

SMI Assignment 1

Andrew Martin

August 10, 2017

Question 1:
 $(Y_i)_{i=1}^n$ are i.i.d $N(\mu, \sigma^2)$ r.vs with sample mean \bar{Y} .

a) Find $E[\bar{Y}^2]$

$$\begin{aligned} \text{var}(\bar{Y}) &= E[\bar{Y}^2] - E[\bar{Y}]^2 \\ \implies E[\bar{Y}^2] &= \text{var}(\bar{Y}) + E[\bar{Y}]^2 \\ &= \text{var}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) + E\left[\frac{1}{n} \sum_{i=1}^n Y_i\right]^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{var}(Y_i) + \frac{1}{n^2} E\left[\sum_{i=1}^n Y_i\right]^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 + \frac{1}{n^2} \left(\sum_{i=1}^n E[Y_i]\right)^2 \\ &= \frac{n}{n^2} \sigma^2 + \frac{1}{n^2} (\sum_{i=1}^n \mu)^2 \\ &= \frac{\sigma^2}{n} + \mu^2 \end{aligned}$$

b) For each i with $1 \leq i \leq n$, prove \bar{Y} and $Y_i - \bar{Y}$ are uncorrelated

$$\begin{aligned} \text{cov}(\bar{Y}, Y_i - \bar{Y}) &= E[(\bar{Y} - E[\bar{Y}])(Y_i - \bar{Y} - E[Y_i - \bar{Y}])] \\ &= E[(\bar{Y} - E[\bar{Y}])(Y_i - \bar{Y} - E[Y_i] + E[\bar{Y}])] \end{aligned}$$

From above $E[\bar{Y}] = \mu$, so:

$$\begin{aligned} &= E[(\bar{Y} - \mu)((Y_i - \bar{Y}) - \mu + \mu)] \\ &= E[(\bar{Y} - \mu)(Y_i - \bar{Y})] \\ &= E[\bar{Y}Y_i - Y_i\mu + \mu\bar{Y} - \bar{Y}^2] \\ &= E[\bar{Y}Y_i] - E[Y_i\mu] + E[\mu\bar{Y}] - E[\bar{Y}^2] \\ &= E[\bar{Y}Y_i] - \mu E[Y_i] + \mu E[\bar{Y}] - \frac{\sigma^2}{n} - \mu^2 \end{aligned}$$

But From Tute 1: $E[\bar{Y}Y_i] = \frac{-\sigma^2}{n} + \mu^2$, so:

$$\begin{aligned} &= \frac{-\sigma^2}{n} + \mu^2 - \mu E[Y_i] + \mu E[\bar{Y}] - \frac{\sigma^2}{n} - \mu^2 \\ &= -\mu\mu + \mu\mu = 0 \end{aligned}$$

Therefore they are uncorrelated.

Question 2:
 $(Z_i)_{i=1}^p$ are i.i.d $N(0, 1)$ random variables with

$$X = \sum_{i=1}^p Z_i^2$$

a) Find the moment generating function and distribution of X , assuming $M_{Z^2}(t) = (1 - 2t)^{-1/2}$, where $Z \sim N(0, 1)$

$$M_X(t) = E[e^{tX}] = E\left[e^{t \sum_{i=1}^p Z_i^2}\right]$$

Note that

$$M_{Z^2}(t) = E[e^{tZ^2}] = E[e^{te^{Z^2}}] = (1 - 2t)^{-1/2}$$

So:

$$\begin{aligned} M_X(t) &= E\left[e^{te^{\sum_{i=1}^p Z_i^2}}\right] \\ &= E\left[e^t \prod_{i=1}^p e^{Z_i^2}\right] \end{aligned}$$

Since Z_i are i.i.d

$$\begin{aligned} &= \prod_{i=1}^p E\left[e^{tZ_i^2}\right] \\ &= \prod_{i=1}^p M_{Z^2}(t) \\ &= \prod_{i=1}^p (1 - 2t)^{-1/2} \end{aligned}$$

b) $(Y_i)_{i=1}^p$ are independent normal random variables with different means and variances
i.e.

