```yaml
1      Course: Statistical_Inference
2      Lesson: Variance
3
4    - Class: text
5      Output: "Variance. (Slides for this and other Data Science courses may be found at
       github https://github.com/DataScienceSpecialization/courses/. If you care to use
       them, they must be downloaded as a zip file and viewed locally. This lesson
       corresponds to 06_Statistical_Inference/05_Variance.)"
6
7    - Class: text
8      Output: In this lesson, we'll discuss variances of distributions which, like means,
       are useful in characterizing them. While the mean characterizes the center of a
       distribution, the variance and its square root, the standard deviation, characterize
       the distribution's spread around the mean. As the sample mean estimates the
       population mean, so the sample variance estimates the population variance.
9
10   - Class: text
11     Output: The variance of a random variable, as a measure of spread or dispersion,  is,
       like a mean, defined as an expected value. It is the expected squared distance of the
       variable from its mean. Squaring the distance makes it positive so values less than
       and greater than the mean are treated the same. In mathematical terms, if X comes
       from a population with mean mu, then
12
13   - Class: text
14     Output: Var(X) = E( (X-mu)^2 ) = E( (X-E(X))^2 ) =  E(X^2)-E(X)^2
15
16   - Class: text
17     Output: So variance is the difference between two expected values. Recall that E(X),
       the expected value of a random variable from the population, is mu, the mean of that
       population.
18
19   - Class: text
20     Output: Higher variance implies more spread around a mean than lower variance.
21
22   - Class: text
23     Output: Finally, it's easy to show from the definition and the linearity of
       expectations that, if a is a constant, then Var(aX)=a^2*Var(X). This will come in
       handy later.
24
25   #- Class: video
26   #  Output: Would you like to see the equation proving this? You'll need an internet
     connection to see it.
27   #  VideoLink: "http://wilcrofter.github.io/slidex/markDown/varAX.html"
28
29   - Class: figure
30     Output: If you're interested, here's the proof. You might have to stretch out your
       plot window to make it clearer.
31     Figure: plotVform.R
32     FigureType: new
33
34   - Class: text
35     Output:  Let's practice computing the variance of a dice roll now. First we need to
       compute E(X^2). From the definition of expected values, this means we'll take a
       weighted sum over all possible values of X^2. The weight is the probability of X
       occurring.
36
37   - Class: cmd_question
38     Output: For convenience, we've defined a 6-long vector for you, dice_sqr, which holds
       the squares of the integers 1 through 6. This will give us the X^2 values. Look at it
       now.
39     CorrectAnswer: dice_sqr
40     AnswerTests: omnitest(correctExpr='dice_sqr')
41     Hint: Type dice_sqr at the command prompt.
42
43   - Class: cmd_question
44     Output: Now we need weights. For these we can use any of the three PDF's, (dice_fair,
       dice_high, and dice_low) we defined in the previous lesson. Using R's ability to
       multiply vectors componentwise and its function 'sum' we can easily compute E(X^2)
       for any of these dice. Simply sum the product dice_sqr * PDF.  Try this now with
```

```
       dice_fair and put the result in a variable ex2_fair.
45     CorrectAnswer: ex2_fair <- sum(dice_fair * dice_sqr)
46     AnswerTests: expr_creates_var('ex2_fair'); ANY_of_exprs('ex2_fair <- sum(dice_fair *
       dice_sqr)','ex2_fair <- sum(dice_sqr * dice_fair)')
47     Hint: Type 'ex2_fair <- sum(dice_fair * dice_sqr)' at the command prompt.
48
49   - Class: cmd_question
50     Output: Recall that the expected value of a fair dice roll is 3.5. Subtract the
       square of that from ex2_fair to compute the sample variance.
