

OkCupid Data for Introductory Statistics and Data Science Courses (Revised)

Albert Y. Kim *

Department of Mathematics
Middlebury College, Middlebury, VT

Adriana Escobedo-Land

Environmental Studies-Biology Program
Reed College, Portland, OR

April 26, 2021

* Address for correspondence: Department of Mathematics, Middlebury College, Warner Hall, 303 College Street, Middlebury, VT 05753. Email: aykim@middlebury.edu.

OkCupid Data for Introductory Statistics and Data Science Courses (Revised)

Abstract

We present a data set consisting of user profile data for 59,946 San Francisco OkCupid users (a free online dating website) from a period in the 2010s. The data set includes typical user information, lifestyle variables, and text responses to 10 essays questions. We present four example analyses suitable for use in undergraduate introductory probability and statistics and data science courses that use R. The statistical and data science concepts covered include basic data visualization, exploratory data analysis, multivariate relationships, text analysis, and logistic regression for prediction.

Keywords: OkCupid, online dating, data science, big data, logistic regression, text mining.

1 Introduction

Given that the field of data science is gaining more prominence in academia and industry, many statisticians are arguing that statistics needs to stake a bigger claim in data science in order to avoid marginalization by other disciplines such as computer science and computer engineering (Davidson, 2013; Yu, 2014). The importance of emphasizing data science concepts in the undergraduate curriculum is stressed in the American Statistical Association’s (ASA) most recent Curriculum Guidelines for Undergraduate Programs in Statistical Science (American Statistical Association Undergraduate Guidelines Workgroup, 2014).

While precise definition of the exact difference between statistics and data science and its implications for statistics education can be debated (Wickham, 2014), one consensus among many in statistics education circles is that at the very least statistics needs to incorporate a heavier computing component and increase the use of technology for both developing conceptual understanding and analyzing data (GAISE College Group, 2005; Nolan and Lang, 2010). Relatedly, in the hopes of making introductory undergraduate statistics courses more relevant, many statistics educators are placing a higher emphasis on the use of real data in the classroom, a practice the ASA’s Guidelines for Assessment and Instruction in Statistics Education (GAISE) project’s report strongly encourages (GAISE College Group, 2005). Of particular importance to the success of such ambitions are the data sets considered, as they provide the context of the analyses and thus will ultimately drive student interest (Gould, 2010).

It is in light of these discussions that we present this paper centering on data from the online dating website OkCupid, specifically a snapshot of San Francisco California users taken during a period in the 2010s. We describe the data set and present a series of example analyses along with corresponding pedagogical discussions. The example analyses presented in this paper were used in a variety of settings at Reed College in Portland, Oregon: a 90 minute introductory tutorial on R, an introductory probability and statistics course, and a follow-up two-hundred level data science course titled “Case Studies in Statistical Analysis.” The statistical and data science concepts covered include basic data visualization, exploratory data analysis, multivariate relationships, text analysis, and logistic regression for prediction. All examples are presented using the R statistical software program and make use of the `mosaic`, `dplyr`, `stringr`, and `ggplot2` packages (Pruim et al., 2014; Wickham, 2009, 2012; Wickham and Francois, 2014).

2 Data

The data consists of the public profiles of 59,946 OkCupid users who were living within 25 miles of San Francisco, had active profiles during a period in the 2010s, were online in the previous year, and had at least one picture in their profile. Using a Python script, data was scraped from users’ public profiles four days later; any non-publicly facing information such as messaging was not accessible.

Variables include typical user information (such as sex, sexual orientation, age, and ethnicity) and lifestyle variables (such as diet, drinking habits, smoking habits). Note that random noise was added to the age variable for de-identification purposes.

We load the `profiles_revised` data as follows:

```
profiles_revised <- read.csv(file="profiles_revised.csv", header=TRUE,
                             stringsAsFactors=FALSE)
n <- nrow(profiles_revised)
```

Furthermore, text responses to the 10 essay questions posed to all OkCupid users are included as well, such as “My Self Summary,” “The first thing people usually notice about me,” and “On a typical Friday night I am...” However, the essay data has been randomized by rows to decouple them from the profiles data. In other words, the user represented in the first row of `profiles_revised` does not necessarily correspond to the user that wrote the responses in the first row of `essays_revised_and_shuffled`.

We load this randomized essays data as follows:

```
essays_revised_and_shuffled <-  
  read.csv(file="essays_revised_and_shuffled.csv", header=TRUE, stringsAsFactors=FALSE)
```

For a complete list of variables and more details, see the accompanying codebook `okcupid_codebook.txt`.

Analyses of similar data has received much press of late, including Amy Webb’s TED talk “How I Hacked Online Dating” (Webb, 2013) and Wired magazine’s “How a Math Genius Hacked OkCupid to Find True Love.” (Poulsen, 2014) OkCupid co-founder Christian Rudder pens periodical analyses of OkCupid data on the blog OkTrends (<http://blog.okcupid.com/>) and has recently published a book “Dataclysm: Who We Are When We Think No One’s Looking” describing similar analyses (Rudder, 2014). Such publicity surrounding data-driven online dating and the salience of dating matters among students makes this data set one with much potential to be of interest to students, hence facilitating the instruction of statistical and data science concepts.