$$Y_i \sim N(\mu_i, \sigma_i^2), i = 1, \dots, p$$

$$\text{Show that } W = \sum_{i=1}^p \frac{(Y_i - \mu_i)^2}{\sigma_i^2} \sim \chi_p^2$$

where χ_p^2 denotes the chi-squared dist with p degrees of freedom.

$$W = \sum_{i=1}^p \frac{(Y_i - \mu_i)^2}{\sigma_i^2}$$

Note that $\frac{(Y_i - \mu_i)}{\sigma_i} \sim Z \sim N(0, 1)$ So from this:

$$W = \sum_{i=1}^p \frac{(Y_i - \mu_i)^2}{\sigma_i^2} = \sum_{i=1}^p Z_i^2$$

For which, the MGF is as above: $= \prod_{i=1}^p (1 - 2t)^{-1/2}$
Which is the form of a χ^2 MGF

Question 3
 Let $X \sim \text{Bin}(n, p)$
 And $\hat{p}^2 = \left(\frac{X}{n}\right)^2$

a) Find $E[\hat{p}^2]$ and state the bias

$$\begin{aligned}
 E[\hat{p}^2] &= E\left[\left(\frac{X}{n}\right)^2\right] \\
 \implies E\left[\left(\frac{X}{n}\right)^2\right] &= \text{var}\left(\frac{X}{n^2}\right) + E\left[\frac{X}{n}\right]^2 \\
 &= \frac{1}{n^2} \text{var}(X) + \frac{1}{n^2} E[X]^2 \\
 &= \frac{1}{n^2} np(1-p) + \frac{1}{n^2} (np)^2 \\
 &= \frac{p(1-p)}{n} + p^2 \\
 &= \frac{p - p^2 + np^2}{n}
 \end{aligned}$$

The bias is:

$$\frac{p - p^2}{n}$$

b) Show that

$$E\left[\frac{\hat{p}(1-\hat{p})}{n-1}\right] = \frac{p(1-p)}{n}$$

Start with the left hand side:

$$\begin{aligned}
 E\left[\frac{\hat{p}(1-\hat{p})}{n-1}\right] &= \frac{1}{n-1} E[\hat{p}(1-\hat{p})] \\
 &= \frac{1}{n-1} E[\hat{p} - \hat{p}^2] \\
 &= \frac{1}{n-1} (E[\hat{p}] - E[\hat{p}^2]) \\
 &= \frac{1}{n-1} \left(E\left[\frac{X}{n}\right] - \left(\frac{p - p^2 + np^2}{n}\right)\right) \\
 &= \frac{1}{n^2 - n} (E[X] - (p - p^2 + np^2)) \\
 &= \frac{1}{n^2 - n} (np - p + p^2 - np^2) \\
 &= \frac{1}{n^2 - n} (p(n-1) + p - np) \\
 &= \frac{n-1}{n^2 - n} (p(1-p)) \\
 &= \frac{p(1-p)}{n}
 \end{aligned}$$

c) Using a) and b) find an unbiased estimator for p^2 .
 If T is an estimator for θ then the bias is: $b_T(\theta) = E[T] - \theta$. For unbiased, set this to zero.
 I.e. In this case, $b_T(p^2) = E[T] - p^2 = 0$
 So aim to find T such that $E[T] = p^2$ By subtracting the expectations in a) and b)

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

So an unbiased estimator, T , for p^2 , would be:

$$T = \hat{p}^2 - \frac{\hat{p}(1-\hat{p})}{n-1}$$

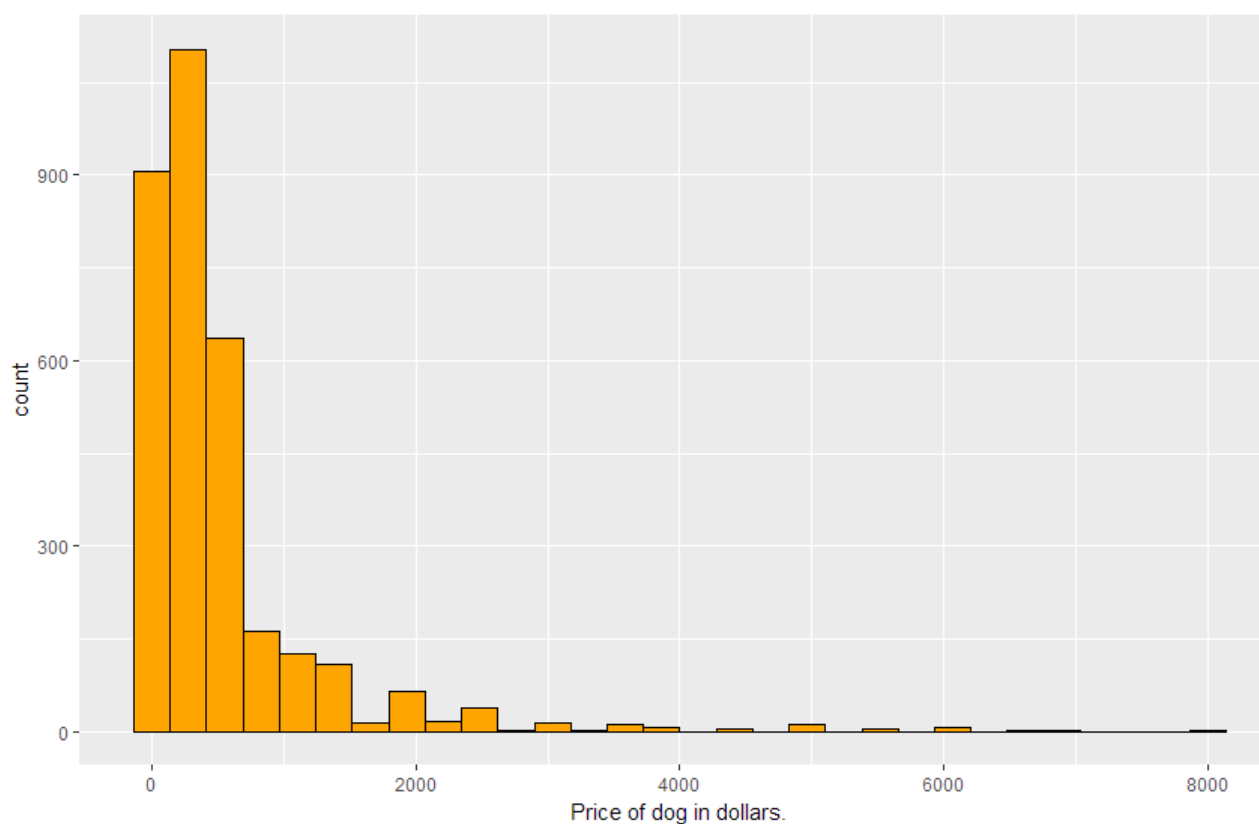
Question 4

- 1) Clean all the variables (include code)
- 2) Produce an appropriate plot for each variable
i.e for categorical - bar chart, for quantitative - histogram.
Label and caption each of these
- 3) For the quantitative variables, identify if they are: unimodal or bimodal, whether it is symmetric, left-skewed or right-skewed. For categorical, identify the most common level