51     CorrectAnswer: ex2_fair-3.5^2
52     AnswerTests:
       ANY_of_exprs('ex2_fair-3.5^2','ex2_fair-3.5*3.5','ex2_fair-(3.5^2)','ex2_fair-(3.5*3.5)
       ')
53     Hint: Type 'ex2_fair-3.5^2' at the command prompt.
54
55   - Class: cmd_question
56     Output: Now use a similar approach to compute the sample variance of dice_high in one
       step. Sum the appropriate product and subtract the square of the mean. Recall that
       edh holds the expected value of dice_high.
57     CorrectAnswer: sum(dice_high * dice_sqr)-edh^2
58     AnswerTests: ANY_of_exprs('sum(dice_high * dice_sqr)-edh^2','sum(dice_sqr *
       dice_high)-edh^2','sum(dice_high * dice_sqr)-edh*edh','sum(dice_sqr *
       dice_high)-edh*edh')
59     Hint: Type 'sum(dice_high * dice_sqr)-edh^2' at the command prompt.
60
61   - Class: text
62     Output: Note that when we talk about variance we're using square units. Because it is
       often more useful to use measurements in the same units as X we define the standard
       deviation of X as the square root of Var(X).
63
64   - Class: figure
65     Output: Here's a figure from the slides. It shows several normal distributions all
       centered around a common mean 0, but with different standard deviations. As you can
       see from the color key on the right, the thinner the bell the smaller the standard
       deviation and the bigger the standard deviation the fatter the bell.
66     Figure: normalVar.R
67     FigureType: new
68
69   - Class: text
70     Output: Just as we distinguished between a population mean and a sample mean we have
       to distinguish between a population variance sigma^2 and a sample variance s^2. They
       are defined similarly but with a slight difference. The sample variance is defined as
       the sum of n squared distances from the sample mean divided by (n-1), where n is the
       number of samples or observations. We divide by n-1 because this is the number of
       degrees of freedom in the system. The first n-1 samples or observations are
       independent given the mean. The last one isn't independent since it can be calculated
       from the sample mean used in the formula.
71
72   - Class: text
73     Output:  In other words, the sample variance is ALMOST the average squared deviation
       from the sample mean.
74
75   - Class: text
76     Output: As with the sample mean, the sample variance is also a random variable with
       an associated population distribution. Its expected value or mean is the population
       variance and its distribution gets more concentrated around the population variance
       with more data. The sample standard deviation is the square root of the sample
       variance.
77
78   - Class: figure
79     Output: To illustrate this point, consider this figure which plots the distribution
       of 10000 variances, Each variance was computed on a sample of standard normals of
       size 10. The vertical line indicates the standard deviation 1.
80     Figure: moreData1.R
81     FigureType: new
82
83   - Class: figure
84     Output: Here we do the same experiment but this time (the taller lump) each of the
       10000 variances is over 20 standard normal samples. We've plotted over the first plot
```

```
        (the shorter lump) and you can see that the distribution of the variances is getting
        tighter and shifting closer to the vertical line.
 85     Figure: moreData2.R
 86     FigureType: new
 87
 88   - Class: figure
 89     Output: Finally, we repeat the experiment using 30 samples for each of the 10000
        variances. You can see that with more data, the distribution gets more concentrated
        around the population variance it is trying to estimate.
 90     Figure: moreData3.R
 91     FigureType: new
 92
 93   - Class: text
 94     Output: Now recall that the means of unbiased estimators equal the values they're
        trying to estimate. We can infer  from the above that the sample variance is an
        unbiased estimator of population variance.
 95
 96   - Class: text
 97     Output: Recall that the average of random samples from a population is itself a
        random variable with a distribution centered around the population mean.
        Specifically, E(X') = mu, where X' represents a sample mean and mu is the population
        mean.
 98
 99   - Class: text
100     Output:  We can show that, if the population is infinite, the variance of the sample
        mean is the population variance divided by the sample size. Specifically,  Var(X') =
        sigma^2 / n. Let's work through this in four short steps.
101
102   - Class: mult_question
103     Output: Which of the following does Var(X') equal? Here X' represents the sample mean
        and 'Sum(X_i)' represents the sum of the n samples X_1,...X_n. Assume these samples
        are independent.
