

Logistic regression

Patrick Breheny

April 17

Introduction

- The general paradigm for penalized regression is

$$Q(\boldsymbol{\beta}|\mathbf{X}, \mathbf{y}) = L(\boldsymbol{\beta}|\mathbf{X}, \mathbf{y}) + P_\lambda(\boldsymbol{\beta})$$

- So far, we have discussed a variety of choices for $P_\lambda(\boldsymbol{\beta})$, but $L(\boldsymbol{\beta}|\mathbf{X}, \mathbf{y})$ has always been the least squares loss function (i.e., linear regression)
- In our next two lectures, we'll explore some different types of loss functions in order to see how penalized regression methods extend to other types of data

Logistic regression

- One of the most common, perhaps even more common than linear regression itself, are studies in which the outcome is binary
- Such studies are particularly common in medical research, where it is common to consider presence/absence of a disease as the outcome
- In this setting, it is natural to model the outcome using a binomial distribution, allowing $\pi_i \equiv \mathbb{P}(Y_i)$ to depend on the features according to

$$\log \frac{\pi_i}{1 - \pi_i} = \mathbf{x}_i^T \boldsymbol{\beta};$$

this is known as the *logistic regression* model

Logistic regression objective function

- As with linear regression, maximum likelihood estimation of β will be problematic if the number of features is large, making penalized likelihood estimation desirable
- Thus, we will estimate β by minimizing the objective function

$$-\frac{1}{n} \sum_{i=1}^n \{y_i \log \pi_i + (1 - y_i) \log (1 - \pi_i)\} + P_\lambda(\beta),$$

where \mathbf{X} and β are included in the likelihood implicitly, as π is a function of $\mathbf{X}\beta$

Similarities and differences

- For the most part, everything we have talked about so far this semester with respect to linear regression carries over to logistic regression: ridge, lasso, MCP, elastic net, etc., penalties have similar effects on regression coefficients as we have seen
- However, two differences are worth discussing:
 - We need new algorithms for model fitting
 - We need new measures of predictive accuracy
- Inference concerning β , of course, would also be different, although this is beyond the scope of a single lecture; some approaches (sample splitting, knockoffs) are straightforward to extend while others (selective inference) are not

Logistic regression: Notation, intercept

- First, let's discuss algorithms for minimizing $Q(\beta|\mathbf{X}, \mathbf{y})$ in the logistic regression case
- We begin by noting that for logistic regression, it is not possible to eliminate the need for an intercept by centering the response variable
- Thus, in the derivations that follow, \mathbf{y} will denote the original vector of 0-1 responses
- The design matrix \mathbf{X} is standardized as before, but now contains an unpenalized column of 1's for the intercept, with corresponding coefficient β_0

Iteratively reweighted least squares algorithm

- Like unpenalized logistic regression itself, algorithms for penalized logistic regression employ Taylor series expansions to produce quadratic approximations to the loss function, thereby allowing us to use our previously derived solutions for linear regression
- This two-step approach is known, generally speaking, as the iteratively reweighted least squares (IRLS) algorithm:
 - (1) Approximate the loss function based on $\beta^{(m)}$
 - (2) Solve for $\beta^{(m+1)}$, the value that minimizes the approximated loss function

with the above two steps alternated until convergence

Taylor series expansion

Approximating the loss function via Taylor series expansion, we have

$$L(\boldsymbol{\beta}|\mathbf{X}, \mathbf{y}) \approx \frac{1}{2n}(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{W}(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}),$$

where

- $\tilde{\mathbf{y}}$, the working response, is defined by

$$\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta}^{(m)} + \mathbf{W}^{-1}(\mathbf{y} - \boldsymbol{\pi})$$

- \mathbf{W} is a diagonal matrix of weights, with elements $w_i = \pi_i(1 - \pi_i)$
- $\boldsymbol{\pi}$ is evaluated at $\boldsymbol{\beta}^{(m)}$

Generalized linear models

- We are focusing today on logistic regression, but the same approach can be applied to fit penalized versions of any generalized linear model (GLM)
- To do so, one simply has to replace $\tilde{\mathbf{y}}$ and \mathbf{W} with the appropriate expressions for the corresponding response distribution and link function used in the model
- For example, both `glmnet` and `ncvreg` have options for fitting Poisson regression models (`family="poisson"`) using the same technique we describe here for logistic regression

