

Understanding Why People Vote

Contents

Introduction	1
Exercices	2
1. <i>Exploration and Logistic Regression</i>	2
Problem 1.1	2
Problem 1.2	3
Problem 1.3	4
Problem 1.4	4
Problem 1.5	5
Problem 1.6	5
2. <i>Trees</i>	6
Problem 2.1	6
Problem 2.2	7
Problem 2.3	7
Problem 2.4	8
3. <i>Interaction Terms</i>	9
Problem 3.1	9
Problem 3.2	9
Problem 3.3	10
Problem 3.4	11
Problem 3.5	12
Problem 3.6	13
Problem 3.7	13

Introduction

In August 2006 three researchers (Alan Gerber and Donald Green of Yale University, and Christopher Larimer of the University of Northern Iowa) carried out a large scale field experiment in Michigan, USA to test the hypothesis that one of the reasons people vote is social, or extrinsic, pressure. To quote the first paragraph of their 2008 research paper:

Among the most striking features of a democratic political system is the participation of millions of voters in elections. Why do large numbers of people vote, despite the fact that ... “the casting of a single vote is of no significance where there is a multitude of electors”? One hypothesis is adherence to social norms. Voting is widely regarded as a citizen duty, and citizens worry that others will think less of them if they fail to participate in elections. Voters’ sense of civic duty has long been a leading explanation of vote turnout...

In this homework problem we will use both logistic regression and classification trees to analyze the data they collected.

The data The researchers grouped about 344,000 voters into different groups randomly - about 191,000 voters were a “control” group, and the rest were categorized into one of four “treatment” groups. These five groups correspond to five binary variables in the dataset.

- **“Civic Duty”** : (variable `civilduty`) group members were sent a letter that simply said “DO YOUR CIVIC DUTY - VOTE!”
- **“Hawthorne Effect”** : (variable `hawthorne`) group members were sent a letter that had the “Civic Duty” message plus the additional message “YOU ARE BEING STUDIED” and they were informed that their voting behavior would be examined by means of public records.
- **“Self”** : (variable `self`) group members received the “Civic Duty” message as well as the recent voting record of everyone in that household and a message stating that another message would be sent after the election with updated records.
- **“Neighbors”** : (variable `neighbors`) group members were given the same message as that for the “Self” group, except the message not only had the household voting records but also that of neighbors - maximizing social pressure.
- **“Control”** : (variable `control`) group members were not sent anything, and represented the typical voting situation.

Additional variables include `sex` (0 for male, 1 for female), `yob` (year of birth), and the dependent variable voting (1 if they voted, 0 otherwise).

Exercices

1. Exploration and Logistic Regression

Problem 1.1 We will first get familiar with the data. Load the CSV file `gerber.csv` into R. What proportion of people in this dataset voted in this election?

```
## 'data.frame': 344084 obs. of 8 variables:
## $ sex      : int  0 1 1 1 0 1 0 0 1 0 ...
## $ yob      : int  1941 1947 1982 1950 1951 1959 1956 1981 1968 1967 ...
## $ voting   : int  0 0 1 1 1 1 1 0 0 0 ...
## $ hawthorne: int  0 0 1 1 1 0 0 0 0 0 ...
## $ civilduty: int  1 1 0 0 0 0 0 0 0 0 ...
## $ neighbors: int  0 0 0 0 0 0 0 0 0 0 ...
## $ self     : int  0 0 0 0 0 0 0 0 0 0 ...
## $ control  : int  0 0 0 0 0 1 1 1 1 1 ...

##      sex      yob      voting      hawthorne
## Min.   :0.0000   Min.   :1900   Min.   :0.0000   Min.   :0.000
```

```
## 1st Qu.:0.0000 1st Qu.:1947 1st Qu.:0.0000 1st Qu.:0.000
## Median :0.0000 Median :1956 Median :0.0000 Median :0.000
## Mean :0.4993 Mean :1956 Mean :0.3159 Mean :0.111
## 3rd Qu.:1.0000 3rd Qu.:1965 3rd Qu.:1.0000 3rd Qu.:0.000
## Max. :1.0000 Max. :1986 Max. :1.0000 Max. :1.000
## civicduty neighbors self control
## Min. :0.0000 Min. :0.000 Min. :0.0000 Min. :0.0000
## 1st Qu.:0.0000 1st Qu.:0.000 1st Qu.:0.0000 1st Qu.:0.0000
## Median :0.0000 Median :0.000 Median :0.0000 Median :1.0000
## Mean :0.1111 Mean :0.111 Mean :0.1111 Mean :0.5558
## 3rd Qu.:0.0000 3rd Qu.:0.000 3rd Qu.:0.0000 3rd Qu.:1.0000
## Max. :1.0000 Max. :1.000 Max. :1.0000 Max. :1.0000
```

