

Topic 3: Classification

Eric B. Laber

Department of Statistical Science, Duke University

Statistics 561



On choices

I'd rather be rich than stupid.
—Pierre-Joseph Proudhon



On having to cover classification in two lectures

I needed to think last night. So I galloped into a wooded glen, and after punch dancing out my rage and suffering an extremely long and very painful fall, I realized what has to be done.

—Rod



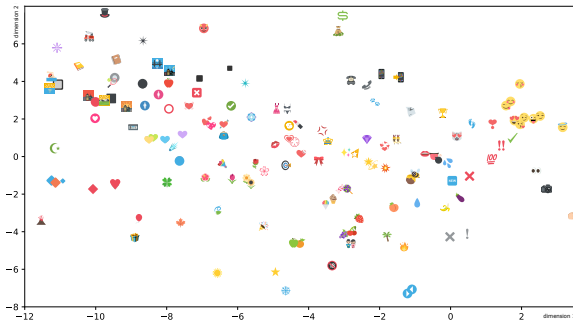
Warm-up (5 minutes)

- ▶ Explain to your group
 - ▶ What is classification? What are some canonical examples?
 - ▶ What is the separability problem?
 - ▶ What is zero-one loss?
- ▶ True or false
 - ▶ Regression is a misnomer in 'logistic regression'
 - ▶ Classification is in the "inner loop" of many RL algorithms
 - ▶ McDonald's is the most prevalent fast-food chain in NC



Classification: quick overview

- ▶ Input-output pairs where output is one of finitely many categories
 - ▶ High-risk v low-risk for complications in surgery
 - ▶ Spam v not spam
 - ▶ Handwritten digit recognition
- ▶ Example from sex-trafficking classification problem



Setup: binary classification

- ▶ Observe $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ comprising i.i.d. draws from P
 - ▶ Inputs: $\mathbf{X} \in \mathbb{R}^p$
 - ▶ Outputs: $Y \in \{-1, 1\}$, aka, label
- ▶ Classifier $c : \mathbb{R}^p \rightarrow \{-1, 1\}$ so that $c(\mathbf{x})$ is the predicted label at input $\mathbf{X} = \mathbf{x}$

0-1 loss

- Natural measure of classification performance is 0-1 loss

$$\ell_0(\mathbf{x}, y; c) \triangleq 1_{y \neq c(\mathbf{x})}$$

so that the expected loss (risk) is

$$\tau(c) \triangleq P\ell_0(\mathbf{X}, Y, c) = P1_{Y \neq c(\mathbf{X})} = P\{Y \neq c(\mathbf{X})\}$$

- E.g., linear classifier $c(\mathbf{x}; \beta) = \text{sign}(\mathbf{X}^\top \beta)$ which has loss

$$\tau(\beta) \triangleq P1_{Y \neq \text{sign}(\mathbf{X}^\top \beta)} = P1_{Y\mathbf{X}^\top \beta < 0}$$

Bayes classifier

- Let \mathcal{C} be class of all (msbl) maps from \mathbb{R}^p into $\{-1, 1\}$ and

$$c^{\text{opt}} = \arg \min_{c \in \mathcal{C}} \tau(c) = \arg \min_{c \in \mathcal{C}} P\{Y \neq c(\mathbf{X})\}$$

then

$$c^{\text{opt}}(\mathbf{x}) = \begin{cases} 1 & \text{if } P(Y = 1|\mathbf{X} = \mathbf{x}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

i.e., $c^{\text{opt}}(\mathbf{x}) = \text{sign}\{2q(\mathbf{x}) - 1\}$ w/ $q(\mathbf{x}) \triangleq P(Y = 1|\mathbf{X} = \mathbf{x})$

Blank page for notes



Blank page for notes



Probabilistic classifiers

- ▶ Natural approach to classification is to estimate construct an estimator $\hat{q}_n(\mathbf{x})$ of $q(\mathbf{x}) = P(Y = 1|\mathbf{X} = \mathbf{x})$
- ▶ Logistic regression posits a model of the form

$$q(\mathbf{x}; \boldsymbol{\beta}^*) = \text{expit}(\mathbf{x}^\top \boldsymbol{\beta}^*) = \frac{\exp(\mathbf{x}^\top \boldsymbol{\beta}^*)}{1 + \exp(\mathbf{x}^\top \boldsymbol{\beta}^*)}$$

Derive log-likelihood for logistic regression



Blank page for notes



In-class exercise (5 min)

- ▶ In your group derive a stochastic gradient descent algorithm for logistic regression: then on to `logistic_regression.R`

Separability problem

- ▶ If classes linearly separable estimates diverge
- ▶ Use ridge with logistic regression¹
- ▶ Penalized negative log-likelihood (0-1 coding for simplicity)

$$\ell_n(\beta) = -\mathbb{P}_n[\mathbf{X}^\top \beta Y + \{1 - q(\mathbf{X}; \beta)\}] + \lambda \|\beta\|^2$$

differentiate wrt β to obtain

$$-\mathbb{P}_n \mathbf{X} \{Y - q(\mathbf{X}; \beta)\} + 2\lambda \beta$$

¹That's right, I'm stating this without qualification of any kind. What's it to you? That's what I thought.



