

Higher criticism thresholding: Optimal feature selection when useful features are rare and weak

David Donoho^{†*} and Jiashun Jin^{§¶1}

[†]Department of Statistics, Stanford University, Stanford, CA 94305; [§]Department of Statistics, Purdue University, West Lafayette, IN 47907; and [¶]Department of Statistics, Carnegie Mellon University, Pittsburgh, PA 15213

Contributed by David Donoho, August 2, 2008 (sent for review July 8, 2008)

In important application fields today—genomics and proteomics are examples—selecting a small subset of useful features is crucial for success of Linear Classification Analysis. We study feature selection by thresholding of feature Z-scores and introduce a principle of threshold selection, based on the notion of *higher criticism* (HC). For $i = 1, 2, \dots, p$, let π_i denote the two-sided P-value associated with the i th feature Z-score and $\pi_{(i)}$ denote the i th order statistic of the collection of P-values. The HC threshold is the absolute Z-score corresponding to the P-value maximizing the HC objective $(i/p - \pi_{(i)})/\sqrt{i/p(1-i/p)}$. We consider a rare/weak (RW) feature model, where the fraction of useful features is small and the useful features are each too weak to be of much use on their own. HC thresholding (HCT) has interesting behavior in this setting, with an intimate link between maximizing the HC objective and minimizing the error rate of the designed classifier, and very different behavior from popular threshold selection procedures such as false discovery rate thresholding (FDR). In the most challenging RW settings, HCT uses an unconventionally low threshold; this keeps the missed-feature detection rate under better control than FDR and yields a classifier with improved misclassification performance. Replacing cross-validated threshold selection in the popular Shrunken Centroid classifier with the computationally less expensive and simpler HCT reduces the variance of the selected threshold and the error rate of the constructed classifier. Results on standard real datasets and in asymptotic theory confirm the advantages of HCT.

false discovery rate | linear classification | threshold selection | rare/weak feature models

The modern era of high-throughput data collection creates data in abundance; however, this data glut poses new challenges. Consider a simple model of linear classifier training. We have a set of labeled training samples (Y_i, X_i) , $i = 1, \dots, n$, where each label Y_i is ± 1 and each feature vector $X_i \in R^p$. For simplicity, we assume the training set contains equal numbers of 1's and -1's and that the feature vectors $X_i \in R^p$ obey $X_i \sim N(Y_i\mu, \Sigma)$, $i = 1, \dots, n$, for an unknown mean contrast vector $\mu \in R^p$; here, Σ denotes the feature covariance matrix and n is the training set size. In this simple setting, one ordinarily uses linear classifiers, taking the general form $L(X) = \sum_{j=1}^p w(j)X(j)$, for a sequence of “feature weights” $w = (w(j): j = 1, \dots, p)$.

Classical theory going back to R. A. Fisher (1) shows that the optimal classifier has feature weights $w \propto \Sigma^{-1}\mu$; at first glance, linear classifier design seems straightforward and settled. However, in many of today's most active application areas, it is a major challenge to construct linear classifiers that work well.

In many ambitious modern applications—genomics and proteomics come to mind—measurements are automatically made on thousands of standard features, but in a given project, the number of observations, n , might be in the dozens or hundreds. In such settings, $p \gg n$, which makes it difficult or impossible to estimate the feature covariance straightforwardly. In such settings one often ignores feature covariances. Working in standardized feature space where individual features have mean zero and variance one, a by-now standard choice uses weights $w(j) \propto$

$\text{Cov}(Y, X(j)) \equiv \mu(j)$ (2, 3). Even when this reduction makes sense, further challenges remain.

When Useful Features Are Rare and Weak

In some important applications, standard measurements generate many features automatically, few of which are likely to be useful in any specific project, but researchers do not know in advance which ones will be useful in a given project. Moreover, reported misclassification rates are relatively high. Hence, the dimension p of the feature vector is very large, and although there may be numerous useful features, they are relatively rare and individually quite weak.

Consider the following *rare/weak feature model* (RW feature model). We suppose the contrast vector μ to be nonzero in only k out of p elements, where $\varepsilon = k/p$ is small, that is, close to zero. As an example, we might have $p = 10,000$, $k = 100$, and so $\varepsilon = k/p = 0.01$. In addition, we suppose that the nonzero elements of μ have *common* amplitude μ_0 . Because the elements $X(j)$ of the feature vector where $\mu(j) = 0$ are entirely uninformative about the value of $Y(j)$, only the k features where $\mu(j) = \mu_0$ are useful. The problem is how to identify and benefit from those rare, weak features. Setting $\tau = \sqrt{n}\mu_0$, we speak of the parameters ε and τ as the sparsity and strength parameters and denote by $RW(\varepsilon, \tau)$ this setting. [Related “sparsity” models are common in estimation settings (4, 5). The RW model includes an additional feature strength parameter τ not present in those estimation models. More closely related to the RW model is work in multiple testing by Ingster and the authors (6–8), although the classification setting gives it a different meaning.]