Answer : 0.3159

Explanation :

Load the dataset into R by using the read.csv command:

Then we can compute the percentage of people who voted by using the table function:

The output tells us that 235,388 people did not vote, and 108,696 people did vote. This means that $108696/(108696+235388) = 0.316$ of all people voted in the election.

Problem 1.2 Which of the four “treatment groups” had the largest percentage of people who actually voted (voting = 1)?

```
## [1] 0.3145377
## [1] 0.3223746
## [1] 0.3451515
## [1] 0.3779482
##           0           1
## 0.3160698 0.3145377
```

Answer :

1. Civic Duty
2. Hawthorne Effect
3. Self
4. **Neighbors**

Explanation :

There are several ways to get this answer. One is to use the tapply function, and compute the mean value of “voting”, sorted by whether or not the people were in each group:

The variable with the largest value in the “1” column has the largest fraction of people voting in their group - this is the Neighbors group.

Problem 1.3 Build a *logistic regression* model for *voting* using the four treatment group variables as the independent variables (civicduty, hawthorne, self, and neighbors). Use all the data to build the model (DO NOT split the data into a training set and testing set).

Which of the following coefficients are significant in the logistic regression model?
Select all that apply.

```
##
## Call:
## glm(formula = voting ~ civicduty + hawthorne + self + neighbors,
##      family = binomial, data = gerber)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9744  -0.8691  -0.8389   1.4586   1.5590
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.863358   0.005006 -172.459 < 2e-16 ***
## civicduty    0.084368   0.012100   6.972 3.12e-12 ***
## hawthorne    0.120477   0.012037  10.009 < 2e-16 ***
## self         0.222937   0.011867  18.786 < 2e-16 ***
## neighbors    0.365092   0.011679  31.260 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 429238  on 344083  degrees of freedom
## Residual deviance: 428090  on 344079  degrees of freedom
## AIC: 428100
##
## Number of Fisher Scoring iterations: 4
```

Answer :

1. Civic Duty
2. Hawthorne Effect
3. Self
4. Neighbors

Explanation :

You can build the logistic regression model with the following command:

If you look at the output of summary(LogModel), you can see that all of the variables are significant.

Problem 1.4 Using a threshold of 0.3, what is the accuracy of the logistic regression model? (When making predictions, you don't need to use the newdata argument since we didn't split our data.)

$$Accuracy = \frac{TruePositive + TrueNegative}{Ntotal}$$

```
##
##      FALSE    TRUE
##    0 134513 100875
##    1  56730  51966
```

```
## [1] 0.5419578
```

Answer : 0.5419578

Explanation :

First compute predictions:

Then, use the table function to make a confusion matrix:

We can compute the accuracy of the sum of the true positives and true negatives, divided by the sum of all numbers in the table:

$$(134513+51966)/(134513+100875+56730+51966) = 0.542$$

Problem 1.5 Using a threshold of 0.5, what is the accuracy of the logistic regression model?

```
##
##      FALSE
##    0 235388
##    1 108696
```

```
## [1] 0.6841004
```

Answer : 0.6841004

Explanation :

First compute predictions:

Then, use the table function to make a confusion matrix:

We can compute the accuracy of the sum of the true positives and true negatives, divided by the sum of all numbers in the table:

$$(235388+0)/(235388+108696) = 0.684$$

Prproblem 1.6 Compare your previous two answers to the percentage of people who did not vote (the baseline accuracy) and compute the AUC of the model.