Warm-up (5 min)

- ▶ What's the definition of a convex function?
- ▶ Show that if $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is convex then $\beta \mapsto P\phi(Y\mathbf{X}^\top\beta)$ is convex (note: we're not assuming ϕ is differentiable)

Large-margin classifiers

- ▶ Decision boundary $\{\mathbf{x} : q(\mathbf{x}) = 1/2\}$
- ▶ Don't need $P(Y|\mathbf{X} = \mathbf{x})$ just $\text{sign}\{2q(\mathbf{x}) - 1\}$, i.e., we don't need to know the boundary just which side we're on
- ▶ Idea! Construct loss functions that penalize distance from correct side of boundary

Surrogate loss functions: laber draws a picture

Surrogate loss functions

- ▶ Common surrogate loss functions include
 - ▶ Squared error loss: $\ell(\mathbf{x}, y; \beta) = (1 - y\mathbf{x}^\top \beta)^2$
 - ▶ Exponential loss: $\ell(\mathbf{x}, y; \beta) = \exp(-y\mathbf{x}^\top \beta)$
 - ▶ Hinge loss: $\ell(\mathbf{x}, y; \beta) = (1 - y\mathbf{x}^\top \beta)_+$
 - ▶ Logistic loss: $\ell(\mathbf{x}, y; \beta) = \log \{1 + \exp(-y\mathbf{x}^\top \beta)\}$
- ▶ Can also consider penalized versions of these loss functions

Nice and relaxing derivation

- Show $\hat{\beta}_n = \arg \min_{\beta} \mathbb{P}_n \log \{1 + \exp(-Y\mathbf{X}^\top \beta)\}$ recovers the logistic regression estimator

Blank page for notes



Price of using surrogates

- ▶ Pro tip: don't google "price of using surrogates," it will lead you down a fascinating (but time consuming) rabbit hole about surrogate mothers and the ethics of paying them
- ▶ We replaced the loss of interest (0-1) with a convex surrogate what are the statistical consequences?
 - ▶ Will we still recover Bayes classifier?
 - ▶ Can any surrogate be used? Are some better than others?
 - ▶ Are there any additional benefits/drawbacks of using a surrogate?

Pointwise consistency aka Fisher Consistency

- ▶ Idea: compare population minimizer of surrogate with minimizer of 0-1 loss, if these agree then we say that the surrogate classifier is Fisher consistent
- ▶ Note all the surrogates we considered are functions of $yf(\mathbf{x})$ for some real-valued function f , e.g., in linear case $y\mathbf{x}^\top\beta$
 - ▶ Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be generic fn indexing $\mathbf{x} \mapsto \text{sign} \{f(\mathbf{x})\}$
 - ▶ Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ denote a surrogate acting on the margin
 - ▶ Squared error loss: $\phi(\alpha) = (1 - \alpha)^2$
 - ▶ Exponential loss: $\phi(\alpha) = \exp(-\alpha)$
 - ▶ Hinge loss: $\phi(\alpha) = (1 - \alpha)_+$
 - ▶ Logistic loss: $\phi(\alpha) = \log \{1 + \exp(-\alpha)\}$

Pointwise consistency aka Fisher Consistency: defns

- ▶ $R(f) \triangleq P1_{Yf(\mathbf{X}) < 0}$ the misclassification error at f
- ▶ R^* the Bayes error and define excess risk at f as $R(f) - R^*$
- ▶ $R_\phi(f) \triangleq P\phi\{Yf(\mathbf{X})\}$ to be the ϕ -risk at f
- ▶ $R_\phi^* \triangleq \inf_f R_\phi(f)$ and define the excess ϕ -risk at f as $R_\phi(f) - R_\phi^*$

Relating excess risks

- Goal is to find a function ψ such that $\psi(\tau) \rightarrow 0$ iff $\tau \rightarrow 0$ and

$$\psi \{R(f) - R^*\} \leq R_\phi(f) - R_\phi^*$$

What does such a result buy us?

Conditional ϕ -risk

- Conditional ϕ -risk

$$\mathbb{E} [\phi \{ Yf(\mathbf{X}) \} | \mathbf{X} = \mathbf{x}] = q(\mathbf{x})\phi \{ f(\mathbf{x}) \} + \{1 - q(\mathbf{x})\} \phi \{ -f(\mathbf{x}) \}$$

- It will be convenient to define

$$C_q(\alpha) \triangleq q\phi(\alpha) + (1 - q)\phi(-\alpha)$$

and

$$H(q) \triangleq \inf_{\alpha} C_q(\alpha) = \inf_{\alpha} \{ q\phi(\alpha) + (1 - q)\phi(-\alpha) \}$$

Conditional ϕ -risk cont'd

- Practice with the notation: show that

$$R_{\phi}^* = \inf_f R_{\phi}(f) = PH\{q(\mathbf{X})\}$$

Blank page for notes



Blank page for notes



Classification-calibration

- ▶ For $q \in [0, 1]$ define

$$H^-(q) \triangleq \inf_{\alpha: \alpha(2q-1) \leq 0} C_q(\alpha)$$

to be optimal conditional ϕ -risk if you disagree with Bayes rule

- ▶ We say that a loss function ϕ is classification calibrated if

$$H^-(q) > H(q)$$

for all $q \neq 1/2$. Thus, ϕ is CC if disagreeing with Bayes rule increases ϕ -risk.