104     AnswerChoices: Var(1/n * Sum(X_i)); E(1/n * Sum(X_i)); mu; sigma
105     CorrectAnswer: Var(1/n * Sum(X_i))
106     AnswerTests: omnitest(correctVal='Var(1/n * Sum(X_i))')
107     Hint: Which of the choices has both Var and the definition of mean in it?
108
109   - Class: mult_question
110     Output: Which of the following does Var(1/n * Sum(X_i)) equal?
111     AnswerChoices: 1/n^2*Var(Sum(X_i)); 1/n^2*E(Sum(X_i)); mu/n^2; sigma/n
112     CorrectAnswer: 1/n^2*Var(Sum(X_i))
113     AnswerTests: omnitest(correctVal='1/n^2*Var(Sum(X_i))')
114     Hint: Remember that fact about Var that we said would be useful before? Now is the
        time to use it.
115
116   - Class: mult_question
117     Output: Recall that Var is an expected value and expected values are linear. Also
        recall that our samples X_1, X_2,...,X_n are independent. What does Var(Sum(X_i))
        equal?
118     AnswerChoices: Sum(Var(X_i)); E(Sum(X_i)); E(mu); Var(sigma)
119     CorrectAnswer: Sum(Var(X_i))
120     AnswerTests: omnitest(correctVal='Sum(Var(X_i))')
121     Hint: By linearity, the variance of the sum equals the sum of the variance.
122
123   - Class: mult_question
124     Output: Finally, each X_i comes from a population with variance sigma^2. What does
        Sum(Var(X_i)) equal? As before, Sum is taken over n values.
125     AnswerChoices: n*(sigma)^2; n*mu; n*E(mu); (n^2)*Var(sigma)
126     CorrectAnswer: n*(sigma)^2
127     AnswerTests: omnitest(correctVal='n*(sigma)^2')
128     Hint: Var(X_i) is the constant value sigma^2 and we're summing over n of them.
129
130   - Class: text
131     Output: So we've shown that
        Var(X')=Var(1/n*Sum(X_i))=(1/n^2)*Var(Sum(X_i))=(1/n^2)*Sum(sigma^2)=sigma^2/n for
        infinite populations when our samples are independent.
132
133   - Class: text
134     Output: The standard deviation of a statistic is called its standard error, so the
```

```
              standard error of the sample mean is the square root of its variance.
135
136    - Class: text
137      Output: We just showed that the variance of a sample mean is sigma^2 / n and we
             estimate it with s^2 / n. It follows that its square root, s / sqrt(n), is the
             standard error of the sample mean.
138
139    - Class: text
140      Output: The sample standard deviation, s, tells us how variable the population is,
             and s/sqrt(n), the standard error, tells us how much averages of random samples of
             size n from the population vary. Let's see this with some simulations.
141
142    - Class: cmd_question
143      Output: The R function rnorm(n,mean,sd) generates n independent (hence uncorrelated)
             random normal samples with the specified mean and standard deviation. The defaults
             for the latter are mean 0 and standard deviation 1. Type the expression
             sd(apply(matrix(rnorm(10000),1000),1,mean)) at the prompt.
144      CorrectAnswer: sd(apply(matrix(rnorm(10000),1000),1,mean))
145      AnswerTests: omnitest(correctExpr='sd(apply(matrix(rnorm(10000),1000),1,mean))')
146      Hint: Type 'sd(apply(matrix(rnorm(10000),1000),1,mean))' at the command prompt.
147
148    - Class: cmd_question
149      Output: This returns the standard deviation of 1000 averages, each of a sample of 10
             random normal numbers with mean 0 and standard deviation 1. The theory tells us that
             the standard error, s/sqrt(n), of the sample means indicates how much averages of
             random samples of size n (in this case 10) vary. Now compute 1/sqrt(10) to see if it
             matches the standard deviation we just computed with our simulation.