What is happening here?

```
## [1] 0.5308461
```

Answer :

1. **Even though all of the variables are significant, this is a weak predictive model.**
2. The model's accuracy doesn't improve over the baseline, but the AUC is high, so this is a strong predictive model.

Explanation :

You can compute the AUC with the following commands (if your model's predictions are called "predictLog"):

Even though all of our variables are significant, our model does not improve over the baseline model of just predicting that someone will not vote, and the AUC is low. So while the treatment groups do make a difference, this is a weak predictive model.

2. *Trees*

Problem 2.1 We will now try out trees. Build a CART tree for *voting* using all data and the same four treatment variables we used before. Don't set the option `method="class"` - we are actually going to create a regression tree here. We are interested in building a tree to explore the fraction of people who vote, or the probability of voting. We'd like CART to split our groups if they have different probabilities of voting. If we used `method='class'`, CART would only split if one of the groups had a probability of voting above 50% and the other had a probability of voting less than 50% (since the predicted outcomes would be different). However, with regression trees, CART will split even if both groups have probability less than 50%.

Leave all the parameters at their default values. You can use the following command in R to build the tree: Plot the tree. **What happens, and if relevant, why?**

0.32

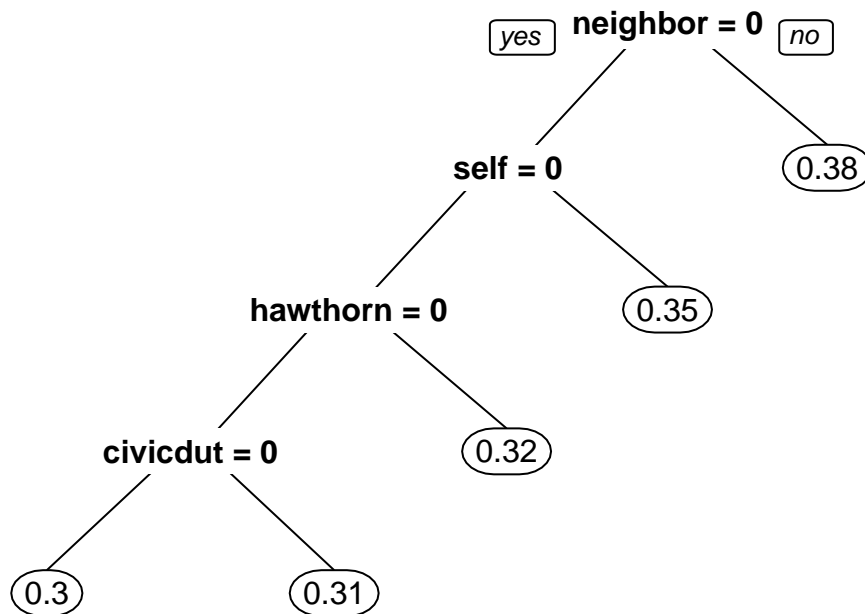
Answer :

1. Only the "Neighbors" variable is used in the tree - it is the only one with a big enough effect.
2. All variables are used - they all make a difference.
3. **No variables are used (the tree is only a root node) - none of the variables make a big enough effect to be split on.**

Explanation :

If you plot the tree, with `prp(CARTmodel)`, you should just see one leaf! There are no splits in the tree, because none of the variables make a big enough effect to be split on.

Problem 2.2 Now build the tree using the command:
to force the complete tree to be built. Then plot the tree.
What do you observe about the order of the splits?



Answer :

1. Civic duty is the first split, neighbor is the last.
2. Neighbor is the first split, civic duty is the last.

Explanation :

You can plot the tree with `prp(CARTmodel2)`.

We saw in Problem 1 that the highest fraction of voters was in the Neighbors group, followed by the Self group, followed by the Hawthorne group, and lastly the Civic Duty group. And we see here that the tree detects this trend.

