

```

1   Course: Regression_Models
2   Lesson: Count_Outcomes
3
4   - Class: text
5   Output: "Count Outcomes. (Slides for this and other Data Science courses may be found
6   at github https://github.com/DataScienceSpecialization/courses. If you care to use
7   them, they must be downloaded as a zip file and viewed locally. This lesson
8   corresponds to Regression_Models/03_03_countOutcomes.)"
9
10  - Class: text
11  Output: "Many data take the form of counts. These might be calls to a call center,
12  number of flu cases in an area, or number of cars that cross a bridge. Data may also
13  be in the form of rates, e.g., percent of children passing a test. In this lesson we
14  will use Poisson regression to analyze daily visits to a web site as the web site's
15  popularity grows, and to analyze the percent of visits which are due to references
16  from a different site."
17
18  - Class: text
19  Output: "Visits to a web site tend to occur independently, one at a time, at a
20  certain average rate. The Poisson distribution describes random processes of this
21  type. A Poisson process is characterized by a single parameter, the expected rate of
22  occurrence, which is usually called lambda. In our case, lambda will be expected
23  visits per day. Of course, as the web site becomes more popular, lambda will grow. In
24  other words, our lambda will depend on time. We will use Poisson regression to model
25  this dependence."
26
27  - Class: cmd_question
28  Output: "Somewhat remarkably, the variance of a Poisson process has the same value as
29  its mean, lambda. You can quickly illustrate this by generating, say, n=1000 samples
30  from a Poisson process using R's rpois(n, lambda) and calculating the sample
31  variance. For example, type var(rpois(1000, 50)). The sample variance won't be
32  exactly equal to the theoretical value, of course, but it will be fairly close."
33  CorrectAnswer: 'var(rpois(1000, 50))'
34  AnswerTests: expr_uses_func('var');expr_uses_func('rpois')
35  Hint: Try typing var(rpois(1000, 50)).
36
37  - Class: mult_question
38  Output: "A famous theorem implies that properly normalized sums of independent,
39  identically distributed random variables will tend to become normally distributed as
40  the number of samples grows large. What is that theorem?"
41  AnswerChoices: The Central Limit Theorem;The Gauss-Markov BLUE Theorem;The
42  Pythagorean Theorem
43  CorrectAnswer: The Central Limit Theorem
44  AnswerTests: omnitest(correctVal= 'The Central Limit Theorem')
45  Hint: It deals with limits.
46
47  - Class: figure
48  Output: "The counts generated by a Poisson process are, strictly speaking, slightly
49  different than the normalized sums of the Central Limit Theorem. However, the counts
50  in a given period of time will represent sums of larger numbers of terms as lambda
51  increases. In fact, it can be formally shown that for large lambda a Poisson
52  distribution is well approximated by a normal. The figure illustrates this effect. It
53  shows progression from a sparse, asymmetric, Poisson probability mass function on the
54  left, to a dense, bell-shaped curve on the right as lambda varies from 2 to 100."
55  Figure: poisson_2_normal.R
56  FigureType: new
57
58  - Class: figure
59  Output: "In a Poisson regression, the log of lambda is assumed to be a linear
60  function of the predictors. Since we will try to model the growth of visits to a web
61  site, the log of lambda will be a linear function of the date:  $\log(\lambda) = b_0 + b_1 \cdot \text{date}$ . This implies that the average number of hits per day, lambda, is exponential
62  in the date:  $\lambda = \exp(b_0) \cdot \exp(b_1)^{\text{date}}$ . Exponential growth is also suggested by
63  the smooth, black curve drawn through the data. Thus  $\exp(b_1)$  would represent the
64  percentage by which visits grow per day."
65  Figure: hits.R
66  FigureType: new
67
68  - Class: video

```

```

37 Output: "If you are connected to the internet right now, would you care to visit the
Leek Group website?"
38 VideoLink: 'http://biostat.jhsph.edu/~jleek/'
39
40 - Class: cmd_question
41 Output: "Our data is in a data frame named hits. Use View(hits), head(hits), or
tail(hits) to examine the data now."