Coordinate descent

- Following the quadratic approximation, the objective function closely resembles the objective we're used to from linear regression, but now with observation weights $\{w_i\}$
- We can still employ coordinate descent, but the presence of these weights changes the form of the updates
- Let

$$\begin{aligned}v_j &= n^{-1} \mathbf{x}_j^T \mathbf{W} \mathbf{x}_j \\ \mathbf{r} &= \mathbf{W}^{-1} (\mathbf{y} - \boldsymbol{\pi}) \\ z_j &= \frac{1}{n} \mathbf{x}_j^T \mathbf{W} (\tilde{\mathbf{y}} - \mathbf{X}_{-j} \boldsymbol{\beta}_{-j}^{(m)}) \\ &= \frac{1}{n} \mathbf{x}_j^T \mathbf{W} \mathbf{r} + v_j \beta_j^{(m)}\end{aligned}$$

Coordinate-wise updating

- Now, the coordinate-descent update for lasso is

$$\beta_j \leftarrow \frac{S(z_j, \lambda)}{v_j}$$

- For the elastic net,

$$\beta_j \leftarrow \frac{S(z_j, \lambda_1)}{v_j + \lambda_2}$$

- For MCP,

$$\beta_j \leftarrow \begin{cases} \frac{S(z_j, \lambda)}{v_j - 1/\gamma} & \text{if } |z_j| \leq v_j \gamma \lambda \\ \frac{z_j}{v_j} & \text{if } |z_j| > v_j \gamma \lambda \end{cases}$$

Remarks

- Note that the coordinate descent portion (2) of the algorithm must also involve the updating of the intercept term:

$$\beta_0 \leftarrow \frac{z_j}{v_j}$$

- Speaking of step (2), we typically do not iterate until convergence in order to obtain $\hat{\beta}^{(m+1)}$, but simply make one pass of coordinate updates, then re-approximate
- As is the case for the traditional IRLS algorithm for GLMs, this algorithm is not guaranteed to converge
- Typically, however, this is not a problem in practice, at least not for pathwise approaches, as the “warm starts” phenomenon provides protection against bad initial values

Model saturation

- One exception to this remark is that convergence does tend to be a problem with saturated models
- To protect against this, the `glmnet` and `ncvreg` packages will terminate the pathwise algorithm early if saturation ($R^2 > .99$) is detected
- Another numerical issue worth mentioning is that weights are typically capped at a minimum value ϵ to prevent fitted probabilities of 0 or 1, which also tends to happen as models become saturated

Reweighting and γ

- Before moving on, it is worth considering how the reweighting affects nonconvex penalties such as MCP and SCAD with respect to convexity and the choice of γ
- In linear regression, the scaling factor by which solutions are adjusted toward their unpenalized solution is a constant ($1 - 1/\gamma$ for MCP) for all values of λ and for each covariate
- Furthermore, for standardized covariates, this constant has a universal interpretation for all linear regression problems, meaning that defaults such as $\gamma = 3$ can be used and will be universally reasonable (though not necessarily optimal, of course)

Reweighting and γ (cont'd)

- In logistic regression, however, this scaling factor ($v_j - 1/\gamma$ for MCP) is different for each data set and for each feature
- This makes choosing an appropriate value for γ considerably more difficult and robs the parameter of a consistent interpretation
- For example, suppose we attempt to use $\gamma = 3$; for logistic regression, w_i cannot exceed 0.25, γ cannot exceed $1/v_j$, and $Q(\beta_j|\beta_{-j})$ is no longer convex and does not have a unique minimum

Adaptive rescaling

- To resolve these difficulties, `ncvreg` takes an approach known as *adaptive rescaling*, which replaces $p_{\lambda,\gamma}(|\beta_j|)$ by $p_{\lambda,\gamma}(|v_j\beta_j|)$
- The consequence is that the updating steps become simple extensions of the linear regression updating steps:

$$\beta_j \leftarrow \frac{F(z_j, \lambda, \gamma)}{v_j} \quad (\text{MCP})$$

$$\beta_j \leftarrow \frac{T_S(z_j, \lambda, \gamma)}{v_j} \quad (\text{SCAD}),$$

where F and T_S are the firm and SCAD thresholding operators

- The purpose of this is to give γ a consistent meaning again and allow simple reasonable defaults such as $\gamma = 3$ to have the same meaning in other GLMs as they do in linear regression

Leukemia data (refresher)