Problem 2.3 Using only the CART tree plot, determine what fraction (a number between 0 and 1) of “Civic Duty” people voted:

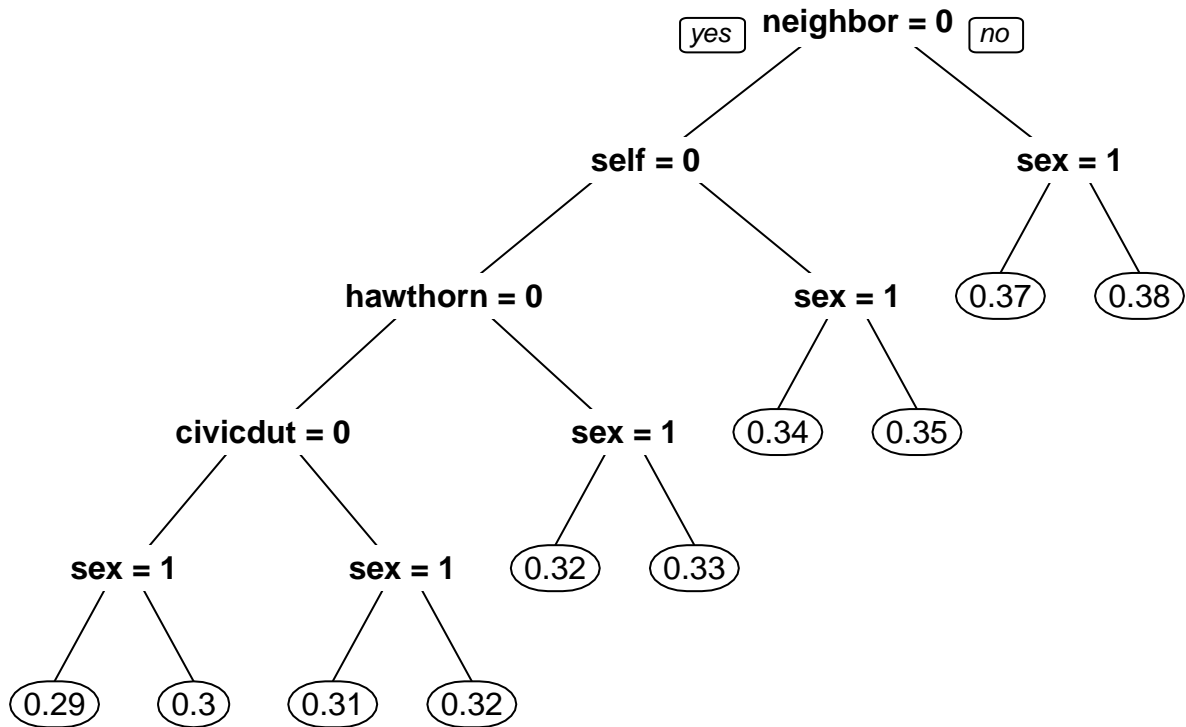
Answer : 0.31

Explanation :

You can find this answer by reading the tree - the people in the civic duty group correspond to the bottom right split, which has value 0.31 in the leaf.

Problem 2.4 Make a new tree that includes the “sex” variable, again with $cp = 0.0$. Notice that sex appears as a split that is of secondary importance to the treatment group.

In the control group, which gender is more likely to vote?



Answer :

1. Men (0)
2. Women (1)

In the “Civic Duty” group, which gender is more likely to vote?

Answer :

1. ** Men (0)**
2. Women (1)

Explanation :

You can generate the new tree using the command:

Then, if you plot the tree with `prp(CARTmodel3)`, you can see that there is a split on the “sex” variable after every treatment variable split. For the control group, which corresponds to the bottom left, $sex = 0$ (male) corresponds to a higher voting percentage.

For the civic duty group, which corresponds to the bottom right, $sex = 0$ (male) corresponds to a higher voting percentage.