Naïve application of the formula $w \propto \text{Cov}(Y, X)$ in the RW setting often leads to very poor results; the vector of empirical covariances $(\widehat{\text{Cov}}_{n,p}(Y, X(j)): j = 1, \dots, p)$ is very high-dimensional and contains mostly “noise” coordinates; the resulting naive classification weights $\hat{w}_{\text{naive}}(j) \propto \widehat{\text{Cov}}_{n,p}(Y, X(j))$ often produce correspondingly noisy decisions. The data glut seriously damages the applicability of such “textbook” approaches.

Feature Selection by Thresholding

Feature selection, that is, working only with an empirically selected subset of features, is a standard response to data glut. Here, and below, we suppose that feature correlations can be ignored and that features are standardized to variance one. We consider subset selectors based on the vector of feature Z-scores with components $Z(j) = n^{-1/2} \sum_i Y_i X_i(j)$, $j = 1, \dots, p$. These are the Z-scores of two-sided tests of $H_{0,j}$: $\text{Cov}(Y, X(j)) = 0$. Under our assumptions $Z \sim N(\theta, I_p)$, where $\theta = \sqrt{n}\mu$ and μ is the feature contrast vector. Features with nonzero $\mu(j)$ typically

Author contributions: D.L.D. and J.J. designed research, performed research, analyzed data, and wrote the paper.

The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

[†]To whom correspondence should be addressed. E-mail: donoho@stat.stanford.edu.

© 2008 by The National Academy of Sciences of the USA

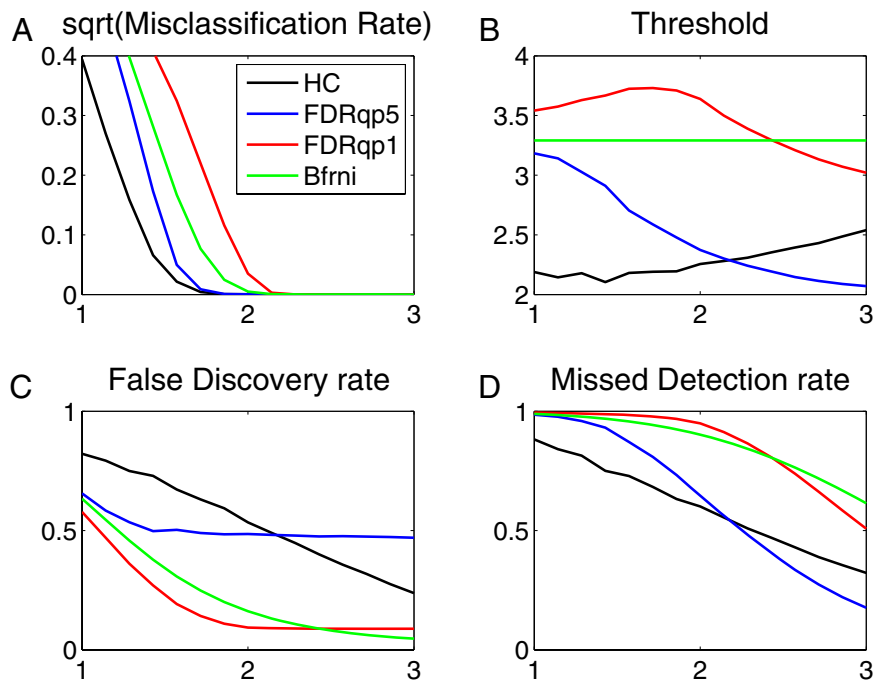


Fig. 2. Monte Carlo performance of thresholding rules in the RW model. (A–D) $P = 1,000$, $\varepsilon = 0.05$, and x axes display τ . (A) $\text{MCR}^{1/2}$. (B) Average threshold. (C) Average FDR. (D) Average MDR. Threshold procedures used: HC (black), Bonferroni (green), FDR ($q = .5$) (blue), FDRT ($q = .1$) (red). Averages from 1,000 Monte Carlo realizations.

features are indeed weak: they have expected Z -scores typically lower than some Z -scores of useless coordinates.