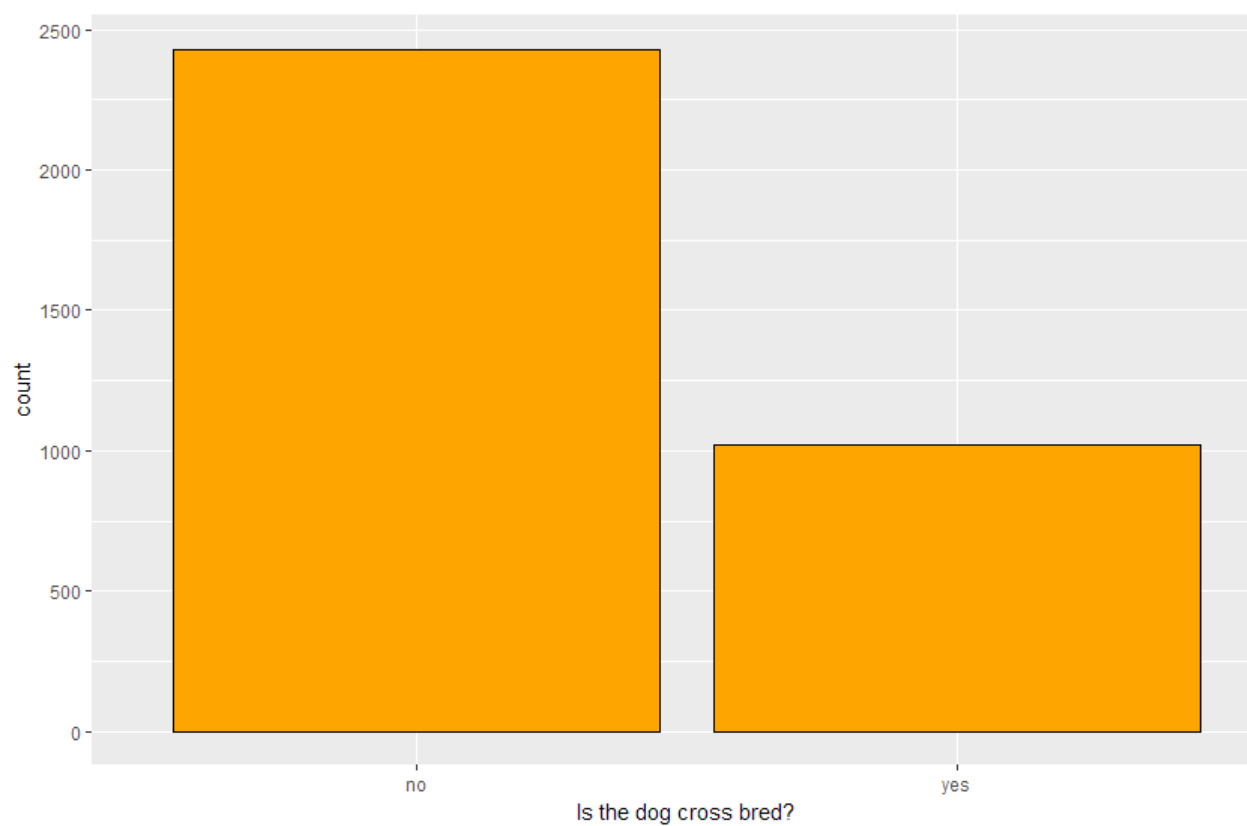
Price
Cross
Pet offered by
Microchip
Vaccination
Desexing status
Relinquished or not

The R code is attached and will be referenced from here, onwards.

