

A CASE STUDY OF LASSO LOGISTIC
REGRESSION FOR PREDICTING
VIROLOGIC RESPONSE USING HIV-1
GENOTYPIC RESISTANCE MUTATIONS

Stefanie Schussler

September 20, 2011

Abstract

We illustrate the use of lasso logistic regression for variable selection among HIV-1 mutations for predicting virologic response to a new treatment regimen for patients with advanced HIV. A control model, using treatment history and baseline covariates, and three models quantifying drug resistance mutations associated with antiretroviral therapy (ART) are considered. The later three models include a genotypic resistance score model, a mutation model with main effects, and a mutation model with interactions. Bootstrapping is used to evaluate the variability of parameter estimates. Prediction accuracy is evaluated based on the area under the ROC curve by means of cross-validation. In this case example, none of the models significantly differ from each other in terms of predictive performance. The set of candidate mutations is expanded beyond those with genotypic scores and the models are re-evaluated. Although significant differences are not found, lasso regression selects mutations that are not selected in the original mutation model. The variable selection varies by fold of cross-validation and by bootstrap sample, though some baseline covariates are commonly selected. The importance of the baseline covariates motivates a partially penalized approach, where all but the baseline covariates are penalized. Within each model, the partially penalized approach is compared to both the typical (fully penalized) lasso approach and to an unpenalized approach. The partially penalized approach does not exhibit significantly higher prediction accuracy, but has some advantages over the other two approaches.

1 INTRODUCTION

While there have been studies that demonstrate logistic regression penalized with the least absolute shrinkage and selection operator (lasso) [22] in the context of high-dimensional data [9], [15], [21], [23], [27], we study its empirical performance using HIV-1 drug resistance mutations. Lasso, a form of L1-penalization, penalizes or shrinks parameter estimates towards zero, in essence performing variable selection while simultaneously estimating the parameters chosen. No more variables are selected than there are observations, the effect of which turns a high-dimensional problem into a low-dimensional one. We illustrate the use of lasso logistic regression for predicting virologic response to ART using HIV-1 drug resistance mutations and treatment history for patients with advanced HIV.

Genotypic sequencing of the HIV-1 retrovirus is used to identify drug resistance mutations associated with loss of susceptibility to ART agents. The interpretation of genotypic resistance tests is complicated because 1) there are many HIV-1 drug resistance mutations; 2) the mutations have varying effects on the virologic response to individual ART drugs; and 3) combinations of mutations may impact the response to the ART regimen. Genotypic drug resistance interpretation systems have been developed to address some of these complications [7], [12], [24]. In particular, the Stanford University HIV Drug Resistance Database (HIVdb) [12] seeks to quantify resistance by assigning a score to each mutation and drug that may be prescribed in the combination regimen. The scores are based on previously reported associations between ART drug activity and specific mutations, in vitro studies aimed at understanding how the mutated virus replicates in the presence of specific drugs, and expert opinion. The

summed score of the weighted mutations quantifies the total resistance for an individual drug. A higher score implies greater resistance.

Although HIVdb simplifies the analysis into an interpretable score, which may be useful for clinical management, the identification of mutations and mutation combinations associated with a lack of susceptibility may provide more details. Further, analysis conducted at the mutation level allows for the inclusion of novel mutations that lack scores. Analysis conducted using pre-determined scores may hamper the ability to discover new mutations or fully take into account the total resistance conferred by a combination of mutations.

This paper illustrates the use of lasso logistic regression for predicting virologic response by comparing predictive performance of varying models based on the results of HIV-1 genotyping. The study population is patients with advanced HIV who have experience with nucleoside reverse transcriptase inhibitors (NRTIs) and who participated in a randomized trial comparing two protease inhibitors (PIs), nelfinavir and ritonavir (NvR study). While ritonavir, as a sole PI, and nelfinavir are no longer commonly used, the data is useful for illustrating the potential utility of lasso. Section 2 provides details on the NvR data, the models, and the validation methods. Section 3 presents the results by examining both variable selection and parameter estimation. The models are compared to each other for prediction accuracy. Within each model, the effect of penalization of the baseline covariates on prediction accuracy is examined by comparing a fully penalized model, a partially penalized model, and a model without penalization. Section 4 discusses the findings of lasso logistic regression applied to the NvR data,

and then discusses advantages and disadvantages found using lasso.

2 MATERIAL AND METHODS

2.1 DATA

The results of the NvR study have been published [16]. Briefly, between January of 1997 and December of 2001, 775 patients with a CD4+ count ≤ 200 cells per mm^3 were randomized to one of two PIs: nelfinavir (NFV) or ritonavir (RTV). The choice of nucleoside and non-nucleoside reverse transcriptase inhibitors (NRTIs and NNRTIs) to be used with the PI was the choice of the clinician/patient. Resistance testing was not used to inform the choice of the regimen. Of the 775 patients enrolled, 610 agreed to participate in a substudy that involved central measures of HIV-RNA sequencing. Following completion of the trial, stored plasma specimens collected prior to randomization were sequenced. Resistance testing was performed at Advanced BioMedical Laboratories (Cinnaminson, NJ), examining the *pro* gene region (encoding the viral protease), which includes 99 amino acids, and over the first 240 amino acids of the *pol* gene region (encoding the viral DNA polymerase reverse transcriptase). Amino acid sequences were evaluated and mutations associated with NRTI, NNRTI, and PI resistance were obtained by comparison with wild-type virus.

Of the 610 patients who agreed to genotyping, 549 had an HIV-RNA viral load of at least 2000 copies/ml at baseline and were successfully genotyped. Patients with previous PI experience, with the exception of saquinavir, were not eligible for the NvR study. We exclude patients who had previously used saquinavir (N=135) or who were

ART-naïve (N=138) to allow for the study of nucleoside resistance. Only five patients were taking zalcitabine (ddC) at randomization and these patients are also excluded. Thus, our analyses are restricted to 228 patients who are NRTI-experienced and PI-naïve at randomization and who had a viral load measured 4 months following randomization. The 4-month time-point was chosen to minimize the effect of changes and/or discontinuation of ART (i.e., non-adherence to the initial regime). As a consequence of the exclusions at baseline, patients in this analysis have little or no PI resistance. Few patients had been prescribed an NNRTI. Thus, only NRTI mutations (a subset of the mutations on reverse transcriptase) are considered.

2.2 VIROLOGIC RESPONSE DEFINITION

HIV-RNA levels were quantified centrally by real-time polymerase chain reaction and sequencing using the Roche Amplicor HIV-1 Monitor (Nutley, New Jersey). The lower limit of detection is 400 copies/ml. Virologic response is defined by an indicator for at least one log₁₀ reduction in RNA (a reduction in RNA by a factor of 10) at 4 months compared to baseline. For example, a reduction in RNA from 10,000 copies/ml at baseline to 1,000 copies/ml at 4 months is a one log₁₀ (log base 10) reduction. There are two patients with a recorded viral load below the detection limit at month 4, who also did not have a one log₁₀ reduction in RNA at month 4 compared to baseline. These two patients are classified as virologic responders.