We compare HC thresholding with three other thresholding rules: (i) FDRT(.5), thresholding with false feature discovery rate (FDR) control parameter $q = 0.5$; (ii) FDRT(.1), thresholding with false feature discovery rate control parameter $q = 0.1$; and (iii) Bonferroni, setting the threshold so that the expected number of false features is 1. These three rules illustrate what we believe to be today's orthodox opinion, which strives to ensure that most features in the classification rule are truly useful, and to strictly control the number of useless features present in the trained classifier. Local false discovery rate control shares the same philosophy. We generated 1,000 Monte Carlo realizations at each choice of parameters. We present results in terms of the dimensionless parameter τ , which is independent of n ; if desired, the reader may choose to translate these results into the form $\mu_0 = \tau/\sqrt{n}$ for a conventional choice of n , such as $n = 40$. Fig. 2 presents the empirical average performance. As compared with traditional approaches, HCT has, in the case of weak signals, a lower threshold, a higher false-feature discovery rate, and lower missed-feature detection rate (MDR); the misclassification rate (MCR) is also improved. In these displays, as the signal strength τ increases, HCT increases, but FDRT decreases (for analysis of this phenomenon, see unpublished work).

HCT Functional and Ideal Thresholding

We now develop connections between HCT and other important notions.

HCT Functional. The *HCT functional* is, informally, the “threshold that HCT is trying to estimate.” More precisely, note that, in the $RW(\varepsilon, \tau)$ model, the empirical distribution function $F_{n,p}(t) = \text{Ave}_j I_{\{Z^{(j)} \leq t\}}$, approximates, for large p and n arbitrary, the theoretical CDF $F_{\varepsilon,\tau}(t) = (1 - \varepsilon) \Phi(t) + \varepsilon \Phi(t - \tau)$, $t \in \mathbf{R}$, where $\Phi(t) = P\{N(0, 1) \leq t\}$ is the standard

normal distribution. The HCT functional is the *result of the HCT recipe on systematically replacing $F_{n,p}(t)$ by $F_{\varepsilon,\tau}(t)$* .

We define the underlying *true positive rate*, $\text{TPR}(t)$; the *false positive rate*, $\text{FPR}(t)$; and the *positive rate*, $\text{PR}(t)$, in the natural way as the expected proportions of, respectively, the useful, the useless, and of all features, having Z -scores above threshold t . The HC objective functional can be rewritten (up to rescaling) as

$$\widetilde{HC} = \frac{\text{PR}(t) - \text{FPR}(t)}{\sqrt{\text{PR}(t)(1 - \text{PR}(t))}} = \frac{\varepsilon(\text{TPR}(t) - \text{FPR}(t))}{\sqrt{\text{PR}(t)(1 - \text{PR}(t))}}. \quad [2]$$

In the $RW(\varepsilon, \tau)$ model, we have $\text{TPR}(t; \varepsilon, \tau) = \Phi(t - \tau) + \Phi(-t - \tau)$, $\text{FPR}(t; \varepsilon, \tau) = 2\Phi(-t)$, and $\text{PR}(t; \varepsilon, \tau) = (1 - \varepsilon) \text{FPR}(t) + \varepsilon \text{TPR}(t)$. Let $t_0 = t_0(\varepsilon, \tau)$ denote the threshold corresponding to the maximization limit α_0 in Definition 2: $\text{PR}(t_0; \varepsilon, \tau) = \alpha_0$. The HCT functional solves a simple maximization in t :

$$T_{HC}(F_{\varepsilon,\tau}) = \text{argmax}_{t \geq t_0} \widetilde{HC}(t; \varepsilon, \tau). \quad [3]$$

Rigorous justification of this formula is supplied in ref. 12, showing that in the $RW(\varepsilon, \tau)$ model, $\hat{I}_{n,p}^{\text{HCT}}$ converges in probability to $T_{HC}(F_{\varepsilon,\tau})$ as p goes to infinity with n either fixed or increasing; so indeed, this is what HCT is “trying to estimate.”

Ideal Threshold. We now study the threshold that (if we only knew it!) would provide optimal classifier performance. Recall that, in our setting, the feature covariance is the identity $\Sigma = I_p$; the quantity $\text{Sep}(w; \mu) = w' \mu / \|w\|_2$ is a fundamental measure of linear classifier performance. The misclassification rate of the trained weights \hat{w} on independent test data with true contrast vector μ obeys

$$P\{\text{Error} | \text{Training Data}, \mu\} = \Phi(-\text{Sep}(\hat{w}; \mu)), \quad [4]$$

where again Φ is the standard normal $N(0, 1)$ CDF. Hence, maximizing Sep is a proxy for minimizing misclassification rate (for more details, see unpublished work).