150      CorrectAnswer: 1/sqrt(10)
151      AnswerTests: omnitest(correctExpr='1/sqrt(10)')
152      Hint: Type '1/sqrt(10)' at the command prompt.
153
154    - Class: mult_question
155      Output: Pretty close, right? Let's try a few more. Standard uniform distributions
             have variance 1/12. The theory tells us the standard error of means of independent
             samples of size n would have which standard error?
156      AnswerChoices: 1/(12*sqrt(n)); 12/sqrt(n); 1/sqrt(12*n); I haven't a clue
157      CorrectAnswer: 1/sqrt(12*n)
158      AnswerTests: omnitest(correctVal='1/sqrt(12*n)')
159      Hint: In this case s is the sqrt(1/12). Divide this by sqrt(n).
160
161    - Class: cmd_question
162      Output:  Compute 1/sqrt(120). This would be the standard error of the means of
             uniform samples of size 10.
163      CorrectAnswer: 1/sqrt(120)
164      AnswerTests: omnitest(correctExpr='1/sqrt(120)')
165      Hint: Type '1/sqrt(120)' at the command prompt.
166
167    - Class: cmd_question
168      Output: Now check it as we did before. Use the expression
             sd(apply(matrix(runif(10000),1000),1,mean)).
169      CorrectAnswer: sd(apply(matrix(runif(10000),1000),1,mean))
170      AnswerTests: omnitest(correctExpr='sd(apply(matrix(runif(10000),1000),1,mean))')
171      Hint: Type 'sd(apply(matrix(runif(10000),1000),1,mean))' at the command prompt.
172
173    - Class: mult_question
174      Output: Pretty close again, right? Poisson(4) are distributions with variance 4; what
             standard error would means of random samples of n Poisson(4) have?
175      AnswerChoices: 2/sqrt(n); 1/sqrt(2*n); 2*sqrt(n); I haven't a clue
176      CorrectAnswer:  2/sqrt(n)
177      AnswerTests: omnitest(correctVal='2/sqrt(n)')
178      Hint: In this case s is 2. Divide this by sqrt(n).
179
180    - Class: cmd_question
181      Output:  We'll do another simulation to test the theory. First, assume you're taking
             averages of 10 Poisson(4) samples and compute the standard error of these means. Use
             the formula you just chose.
182      CorrectAnswer: 2/sqrt(10)
183      AnswerTests: omnitest(correctExpr='2/sqrt(10)')
184      Hint: Type '2/sqrt(10)' at the command prompt.
```

```
185
186    - Class: cmd_question
187      Output: Now check it as we did before. Use the expression
             sd(apply(matrix(rpois(10000,4),1000),1,mean)).
188      CorrectAnswer: sd(apply(matrix(rpois(10000,4),1000),1,mean))
189      AnswerTests: omnitest(correctExpr='sd(apply(matrix(rpois(10000,4),1000),1,mean))')
190      Hint: Type 'sd(apply(matrix(rpois(10000,4),1000),1,mean))' at the command prompt.
191
192    - Class: mult_question
193      Output: Like magic, right? One final test. Fair coin flips have variance 0.25; means
             of random samples of n coin flips have  what standard error?
194      AnswerChoices: 2/sqrt(n); 1/sqrt(2*n); 2*sqrt(n); 1/(2*sqrt(n)); I haven't a clue
195      CorrectAnswer:  1/(2*sqrt(n))
196      AnswerTests: omnitest(correctVal='1/(2*sqrt(n))')
197      Hint: In this case s is 1/2 which is the sqrt of 1/4, the variance. Divide this by
             sqrt(n).
198
199    - Class: cmd_question
200      Output:  You know the drill. Assume you're taking averages of 10 coin flips and
             compute the standard error of these means with the theoretical formula you just picked.