2.3 GENOTYPIC RESISTANCE DATA

Two-hundred and seventy-two mutations are observed on reverse transcriptase in this patient sample (different amino acid substitutions at the same position are treated as

separate mutations). As a starting point, a consensus list given by the International AIDS Society – USA (IAS, December 2008 edition) [10] is used, containing 44 individual established resistance mutations associated with reverse transcriptase. Of these, only 19 mutations are associated with NRTIs. Seventeen of the 19 mutations are observed in the sample population. These 17 mutations are present in the IAS list and have a score in the Stanford HIVdb to at least one of the four prescribed NRTIs (3TC, AZT, DDI, D4T). Hence, the mutations considered in the analysis are well known and recognized as affecting virologic response.

Hereafter, any amino acid substitution from “*A*” to “*S*” at position “*XX*” on reverse transcriptase will be denoted “*AXXS*”. For example, *M41L* denotes an amino acid substitution of leucine *L* (from wildtype methionine *M*) at position 41 in reverse transcriptase.

The set of mutations considered in the mutation model are used to compute the genotype resistance score. Scores are determined using the Stanford HIVdb (version 5.1.2). A regimen-specific genotypic resistance score (GRS) is computed as the sum of the mutation-drug scores of all the NRTI drugs in the regimen. Hereafter, the GRS will refer to the regimen-specific score for the NRTI class.

As an alternate to the GRS, the number of susceptible drugs within a regimen was investigated. HIVdb provides a guideline for partitioning the resistance score into susceptibility categories. We used these guidelines in determining a cut-off value for dividing a drug as susceptible versus resistant to create a drug susceptibility indica-

tor. That is, when a drug-specific score is less than 15, then the drug is considered susceptible. The regimen-specific genotypic susceptibility score (GSS) is computed as the count of susceptible drugs in the regimen by summing all the drug susceptibility indicators in the NRTI regimen at randomization.

For example, suppose a patient had received both AZT and 3TC at randomization and observed mutations M184V and M41L at baseline. According to Stanford HIVdb, mutations M184V and M41L confer scores of -8 and 15 to AZT, respectively, and 60 and 4 to 3TC, respectively. The individual drug scores are 7 for AZT and 64 for 3TC. Using these scores for this patient, AZT is deemed a susceptible drug, while 3TC is deemed a resistant drug. The GRS is the sum of the individual drug scores, yielding a $GRS=71$, while the GSS is the sum of susceptible drugs, yielding a $GSS=1$. The GRS is sensitive to the scores assigned to the mutation-drug combinations and analysis based on the GRS provides more detail as to the effect of a unit change in a score, while the GSS, being less sensitive to these scores, and has the advantage of being easily interpretable as the count of susceptible drugs.

2.4 MODELS

Four models are considered. The first model (control model) includes the following covariates assessed at baseline: CD4+ count (cells/ mm^3), log-10 (log base 10) HIV-RNA (copies/ml), an indicator for a prior AIDS opportunistic condition, gender (female versus male), race (white versus non-white), age (years), randomly assigned PI (NFV versus RTV), the number of prior NRTI/NNRTI drugs (1 to 6), and the number of new NRTI/NNRTI drugs in the regimen (0, 1, or 2). This model does not include

any information from the genotypic resistance tests. The second model includes the regimen-specific GRS (i.e., the total scores for all mutational effects in the nucleoside regimen) as well as the covariates in the control model. The third model includes the covariates in the control model plus indicators for the presence of individual mutations and interactions between mutations and the NRTIs in the regimen prescribed at randomization, for congruency with the GRS model. The fourth model adds two types of interaction terms to the third model: those between the NRTIs prescribed at randomization and those between the mutations. There are 70 mutation-mutation, 59 mutation-drug, and 5 drug-drug interaction terms considered. A hierarchy is not enforced to accommodate variable selection. That is, if an interaction is used, their main effects are not forced into the model. All variables considered in each model are candidates for inclusion in the resulting model after lasso penalization, in which variable selection is conducted.

Let I denote the set of indicators for each ART agent at randomization and P denote their indices. Let W denote the set of potential covariates and K denote their indices. Let J denote the set of indices for each mutation. Define $\pi = \text{pr}(\text{virologic response})$. Therefore, the four models are:

1. Control model ($p=9$):

$$\text{logit}(\pi) = \beta_0 + \sum_{k \in K} \beta_k W_k$$

2. GRS model ($p=10$): the control model + the GRS for the NRTI drug class

$$\begin{aligned} \text{logit}(\pi) &= \beta_0 + \beta_{NRTI} GRS_{NRTI} + \sum_{k \in K} \beta_k W_k \\ &= \beta_0 + \beta_{NRTI} \sum_{j \in J} \sum_{p \in P} GRS_{Mut_j, ART_p} \cdot I_{ART_p} \cdot Mut_j + \sum_{k \in K} \beta_k W_k \end{aligned}$$

3. Mutation model ($p=85$): the control model + individual mutations + mutation-drug interactions

$$\text{logit}(\pi) = \beta_0 + \sum_{k \in K} \beta_k W_k + \sum_{j \in J} \beta_j Mut_j + \sum_{j \in J} \sum_{p \in P} \beta_{j,p} Mut_j \cdot I_{ART_p}$$

4. Mutation interaction model ($p=160$): the control model + individual mutations + mutation-drug interactions + mutation-mutation interactions + drug-drug interactions

$$\begin{aligned} \text{logit}(\pi) = & \beta_0 + \sum_{k \in K} \beta_k W_k + \sum_{j \in J} \beta_j Mut_j + \sum_{j \in J} \sum_{p \in P} \beta_{j,p} Mut_j \cdot I_{ART_p} + \\ & \sum_{j_1 \in J} \sum_{j_2 > j_1 \in J} \beta_{j_1, j_2} Mut_{j_1} \cdot Mut_{j_2} + \sum_{p_1 \in P} \sum_{p_2 > p_1 \in P} \beta_{p_1, p_2} I_{ART_{p_1}} \cdot I_{ART_{p_2}} \end{aligned}$$

Models 2 – 4 may be extended to be based on a set of candidate mutations beyond those limited by using the IAS consensus list or HIVdb, allowing for the possible discovery of novel mutations. These models are additionally analyzed using the top 50 most frequently observed polymorphisms and resistance mutations, instead of the previous set of 17 mutations; although, nine of the previous 17 mutations are included in this set of 50. These 50 polymorphisms/mutations occur in at least 10 patients; hereafter, loosely referenced as mutations. The GRS is recomputed but only 11 of the top 50 mutations have scores given by HIVdb. For this expanded mutation model, there are 1006 mutation-mutation interactions, 200 mutation-drug interactions, and 5 drug-drug interactions.

