

[Courseware \(/courses/MITx/15.071x/1T2014/courseware/\)](/courses/MITx/15.071x/1T2014/courseware/)[Course Info \(/courses/MITx/15.071x/1T2014/info/\)](/courses/MITx/15.071x/1T2014/info/)[Discussion \(/courses/MITx/15.071x/1T2014/discussion/forum/\)](/courses/MITx/15.071x/1T2014/discussion/forum/)[Progress \(/courses/MITx/15.071x/1T2014/progress/\)](/courses/MITx/15.071x/1T2014/progress/)[Syllabus \(/courses/MITx/15.071x/1T2014/4264e68418f34d839cf0b33a5da644b2/\)](/courses/MITx/15.071x/1T2014/4264e68418f34d839cf0b33a5da644b2/)[Schedule \(/courses/MITx/15.071x/1T2014/2891f8bf120945b9aa12e6601739c3e6/\)](/courses/MITx/15.071x/1T2014/2891f8bf120945b9aa12e6601739c3e6/)

DETECTING VANDALISM ON WIKIPEDIA

Wikipedia (<http://en.wikipedia.org/wiki/Wikipedia>) is a free online encyclopedia that anyone can edit and contribute to. It is available in many languages and is growing all the time. On the English language version of Wikipedia:

- There are currently 4.3 million pages (http://en.wikipedia.org/wiki/Wikipedia:Size_of_Wikipedia).
- There have been a total of 653 million edits (http://en.wikipedia.org/wiki/Wikipedia:Pruning_article_revisions) (also called revisions) over its lifetime.
- There are approximately 130,000 edits per day (http://en.wikipedia.org/wiki/Wikipedia:WikiProject_Editing_trends/Raw_data/Revisions_per_day).

One of the consequences of being editable by anyone is that some people *vandalize* pages. This can take the form of removing content, adding promotional or inappropriate content, or more subtle shifts that change the meaning of the article. With this many articles and edits per day it is difficult for humans to detect all instances of vandalism and *revert* (undo) them. As a result, Wikipedia uses *bots* - computer programs that automatically revert edits that look like vandalism. In this assignment we will attempt to develop a vandalism detector that uses machine learning to distinguish between a valid edit and vandalism.

The data for this problem is based on the revision history of the page Language (<http://en.wikipedia.org/wiki/Language>). Wikipedia provides a history for each page that consists of the state of the page at each revision. Rather than manually considering each revision, a script was run that checked whether edits stayed or were reverted. If a change was eventually reverted then that revision is marked as vandalism. This may result in some misclassifications, but the script performs well enough for our needs.

As a result of this preprocessing, some common processing tasks have already been done, including lower-casing and punctuation removal. The columns in the dataset are:

- Vandal = 1 if this edit was vandalism, 0 if not.
- Minor = 1 if the user marked this edit as a "minor edit", 0 if not.
- Loggedin = 1 if the user made this edit while using a Wikipedia account, 0 if they were not.
- Added = The *unique* words added.
- Removed = The *unique* words removed.

Notice the repeated use of *unique*. The data we have available is not the bag of words - rather it is the set of words that were removed or added. For example, if a word was removed multiple times in a revision it will only appear one time in the "Removed" column.



PROBLEM 1.1 - BAGS OF WORDS (1 point possible)

Load the data `wiki.csv` (`/c4x/MITx/15.071x/asset/wiki.csv`) with the option `stringsAsFactors=FALSE`, calling the data frame "wiki". Convert the "Vandal" column to a factor using the command `wiki$Vandal = as.factor(wiki$Vandal)`.

How many cases of vandalism were detected in the history of this page?

[Show Answer](#)

You have used 0 of 3 submissions

PROBLEM 1.2 - BAGS OF WORDS (1 point possible)

We will now use the bag of words approach to see if we can improve our model further. We now have two columns of textual data, with different meanings. For example, adding rude words has a different meaning to removing rude words. We'll start like we did in class by building a document term matrix from the Added column. The text already is lowercase and stripped of punctuation. So to pre-process the data, just complete the following four steps:

- 1) Create the corpus for the Added column, and call it "corpusAdded".
- 2) Remove the English-language stopwords.
- 3) Stem the words.
- 4) Build the DocumentTermMatrix, and call it dtmAdded.