42 CorrectAnswer: 'View(hits)'
43 AnswerTests: ANY_of_exprs('View(hits)', 'head(hits)', 'tail(hits)')
44 Hint: "Type View(hits), head(hits), or tail(hits). Any of these will do."
45
46 - Class: text
47 Output: "There are three columns of data labeled date, visits, and simplystats
respectively. The simplystats column records the number of visits which are due to
references from another site, the Simply Statistics blog. We'll come back to that
column later. For now, we are interested in the date and visits columns. The date
will be our predictor."
48
49 - Class: cmd_question
50 Output: "Our dates are represented in terms of R's class, Date. Verify this by typing
class(hits[, 'date']), or something equivalent."
51 CorrectAnswer: class(hits[, 'date'])
52 AnswerTests: ANY_of_exprs("class(hits[, 'date'])", 'class(hits[, "date"])',
'class(hits[, 1])', 'class(hits$date)')
53 Hint: Type class(hits[, 'date']), or something equivalent.
54
55 - Class: cmd_question
56 Output: "R's Date class represents dates as days since or prior to January 1, 1970.
They are essentially numbers, and to some extent can be treated as such. Dates can,
for example, be added or subtracted, or easily converted to numbers. Type
as.integer(head(hits[, 'date'])) to see what I mean."
57 CorrectAnswer: class(hits[, 'date'])
58 AnswerTests: ANY_of_exprs("as.integer(head(hits[, 'date']))",
'as.integer(head(hits[, "date"]))', 'as.integer(head(hits[, 1]))',
'as.integer(head(hits$date))')
59 Hint: Type as.integer(head(hits[, 'date']), or something equivalent.
60
61 - Class: cmd_question
62 Output: "The arithmetic properties of Dates allow us to use them as predictors. We'll
use Poisson regression to predict log(lambda) as a linear function of date in a way
which maximizes the likelihood of the counts we actually see. Our formula will be
visits ~ date. Since our outcomes (visits) are counts, our family will be 'poisson',
and our third argument will be the data, hits. Create such a model and store it in a
variable called mdl using the following expression or something equivalent, mdl <-
glm(visits ~ date, poisson, hits)."
63 CorrectAnswer: mdl <- glm(visits ~ date, poisson, hits)
64 AnswerTests: creates_glm_model('mdl <- glm(visits ~ date, poisson, hits)')
65 Hint: Type mdl <- glm(visits ~ date, poisson, hits) or something equivalent.
66
67 - Class: figure
68 Output: "The figure suggests that our Poisson regression fits the data very well. The
black line is the estimated lambda, or mean number of visits per day. We see that
mean visits per day increased from around 5 in early 2011 to around 10 by 2012, and
to around 20 by late 2013. It is approximately doubling every year."
69 Figure: model_1.R
70 FigureType: new
71
72 - Class: cmd_question
73 Output: "Type summary(mdl) to examine the estimated coefficients and their
significance."
74 CorrectAnswer: summary(mdl)
75 AnswerTests: omnitest('summary(mdl)')
76 Hint: Just type summary(mdl)
77
78 - Class: text
79 Output: "Both coefficients are significant, being far more than two standard errors
from zero. The Residual deviance is also very significantly less than the Null,
indicating a strong effect. (Recall that the difference between Null and Residual
deviance is approximately chi-square with 1 degree of freedom.) The Intercept

```

coefficient, b0, just represents log average hits on R's Date 0, namely January 1, 1970. We will ignore it and focus on the coefficient of date, b1, since $\exp(b1)$ will estimate the percentage at which average visits increase per day of the site's life."

```
80
81 - Class: cmd_question
82 Output: "Get the 95% confidence interval for exp(b1) by exponentiating confint(mdl,
83 'date')"
84 CorrectAnswer: exp(confint(mdl, 'date'))
85 AnswerTests: ANY_of_exprs("exp(confint(mdl, 'date'))", 'exp(confint(mdl, "date"))',
86 "exp(confint(mdl, 2))")
87 Hint: Just type exp(confint(mdl, 'date')) or exp(confint(mdl, 2)).
88
89 - Class: text
90 Output: "Visits are estimated to increase by a factor of between 1.002192 and
91 1.002399 per day. That is, between 0.2192% and 0.2399% per day. This actually
92 represents more than a doubling every year."