Each model is fit with lasso logistic regression. Lasso constrains the sum of the absolute

values of the coefficients. The amount of shrinkage to parameter estimates is controlled by a tuning parameter, λ . Let $\ell(\beta)$ denote the log-likelihood for lasso logistic regression. Then the maximum penalized likelihood estimates (MPLEs) obtained for lasso are:

$$\hat{\beta}_{lasso} = \underset{\beta}{\operatorname{argmax}} \{l(\beta) - \lambda \sum_{j=1}^p |\beta_j|\}$$

The parameter estimates determined by lasso (MPLEs) are compared to the maximum likelihood estimates (MLEs) using the same set of variables, to contrast the amount of shrinkage of individual variables. Lasso does not provide estimates of parameter variability for the MPLEs. Three-thousand bootstrap samples are used to produce bootstrap confidence intervals for the parameters of the variables selected by lasso. In addition, the estimates of parameter variability for the MLEs are examined.

2.5 VALIDATION

Predicted probabilities of virologic response are used to classify patients as predicted responders or non-responders using the observed fraction of virologic responders. This allows for a supervised approach of comparing the observed virologic response with the predicted virologic response. Two nested five-fold cross-validations are performed to find the optimal tuning parameter. The outer-round partitions the data creating five corresponding training (80% of observations) and test (20% of observations) data sets. The inner-round cross-validations are performed using a predefined set of tuning parameters and then pooled back together to form the outer-round training sets. The tuning parameter which minimizes misclassification across the outer-round training set is deemed optimal. Using the optimized tuning parameter and the outer-round training set, the model is refit to obtain parameter estimates. The model is then fit with these

parameter estimates to corresponding outer-round test set using with the optimized tuning parameter. The misclassification error rates, residual-sum-of-squares based on the predicted probabilities of virologic response, and the area under the receiver operating characteristic curve (AUC) are computed using the outer-round test sets.

A typical lasso approach penalizes all variables in the model. A concern of this fully penalized approach is that some variables may be over-shrunk (underestimated) and may be dropped from the model. One way to circumvent this is to penalize only certain variables. In this case example, the mutations are all naturally on the same scale (binary), but the scaling of the control variables differ. Despite standardizing, we believe all the control variables are important based on previous studies, but we do not know which genotype variables are important. Thus, a partially penalized approach is additionally taken where all of the control variables are not penalized and remain in the model, while the remaining variables (i.e., the mutations/GRS, and interaction terms) are penalized by lasso. In this approach, the unpenalized control variables are first estimated in an unpenalized model by themselves, and then the resulting fitted values are treated as an offset in the penalized models. Within each of the four models, the fully penalized and partially penalized approaches are both compared for prediction accuracy against an unpenalized approach of logistic regression without variable selection.

The results of lasso logistic regression strongly depend on the choice of the tuning parameter. Each fold of cross-validation may produce different results. Partitioning the data with two nested cross-validations increases the chance of having a cross-validation training or test set with no variability in response (i.e., all responders or

all non-responders). To circumvent this issue, each fold of cross-validation has equal numbers of responders, as well as equal numbers of non-responders. For stability, this entire process is repeated 100 times, each time randomly partitioning the data for cross-validation. All evaluation measures are estimated by the average of these 100 iterations.

The receiver operating characteristic curve accounts for true and false positives. The AUC discriminator is tested for a significant difference from chance alone (H_0 : $AUC=0.5$) and is evaluated based on the corresponding two-tailed p-value. Significant differences are examined between the AUCs within each of the 100 iterations only. The data is partitioned the same for cross-validation within each iteration; consequently the AUCs are dependent [8]. Every $\binom{4}{2} = 6$ pairs of models, within each of the 100 iterations, are checked for a significant difference between their AUCs, using a Bonferroni multiple testing correction (FWER $\alpha=0.05$).

Analyses are conducted using R [1]. R package *glmnet* (version 1.1-4) performs lasso logistic regression, package *verification* (version 1.31) computes the AUC, and package *caTools* (version 1.10) computes its associated p-value.

3 RESULTS

3.1 POPULATION PROFILE

Table 1 describes the characteristics of the 228 patients at study entry. At month 4, 113 patients (49.6%) have at least a one log-10 reduction in HIV-RNA level. Almost

half of the patients in this sample had a history of an AIDS clinical event. Less than half of the patients in this sample had multiple NRTI mutations at the time of randomization. One patient, who is a virologic non-responder, is missing a CD4+ cell count at randomization and is excluded from some analyses.

Table 1: Patient characteristics at the time of randomization [n=228]

Age [years, median(Q1, Q3)]	38 (32,43)
Female [n(%)]	41 (18.0)
White [n(%)]	66 (28.9)
HIV RNA [log10 copies/ml, median(Q1, Q3)]	4.95 (4.27,5.43)
CD4+ count* [cells/mm ³ , median(Q1, Q3)]	48 (15,84.5)
History of AIDS clinical event [n(%)]	113 (49.6)
Number of previous ART agents at or prior to randomization (NRTI/NNRTIs) [median(Q1, Q3)]	2 (2,4)
Number of mutations [median(min, max)] out of 43	3 (0,9)
Number of NRTI mutations at baseline [median(min, max)] out of 17	1 (0,8)
Number of new NRTIs to be given in planned ART regimen [n(%)]	
0	191 (83.8%)
1	28 (12.3%)
2	9 (3.9%)
Number of new NNRTIs to be given in planned ART regimen [n(%)]	
0	220 (96.5%)
1	8 (3.5%)
Randomized to NFV [n(%)]	119 (52.2)
NRTI GRS [median(Q1, Q3)]	51 (0, 91)

* 1 patient with a missing value who is a virologic non-responder

Q1 = 1st quartile; Q3 = 3rd quartile

Prior to randomization, all patients had exposure to between one and five NRTIs, while 19 patients had exposure to either one or two NNRTIs. Following randomization, in terms of *new* use of NRTIs or NNRTIs, 43 of the 228 patients are prescribed at least one *new* NRTI and/or *new* NNRTI. Of these, 6 patients are only prescribed a new NNRTI, 2 are prescribed both a new NNRTI and new NRTI, and 35 are only prescribed one or more new NRTIs. Following randomization, in terms of NRTIs (new or

continued use), 3 of the 228 patients are not prescribed any NRTI, 16 are prescribed one NRTI (new or continued), and 209 are prescribed two NRTIs (new or continued). The most common NRTI combinations are AZT+3TC (88 patients) and d4T+3TC (87 patients). Fifteen patients have new or continued to use of an NNRTI following randomization.

Figure 1 shows the distribution of the GRS for the NRTI class. The GRS ranges from -5 to 194 and 94 patients have a score of zero. The median GRS is 51; the median non-zero score for the NRTI class is 86. One-hundred and thirty patients (57.0%) have an NRTI GRS of 15 or greater, suggesting moderate to high levels of resistance to the NRTI class.