Before we continue we note that even though this data consists of publicly facing material, one should proceed with caution before scraping and using data in fashion similar to ours, as the Computer Fraud and Abuse Act (CFAA) makes it a federal crime to access a computer without authorization from the owner (Penenberg, 2014). In our case, permission to use and disseminate the data was given by its owners (See Acknowledgements).

3 Example Analyses

We present example analyses that address the following questions:

1. How do the heights of male and female OkCupid users compare?
2. What does the San Francisco online dating landscape look like? Or more specifically, what is the relationship between users’ sex and sexual orientation?
3. How accurately can we predict a user’s sex using their listed height?

For each question, we present an exercise as would be given to students in a lab setting, followed by a pedagogical discussion.

3.1 Male and Female Heights

3.1.1 Exercise

We compare the distribution of male and female OkCupid users’ heights. Height is one of 3 numerical variables in this data set (the others being age and income). This provides us an opportunity to investigate numerical summaries using the `favstats()` function from the `mosaic` package:

```
require(mosaic)  
favstats(~height, data=profiles_revised)  
  
##   min Q1 median Q3 max mean sd      n missing  
##    1 66    68 71  95   68  4 59943         3
```

We observe that some of the heights are nonsensical, including heights of 1 inch and 95 inches (equaling 7’11”). We deem heights between 55 and 80 inches to be reasonable and remove the rest. While there is potential bias in discarding users with what we deem non-reasonable heights, since out of the 59946 users there are only 117 who would be discarded, the effect would not be substantial. Therefore we keep only those users with heights between 55 and 80 inches using the `filter()` function from the `dplyr` package:

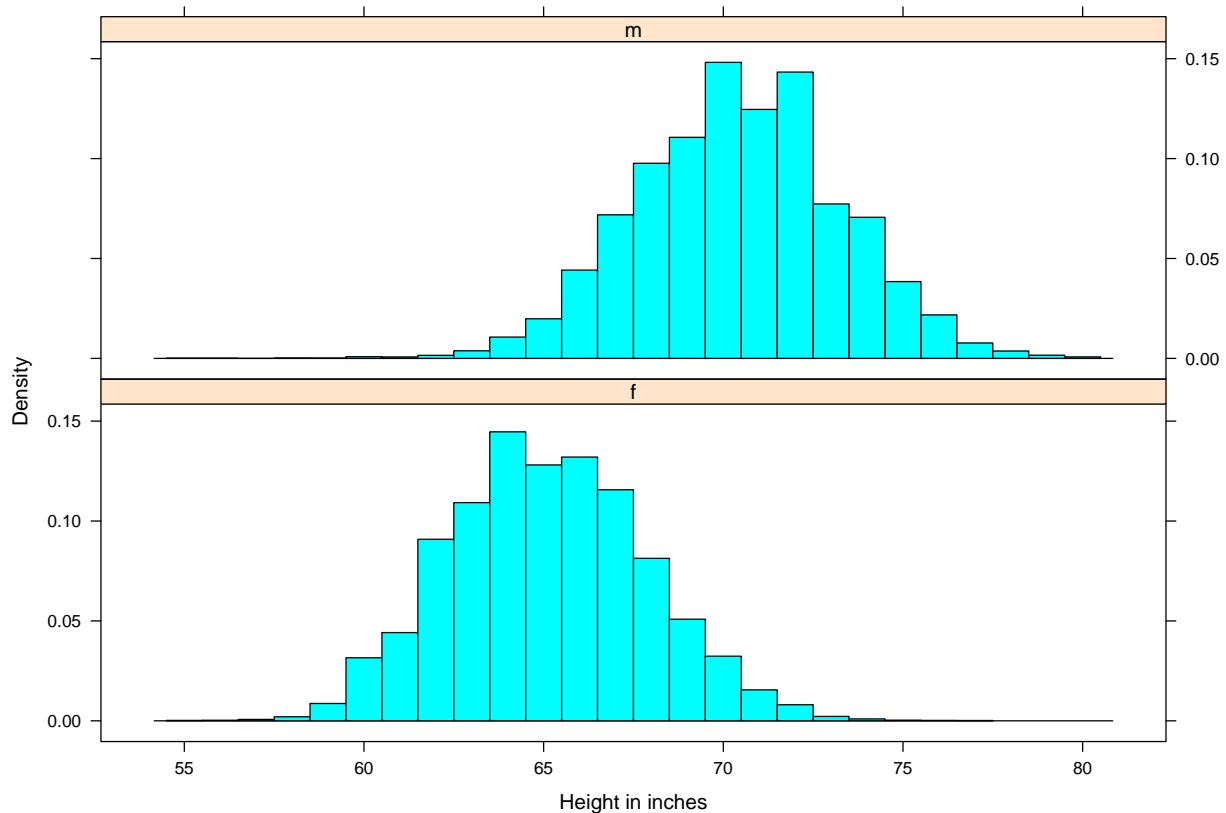


Figure 1: Histograms of user heights split by sex.

```
require(dplyr)
profiles_revised.subset <- filter(profiles_revised, height>=55 & height <=80)
```