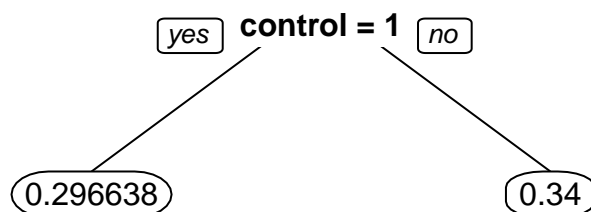
3. Interaction Terms

Problem 3.1 We know trees can handle “nonlinear” relationships, e.g. “in the ‘Civic Duty’ group and female”, but as we will see in the next few questions, it is possible to do the same for logistic regression. First, let’s explore what trees can tell us some more.

Let’s just focus on the “Control” treatment group. Create a regression tree using just the “control” variable, then create another tree with the “control” and “sex” variables, both with $cp=0.0$.

In the “control” only tree, **what is the absolute value of the difference in the predicted probability of voting between being in the control group versus being in a different group?**

You can use the absolute value function to get answer, i.e. $\text{abs}(\text{Control Prediction} - \text{Non-Control Prediction})$. Add the argument “digits = 6” to the prp command to get a more accurate estimate.



```
## [1] 0.043362
```

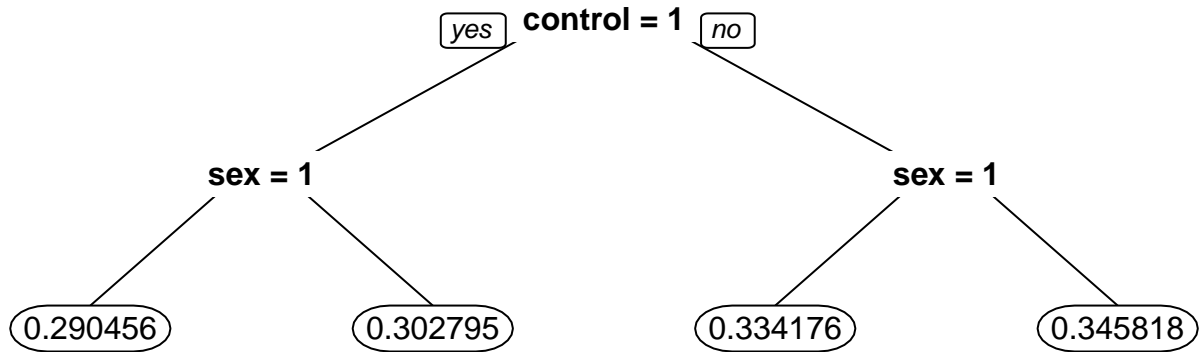
Answer : 0.043362

Explanation :

You can build the two trees with the following two commands:

The split says that if $\text{control} = 1$, predict 0.296638, and if $\text{control} = 0$, predict 0.34. The absolute difference between these is 0.043362.

Problem 3.2 Now, using the second tree (with control and sex), determine **who is affected more by NOT being in the control group** (being in any of the four treatment groups):



```
## [1] 0.04372
```

```
## [1] 0.043023
```

Answer :

1. Men, by a margin of more than 0.001
2. Women, by a margin of more than 0.001
3. They are affected about the same (change in probability within 0.001 of each other).

Explanation :

You can plot the second tree using the command:

The first split says that if control = 1, go left. Then, if sex = 1 (female) predict 0.290456, and if sex = 0 (male) predict 0.302795. On the other side of the tree, where control = 0, if sex = 1 (female) predict 0.334176, and if sex = 0 (male) predict 0.345818. So for women, not being in the control group increases the fraction voting by 0.04372. For men, not being in the control group increases the fraction voting by 0.04302. So men and women are affected about the same.

Problem 3.3 Going back to logistic regression now, create a model using “sex” and “control”. Interpret the coefficient for “sex”:

```
##
## Call:
## glm(formula = voting ~ control + sex, family = binomial, data = gerber)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9220  -0.9012  -0.8290   1.4564   1.5717
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.635538   0.006511 -97.616  < 2e-16 ***
## control      -0.200142   0.007364 -27.179  < 2e-16 ***
## sex          -0.055791   0.007343  -7.597 3.02e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 429238  on 344083  degrees of freedom
## Residual deviance: 428443  on 344081  degrees of freedom
## AIC: 428449
##
## Number of Fisher Scoring iterations: 4
```

Answer :

1. Coefficient is negative, reflecting that women are less likely to vote
2. Coefficient is negative, reflecting that women are more likely to vote
3. Coefficient is positive, reflecting that women are less likely to vote
4. Coefficient is positive, reflecting that women are more likely to vote

Explanation :

You can create the logistic regression model by using the following command:

If you look at the summary of the model, you can see that the coefficient for the “sex” variable is -0.055791. This means that women are less likely to vote, since women have a larger value in the sex variable, and a negative coefficient means that larger values are predictive of 0.