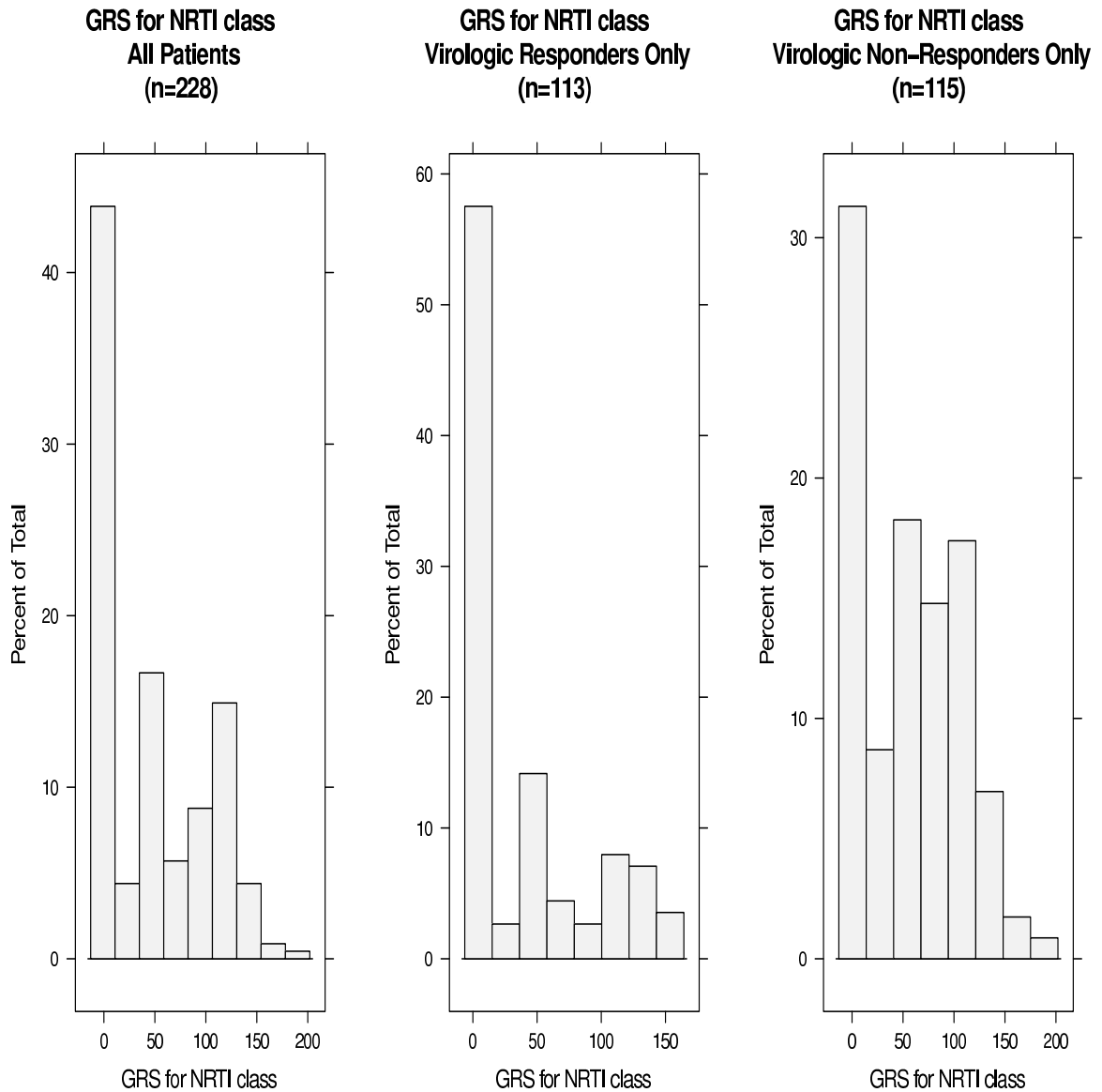
Unadjusted, univariate associations of baseline covariates are summarized in Table 2. When present, 5 of 17 NRTI mutations occurred once or twice in patient viruses. Three of the remaining 12 mutations (M184I, L210W, T215Y) are associated with a reduced likelihood of a virologic response ($p < 0.05$). As resistance towards the NRTI class (the NRTI GRS) increases, the log odds of virologic response decreases ($p = 0.003$). For example, a 50-unit higher GRS is associated with 66% lower odds ($e^{(-0.0084 \times 50)} = 0.66$) of virologic response. Use of nelfinavir ($p = 0.035$) is significantly associated with a poorer virologic response than use of ritonavir, while older age ($p = 0.0002$) is significantly associated with virologic response.

Additional models were considered (results not shown) which varied in how prior and new uses of any NRTI or NNRTI were parameterized. Regardless of how ART is

Table 2: Univariate logistic regression results, testing for significant differences in virologic response for mutations, scores, drug exposure, and patient characteristics

n=228		Virologic Responders	Virologic Non-responders	Log Odds Ratio of VR	95% CI for Log OR	Logistic P-value
Discrete	n	n (% of 113)	n (% of 115)			
NRTI mutations at baseline						
M41L	57	23 (20.4%)	34 (29.6%)	-0.50	-1.10, 0.11	0.110
A62V	2	1 (0.9%)	1 (0.9%)	-	-	-
K65R	1	1 (0.9%)	0 (0.0%)	-	-	-
D67N	38	16 (14.2%)	22 (19.1%)	-0.36	-1.06, 0.34	0.315
K70R	36	18 (15.9%)	18 (15.7%)	0.02	-0.69, 0.73	0.954
K70E	1	1 (0.9%)	0 (0.0%)	-	-	-
L74V	6	3 (2.7%)	3 (2.6%)	0.02	-1.60, 1.64	0.983
V75I	1	0 (0.0%)	1 (0.9%)	-	-	-
F116Y	1	0 (0.0%)	1 (0.9%)	-	-	-
M184V	91	40 (35.4%)	51 (44.3%)	-0.37	-0.91, 0.16	0.168
M184I	10	1 (0.9%)	9 (7.8%)	-2.25	-4.34, -0.17	0.034
L210W	32	9 (8.0%)	23 (20.0%)	-1.06	-1.88, -0.24	0.011
T215F	15	8 (7.1%)	7 (6.1%)	0.16	-0.89, 1.21	0.763
T215Y	70	23 (20.4%)	47 (40.9%)	-0.99	-1.58, -0.41	0.001
K219E	8	4 (3.5%)	4 (3.5%)	0.02	-1.39, 1.43	0.980
K219Q	15	9 (8.0%)	6 (5.2%)	0.45	-0.61, 1.52	0.406
Y181C	7	5 (4.4%)	2 (1.7%)	0.96	-0.70, 2.62	0.257
NRTI use prior to randomization						
AZT	190	89 (78.8%)	101 (87.8%)	-0.67	-1.38, 0.05	0.069
ddI	72	41 (36.3%)	31 (27.0%)	0.43	-0.13, 1.00	0.131
ddC	39	15 (13.3%)	24 (20.9%)	-0.54	-1.25, 0.16	0.131
d4T	109	52 (46.0%)	57 (49.6%)	-0.14	-0.66, 0.38	0.592
3TC	203	99 (87.6%)	104 (90.4%)	-0.29	-1.13, 0.55	0.496
NRTI use at randomization						
AZT	92	41 (36.3%)	51 (44.3%)	-0.34	-0.87, 0.20	0.215
ddI	39	22 (19.5%)	17 (14.8%)	0.33	-0.36, 1.03	0.349
ddC	0	0 (0.0%)	0 (0.0%)	-	-	-
d4T	121	63 (55.8%)	58 (50.4%)	0.21	-0.31, 0.73	0.421
3TC	183	88 (77.9%)	95 (82.6%)	-0.30	-0.96, 0.36	0.370
NNRTI use prior to randomization						
NEV	17	7 (6.2%)	10 (8.7%)	-0.37	-1.37, 0.64	0.474
DLV	3	2 (1.8%)	1 (0.9%)	-	-	-
NNRTI use at randomization						
NEV	14	5 (4.4%)	9 (7.8%)	-0.61	-1.73, 0.52	0.291
EFV	1	1 (0.9%)	0 (0.0%)	-	-	-
Randomization Drug (1=NFV; 0=RTV)	119	51 (45.1%)	68 (59.1%)	-0.56	-1.09, -0.04	0.035
AIDS at BL (1=yes; 0=no)	112	57 (50.4%)	55 (47.8%)	0.10	-0.41, 0.62	0.693
Gender (1=female; 0=male)	41	25 (22.1%)	16 (13.9%)	0.56	-0.13, 1.25	0.109
Race (1=white; 0=non-white)	66	39 (34.5%)	27 (23.5%)	0.54	-0.04, 1.12	0.068
n=228		Virologic Responders†	Virologic Non-responders			
Continuous	Median (Q1, Q3)	Median (Q1, Q3)	Median (Q1, Q3)	Log Odds Ratio	95% CI for Log OR	Logistic p-value
# of prior NRTI\NNRTIs	2 (2, 4)	2 (2, 3)	3 (2, 4)	-0.15	-0.41, 0.11	0.255
# of new NRTI\NNRTIs	0 (0, 0)	0 (0, 0)	0 (0, 0)	0.34	-0.17, 0.84	0.192
CD4+ count at BL (cells/mm ³)	48 (15, 84.5)	54 (16, 85)	45 (15, 83.75)	0.0006	-0.003, 0.005	0.749
AGE (years)	38 (32, 43)	40 (35, 46)	36 (31.5, 41)	0.07	0.03, 0.11	0.0002
log-10 RNA at BL (log-10 copies/mL)	4.95 (4.27, 5.43)	5 (4.31, 5.42)	4.93 (4.25, 5.45)	0.17	-0.15, 0.48	0.301
GRS for NRTI class	51 (0, 91)	0 (0, 67)	55 (0, 97.5)	-0.01	-0.014, -0.003	0.003