We compare the distributions of male and female heights using histograms. While we could plot two separate histograms without regard to the scale of the two axes, in Figure 1 we instead use the `histogram()` function from the `mosaic` package to:

1. Plot heights given sex by defining the formula: `~ height | sex`.
2. Plot them simultaneously in a *lattice* consisting of two rows and one column of plots by setting `layout=c(1,2)`
3. Plot them with bin widths matching the granularity of the observations (inches) by setting `width=1`. The `histogram()` function automatically matches the scales of the axes for both plots.

```
histogram(~height | sex, width=1, layout=c(1,2), xlab="Height in inches",
          data=profiles_revised.subset)
```

3.1.2 Pedagogical Discussion

This first exercise stresses many important considerations students should keep in mind when working with real data. Firstly, it emphasizes the importance of performing an exploratory data analysis to identify anomalous observations and confronts students with the question of what to do with them. For example, while a height of 1 inch is clearly an outlier that needs to be removed, at what point does a height no longer become reasonable and what impact does the removal of unreasonable heights have on the conclusions? In our case, since only a small number of observations are removed, the impact is minimal.

Secondly, this exercise demonstrates the power of simple data visualizations such as histograms to convey insight and hence emphasizes the importance of putting careful thought into their construction. In our case, while having students plot two histograms simultaneously in order to demonstrate that males have on average greater height may seem to be a pedantic goal at first, we encouraged students to take a closer look at the histograms and steered their focus towards the unusual peaks at 72 inches (6 feet) for males and 64 inches (5'4") for females. Many of the students could explain the phenomena of the peak at 72 inches for men: sociological perceptions of the rounded height of 6 feet. On the other hand, consensus was not as strong about perceptions of the height of 5'4" for women. Instructors can then refer students to the entry on OkCupid's blog OkTrends "The Biggest Lies in Online Data" (Rudder, 2010) to show they have replicated (on a smaller scale) a previous analysis and then show other analyses conducted by OkCupid.

Further questions that can be pursued from this exercise include "How can we question if those peaks are significant or due to chance?," "Are we only observing men who are just under 6 feet rounding up, or are men just over 6 feet rounding down as well?," or "How can we compare the distribution of listed heights on OkCupid to the actual San Francisco population's heights?"

3.2 Relationship Between Sex and Sexual Orientation

3.2.1 Exercise

Since among the most important considerations in assessing a potential mate are their sex and sexual orientation, in this exercise we investigate the relationship between these two variables. At the time, OkCupid allowed for two possible sex choices (male or female) and three possible sexual orientation choices (gay, bisexual, or straight)¹. First, we perform a basic exploratory data analysis on these variables using barcharts in Figure 2:

```
par(mfrow=c(1, 2))
barplot(table(profiles_revised$sex)/n, xlab="sex", ylab="proportion")
barplot(table(profiles_revised$orientation)/n, xlab="orientation", ylab="proportion")
```

However, in order to accurately portray the dating landscape we can't just consider the **marginal distributions** of these variables, we must consider their **joint** and **conditional distributions** i.e. the cross-classification of the two variables. We describe the distribution of sexual orientation conditional on sex. For example, we can ask of the female population, what proportion are bisexual? We do this using the `tally()` function from the `mosaic` package and ensure both columns sum to 1 by setting `format='proportion'`. Furthermore, we visualize their joint distribution, as represented by their contingency table, via the `mosaicplot` shown in Figure 3.

```
tally(orientation ~ sex, data=profiles_revised, format='proportion')

##           sex
## orientation  f    m
##   bisexual 0.083 0.022
##    gay     0.066 0.111
##   straight 0.851 0.867
```

¹OkCupid has since relaxed these categorizations to allow for a broader range of choices for both sex and sexual orientation.