- For the second half of this lecture, I'd like to (a) discuss various measures of predictive accuracy (b) demonstrate how to fit penalized logistic regression models with `glmnet` and `ncvreg` and (c) illustrate using a real data set
- For our case study data set, we will use the Leukemia data that we previously analyzed back at the beginning of the course using a multiple testing approach
- To refresh your memory, the data involved expression levels for 7,129 genes and 72 patients, of whom 47 patients had acute lymphoblastic leukemia (ALL) and the other 25 patients had acute myeloid leukemia (AML)

Deviance

- The most natural measure of predictive accuracy is to simply use the loss function (i.e., the log-likelihood); this is equivalent to the idea of deviance in GLM theory
- Specifically, let

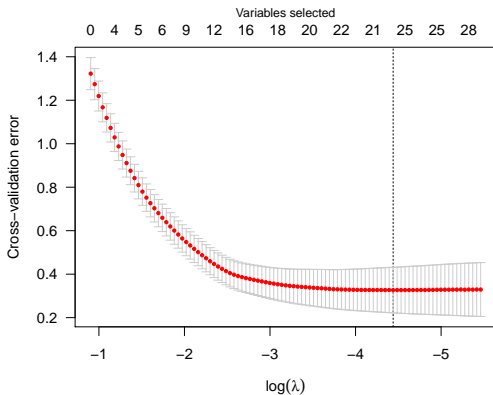
$$d_i^2 = -2\{y_i \log \hat{\pi}_{(i)} + (1 - y_i) \log (1 - \hat{\pi}_{(i)})\},$$

where $\hat{\pi}_{(i)}$ denotes the predicted value of π_i based on the cross-validation fold from which observation i was excluded

- Then $D = \sum_i d_i^2$ is known as the *deviance*; the factor of 2 is so that the difference in deviance between two nested models will follow a χ^2 distribution in classical likelihood theory, but is essentially arbitrary for our purposes

Deviance for the Leukemia data (lasso)

```
cvfit <- cv.glmnet(X, y, family="binomial", gamma=4)
```

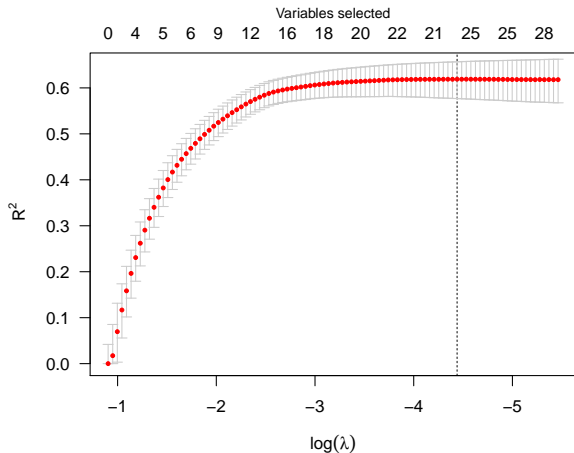


R^2

- As with linear regression, it is often more interpretable to look at the fraction of RSS/deviance explained, rather than the total (although of course the two are equivalent in terms of choosing an optimal value of λ)
- Letting $D(\lambda)$ denote the (cross-validated) deviance at a given value of λ and D_0 denote the deviance of the null model, the coefficient of determination is given by

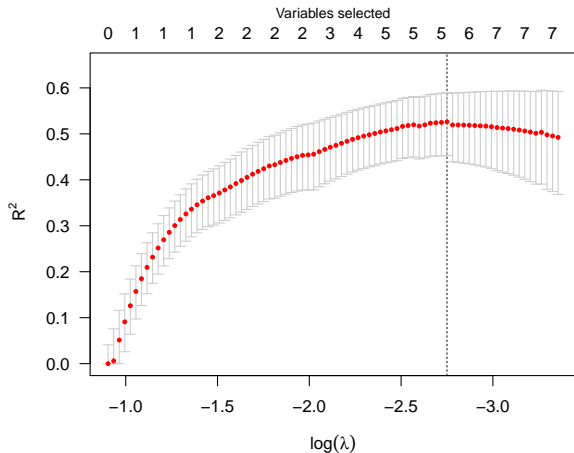
$$R^2 = \frac{D_0 - D(\lambda)}{D_0}$$

R^2 for the Leukemia data (lasso)



$$R_{\max}^2 = 0.62$$

R^2 for the Leukemia data (MCP)