Figure 1: Distribution of the GRS for the NRTI class by virologic response



parameterized, when the GRS is included in the model, it remains significant. ART history and new ART use (either NRTI or NNRTIs) were not significantly associated with virologic response despite varying parameterizations with or without analysis of the GRS.

3.2 COMPARISON OF MPLES AND MLES OF THE FULLY PENALIZED LASSO APPROACH

Lasso shrinks all estimates toward zero. Shrinkage can be seen by comparing the MPLEs with their corresponding MLEs. The MPLEs are estimated based on the tuning parameter that minimizes misclassification of virologic response using lasso (fully penalized) logistic regression. The corresponding MLEs are obtained by using logistic regression on the corresponding set of variables. The parameter estimates are given as the log of the odds ratio of virologic response.

In the control model (Table 3), 7 of the 9 variables are selected by lasso. The larger magnitude parameter estimates tend to be shrunk more than the smaller ones. For example, the parameter estimate for white race has an MPLE of 0.559 and corresponding MLE of 0.834 (a reduction from the MLE to the MPLE of 33%), while the parameter estimate for age has an MPLE of 0.070 and corresponding MLE of 0.092 (a reduction of 24%). The estimates for the log of the odds ratio are lower for a patient assigned nelfinavir compared to ritonavir. The control model is the only model (based on the full data) that selects prior ART use.

In the GRS model (Table 4), 7 of the 10 variables are selected by lasso. The NRTI GRS is selected by lasso with a negative parameter estimate indicating the likelihood of virologic response decreases as the NRTI resistance score increases. The associated p-value is 0.002 based on logistic regression.

In the mutation model (Table 5), 9 of the 85 variables are selected by lasso. Two base-

line mutations (M184I and T215Y) and two mutation-drug interactions ($L210W \times 3TC$ and $T215Y \times 3TC$) are selected by lasso. Despite 3TC being in both mutation-drug interactions, it is not selected by lasso. These four selected variables have negative parameter estimates, indicating the likelihood of virologic response decreases with their presence. Standard errors of the point estimates for these four parameters, corresponding to the logistic regression, are large. Consider mutation M184I, which has the largest standard error of 1.111 point estimate of -1.844, and p-value of 0.097 using the logistic regression. Only one of the 10 patients who have this mutation at baseline is a virologic responder.

In the mutation interaction model (results not shown), 10 of the 160 variables are selected by lasso. The mutation interaction model has the same selected variables as those in the mutation model, as well as the interaction between mutations M184V and Y181C which has a positive point estimate.

In the expanded mutation model using the top 50 most frequently observed mutations, 20 of the 259 variables are selected by lasso. There are five mutations selected by lasso, with their parameter estimates given in Table 6. The other variables selected by lasso are not shown. Both mutations (M184I, T215Y) selected in the previous mutation model (based on the 17 candidate mutations) are also selected in this expanded model. Mutation M184I and T215Y have MPLEs of -0.289 and -0.263, respectively, in the expanded mutation model compared to MPLEs of -0.478 and -0.322, respectively, in the previous mutation model. Neither mutation has a significant p-value based on logistic regression. Again M184I has a large standard error (1.331). Of the remaining

three selected mutations, K173E ($n=13$, 5.7%; 11 are virologic responders) and T200A ($n=56$, 24.6%; 36 are virologic responders) are the only two mutations significantly associated with a favorable virologic response ($p=0.033$ and $p=0.003$, respectively from the logistic regression). These two mutations neither occur in HIVdb nor in the IAS consensus list. Nearly all patients with mutation K173E had taken multiple NRTIs prior to randomization; likewise, for mutation T200A. Patients observing mutation K173E also observe between 3 and 9 additional mutations (median value of 5) out of 50, whereas patients observing mutation T200A observe between 2 and 12 additional mutations (median value of 6) out of 50. Among these additional mutations, common mutations include: M184V, R211K, D177E, and K122E, M41L, and T215Y. In addition, patients with either of these two mutations, who are virologic responders, were commonly prescribed 3TC and AZT prior to randomization ($\geq 70\%$ were taking at least these two drugs) and are taking either 3TC and AZT or 3TC and D4T at randomization.

In summary, the signs of the parameter estimates are in agreement between the MPLEs (for the fully-penalized lasso logistic regression) and the MLEs. For the fully penalized lasso logistic regression, randomization PI, which is significant at predicting virologic response in the univariate analysis, is selected across all four models. In addition, the baseline variables for (older) age, female gender, HIV-RNA, and white race are selected across all four models and are associated with an increased likelihood of virologic response. The signs of these coefficients are in agreement across the models that select these control variables. When the GRS is modeled with treatment history (prior NRTI use), treatment history is not selected; however, without including the

genotype (i.e., the control model), treatment history is selected.