Table 1. Error rates of standard classifiers on standard examples from Dettling (16)

| Method | ALL/reg | Col/reg | Pro/reg | m-reg | R |
|--------|-----------|------------|------------|-------|-----|
| Bagboo | 4.08/0.59 | 16.10/0.52 | 7.53/0 | 0.59 | 6 |
| Boost | 5.67/1 | 19.14/1 | 8.71/0.18 | 1 | 7.5 |
| RanFor | 1.92/0.02 | 14.86/0.32 | 9.00/0.22 | 0.32 | 2 |
| SVM | 1.83/0 | 15.05/0.35 | 7.88/0.05 | 0.35 | 3 |
| DLDA | 2.92/0.28 | 12.86/0 | 14.18/1 | 1 | 7.5 |
| KNN | 3.83/0.52 | 16.38/0.56 | 10.59/0.46 | 0.56 | 5 |
| PAM | 3.55/0.45 | 13.53/0.11 | 8.87/0.20 | 0.45 | 4 |
| HCT | 2.86/0.27 | 13.77/0.14 | 9.47/0.29 | 0.29 | 1 |

reg, regret; col, colon; Pro, prostate; m-reg, maximum regret; R, rank based on m-reg.

$\varepsilon \cdot (1 - \text{TPR})(t)$ is a Lipschitz function of the ROC coordinates, all three maximizers must offer similar performance.

The maximizer of Proxy_2 has a very elegant characterization, as the point in t where the *secant* to the ROC curve is double the *tangent* to the ROC curve, $\frac{\text{TPR}'}{\text{FPR}'} = \frac{\text{TPR}}{2\text{FPR}}$ at $t = t_{\text{Proxy}_2}$. The maximizer of Proxy_1 obeys a slightly more complex relationship $\frac{\text{TPR}'}{\text{FPR}'} = (2 \frac{\text{FPR}}{\text{TPR}} (1 - \varepsilon/2)(1 - \varepsilon) + \varepsilon)^{-1}$ at $t = t_{\text{Proxy}_1}$. For small enough ε , this nearly follows the same rule: secant $\approx 2 \times$ tangent.

For comparative purposes, FDR thresholding finds a point on the ROC curve with prescribed *secant*: $\frac{\text{TPR}}{\text{FPR}} = \frac{1 - \varepsilon}{\varepsilon} \varepsilon (q^{-1} - 1)$ at $t = -\varepsilon \varepsilon (q^{-1} - 1)$ at $t = t_{\text{FDR},q}$. Further, a *local* false discovery rate threshold yields a point on the ROC curve with prescribed *tangent* $\frac{\text{TPR}'}{\text{FPR}'} = \frac{1 - \varepsilon}{\varepsilon} (q^{-1} - 1)$ at $t = t_{\text{localFDR},q}$. Defining the *true discovery rate* $\text{TDR} \equiv 1 - \text{FDR}$, we see that HCT obeys $\frac{\text{FDR}}{\text{TDR}} \approx \frac{1}{2} \frac{\text{local FDR}}{\text{local TDR}}$ at $t = t_{\text{Proxy}_2}$. HCT and its proxies are thus visibly quite different from prescribing FDR or local FDR, which again underscores the distinction between avoiding false-feature selection and maximizing classifier performance.

Complements

Performance on Standard Datasets. In recent literature on classification methodology, a collection of six datasets has been used frequently for illustrating empirical classifier performance (16). We have reservations about the use of such data to illustrate HCT, because no one can say whether any specific such dataset is an example of rare/weak feature model. However, such comparisons are sure to be requested, so we report them here.

Of the standard datasets reported in ref. 16, three involve two-class problems of the kind considered here; these are the ALL (10), Colon (17), and Prostate (18) datasets. In ref. 17, 3-fold random training test splits of these datasets were considered, and seven well known classification procedures were implemented: Bagboost (16), LogitBoost (19), SVM (20), Random Forests (21), PAM (9), and the classical methods *DLDA* and *KNN*. We applied HCT in a completely out-of-the-box way by using definitions standard in the literature. HCT-hard, which uses feature weights based on hard thresholding of feature Z -scores, gave quite acceptable performance. For comparison, introduce the relative regret measure $\text{Regret}(A) = [\text{err}(A) - \min_{A'} \text{err}(A')]/[\max_{A'} \text{err}(A') - \min_{A'} \text{err}(A')]$. This compares the achieved error rate with the best and worst performance seen across algorithms. We report error rates and regrets side by side in Table 1, where rows 2–7 are from Dettling (16), row 8 is provided by Tibshirani, and row 9 is the result of HCT-hard.