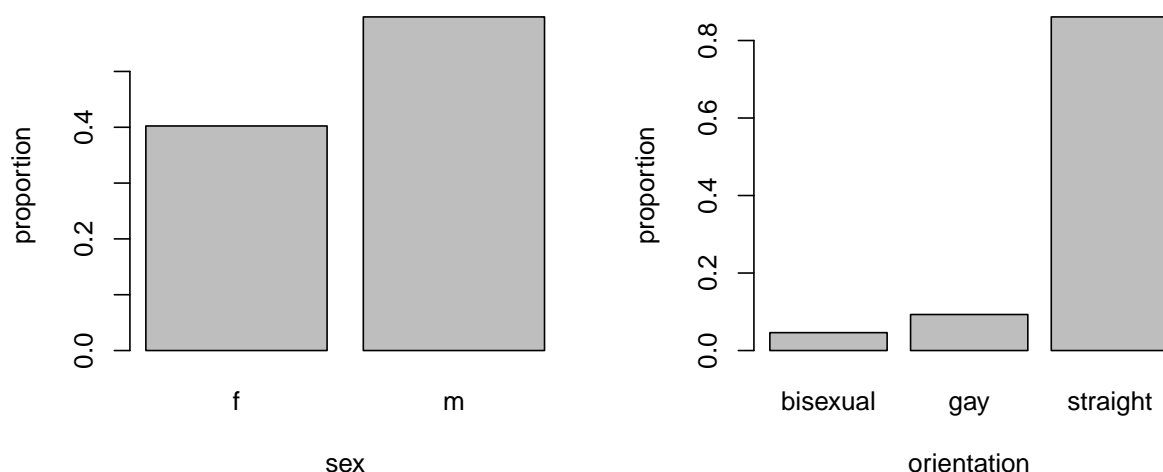


Figure 2: Distributions of sex and sexual orientation.

```
sex.by.orientation <- tally(~sex + orientation, data=profiles_revised)
sex.by.orientation

##      orientation
## sex bisexual   gay straight
## f      1996   1588  20533
## m       771   3985  31073

mosaicplot(sex.by.orientation, main="Sex vs Orientation", las=1)
```

Do these results generalize to the entire San Francisco online dating population?

3.2.2 Pedagogical Discussion

This exercise is an opportunity to reinforce statistical notions such as marginal/joint/conditional distributions and sampling bias. The data indicate that the San Francisco OkCupid dating population skews male and while the proportions of males and females who list themselves as straight are similar, a higher proportion of males list themselves as gay while a higher proportion of females list themselves as bisexual. Many students were not surprised by these facts as they were well aware of the gender imbalance issues in the large technology sector in the San Francisco Bay Area and San Francisco's history of being a bastion for the gay community.

The question of generalizability was presented in an introductory probability and statistics assignment. Almost all students were able to recognize the selection biases of who signs up for this particular site and hence the non-generalizability of the results. For example, some recognized that OkCupid's demographic is most likely different than other dating websites' demographics such as [match.com](https://www.match.com) (which is not free) or [christiansingles.com](https://www.christiansingles.com) (which is targeted towards Christians). So while 59946 users may initially seem like a large sample, we emphasized to students that bigger isn't always better when it comes to obtaining accurate inference. This proved an excellent segue to Kate Crawford of Microsoft Research's YouTube

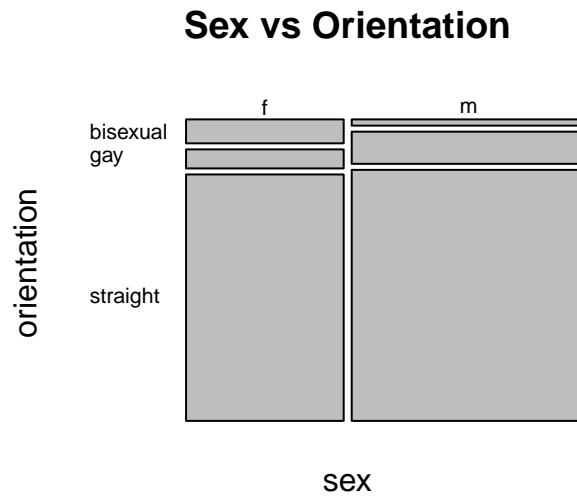


Figure 3: Joint distribution of sex and sexual orientation.

talk “Algorithmic Illusions: Hidden Biases of Big Data” ([Crawford, 2013](#)) where she discusses examples of sampling bias in the era of “Big Data.”

Further questions one can pose to students include “Which dating demographic would you say has it the best and worst in terms of our simplified categorization?” and “What variable do you think should be incorporated next in order to represent the OkCupid dating pool as faithfully as possible?”

3.3 Predictors of Sex

3.3.1 Exercise

This exercise provides an opportunity to fit a predictive model for sex using logistic regression. In order to reinforce the concepts of logistic regression, we keep things simple and consider only one predictor variable in the logistic model: height. We restrict consideration to users whose heights are “reasonable” as defined in Section 3.1 and in order to speed up computation and improve graphical outputs, we only consider a random sample of 5995 users (10% of the data).

However, to ensure the replicability of our results (in other words ensuring the same 5995 users are “randomly” selected each time we run the code), we demonstrate the use of the `set.seed()` function. R’s random number generator is not completely random, but rather is *pseudorandom* in that it generates values that are statistically indistinguishable from a truly random sequence of values, but are generated by a deterministic process. This deterministic process takes in a *seed* value and for the same seed value, R will generate the same sequence of values. For example, consider generating a random sequence of the numbers 1 through 10 using the `sample()` function for various seed values. We see that setting the seed to the same (arbitrarily chosen) value 76 yields the same sequence, whereas changing the seed value to 79 yields a difference sequence. Play around with this function to get a feel for it.

```
set.seed(76)
sample(1:10)
```



```
## [1] 5 1 10 4 2 6 7 3 9 8

set.seed(76)
sample(1:10)

## [1] 5 1 10 4 2 6 7 3 9 8

set.seed(79)
sample(1:10)

## [1] 3 8 7 2 6 10 1 9 4 5
```

We proceed by setting the seed value to the value 76 and sample 5995 users at random by using the `sample_n()` function from the `dplyr` package

```
profiles_revised <- filter(profiles_revised, height>=55 & height <=80)
set.seed(76)
profiles_revised <- sample_n(profiles_revised, 5995)
```