3.3 COMPARISON OF MODEL PREDICTIONS

The different measures of prediction accuracy computed indicate that incorporating genotypic information does not yield a better model in terms of prediction (see Table 7). The AUCs range from 0.551 to 0.699 across all 100 iterations and all four models. The average AUC of each model over the 100 iterations ranges from 0.6101 to 0.640. The control model produces an average test error rate of 41.3%, average residual-sum-of-squares for predicted probability of 54.38, and average AUC of 0.623. Comparing the AUCs between each pair of models, within iteration, only one of the 100 iterations shows a significant difference between the GRS model and the mutation interaction model. There are no other significant differences found between any two models using a Bonferroni adjustment.

As an alternative to the GRS model, a GSS model was investigated based on the estimated number of susceptible drugs in the NRTI drug class at randomization. This model and the GRS model describe the data similarly, except in opposing ways (resistance versus susceptibility). The conclusions were similar (results not shown).

The top 50 most frequently observed mutations in reverse transcriptase, associated with NRTI and NNRTIs, were examined instead of the 17 previously described that are associated with NRTIs. The inclusion of these candidate variables allows for mutations to be selected which do not have scores in HIVdb and which may not be well studied in literature. Their inclusion sharply increases the number of candidate variables in the

mutation interaction model such that there were more candidate parameters than observations ($p \gg n$) (data not shown). None of the models were significantly different from each other based on the AUCs using a Bonferroni adjustment.

Finally, the three approaches to penalization (full penalization, partial penalization, and without penalization) are compared across the four models based on the point estimates of the AUC (results not shown). We find no statistically significant differences when comparing these three approaches within each model. However, we do observe patterns of consistency where the partially penalized approach (penalizing all variables except the control variables) consistently performs better across the four models compared to the other two penalization approaches. For instance, in the mutation interaction model, the average AUC of the partial penalized approach is 0.647 as compared to the average AUCs of the full penalized and unpenalized approaches (0.610 and 0.555, respectively). Each has an average standard error of approximately 0.04. Further, we compare the four models within the three penalization approaches. In the approach without penalization: 8 of 100 iterations found a significant difference between the control and mutation model; 10 iterations found a significant difference between the control and the mutation interaction model; 25 iterations found a significant difference between the GRS and the mutation model; 1 iteration found a significant difference between the mutation and mutations interaction model; and, 37 iterations found a significant difference between the GRS and the mutation interaction model. No statistically significant differences are found in the fully penalized and partially penalized approaches, although the GRS model consistently performs better than the other three models based on the point estimates of the AUCs.

4 DISCUSSION

Based on the results of several prospective studies, treatment guidelines recommend genotypic resistance testing to guide the choice of a new treatment regimen for patients with virologic failure [2], [5], [11]. Following early work that used simple schemes such as counts or weighted counts of the number of mutations to predict virologic response to a new regimen [3], [17], [26], [30], a number of genotypic resistance test interpretation algorithms were developed and compared [4], [14], [20]. Most of these algorithms were aimed at reducing the dimension of the genotypic test results to a single score for each drug within a class. In this paper, we use lasso logistic regression for variable selection and estimation to compare models utilizing treatment history (number of prior NRTI/NNRTI drugs) to treatment history plus the genotype, adjusting for the eight control variables studied herein (CD4+ cell count, log-10 HIV-RNA, prior AID opportunistic condition, gender, age, race, randomly assigned PI, number of new NRTI/NNRTI drugs in the regimen) expressed either as mutations or the GRS for the NRTI regimen. Mutation by drug interactions are employed in order to address potential synergistic and antagonistic effects [13] that could be missed with current scoring systems.

Our main finding regarding the four different models and three different approaches to penalization that we compared are: 1) predictive accuracy does not vary significantly among models using mutations, GRS, and treatment history, after adjusting for the eight control variables mentioned above; 2) in an expanded mutation model, two polymorphisms/mutations that do not occur in HIVdb or IAS consensus list, and which are therefore not part of the GRS, are identified with lasso logistic regression; 3) a poorer

virologic outcome is observed with M184I compared to M184V, despite M184I being a transitional mutation to the more mature M184V mutation; and 4) partially penalized models in which treatment history and the other eight control variables listed above are forced into the model tends to perform better than fully penalized models. Below we briefly discuss each of these findings and follow that discussion with some comments on lasso based on our experience with this case example.

Consistent with our findings, based on the AUCs, two studies by Prosperi et al. found 1) no statistical differences between models using a genotypic score compared to models using mutation or mutation interactions, after adjusting for treatment history and baseline covariates [18]; and 2) modeling the genotype alone did not significantly differ from modeling treatment history alone [19]. In the GRS model presented, treatment history is a candidate for selection but is not selected by lasso. The lack of statistical differences between our models using treatment history (the control model) compared to the HIV genotype may be attributed in part to the sample population being too homogeneous among their treatment histories, which consequently may have allowed for similar resistance patterns among the mutations to arise at the time of baseline.

We expanded the set of baseline mutations to widen the scope of the mutation analysis. In doing so, we found two polymorphism/mutations (K173E, T200A) that are selected by lasso and are significant using a corresponding logistic regression. These polymorphisms could be compensatory when they occur with known NRTI mutations, restoring or enhancing activity to NRTIs. Other mutations in HIV-1 have been described which confer either enhanced susceptibility or improved virologic response to

certain ART drugs. This has been described for the M184V mutation, which is associated with improved activity to tenofovir [12]. In addition, some NRTI mutations, such as T215Y, have been associated with phenotypic hypersusceptibility to NNRTIs [11].

Mutation M184I is considered a transitional mutation, occurring for a short time while the viral population evolves to M184V. Although mutation M184I is observed in only ten patients, nine of them are virologic non-responders. Mutations M184I alone is significantly associated with a poorer virologic response, whereas M184V is not. Further, the point estimate of M184I suggests a worse outcome than the point estimate of M184V, although the Stanford HIVdb yields the same score with either mutation. Mutation M184I is also selected in the mutation model, the mutation interaction model, and the expanded mutation model, while M184V is not selected. According to HIVdb, both of these mutations individually have large resistance scores with lamivudine (3TC), but allow zidovudine (AZT) and tenofovir to become more active.

Treatment history and the other eight control variables studied herein have been shown to be predictive of virologic response [3], [18], [19], [20], motivating our partially penalized approach which retains these variables in all four of our models, while allowing the genotype variables to incur shrinkage. We find the partially penalized approach tends to have better predictions, based on the point estimates of the AUC and the test error rates, compared to both the fully penalized lasso logistic regression and the unpenalized logistic regression.

