

STAT460 – Homework 7

Due: Mar. 11 at the start of class.

A major issue with antiretroviral drugs is the mutation of the virus' genes. Because of its high rate of replication (10^9 to 10^{10} virus per person per day) and error-prone polymerase¹, HIV can easily develop mutations that alter susceptibility to antiretroviral drugs. The emergence of resistance to one or more antiretroviral drugs is one of the more common reasons for therapeutic failure in the treatment of HIV.

In the following paper², a sample of *in vitro*³ HIV viruses were grown and exposed to a particular antiretroviral therapy. The susceptibility of the virus to treatment and the number of genetic mutations of each virus were recorded.



Genotypic predictors of human immunodeficiency virus type 1 drug resistance

Soo-Yon Rhee*, Jonathan Taylor†, Gauhar Wadhera*, Asa Ben-Hur*[§], Douglas L. Brutlag‡, and Robert W. Shafer*[¶]

Division of Infectious Diseases, Departments of *Medicine, †Statistics, and ‡Biochemistry, Stanford University, Stanford, CA 94305

Communicated by Bradley Efron, Stanford University, Stanford, CA, August 28, 2006 (received for review December 5, 2005)

Understanding the genetic basis of HIV-1 drug resistance is essential to developing new antiretroviral drugs and optimizing the use of existing drugs. This understanding, however, is hampered by the large numbers of mutation patterns associated with cross-resistance within each antiretroviral drug class. We used five statistical learning methods (decision trees, neural networks, support vector regression, least-squares regression, and least angle regression) to relate HIV-1 protease and reverse transcriptase mutations to *in vitro* susceptibility to 16 antiretroviral drugs. Learning methods were trained and tested on a public data set of genotype-phenotype correlations by 5-fold cross-validation. For each learning method, four mutation sets were used as input

Results

Drug Susceptibility Results, Input Mutations, and Learning Methods.

For each of the three drug classes, we created four mutation sets that included (i) a complete set of all mutations present in ≥ 2 sequences, (ii) an expert panel mutation set (9), and (iii) a set of nonpolymorphic treatment-selected mutations (TSMs) derived from a database linking protease and RT sequences to the treatment histories of persons from whom the sequenced viruses were obtained (10) (Table 1). A control set of the 30 most common mutations in the data set was also created (see *Supporting Text*, which is published as supporting information on the PNAS web site). Predictions using these 30 mutations were consistently inferior

Anytime I ask for 'test set prediction error' for a method, use the following

```
X = hiv.train$x
Y = hiv.train$y

X_0 = hiv.test$x
Y_0 = hiv.test$y

p = ncol(X)
n = nrow(X)

Y.hat = ##### prediction of method on X_0 #####

print(mean(Yhat == Y_0))
```

¹An enzyme that 'stitches' back together DNA or RNA after replication

²The entire paper is on the website. Try to see what you can get out of it.

³Latin for *in glass*, sometimes known colloquially as a test tube

1. Recall from a previous homework that we made the following plot of the log transformed susceptibility of a virus to the considered treatment, with large values indicating the virus is resistant (that is, not susceptible). Run

```
hist(Y)
```

Divide the response Y into two classes based on the apparent grouping:

```
thresh = ????  
Y_class = rep(0,n)  
Y_class[Y<thresh] = 1  
  
Y_0_class = rep(0,nrow(X_0))  
Y_0_class[Y_0 < thresh] = 1
```

Be sure to report the value you use for `thresh`. Also, report `table(Y_class)`. In terms of log transformed susceptibility, what does `Y_class = 1` correspond to?

2. Let's do some predictions on the test set using LDA.
 - (a) First, attempt to run LDA using `Y_class` as the response (remember that `X` must be a data frame). What happens? What do you think this means?
 - (b) We can correct for this problem by detecting and deleting the miscreant columns

```
out0 = apply(X[Y_class==0,],2,sd) > 1e-16  
out1 = apply(X[Y_class==1,],2,sd) > 1e-16  
  
nonConstantVars = out0*out1
```

```
X.lda = X[,nonConstantVars]  
X_0.lda = X_0[,nonConstantVars]
```

Now, run LDA using `X.lda` as the design matrix. What happens? What do you think this means?

- (c) It turns out the previous issue isn't too important; only a warning. Let's predict the test set and get the prediction error

```
out.lda = lda(Y_class~.,data=data.frame(X.lda))  
  
Yhat.lda = predict(out.lda,data.frame(X_0.lda))  
  
print(mean(Yhat.lda$class == Y_0_class ) )
```

What do you get?

3. Now, let's do logistic lasso

```
out = cv.glmnet(X,Y_class,family='binomial',alpha=1,standardize=F)  
Yhat.glmnet = predict(out,X_0,s='lambda.min',type='class')  
print(mean(Yhat.glmnet == Y_0_class ) )
```

4. An inferential question would be: what gene mutations are most related to producing a susceptible virus?

- (a) What gene mutations are related to susceptibility?

```
betaHat.glmnet = coef(lasso.cv.glmnet,s='lambda.min')  
which(abs(betaHat.glmnet) > 0)
```

- (b) **Harder question:** Which gene mutations lead to an increase in the estimated probability of a virus being susceptible to this particular drug? *Hint:* Consider the signs of the coefficients.