If the code `length(stopwords("english"))` does not return 174 for you, then please run the line of code in this file (`/c4x/MITx/15.071x/asset/stopwords.txt`), which will store the standard stop words in a variable called `sw`. When removing stop words, use `tm_map(corpusAdded, removeWords, sw)` instead of `tm_map(corpusAdded, removeWords, stopwords("english"))`.

How many terms appear in dtmAdded?

[Show Answer](#)

You have used 0 of 5 submissions

PROBLEM 1.3 - BAGS OF WORDS (1 point possible)

Filter out sparse terms by keeping only terms that appear in 0.3% or more of the revisions, and call the new matrix `sparseAdded`. How many terms appear in `sparseAdded`?

[Show Answer](#)

You have used 0 of 3 submissions

PROBLEM 1.4 - BAGS OF WORDS (1 point possible)

Convert `sparseAdded` to a data frame called `wordsAdded`, and then prepend all the words with the letter A, by using the command:

```
colnames(wordsAdded) = paste("A", colnames(wordsAdded))
```

Now repeat all of the steps we've done so far (create a corpus, remove stop words, stem the document, create a sparse document term matrix, and convert it to a data frame) to create a Removed bag-of-words dataframe, called wordsRemoved, except this time, prepend all of the words with the letter R:

```
colnames(wordsRemoved) = paste("R", colnames(wordsRemoved))
```

How many words are in the wordsRemoved data frame?

Show Answer

You have used 0 of 5 submissions

PROBLEM 1.5 - BAGS OF WORDS (1 point possible)

Combine the two dataframes (using cbind) into a data frame called wikiWords then add the Vandal column (HINT: remember how we added the dependent variable back into our data frame in the Twitter lecture). Set the random seed to 123 and then split the data set using sample.split from the "caTools" package to put 70% in the training set.

What is the accuracy on the test set of a baseline method that always predicts "not vandalism" (the most frequent outcome)?

Show Answer

You have used 0 of 5 submissions

PROBLEM 1.6 - BAGS OF WORDS (1 point possible)

Build a CART model to predict Vandal, using all of the other variables as independent variables. Use the training set to build the model and the default parameters (don't set values for minbucket or cp).

What is the accuracy of the model on the test set, using a threshold of 0.5? (Remember that if you add the argument type="class" when making predictions, the output of predict will automatically use a threshold of 0.5.)

Show Answer

You have used 0 of 5 submissions

PROBLEM 1.7 - BAGS OF WORDS (1 point possible)

Plot the CART tree. How many word stems does the CART model use?

[Show Answer](#)*You have used 0 of 3 submissions*

PROBLEM 1.8 - BAGS OF WORDS (1 point possible)

Given the performance of the CART model relative to the baseline, what is the best explanation of these results?

- ☐ We have a bad testing/training split.
- ☐ The CART model overfits to the training set.
- ☐ Although it beats the baseline, bag of words is not very predictive for this problem.
- ☐ We over-sparsified the document-term matrix.

[Show Answer](#)*You have used 0 of 1 submissions*

PROBLEM 2.1 - PROBLEM-SPECIFIC KNOWLEDGE (1 point possible)

We weren't able to improve on the baseline using the raw textual information. More specifically, the words themselves were not useful. There are other options though, and in this section we will try two techniques - identifying a key class of words, and counting words.

The key class of words we will use are website addresses. "Website addresses" (also known as URLs - Uniform Resource Locators) are comprised of two main parts. An example would be "http://www.google.com". The first part is the protocol, which is usually "http" (HyperText Transfer Protocol). The second part is the address of the site, e.g. "www.google.com". We have stripped all punctuation so links to websites appear in the data as one word, e.g. "httpwwwgooglecom". We hypothesize that given that a lot of vandalism seems to be adding links to promotional or irrelevant websites, the presence of address is a sign of vandalism.