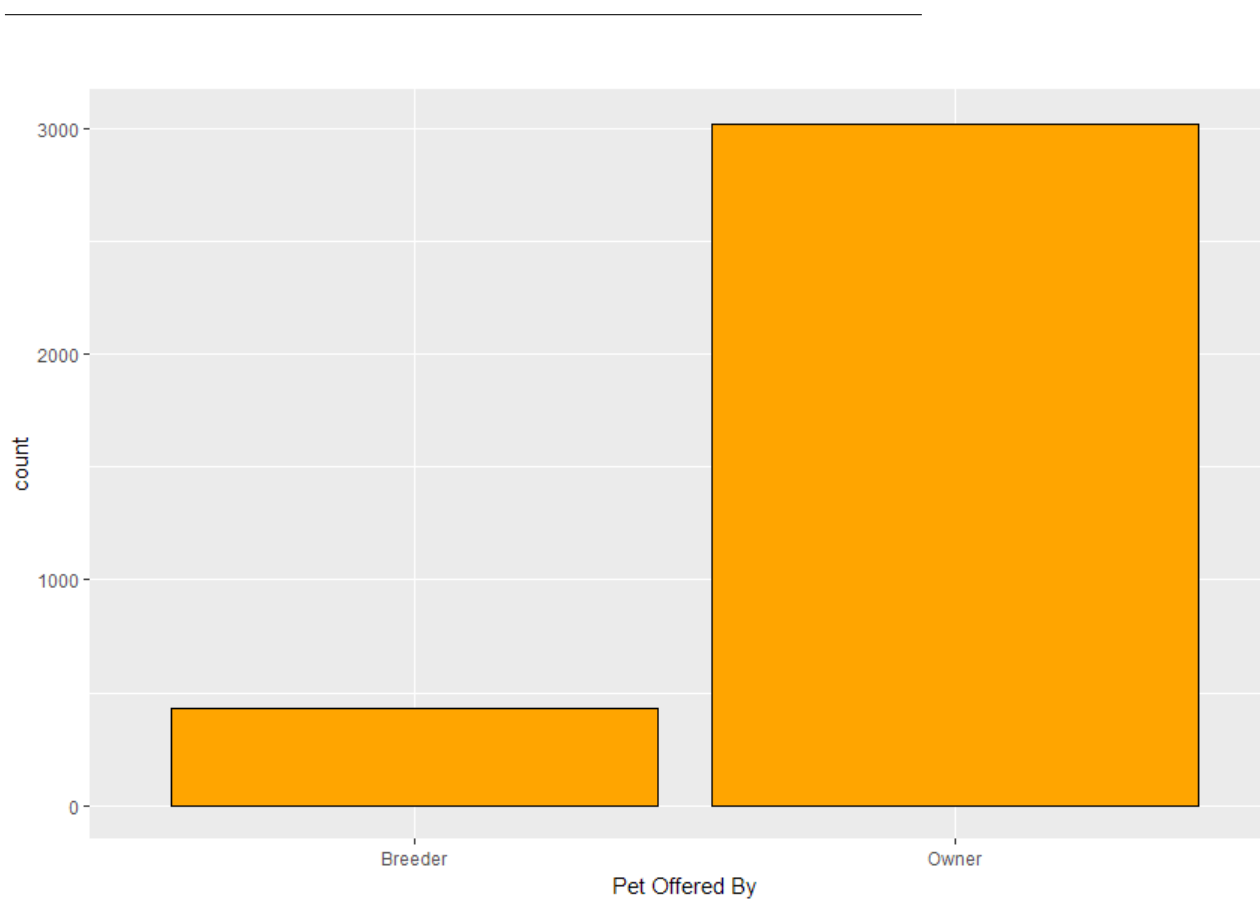
```
1 #R Code for SMI Assignment 1
2 #Andrew Martin
3 #10/08/2017
4 library('readxl')
5 library('tidyverse')
6 library('rmarkdown')
7 library('ggplot2')
8 setwd("F:/Documents/Uni/SMI")
9
10
11 gumtree = read_xlsx("Gumtree_dogs.xlsx")
12
13
14 #####Price ----
15 ##Cleaning -
16 #this should be a number so "NA" chars are removed
17 gumtree$price[gumtree$price=="NA"]=NA
18 class(gumtree$price)
19
20 #class is changed to numeric
21 gumtree$price=as.numeric(gumtree$price)
22
23 #since price is a quantitative, continuous variable, a histogram is appropriate
24
25 ggplot(gumtree,aes(x = price)) +
26   geom_histogram(col = "black", fill = "orange") +
27   labs(x = "Price of dog in dollars.")
28 #This graph is right - skewed, which suggests price is right skewed
29 #it has a single peak which suggests it is unimodal
30
```