$$R_{\max}^2 = 0.53$$

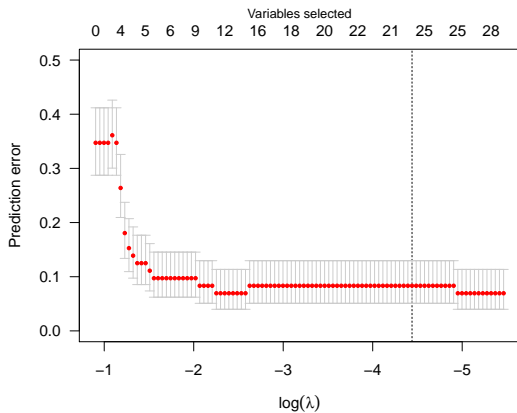
Misclassification error

- A more black-and-white measure of predictive accuracy is the ability of the model to predict the outcome (here, leukemia type) correctly, in the sense that $\hat{\pi}_{(i)} > 0.5$ implies a prediction that $Y_i = 1$
- This is generally referred to as the *misclassification error* of the model
- The advantage of misclassification error is ease of interpretation: everyone know what it means to say that the model got 78% of its predictions correct
- The disadvantage is that it is less stable and well-behaved than the deviance; small changes in λ can lead to large changes in the misclassification error as $\hat{\pi}_{(i)}$ moves from one side of 0.5 to the other

Standard errors

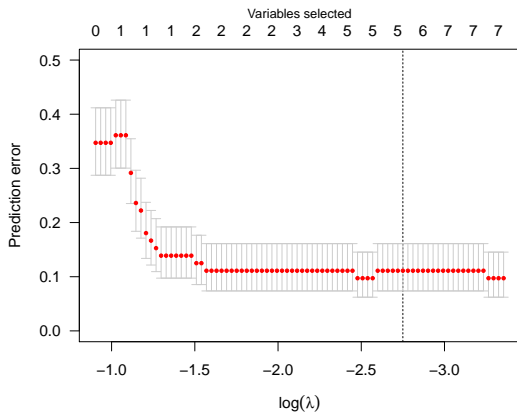
- We've been showing standard errors, but I have not yet discussed how they are estimated
- For deviance and R^2 , `glmnet` and `ncvreg` estimate the standard error of \bar{d}^2 through the simple formula
$$SE = SD(d_i^2) / \sqrt{n}$$
- `glmnet` uses the same approach for misclassification error, though `ncvreg` calculates exact (Clopper-Pearson) confidence intervals based on the observed number of misclassifications and trials

Misclassification error (lasso)



$$\text{MCE}(\lambda_{CV}) = 8.3\%$$

Misclassification error (MCP)

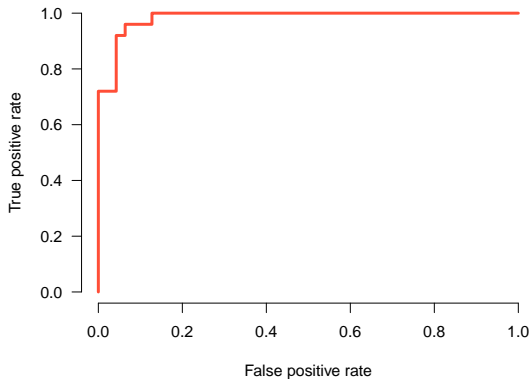


$$\text{MCE}(\lambda_{\text{CV}}) = 11.1\%$$

Area under the curve

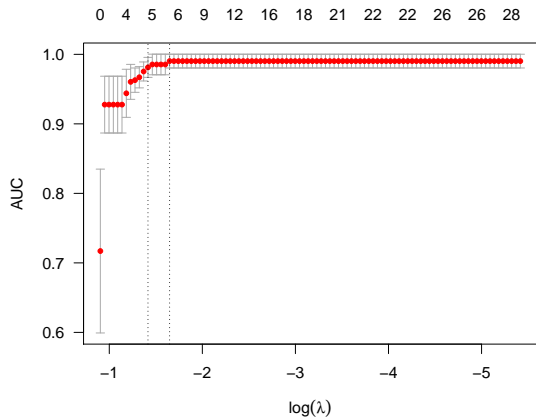
- Finally, along the same lines, one measure that some people find appealing is the area under the curve (AUC), which reflects the concordance between the predicted values and the truth
- The idea here is, rather than use $\hat{\pi}_{(i)} > 0.5$ as a cutoff, we could instead use a variable cutoff $\hat{\pi}_{(i)} > c$
- This would produce false positive and false negative rates for each value of c ; we could integrate the resulting curve to obtain an overall measure of accuracy over all possible cutoffs

ROC curve at λ_{CV}



AUC = 98.4%

AUC (lasso)



$$AUC_{\max} = 99.0\%$$