We convert the `sex` variable to a binary `is.female` variable, whose value is 1 if the user is female and 0 if the user is male, using the `ifelse()` function. Alternatively, we could have coded `is.female` with TRUE/FALSE values, but for plotting purposes we code this variable using 1/0 numerical values. We create the `is.female` variable using the `mutate()` function from the `dplyr` package, which allows us create new variables from existing ones. We plot the points as in Figure 4, making use of the `ggplot2` package and defining an initial base plot.

```
require(ggplot2)
profiles_revised <- mutate(profiles_revised, is.female = ifelse(sex=="f", 1, 0))
base.plot <- ggplot(data=profiles_revised, aes(x=height, y=is.female)) +
  scale_y_continuous(breaks=0:1) +
  theme(panel.grid.minor.y = element_blank()) +
  xlab("Height in inches") +
  ylab("Is female?")
```

We modify this base plot as we go:

```
base.plot + geom_point()
```

This plot is not very useful, as the overlap of the points makes it difficult for determine how many points are involved. We use the `geom_jitter()` function to add a little random noise to each clump of points both along the x and y axes as shown in Figure 5. We observe, for example, there are much fewer males with height 63 inches than 70 inches.

```
base.plot + geom_jitter(position = position_jitter(width = .2, height=.2))
```

We fit both linear and logistic regression models using height as the sole predictor. In order to summarize the results, we use the `msummary()` function from the `mosaic` package rather than the standard `summary()` function, as its output is much more digestible. Furthermore, we extract the coefficients of the linear model using the `coef()` function.

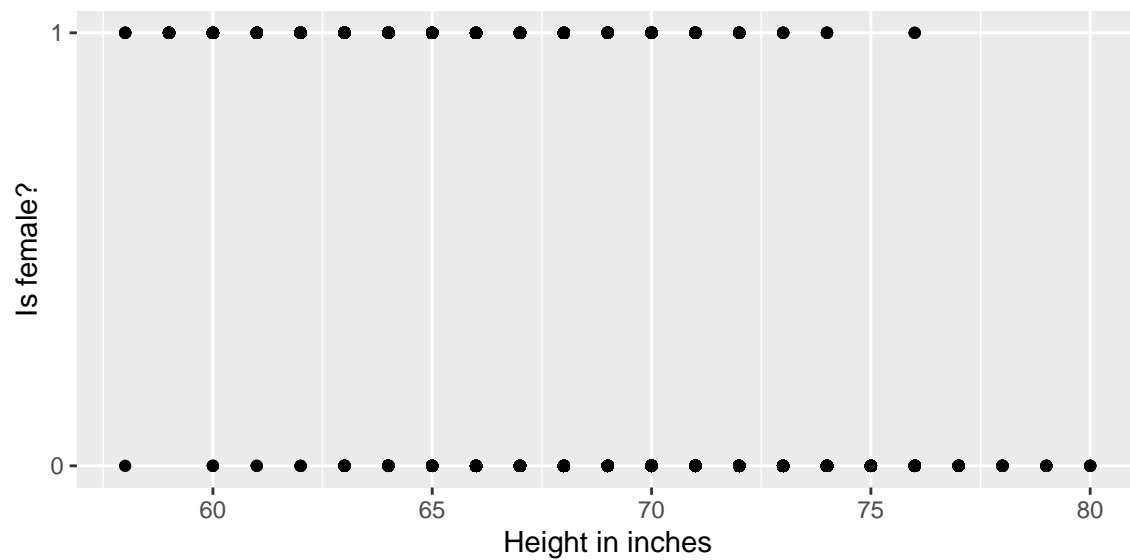


Figure 4: Female indicator vs height.