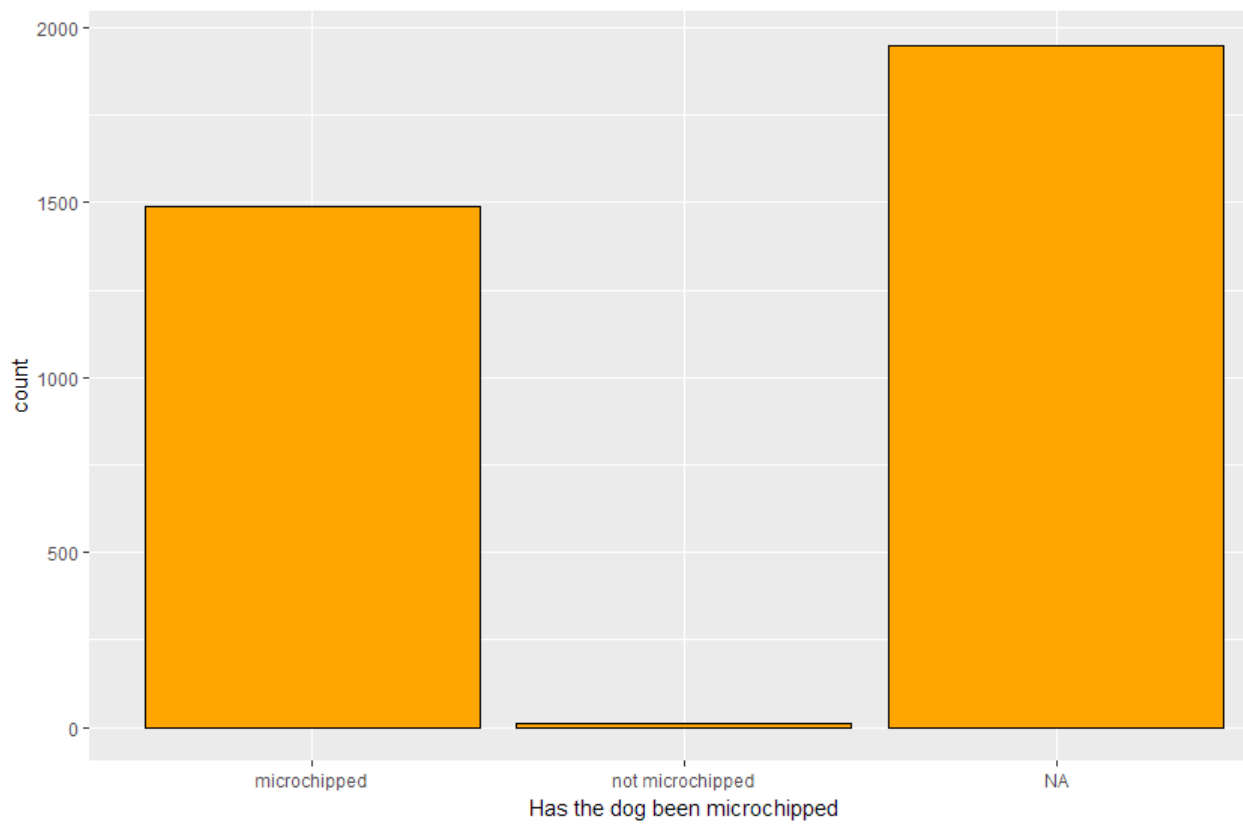
```
33 ##cross ----
34 #told to ignore the "Cross" (uppercase) section and use the "cross" (lowercase)
35 #interesting note:
36 length(which(gumtree$Cross!=gumtree$cross)) #gives 635 -> that means 635 values are different!
37
38 #cross is chr which is good
39 class(gumtree$cross)
40 #nas are changed so R can read them
41 gumtree$cross[gumtree$cross=="NA"]=NA
42
43 #since cross has two states, yes and no, a bar chart is more appropriate
44 ggplot(gumtree,aes(x = cross)) +
45   geom_bar(col = "black", fill = "orange") +
46   labs(x = "Is the dog cross bred?")
47
48 #clearly a greater portion is not cross bred
```



