

[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/)

Help

QUICK QUESTION 7 (1/1 point)

In the previous video, we used CART and Random Forest to predict sentiment. Let's see how well logistic regression does. Build a logistic regression model (using the training set) to predict "Negative" using all of the independent variables. You may get a warning message after building your model - don't worry (we explain what it means in the explanation).

Now, make predictions using the logistic regression model:

```
predictions = predict(tweetLog, newdata=testSparse, type="response")
```

where "tweetLog" should be the name of your logistic regression model. You might also get a warning message after this command, but don't worry - it is due to the same problem as the previous warning message.

Build a confusion matrix (with a threshold of 0.5) and compute the accuracy of the model. What is the accuracy?

Answer: 0.8197183

EXPLANATION

You can build a logistic regression model in R by using the command:

```
tweetLog = glm(Negative ~ ., data=trainSparse, family="binomial")
```

Then you can make predictions and build a confusion matrix with the following commands:

```
predictLog = predict(tweetLog, newdata=testSparse, type="response")
```

```
table(testSparse$Negative, predictLog >= 0.5)
```

The accuracy is $(254+37)/(254+46+18+37) = 0.8197183$, which is worse than the baseline. If you were to compute the accuracy on the training set instead, you would see that the model does really well on the training set - this is an example of over-fitting. The model fits the training set really well, but does not perform well on the test set. A logistic regression model with a large number of variables is particularly at risk for overfitting.

Note that you might have gotten a different answer than us, because the glm function struggles with this many variables. The warning messages that you might have seen in this problem have to do with the number of variables, and the fact that the model is overfitting to the training set. We'll discuss this in more detail in the Homework Assignment.

Is this worse or better than the baseline model accuracy of 84.5%? Think about the properties of logistic regression that might make this the case!

You have used 1 of 5 submissions



About (<https://www.edx.org/about-us>) Jobs (<https://www.edx.org/jobs>)
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/1082353830440950827>)



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

© 2014 edX, some rights reserved.

[Terms of Service and Honor Code](#) -
[Privacy Policy \(https://www.edx.org/edx-privacy-policy\)](https://www.edx.org/edx-privacy-policy)