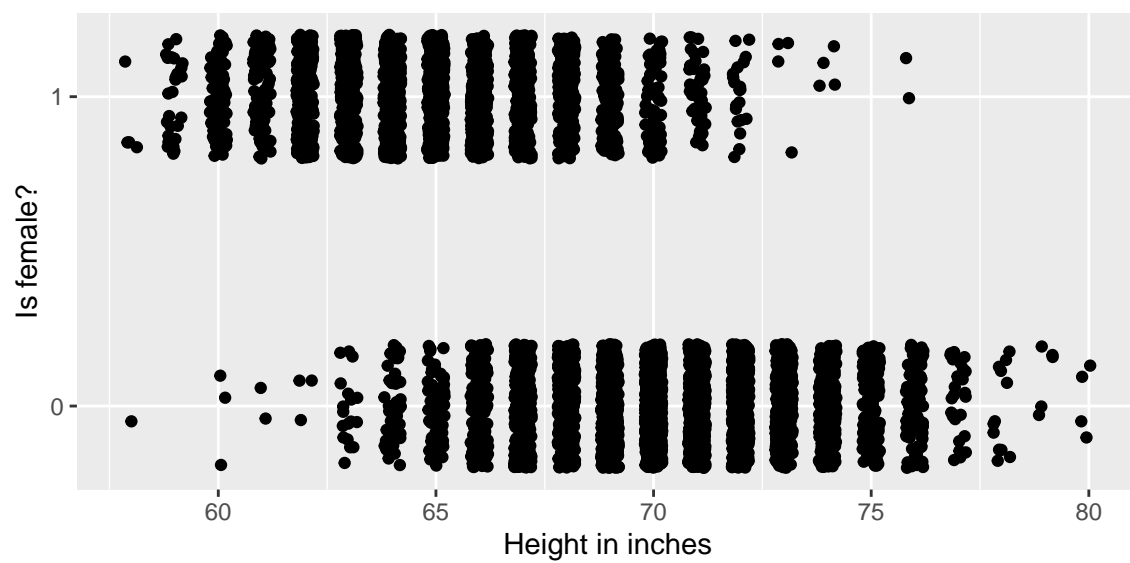


Figure 5: Female indicator vs height (jittered).

```
linear.model <- lm(is.female ~ height, data=profiles_revised)
msummary(linear.model)
```

```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  6.28801    0.08267   76.1  <2e-16 ***
## height      -0.08631    0.00121  -71.4  <2e-16 ***
##
## Residual standard error: 0.36 on 5993 degrees of freedom
## Multiple R-squared:  0.46, Adjusted R-squared:  0.46
## F-statistic: 5.1e+03 on 1 and 5993 DF,  p-value: <2e-16
```

```
b1 <- coef(linear.model)
b1
```

```
## (Intercept)      height
##          6.288      -0.086
```

```
logistic.model <- glm(is.female ~ height, family=binomial, data=profiles_revised)
msummary(logistic.model)
```

```
## Coefficients:
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)  44.0638    1.1267   39.1  <2e-16 ***
## height      -0.6572    0.0167  -39.4  <2e-16 ***
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 8030.6  on 5994  degrees of freedom
## Residual deviance: 4466.5  on 5993  degrees of freedom
## AIC: 4471
##
## Number of Fisher Scoring iterations: 6
```

```
b2 <- coefficients(logistic.model)
b2
```

```
## (Intercept)      height
##          44.06      -0.66
```

In both cases, we observe that the coefficient associated with height is negative (-0.09 and -0.66 for the linear and logistic regressions respectively). In other words, as height increases, the fitted probability of being female decreases as is expected. We plot both regression lines in Figure 6, with the linear regression in red and the logistic regression in blue. The latter necessitates the function `inverse.logit()` in order to compute the inverse logit of the linear equation to obtain the fitted probabilities \hat{p}_i :

$$\hat{p}_i = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 \times \text{height}_i)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 \times \text{height}_i)} = \frac{1}{1 + \exp(-(\hat{\beta}_0 + \hat{\beta}_1 \times \text{height}_i))}$$

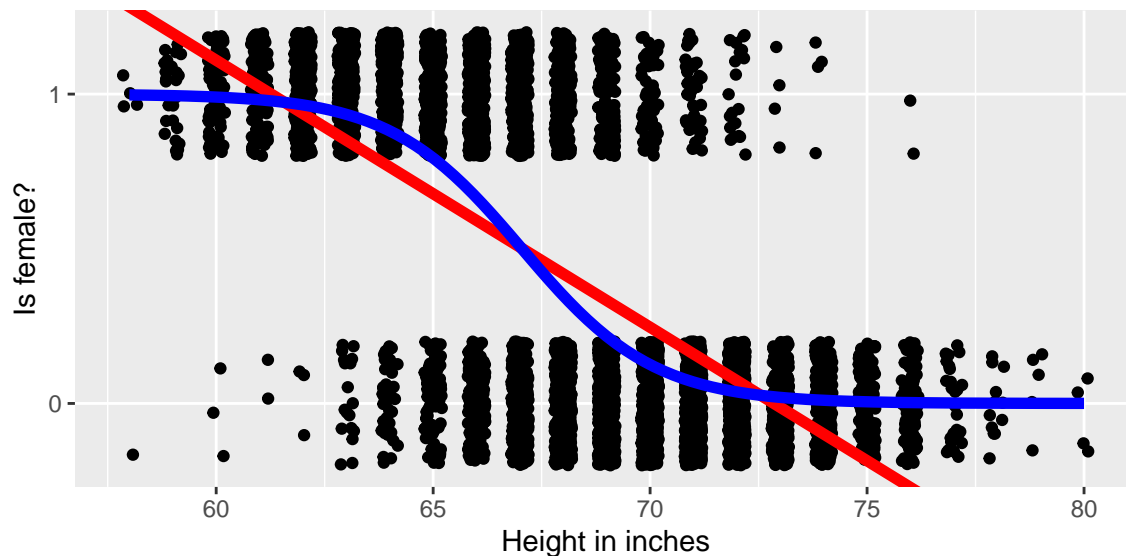


Figure 6: Predicted linear (red) and logistic (blue) regression curves.

```
inverse.logit <- function(x, b){
  linear.equation <- b[1] + b[2]*x
  1/(1+exp(-linear.equation))
}
base.plot + geom_jitter(position = position_jitter(width = .2, height=.2)) +
  geom_abline(intercept=b1[1], slope=b1[2], col="red", size=2) +
  stat_function(fun = inverse.logit, args=list(b=b2), color="blue", size=2)
```

We observe that linear regression (red curve) yields fitted probabilities greater than 1 for heights less than 61 inches and less than 0 for heights over 73 inches, which do not make sense. This is not a problem with logistic regression as the shape of the logistic curve ensures that all fitted probabilities are between 0 and 1. We therefore deem logistic regression to be a more appropriate technique for this data than linear regression.