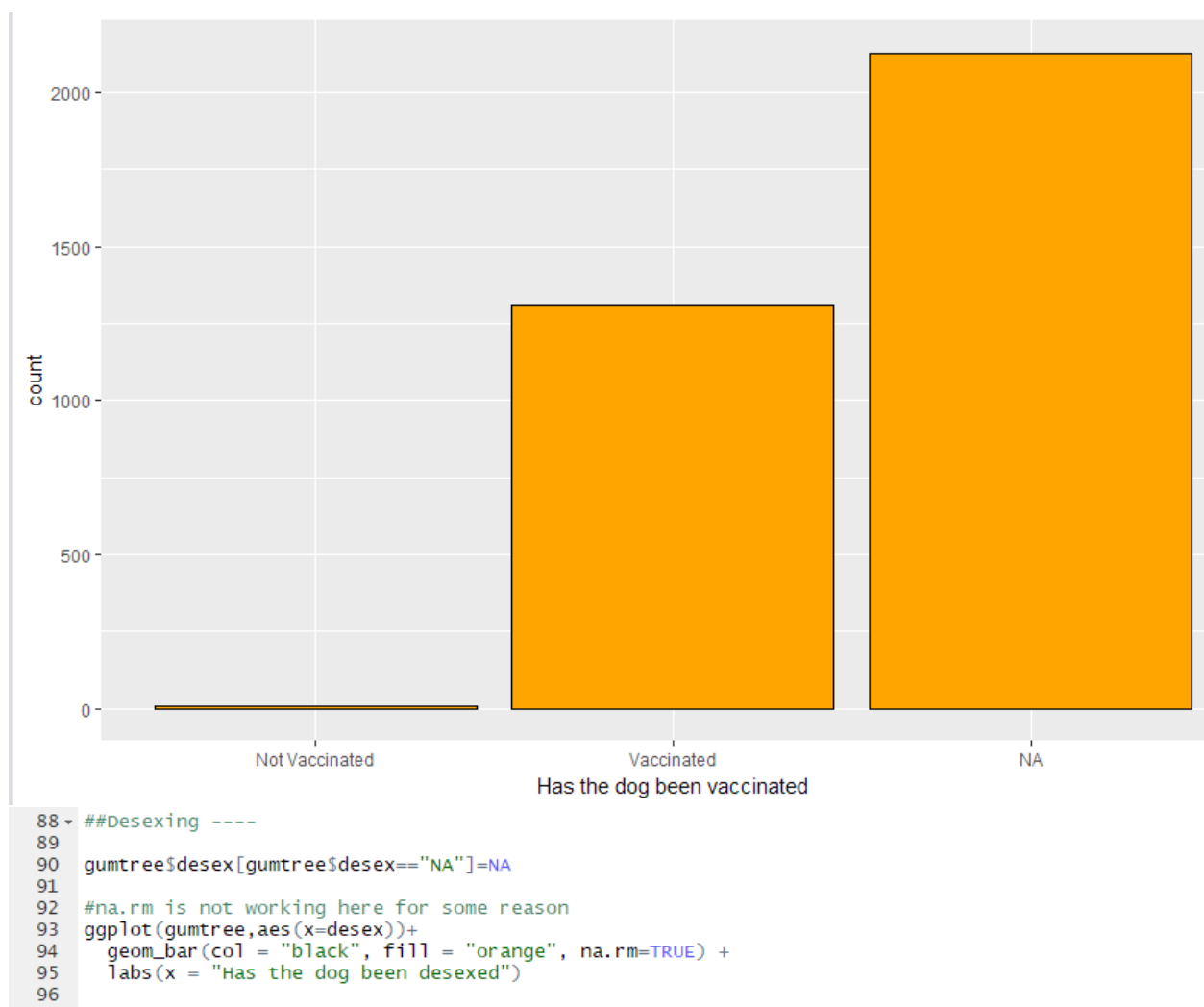
```
53 ▾ ##Pet offered by ----  
54 gumtree$"Pet Offered By:"[gumtree$"Pet Offered By:"=="NA"]=NA  
55  
56  
57 ggplot(gumtree,aes(x = gumtree$'Pet Offered By:')) +  
58   geom_bar(col = "black", fill = "orange") +  
59   labs(x = "Pet Offered By")  
60  
61 #a majority of the pets were offered by an owner.  
62
```

