0, nsamples)
for (i in 1:nsamples){
  samp<-sample(pop1, 10, replace=TRUE) #taking a new sample of size 10 from population
  sd.sample[i]<-samp.sd(samp) #calculating and storing sd using sample formula on new sample
  sd.pop[i]<-pop.sd(samp) #calculating and storing sd using pop formula on new sample
}
#Displaying histograms of the 100000 standard deviations we calculated used
#the population and sample equations;
#adding true population standard deviation in red and
#mean of the simulated standard deviations in blue.
par(mfrow=c(1,2))
hist(sd.sample, freq=FALSE, xlim=c(0.3, 3.75)); #histogram of generated sample sds
abline(v=pop.sd(pop1), col="red"); abline(v=mean(sd.sample), col="blue")
hist(sd.pop, freq=FALSE, xlim=c(0.3, 3.75)); #histogram of generated pop sds
abline(v=pop.sd(pop1), col="red"); abline(v=mean(sd.pop), col="blue")

```

Histogram of sd.sample



Histogram of sd.pop



```

par(mfrow=c(1,1))

```