The partially penalized approach is similar to the adaptive group lasso [25] with two

groups, one for the control variables and one for the remaining variables. The “group” reference implies shrinkage to the groups of variables towards zero, while the “adaptive” reference implies different weights to be applied to each group. The effect is that each group has a weighted tuning parameter. In our case, the control variables group is left unpenalized, which is equivalent to having a weight of zero with the adaptive group lasso. One difference with our example is the control variable group is treated as an offset in our approach, so does not have a parameter estimate. An advantage of penalizing the mutations as one group is that the mutations may share sparsity patterns. Furthermore, to improve interpretability, the hierarchy may be enforced in the mutation and mutation interaction models by grouping all higher-order terms together with all lower-order terms using the adaptive group lasso [28]. Although this idea is similar to that of adaptive group lasso, it differs from another perspective. That is, that the tuning parameter for each group is not an adjustment made to the overall tuning parameter by applying a group weight; rather, each group would have its own derived tuning parameters based only on the variables that are members of its group.

Future work is needed to examine other partially penalized approaches. For instance, the mutation model is composed of treatment history and the other eight studied control variables, the individual mutations, and the mutation-drug interactions. Perhaps using a different tuning parameter for the set of individual mutations versus the set of mutation-drug interactions would provide better predictions. Determining the number of groups of variables to penalize could be a tuning parameter itself. Further, an alternative partially penalized approach can be conducted, where instead of modeling treatment history and the other eight control variables as an offset; they can be

modeled together with the mutations, where only the mutations are penalized. This approach would be a special case of a general penalty function by Fan and Li (2006) [6] where each variable has its own tuning parameter and the baseline covariates would have their tuning parameters set at zero.

In summary, in this case study we find little difference in prediction between models using treatment history, mutations, or a genotypic summary score, after adjusting for the control variables. Notwithstanding the lack of differences found between the AUCs of each model, the genotypic summary score is itself a significant predictor of virologic response whereas treatment history is not. One useful finding is that prediction by lasso shows a tendency for improvement when selective shrinkage is used, allowing a model to retain select variables in the form of an offset while penalizing the rest. Future work on partially penalized lasso models may be warranted.

5 ACKNOWLEDGEMENTS

This research was supported by Terry Bein Community Programs for Clinical Research on AIDS and National Institutes of Health grants 5-U01-A1-42170 and T32-AI007432.

We thank Deborah Wentworth for explanation of the NvR data set.

References

- [1] Comprehensive R Archive Network (CRAN), <http://cran.r-project.org/web/packages/>.

- [2] Baxter JD, Mayers DL, Wentworth DN, Neaton JD, et al. 2000. A Randomized Study of Antiretroviral Management Based on Plasma Genotypic Antiretroviral Resistance Testing in Patients Failing Therapy, *AIDS*; 14(9): F83-F93.
- [3] DeGruttola V, Dix L, D'Aquila R, et al. 2000. The Relation Between Baseline HIV Drug Resistance and Response to Antiretroviral Therapy: Re-Analysis of Retrospective and Prospective Studies Using a Standardized Data Analysis Plan, *Antivir Ther*; 5:41-8.
- [4] De Luca A, Cinngolani A, Di Giambenedetto S, et al. 2003. Variable Prediction of Antiretroviral Treatment Outcome by Different Systems for Interpreting Genotypic Human Immunodeficiency Virus Type 1 Drug Resistance, *J Infect Dis*, 187:1934-43.
- [5] Department of Health and Human Services Panel of Antiretroviral Guidelines for Adults and Adolescents. Jan 10, 2011. Guidelines for the Use of Antiretroviral Agents in HIV-1 Infected Adults and Adolescents, Available at <http://www.aidsinfo.nih.gov/ContentFiles/AdultsandAdolescentGL.pdf>. Accessed 1/23/11.
- [6] Fan J, Li R. 2006. Statistical Challenges with High Dimensionality: Feature Selection in Knowledge Discovery, *Proceedings of the International Congress of Mathematicians*, Vol. III, European Mathematical Society, Zurich, 595-622.
- [7] French ANRS (National Agency for AIDS Research), Available online at <http://www.hivfrenchresistance.org/index.html>.
- [8] Hanley JA, McNeil BA. 1983. A Method of Comparing the Areas under Receiver Operating Characteristic Curves Derived from the Same Cases, *Radiology*; 148: 839-43.
- [9] Huang J, Ma S, Zhang CH. 2007. Adaptive Lasso for Sparse High-Dimensional Regression Models, University of Iowa, Dept of Statistics and Actuarial Sciences, Technical Report No. 374, Revision 1.
- [10] International AIDS Society – USA. Johnson VA, Brun-Vézinet F, Clotet B, et al. 2008. Update of the Drug Resistance Mutations in HIV-1: December 2008, *Top HIV Med*; 16(5): 138–45.
- [11] International AIDS Society – USA. Thompson MA, Aberg JA, Cahn P et al. 2010. Antiretroviral Treatment of Adult HIV Infection, 2010 Recommendations of the International AIDS Society, USA Panel, *JAMA*; 304(3): 321–33.
- [12] Liu TF, Shafer RW. 2006. Web Resources for HIV type 1 Genotypic-Resistance Test Interpretation, *Clin Infect Dis*; 42(11):1608–18. Epub 2006 Apr 28.
- [13] Larder BA. 1994. Interactions Between Drug Resistance Mutations in Human Immunodeficiency Virus Type 1 Reverse Transcriptase, *J Gen Vir*; 75:951-7.
- [14] Maggiolo F, Airoidi M, Callegaro A, et al. 2007. Prediction of Virologic Outcome of Salvage Antiretroviral Treatment by Different Systems for Interpreting Genotypic HIV Drug Resistance, *J Int Assoc Physicicans AIDS CARE*; 6(2):87-93.