We can search for the presence of a web address in the words added by searching for "http" in the Added column. The grepl function returns TRUE if a string is found in another string, e.g.

```
grepl("cat","dogs and cats",fixed=TRUE) # TRUE
```

```
grepl("cat","dogs and rats",fixed=TRUE) # FALSE
```

Create a copy of your dataframe from the previous question:

```
wikiWords2 = wikiWords
```

Make a new column in wikiWords2 that is 1 if "http" was in Added:

```
wikiWords2$HTTP = ifelse(grepl("http",wiki$Added,fixed=TRUE), 1, 0)
```

Based on this new column, how many revisions added a link?

[Show Answer](#)*You have used 0 of 3 submissions*

PROBLEM 2.2 - PROBLEM-SPECIFIC KNOWLEDGE (1 point possible)

In problem 1.5, you computed a vector called "spl" that identified the observations for the training and testing set. Use that variable (do not recompute it with sample.split) to make new training and testing sets:

```
wikiTrain2 = subset(wikiWords2, spl==TRUE)
```

```
wikiTest2 = subset(wikiWords2, spl==FALSE)
```

Then create a new CART model using this new variable as one of the independent variables.

What is the new accuracy of the CART model on the test set, using a threshold of 0.5?

Show Answer

You have used 0 of 5 submissions

PROBLEM 2.3 - PROBLEM-SPECIFIC KNOWLEDGE (1 point possible)

Another possibility is that the number of words added and removed is predictive, perhaps more so than the actual words themselves. We already have a word count available in the form of the document-term matrices (DTMs).

Sum the rows of dtmAdded and dtmRemoved and add them as new variables in your data frame wikiWords2 (called NumWordsAdded and NumWordsRemoved) by using the following commands:

```
wikiWords2$NumWordsAdded = rowSums(as.matrix(dtmAdded))
```

```
wikiWords2$NumWordsRemoved = rowSums(as.matrix(dtmRemoved))
```

What is the average number of words added?

Show Answer

You have used 0 of 3 submissions

PROBLEM 2.4 - PROBLEM-SPECIFIC KNOWLEDGE (1 point possible)

In problem 1.5, you computed a vector called "spl" that identified the observations for the training and testing set. Use that variable (do not recompute it with sample.split) to make new training and testing sets with wikiWords2. Create the CART model again (using the training set and the default parameters).

What is the new accuracy of the CART model on the test set?

Show Answer

You have used 0 of 5 submissions

PROBLEM 3.1 - USING NON-TEXTUAL DATA (1 point possible)

We have two pieces of "metadata" (data about data) that we haven't yet used. Make a copy of wikiWords2, and call it wikiWords3:

```
wikiWords3 = wikiWords2
```

Then add the two original variables Minor and Loggedin to this new data frame:

```
wikiWords3$Minor = wiki$Minor
```

```
wikiWords3$Loggedin = wiki$Loggedin
```

In problem 1.5, you computed a vector called "spl" that identified the observations for the training and testing set. Use that variable (do not recompute it with sample.split) to make new training and testing sets with wikiWords3.

Build a CART model using all the training data. What is the accuracy of the model on the test set?

Show Answer

You have used 0 of 5 submissions

PROBLEM 3.2 - USING NON-TEXTUAL DATA (1 point possible)

There is a substantial difference in the accuracy of the model using the meta data. Is this because we made a more complicated model?

Plot the CART tree. How many splits are there in the tree?

Show Answer

You have used 0 of 3 submissions

Please remember not to ask for or post complete answers to homework questions in this discussion forum.

Show Discussion

 New Post



Press (<https://www.edx.org/press>) FAQ (<https://www.edx.org/student-faq>)
Contact (<https://www.edx.org/contact>)



EdX is a non-profit created by founding partners Harvard and MIT whose mission is to bring the best of higher education to students of all ages anywhere in the world, wherever there is Internet access. EdX's free online MOOCs are interactive and subjects include computer science, public health, and artificial intelligence.



(<http://www.meetup.com/edX-Global-Community/>)



(<http://www.facebook.com/EdxOnline>)



(<https://twitter.com/edXOnline>)



(<https://plus.google.com/108235383044095082>)



(<http://youtube.com/user/edxonline>)

© 2014 edX, some rights reserved.

Terms of Service and Honor Code -
Privacy Policy (<https://www.edx.org/edx-privacy-policy>)