Additionally, column 5 is the maximum regret across three different datasets, and column 6 is the rank based on the maximum regret. In the random-split test, HCT-hard was the minimax regret procedure, always being within 29% of the best known performance, whereas every other procedure was worse in relative performance in at least some cases.

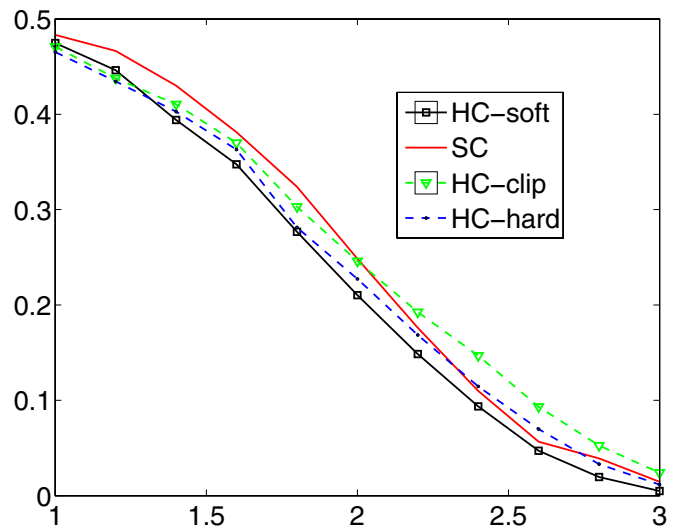


Fig. 5. Comparison of error rates by using Shrunken Centroids, threshold choice by cross-validation, and linear classifiers by using HCT-based threshold selection. Simulation assuming the RW model. Black, HCT-soft; red, Shrunken Centroids; green, HCT-clip; blue, HCT-hard. x axis displays τ .

It is worth remarking that HCT-based feature selection classifiers are radically simpler than all of the other methods being considered in this competition, requiring no tuning or cross-validation to achieve the presented results.

Comparison to Shrunken Centroids. The well known “Shrunken Centroids” (SC) algorithm (9) bears an interesting comparison to the procedures discussed here. In the two-class setting, SC amounts to linear classification with feature weights obtained from soft thresholding of feature Z -scores. Consequently, HCT-soft can be viewed as a modification to SC, choosing thresholds by HCT rather than cross-validation. We made a simulation study contrasting the performance of SC with HCT-hard, HCT-soft, and HCT-clip in the rare/weak features model. We conducted 100 Monte Carlo simulations, where we chose $p = 10,000$, $k = 100$ (so $\varepsilon = k/p = 0.01$), $n = 40$, and $\tau \in [1, 3]$. Over this range, the best classification error rate ranged from nearly 50%—scarcely better than ignorant guessing—to $<3\%$. Fig. 5 shows the results. Apparently, HCT-soft and SC behave similarly—with HCT-soft consistently better (here SC is implemented with a threshold picked by 10-fold cross-validations). However, HCT-soft and SC are not at all similar in computational cost at the training stage, as HCT-soft requires no cross-validation or tuning. The similarity of the two classifiers is, of course, explainable by using discussions above. Cross-validation is “trying” to estimate the ideal threshold, which the HCT functional also approximates. In Table 2, we tabulated the mean and standard deviation (SD) of HCT and cross-validated threshold selection (CVT). We see that CVT is on average larger than the HCT in this range of parameters. We also see that CVT has a signifi-

Table 2. Comparison of HCT and CVT

| τ | HCT mean | CVT mean | HCT SD | CVT SD |
|--------|----------|----------|--------|--------|
| 1.0 | 2.2863 | 3.8192 | 0.3746 | 1.9750 |
| 1.4 | 2.2599 | 3.3255 | 0.3401 | 1.7764 |
| 1.8 | 2.2925 | 3.0943 | 0.3400 | 1.3788 |
| 2.2 | 2.3660 | 2.6007 | 0.2921 | 0.8727 |
| 2.6 | 2.5149 | 2.5929 | 0.2644 | 0.5183 |
| 3.0 | 2.6090 | 2.9904 | 0.2698 | 0.5971 |