201      CorrectAnswer: 1/(2*sqrt(10))
202      AnswerTests: omnitest(correctExpr=' 1/(2*sqrt(10))')
203      Hint: Type ' 1/(2*sqrt(10))' at the command prompt.
204
205    - Class: cmd_question
206      Output: Now check it as we did before. Use the expression
             sd(apply(matrix(sample(0:1,10000,TRUE),1000),1,mean)).
207      CorrectAnswer: sd(apply(matrix(sample(0:1,10000,TRUE),1000),1,mean))
208      AnswerTests:
             omnitest(correctExpr='sd(apply(matrix(sample(0:1,10000,TRUE),1000),1,mean))')
209      Hint: Type 'sd(apply(matrix(sample(0:1,10000,TRUE),1000),1,mean))' at the command
             prompt.
210
211    - Class: text
212      Output: Finally, here's something interesting. Chebyshev's inequality helps interpret
             variances. It states that the probability that a random variable X is at least k
             standard deviations from its mean is less than 1/(k^2). In other words, the
             probability that X is at least 2 standard deviations from the mean is less than 1/4,
             3 standard deviations 1/9, 4 standard deviations 1/16, etc.
213
214    - Class: text
215      Output: However this estimate is quite conservative for random variables that are
             normally distributed, that is, with bell-curve distributions. In these cases, the
             probability of being at least 2 standard deviations from the mean is about 5% (as
             compared to Chebyshev's upper bound of 25%) and the probability of being at least 3
             standard deviations from the mean is roughly .2%.
216
217    - Class: mult_question
218      Output: Suppose you had a measurement that was 4 standard deviations from the
             distribution's mean. What would be the upper bound  of the probability of this
             happening using Chebyshev's inequality?
219      AnswerChoices: 6%; 0%; 11%; 25%; 96%
220      CorrectAnswer: 6%
221      AnswerTests: omnitest(correctVal='6%')
222      Hint: Chebyshev's inequality estimates that probability as 1/16. Convert this to a
             probability.
223
224
225    - Class: mult_question
226      Output: Now to review. The sample variance estimates what?
227      AnswerChoices: population variance; sample mean; sample standard deviation; population
228      CorrectAnswer: population variance
229      AnswerTests: omnitest(correctVal='population variance')
230      Hint: Which choice has the word variance in it?
231
232    - Class: mult_question
233      Output: The distribution of the sample variance is centered at what?
234      AnswerChoices: population variance; sample mean; sample standard deviation; population
235      CorrectAnswer: population variance
```

```
236        AnswerTests: omnitest(correctVal='population variance')
237        Hint: What is the sample variance estimating?
238
239    - Class: mult_question
240      Output: True or False - The sample variance gets more concentrated around the
           population variance with larger sample sizes
241      AnswerChoices: True; False
242      CorrectAnswer: True
243      AnswerTests: omnitest(correctVal='True')
244      Hint: Is more data better than less data?
245
246    - Class: mult_question
247      Output: The variance of the sample mean is the population variance divided by ?
248      AnswerChoices: n; n^2; sqrt(n); I haven't a clue
249      CorrectAnswer: n
250      AnswerTests: omnitest(correctVal='n')
251      Hint: Remember the 4 step proof starting with Var(X')=...? The last step had an n
           divided by an n^2.
252
253    - Class: mult_question
254      Output: The standard error of the sample mean is the sample standard deviation s
           divided by ?
255      AnswerChoices: n; n^2; sqrt(n); I haven't a clue
256      CorrectAnswer: sqrt(n)
257      AnswerTests: omnitest(correctVal='sqrt(n)')
258      Hint: Remember the many many examples we went through. The sqrt(n) figured
           prominently in them.
259
260    - Class: text
261      Output: Congrats! You've concluded this vary long lesson on variance. We hope you
           liked it vary much.
262
263    - Class: mult_question
264      Output: "Would you like to receive credit for completing this course on
265        Coursera.org?"
266      CorrectAnswer: NULL
267      AnswerChoices: Yes;No
268      AnswerTests: coursera_on_demand()
269      Hint: ""
270
```