- [15] Meier L, Van De Geer SA, Buhlmann P. 2008. The Group Lasso for Logistic Regression, *Royal Statist Soc B*; 70(1): 53-71.
- [16] Perez G, MacArthur RD, Walmsley S, et al. 2004. A Randomized Clinical Trial Comparing Nelfinavir and Ritonavir in Patients with Advanced HIV Disease (CPCRA 042/CTN 102), *HIV Clin Trials*; 5(1):7-18.
- [17] Ormaasen V, Sandvik L, Asjo B, et al. 2004. An Algorithm-Based Genotypic Resistance Score is Associated with Clinical Outcome in HIV-1 Infected Adults on Antiretroviral Therapy, *HIV Med*; 5:400-6.
- [18] Prosperi M, Altmann A, Rosen-Zvi M, et al. 2009. Investigation of Expert Rule Bases, Logistic Regression, and Non-Linear Machine Learning Techniques for Predicting Response to Antiretroviral Treatment, *Antivir Ther*; 14: 433-42.
- [19] Prosperi M, Rosen-Zvi M, Altmann A, et al. 2010. Antiretroviral Therapy Optimization without Genotype Resistance Testing: A Perspective on Treatment History Based Models, *PLoS ONE*, 5(10): e13753 1-8.
- [20] Rhee S, Fessel W, Liu T, et al. 2009. Predictive Value of HIV-1 Genotypic Resistance Test Interpretation Algorithms, *J Infect Dis*; 200(3): 453-63.
- [21] Steyerberg EW, Eijkermans MJ, Harrell FE Jr, Habbema JD, 2000. Prognostic Modelling with Logistic Regression Analysis: A Comparison of Selection and Estimation Methods in Small Data Sets, *Stat Med*, 19(8):1059-79.
- [22] Tibshirani R. 1996. Regression Shrinkage and Selection via the Lasso, *Royal Statist Soc B*, 58: 267-88
- [23] Van De Geer SA. 2008. High-Dimensional Generalized Linear Models and the Lasso, *Annals of Stat.*, 36(2):614-45.
- [24] Van Laethem K, De Luca A, Antinori A, et al. 2002. A Genotypic Drug Resistance Interpretation Algorithm that Significantly Predicts Therapy Response in HIV-1-Infected Patients, *Antivir Ther*; Jun, 7(2):123-9. Available online at <http://www.rega.kuleuven.be/cev/index.php?id=30>.
- [25] Wang H, Leng C. 2008. A Note on Adaptive Group Lasso, *Computational Statistics and Data Analysis*, 52: 5277-86.
- [26] Wittkop L, Commenges S, Pellegrin I, et al. 2008. Alternative Methods to Analyse the Impact of HIV Mutations on Virologic Response to Antiviral Therapy, *BMC Med Res Methodol*; 8(68):1-9
- [27] Wu TT, Chen YF, Hastie T, et al. 2009. Genome-Wide Association Analysis by Lasso Penalized Logistic Regression, *Bioinformatics*, 25(6): 714-21.
- [28] Yuan M, Lin Y. 2006. Model Selection and Estimation in Regression with Grouped Variables, *Royal Statist Soc B*, 68(1): 49-67.
- [29] Zhao P, Yu B. 2007. On Model Selection Consistency of Lasso, *JMLR*, 7(2): 2541-64.

- [30] Zolopa AR, Shafer RW, Warford A, et al. 1999. HIV-1 Genotypic Resistance Patterns Predict Response to Saquinavir-Ritonavir Therapy in Patients in Whom Previous Protease Inhibitor Therapy Had Failed, *Ann Intern Med*, 131:813-21.

Table 3: Lasso logistic regression and ordinary logistic regression summaries using the full data with the control model (model 1). Lasso selects 7 of the 9 potential variables. The MPLEs and MLEs are estimates of the log odds ratios of virologic response and the 95% confidence intervals are based on 3000 bootstrap samples.

n=227†	Lasso Logistic Regression			Ordinary Logistic Regression		
	MPLEs	Low 95% CI	Upper 95% CI	MLEs	Std. Error	P-value
(Intercept)	-3.690	-5.126	0.032	-5.118	1.387	<0.001
new NRTI/NNRTI use	0.009	-0.429	0.368	0.115	0.286	0.687
prior NRTI/NNRTI use	-0.043	-0.335	0.071	-0.124	0.144	0.390
randomization PI (NFV)	-0.343	-0.919	0.000	-0.509	0.289	0.079
log10 RNA at BL	0.222	-0.059	0.426	0.384	0.181	0.034
age at entry	0.070	0.000	0.099	0.092	0.021	<0.001
female	0.410	0.000	1.089	0.667	0.377	0.077
white race	0.559	0.000	1.258	0.834	0.327	0.011

†One observation with a missing CD4+ cell count is excluded.

Table 4: Lasso logistic regression and ordinary logistic regression summaries for the full data with the GRS model (model 2). Lasso selects 7 of the 10 potential variables. The MPLEs and MLEs are estimates of the log odds ratios of virologic response and the 95% confidence intervals are based on 3000 bootstrap samples.

n=227†	Lasso Logistic Regression			Ordinary Logistic Regression		
	MPLEs	Low 95% CI	Upper 95% CI	MLEs	Std. Error	P-value
(Intercept)	-4.123	-5.331	0.107	-5.412	1.388	<0.001
new NRTI/NNRTI use	0.088	-0.378	0.401	0.196	0.289	0.498
randomization PI (NFV)	-0.357	-0.925	0.000	-0.498	0.294	0.090
log10 RNA at BL	0.307	-0.033	0.475	0.441	0.185	0.017
age at entry	0.077	0.000	0.101	0.095	0.022	<0.001
female	0.410	0.000	1.046	0.601	0.385	0.118
white race	0.580	0.000	1.262	0.787	0.332	0.018
NRTI GRS	-0.007	-0.011	0.000	-0.009	0.003	0.002

†One observation with a missing CD4+ cell count is excluded.

Table 5: Lasso logistic regression and ordinary logistic regression summaries for the full data with the mutation model (Model 3). Lasso selects 9 of the 85 potential variables. The MPLEs and MLEs are estimates of the log odds ratios of virologic response and the 95% confidence intervals are based on 3000 bootstrap samples.

n=227†	Lasso Logistic Regression			Ordinary Logistic Regression		
	MPLEs	Low 95% CI	Upper 95% CI	MLEs	Std. Error	P-value
(Intercept)	-1.731	-5.700	0.089	-5.716	1.407	<0.001
randomization PI (NFV)	-0.132	-1.174	0.000	-0.587	0.302	0.052
log10 RNA at BL	0.070	-0.237	0.488	0.541	0.191	0.005
age at entry	0.042	0.018	0.129	0.092	0.022	<0.001
female	0.018	-0.068	1.377	0.636	0.391	0.103
white race	0.147	0.000	1.411	0.750	0.339	0.027
M184I	-0.478	-3.014	0.000	-1.844	1.111	0.097
T215Y	-0.322	-0.851	0.000	-0.678	0.642	0.291
L210W \times 3TC	-0.337	-3.795	0.000	-1.067	0.616	0.083
T215Y \times 3TC	-0.107	-0.564	1.392	-0.159	0.721	0.826

†One observation with a missing CD4+ cell count is excluded.

Table 6: Lasso logistic regression and ordinary logistic regression summaries for the full data with the expanded mutation model based on top 50 most frequently occurring mutations. Lasso selects 20 of the 259 potential variables. The MPLEs and MLEs are estimates of the log odds ratios of virologic response. Only the individual mutations are shown.

Mutation	Lasso	Unpenalized Logistic Regression		
	MPLEs	MLEs	Std. Error	P-value
K173E	0.359	1.837	0.861	0.033
M184I	-0.289	0.695	1.331	0.601
T200A	0.243	1.457	0.489	0.003
T215Y	-0.263	-0.927	0.774	0.231
V245M	0.104	0.745	0.501	0.137

Table 7: Model predictions using cross-validation [n=228]

POOLED TEST DATA: <i>AVERAGE</i> VALUES OVER 100 ITERATIONS PENALTY CHOSEN TO MINIMIZE MISCLASSIFICATION AND MAXIMIZE VARIABLE SELECTION					
			AUC CURVE		
Model	Error Rates	RSS for Pred Probs	Area	SE	Avg p-value
CONTROL	41.3%	54.38	0.623	0.037	0.0019
MUTATION	41.4%	55.81	0.618	0.037	0.0049
GRS	40.6%	53.43	0.640	0.036	0.0010
MUTATION INTERACTION	41.7%	56.39	0.610	0.037	0.0095