Problem 3.4 The regression tree calculated the percentage voting exactly for every one of the four possibilities (Man, Not Control), (Man, Control), (Woman, Not Control), (Woman, Control). However, logistic regression on the “sex” and “control” variables considers these variables separately, not jointly, and therefore did not do as well.

We can quantify this precisely. Create the following dataframe (this contains all of the possible values of sex and control), and evaluate your logistic regression using the predict function (where “LogModelSex” is the name of your logistic regression model that uses both control and sex):

The four values in the results correspond to the four possibilities in the order they are stated above ((Man, Not Control), (Man, Control), (Woman, Not Control), (Woman, Control)).

What is the absolute difference between the tree and the logistic regression for the (Woman, Control) case?

Give an answer with five numbers after the decimal point.

```
##           1           2           3           4
## 0.3462559 0.3024455 0.3337375 0.2908065

## [1] 0.0003505
```

Answer : 0.0003505

Explanation :

The CART tree predicts 0.290456 for the (Woman, Control) case, and the logistic regression model predicts 0.2908065. So the absolute difference, to five decimal places, is 0.00035.

Problem 3.5 So the difference is not too big for this dataset, but it is there. We're going to add a new term to our logistic regression now, *that is the combination of the "sex" and "control" variables* - so if this new variable is 1, that means the person is a woman AND in the control group. We can do that with the following command:

How do you interpret the coefficient for the new variable in isolation? That is, how does it relate to the dependent variable?

```
##
## Call:
## glm(formula = voting ~ sex + control + sex:control, family = binomial,
##      data = gerber)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9213  -0.9019  -0.8284   1.4573   1.5724
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.637471   0.007603 -83.843  < 2e-16 ***
## sex          -0.051888   0.010801  -4.804 1.55e-06 ***
## control      -0.196553   0.010356 -18.980 < 2e-16 ***
## sex:control  -0.007259   0.014729  -0.493   0.622
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 429238  on 344083  degrees of freedom
## Residual deviance: 428442  on 344080  degrees of freedom
## AIC: 428450
##
## Number of Fisher Scoring iterations: 4
```

Answer :

1. If a person is a woman or in the control group, the chance that she voted goes up.
2. If a person is a woman and in the control group, the chance that she voted goes up.
3. If a person is a woman or in the control group, the chance that she voted goes down.

4. If a person is a woman and in the control group, the chance that she voted goes down.

Explanation :

This coefficient is negative, so that means that a value of 1 in this variable decreases the chance of voting. This variable will have variable 1 if the person is a woman and in the control group.

Problem 3.6 Run the same code as before to calculate the average for each group:

Now what is the difference between the logistic regression model and the CART model for the (Woman, Control) case?

Again, give your answer with five numbers after the decimal point.

```
##           1           2           3           4
## 0.3458183 0.3027947 0.3341757 0.2904558
```

```
## [1] 2e-07
```

Answer : 2e-07 so 0

Explanation :

The logistic regression model now predicts 0.2904558 for the (Woman, Control) case, so there is now a very small difference (practically zero) between CART and logistic regression.

Problem 3.7 This example has shown that trees can capture nonlinear relationships that logistic regression can not, but that we can get around this sometimes by using variables that are the combination of two variables.

Should we always include all possible interaction terms of the independent variables when building a logistic regression model?

Answer :

1. Yes
2. No

Explanation :

We should not use all possible interaction terms in a logistic regression model due to overfitting. Even in this simple problem, we have four treatment groups and two values for sex. If we have an interaction term for every treatment variable with sex, we will double the number of variables. In smaller data sets, this could quickly lead to overfitting.