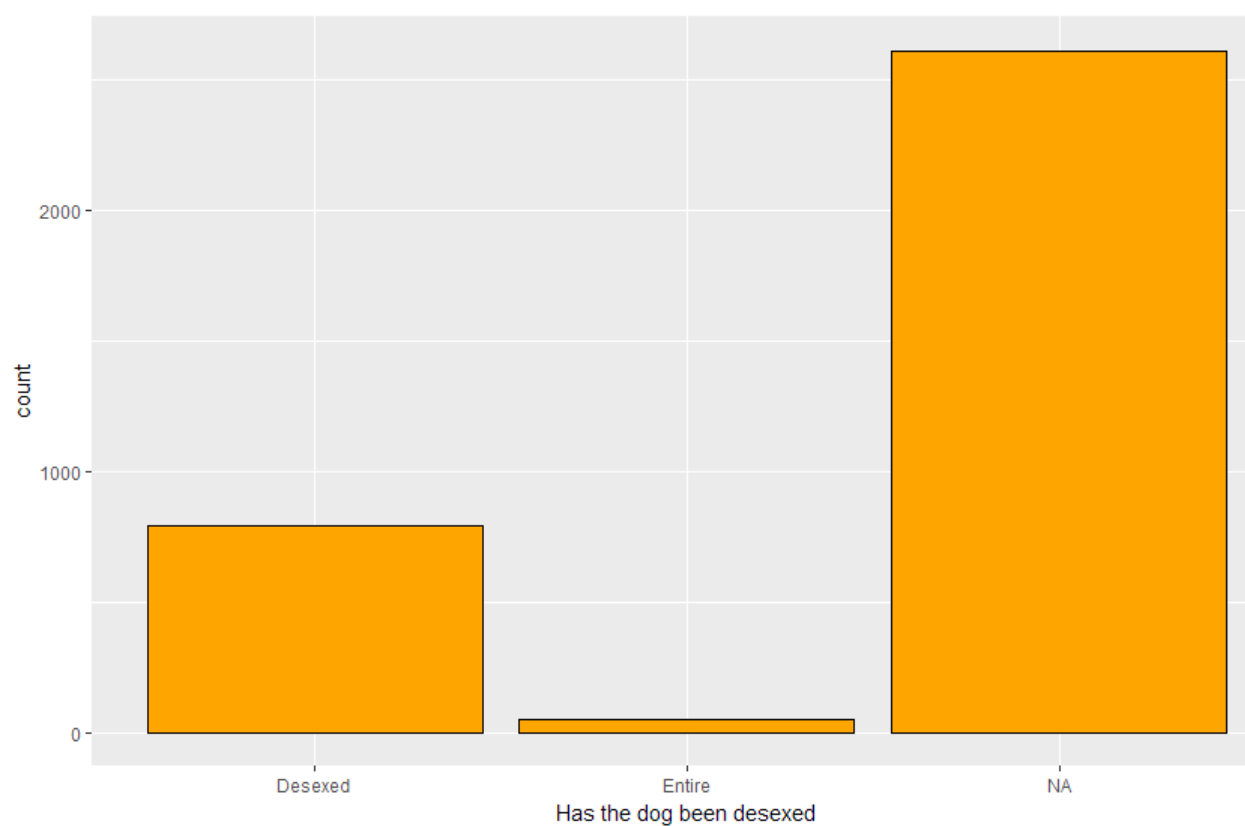



```
63 ▾ ##Microchip ----
64   gumtree$micro[gumtree$micro=="NA"]=NA
65
66   #na.rm is not working here for some reason
67   ggplot(gumtree,aes(x = micro)) +
68     geom_bar(col = "black", fill = "orange", na.rm = TRUE) +
69     labs(x = "Has the dog been microchipped")
70
71
72   #most fields are NA, but few to no dogs are not microchipped
73
```



```
75 #Vaccination ----
76
77
78 gumtree$ vacc[gumtree$ vacc=="NA"]=NA
79
80 #na.rm is not working here for some reason
81 ggplot(gumtree,aes(x=vacc))+
82   geom_bar(col = "black", fill = "orange", na.rm=TRUE) +
83   labs(x = "Has the dog been vaccinated")
84
85 #A very small percentage are labelled as not vaccinated, and most are vaccinated.
86
87
```





```
98 ##Relinquished ----
99 gumtree$relinquished[gumtree$relinquished=="NA"]=NA
100
101 ggplot(gumtree,aes(x = relinquished)) +
102   geom_bar(col = "black", fill = "orange", na.rm=TRUE) +
103   labs(x = "Is the dog being relinquished?")
104 #Most dogs have been relinquished
105
```