However, when predicting a user's gender, just using the fitted probabilities \hat{p}_i is insufficient; a decision threshold is necessary. In other words, a point at which if the fitted probability of a user being female is exceeded, we *predict* that user to be female. Looking at the histogram of fitted probabilities, we pick a decision threshold p^* such that for all users with $\hat{p}_i > p^*$, we predict those users to be female. We opt for $p^* = 0.5$ since it splits the values somewhat nicely and highlight this value in red in Figure 7. In order to evaluate the performance of our model and our decision threshold, we produce a contingency table comparing the true (`is.female`) and predicted (`predicted.female`) values:

```
profiles_revised$p.hat <- fitted(logistic.model)
ggplot(data=profiles_revised, aes(x=p.hat)) +
  geom_histogram(binwidth=0.1) +
  xlab(expression(hat(p))) +
  ylab("Frequency") +
  xlim(c(0,1)) +
  geom_vline(xintercept=0.5, col="red", size=1.2)
```

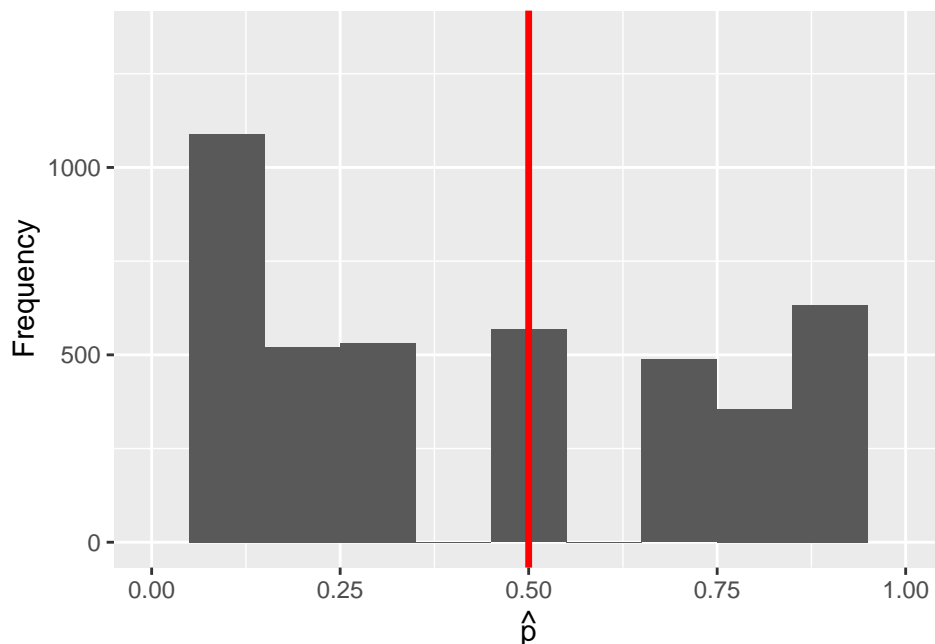


Figure 7: Fitted probabilities of being female and decision threshold (in red).

```
profiles_revised <- mutate(profiles_revised, predicted.female = p.hat >= 0.5)
tally(~is.female + predicted.female, data=profiles_revised)

##      predicted.female
## is.female TRUE FALSE
##      0   593   3050
##      1  1910   442
```

How did our predictions fare?

3.3.2 Pedagogical Discussion

We find that the jump from linear to logistic regression is hard for many students to grasp at first. For example, students often ask “Why the log and exp functions?” and “So we are not modelling the outcome variable Y_i , we’re modeling the probability p_i that Y_i equals 1?” This exercise allows students to build up to the notion of logistic regression from the ground up using visual tools. We also argue that on top of fitting models and interpreting any results, students should also use the results to make explicit predictions and evaluate any model’s predictive power. We asked the students “For what proportion of people did the model guess wrong?” referring to the misclassification error rate, in this case 17.26%. Also solving for height using $p^* = 0.5$ yields a height of 67.05 inches, corresponding to 5 foot 7 inches, which is the height in Figure 1 at which the proportion of males starts to exceed the proportion of females. This point can be highlighted to students, tying together this exercise with the exercise in Section 3.1.

Further questions to ask of students include building a model with more than one predictor, evaluating the *false positive rate* (the proportion of users who were predicted to be female who were actually male), evaluating the *false negative rate* (the proportion of users who were predicted to be male who were actually female), the impact of varying the decision threshold p^* , and asking questions about out-of-sample predictions (using different data to fit and evaluate the model).