93
94 - Class: figure
95 Output: "Our model looks like a pretty good description of the data, but no model is
96 perfect and we can often learn about a data generation process by looking for a
97 model's shortcomings. As shown in the figure, one thing about our model is 'zero
98 inflation' in the first two weeks of January 2011, before the site had any visits.
99 The model systematically overestimates the number of visits during this time. A less
100 obvious thing is that the standard deviation of the data may be increasing with
101 lambda faster than a Poisson model allows. This possibility can be seen in the
102 rightmost plot by visually comparing the spread of green dots with the standard
103 deviation predicted by the model (black dashes.) Also, there are four or five bursts
104 of popularity during which the number of visits far exceeds two standard deviations
105 over average. Perhaps these are due to mentions on another site."
106 Figure: shortcomings.R
107 FigureType: new
108
109 - Class: figure
110 Output: "It seems that at least some of them are. The simplystats column of our data
111 records the number of visits to the Leek Group site which come from the related site,
112 Simply Statistics. (I.e., visits due to clicks on a link to the Leek Group which
113 appeared in a Simply Statistics post.) "
114 Figure: bursts.R
115 FigureType: new
116
117 - Class: cmd_question
118 Output: "In the figure, the maximum number of visits occurred in late 2012. Visits
119 from the Simply Statistics blog were also at their maximum that day. To find the
120 exact date we can use which.max(hits[, 'visits']). Do this now."
121 CorrectAnswer: which.max(hits[, 'visits'])
122 AnswerTests: omnitest("which.max(hits[, 'visits'])", 704)
123 Hint: Type which.max(hits[, 'visits']) or something equivalent.
124
125 - Class: cmd_question
126 Output: "The maximum number of visits is recorded in row 704 of our data frame. Print
127 that row by typing hits[704,]."
128 CorrectAnswer: hits[704,]
129 AnswerTests: omnitest(correctExpr='hits[704,]')
130 Hint: Just type hits[704,].
131
132 - Class: cmd_question
133 Output: "The maximum number of visits, 94, occurred on December 4, 2012, of which 64
134 came from the Simply Statistics blog. We might consider the 64 visits to be a special
135 event, over and above normal. Can the difference, 94-64=30 visits, be attributed to
136 normal traffic as estimated by our model? To check, we will need the value of lambda
137 on December 4, 2012. This will be entry 704 of the fitted.values element of our
138 model. Extract mdl$fitted.values[704] and store it in a variable named lambda."
139 CorrectAnswer: lambda <- mdl$fitted.values[704]
140 AnswerTests: omnitest(correctExpr='lambda <-
141 mdl$fitted.values[704]');expr_creates_var('lambda')
142 Hint: Just type lambda <- mdl$fitted.values[704].
143
144 - Class: cmd_question
145 Output: "The number of visits explained by our model on December 4, 2012 are those of
```

```

a Poisson random variable with mean lambda. We can find the 95th percentile of this
distribution using qpois(.95, lambda). Try this now."
120 CorrectAnswer: qpois(.95, lambda)
121 AnswerTests: ANY_of_exprs('qpois(.95, lambda)', 'qpois(0.95, lambda)')
122 Hint: Type qpois(.95, lambda) or qpois(0.95, lambda).
123
124 - Class: text
125 Output: "So, 95% of the time we would see 33 or fewer visits, hence 30 visits would
not be rare according to our model. It would seem that on December 4, 2012, the very
high number of visits was due to references from Simply Statistics. To gauge the
importance of references from Simply Statistics we may wish to model the proportion
of traffic such references represent. Doing so will also illustrate the use of glm's
parameter, offset, to model frequencies and proportions."