4 Conclusions

We present a data set consisting of actual San Francisco OkCupid users' profiles during a period in the 2010s and present example analyses of different levels of sophistication for direct use in the classroom in a similar fashion to [Horton et al. \(2015\)](#). We feel that this data set is ideal for use in introductory statistics and data science courses as the salience of the data set provides students with an interesting vehicle for learning important concepts. By presenting questions to students that allow for the use of their background knowledge, whether it be from the news, stereotypes, or sociological knowledge, students are much better primed to absorb statistical lessons. Furthermore,

1. The data consists of a rich array of categorical, ordinal, numerical, and text variables.
2. This is an instance of real data that is messy, has many suspicious values that need to be accounted for, and includes categorical variables of a complicated nature (for instance, there are 218 unique responses to the ethnicity variable). This reinforces to students that time and energy must be often invested into preparing data for analysis.
3. The data set is of modest size. While $n = 59946$ is not an overwhelmingly large number of observations, it is still much larger than typical data sets used in many introductory probability and statistics courses.

All the files, including the original data and the R Sweave `JSE.Rnw` file used to create this document, can be found at https://github.com/rudeboybert/JSE_OkCupid. Note that the file `profiles_revised.csv.zip` must be unzipped first. All R code used in this document can be outputted into an R script file by using the `purl()` function in the `knitr` package on `JSE.Rnw`:

```
library(knitr)
purl(input="JSE.Rnw", output="JSE.R", quiet=TRUE)
```

Acknowledgements

First, we thank OkCupid president and co-founder Christian Rudder for agreeing to our use of this data set (under the condition that the data set remains public). Second, we thank Everett Wetchler for providing the data. Finally, we thank the reviewers for their helpful comments.

References

- American Statistical Association Undergraduate Guidelines Workgroup (2014). “2014 Curriculum Guidelines for Undergraduate Programs in Statistical Science”, *Technical report*, American Statistical Association, Alexandria, VA.
URL: <http://www.amstat.org/education/curriculumguidelines.cfm>, last accessed April 11, 2015
- Crawford, K. (2013). “Algorithmic Illusions: Hidden Biases of Big Data”, *Strata Conference* .
URL: <https://www.youtube.com/watch?v=irP5RCdpilc>, last accessed April 11, 2015
- Davidson, M. (2013). “Aren’t We Data Science?”, *AMSTAT News* .
URL: <http://magazine.amstat.org/blog/2013/07/01/datascience/>, last accessed April 11, 2015
- GAISE College Group (2005). “Guidelines for Assessment and Instruction in Statistics Education”, *Technical report*, American Statistical Association, Alexandria, VA.
URL: <http://www.amstat.org/education/gaise>, last accessed April 11, 2015
- Gould, R. (2010). “Statistics and the Modern Student”, *International Statistics Review* **78**(2): 297–315.
- Horton, N. J., Baumer, B. and Wickham, H. (2015). “Setting the stage for data science: integration of data management skills in introductory and second courses in statistics”, *CHANCE* **28**(2): 40–50.
- Nolan, D. and Lang, D. T. (2010). “Computing in the Statistics Curricula”, *The American Statistician* **64**(2): 97–107.
- Penenberg, A. L. (2014). “Did the mathematician who hacked OkCupid violate federal computer laws?”, *Pando Daily* .
URL: <http://pando.com/2014/01/22/did-the-mathematician-who-hacked-okcupid-violate-federal-computer-laws/>, last accessed April 11, 2015
- Poulsen, K. (2014). “How a Math Genius Hacked OkCupid to Find True Love”, *WIRED* .
URL: <http://www.wired.com/wiredscience/2014/01/how-to-hack-okcupid/>, last accessed April 11, 2015
- Pruim, R., Kaplan, D. and Horton, N. (2014). “*mosaic: Project MOSAIC (mosaic-web.org) statistics and mathematics teaching utilities*”. R package version 0.9.1-3.
URL: <http://CRAN.R-project.org/package=mosaic>, last accessed April 11, 2015
- Rudder, C. (2010). “The Biggest Lies in Online Data”, *OkTrends: dating research from OkCupid* .
URL: <http://blog.okcupid.com/index.php/the-biggest-lies-in-online-dating>, last accessed April 11, 2015
- Rudder, C. (2014). *Dataclism: Who We Are When We Think No One’s Looking*, Crown.
- Webb, A. (2013). “How I Hacked Online Dating”, *TED Talks* .
- Wickham, H. (2009). *ggplot2: elegant graphics for data analysis*, Springer New York.
URL: <http://had.co.nz/ggplot2/book>
- Wickham, H. (2012). “*stringr: Make it easier to work with strings*”. R package version 0.6.2.
URL: <http://CRAN.R-project.org/package=stringr>, last accessed April 11, 2015
- Wickham, H. (2014). “How are Data Science and Statistics different?”, *IMS Bulletin* **43**(6).
URL: <http://bulletin.imstat.org/2014/09/data-science-how-is-it-different-to-statistics/>, last accessed April 11, 2015
- Wickham, H. and Francois, R. (2014). “*dplyr: a grammar of data manipulation*”. R package version 0.2.
URL: <http://CRAN.R-project.org/package=dplyr>, last accessed April 11, 2015
- Yu, B. (2014). “IMS Presidential Address: Let Us Own Data Science”, *IMS Bulletin* **43**(7).
URL: <http://bulletin.imstat.org/2013/10/president%E2%80%99s-welcome-bin-yu/>, last accessed April 11, 2015