126
127 - Class: text
128 Output: "A Poisson process generates counts, and counts are whole numbers, 0, 1, 2,
3, etc. A proportion is a fraction. So how can a Poisson process model a proportion?
The trick is to include the denominator of the fraction, or more precisely its log,
as an offset. Recall that in our data set, 'simplystats' is the visits from Simply
Statistics, and 'visits' is the total number of visits. We would like to model the
fraction simplystats/visits, but to avoid division by zero we'll actually use
simplystats/(visits+1). A Poisson model assumes that log(lambda) is a linear
combination of predictors. Suppose we assume that log(lambda) = log(visits+1) + b0 +
b1*date. In other words, if we insist that the coefficient of log(visits+1) be equal
to 1, we are predicting the log of mean visits from Simply Statistics as a proportion
of total visits: log(lambda/(visits+1)) = b0 + b1*date."
129
130 - Class: cmd_question
131 Output: "glm's parameter, offset, has precisely this effect. It fixes the coefficient
of the offset to 1. To create a model for the proportion of visits from Simply
Statistics, we let offset=log(visits+1). Create such a Poisson model now and store it
as a variable called mdl2."
132 CorrectAnswer: mdl2 <- glm(simplestats ~ date, poisson, hits, offset=log(visits+1))
133 AnswerTests: creates_glm_model('mdl2 <- glm(simplestats ~ date, poisson, hits,
offset=log(visits+1))')
134 Hint: "Enter mdl2 <- glm(formula = simplestats ~ date, family = poisson, data = hits,
offset = log(visits + 1)), or something equivalent."
135
136 - Class: cmd_question
137 Output: "Although summary(mdl2) will show that the estimated coefficients are
significantly different than zero, the model is actually not impressive. We can
illustrate why by looking at December 4, 2012, once again. On that day there were 64
actual visits from Simply Statistics. However, according to mdl2, 64 visits would be
extremely unlikely. You can verify this weakness in the model by finding mdl2's 95th
percentile for that day. Recalling that December 4, 2012 was sample 704, find
qpois(.95, mdl2$fitted.values[704])."
138 CorrectAnswer: qpois(.95, mdl2$fitted.values[704])
139 AnswerTests: ANY_of_exprs('qpois(.95, mdl2$fitted.values[704])', 'qpois(0.95,
mdl2$fitted.values[704])')
140 Hint: Just type qpois(.95, mdl2$fitted.values[704]).
141
142 - Class: mult_question
143 Output: "A Poisson distribution with lambda=1000 will be well approximated by a
normal distribution. What will be the variance of that normal distribution?"
144 AnswerChoices: lambda;lambda squared;the square root of lambda.
145 CorrectAnswer: lambda
146 AnswerTests: omnitest(correctVal= 'lambda')
147 Hint: The mean and the variance of a Poisson distribution are equal.
148
149 - Class: mult_question
150 Output: "When modeling count outcomes as a Poisson process, what is modeled as a
linear combination of the predictors?"
151 AnswerChoices: The log of the mean;The mean;The counts
152 CorrectAnswer: The log of the mean
153 AnswerTests: omnitest(correctVal= 'The log of the mean')
154 Hint: Count outcomes and their means are never negative, but linear combinations of
predictors may be.
155
156 - Class: mult_question

```

```
157 Output: "What parameter of the glm function allows you to include a predictor whose
158 coefficient is fixed to the value 1?"
159 AnswerChoices: offset;data;b0;family;formula
159 CorrectAnswer: offset
160 AnswerTests: omnitest(correctVal= 'offset')
161 Hint: We used this parameter to include log(visits+1) when computing mdl2.
162
163 - Class: text
164 Output: "That completes the Poisson GLM example. Thanks for sticking with it. I hope
165 we've made it count."
166
166 - Class: mult_question
167 Output: "Would you like to receive credit for completing this course on
168 Coursera.org?"
169 CorrectAnswer: NULL
170 AnswerChoices: Yes;No
171 AnswerTests: coursera_on_demand()
172 Hint: ""
173
```