Ex. classification-calibrated loss

- Show exponential loss is classification-calibrated

Blank page for notes



Blank page for notes



Bounding the excess ϕ -risk

Theorem

Let ϕ be convex. Then ϕ is classification calibrated iff it is differentiable at 0 and $\phi'(0) < 0$. Furthermore, if ϕ is classification calibrated then

$$\psi(\theta) = \phi(0) - H \left(\frac{1 + \theta}{2} \right)$$

satisfies

$$\psi \{R(f) - R^*\} \leq R_\phi(f) - R_\phi^*,$$

where $\psi(\theta) \rightarrow 0$ iff $\theta \rightarrow 0$.

Statistical inference

► Convexity buys you a lot, many derivations much simpler

► High-level overview of assumptions

(A1) $\ell(\mathbf{x}, y; \beta)$ is convex in β for each (\mathbf{x}, y)

(A2) $Q(\beta) = P\ell(\mathbf{X}, Y; \beta)$ exists and is finite for all β

(A3) $\beta^* = \arg \min Q(\beta)$ exists and is unique

(A4) $Q(\beta)$ is twice continuously differentiable in a nbrhd of β^* and $H = \nabla^2 Q(\beta^*)$ is positive definite

Theorem

Assume (A1)-(A4) and let $\Omega = P\nabla\ell(\mathbf{X}, Y; \beta^*)\nabla\ell(\mathbf{X}, Y; \beta^*)^\top$ then

$$\sqrt{n}(\hat{\beta}_n - \beta^*) \rightsquigarrow \text{Normal}(0, H^{-1}\Omega H^{-1}).$$



Statistical inference cont'd

- ▶ Asymptotic normality \Rightarrow we can use standard methods of inference for β^* as in the regression case, e.g., Wald-type intervals etc.
- ▶ Unlike regression case, measures of performance are not well-behaved
- ▶ Recall our measures of performance
 - ▶ Population-level error: $R \triangleq P1_{Y\mathbf{X}^\top\beta^* < 0}$
 - ▶ Conditional error: $C(\hat{\beta}_n) \triangleq P1_{Y\mathbf{X}^\top\hat{\beta}_n < 0}$
 - ▶ Average error: $A_n \triangleq \mathbb{E}C(\hat{\beta}_n)$

Comparing measures of performance

- ▶ Fact: the three measures of performance need not coincide even in infinite samples

Blank page for notes



Blank page for notes



Potpourri: local linear models

- ▶ Let $K_\sigma : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}_+$ be a kernel fn, e.g.,
 $K_\sigma(\mathbf{x}, \mathbf{x}') = \exp \{ -\|\mathbf{x} - \mathbf{x}'\|^2 / \sigma^2 \}$
- ▶ Local linear model: estimate coefficient fn $\hat{\beta}_n(\mathbf{x})$ for each \mathbf{x} via

$$\hat{\beta}_n(\mathbf{x}) = \arg \min_{\beta} \mathbb{P}_n (Y - \mathbf{X}^\top \beta)^2 K_\sigma(\mathbf{X}, \mathbf{x}),$$

tune σ using CV etc.

- ▶ Natural (and easy!) extension of linear model
- ▶ Studying coeff fn can generate insights about how the mean of Y changes in different regions of the input space (but this is not trivial if p is large)

Potpourri: local large-margin classifiers

- ▶ Let K_σ be a kernel and ϕ a surrogate loss function, define

$$\hat{\beta}_n(\mathbf{x}) = \arg \min_{\beta} \mathbb{P}_n \phi(Y \mathbf{X}^\top \beta) K_\sigma(\mathbf{X}, \mathbf{x})$$

- ▶ E.g., local logistic regression with ridge penalty

$$\hat{\beta}_n^\lambda(\mathbf{x}) = \arg \min_{\beta} \mathbb{P}_n \log \{1 + \exp(-Y \mathbf{X}^\top \beta)\} K_\sigma(\mathbf{X}, \mathbf{x}) + \lambda \|\beta\|^2$$

Potpourri: trees

- ▶ Canonical local model: classification and regression trees
- ▶ laber draws a tree:

- ▶ Tree is an additive model with form $f(\mathbf{x}) = \sum_{m=1}^M \beta_m \mathbf{1}_{\mathbf{x} \in R_m}$
where R_1, \dots, R_M partition the input space

Potpourri: trees cont'd

- Note that if regions R_1, \dots, R_m were given and Y were continuous then

$$\hat{\beta}_n = \arg \min_{\beta \in \mathbb{R}^M} \mathbb{P}_n \{Y - \mathcal{I}(\mathbf{X})^\top \beta\}^2,$$

where $\mathcal{I}(\mathbf{x}) = (1_{\mathbf{x} \in R_1}, \dots, 1_{\mathbf{x} \in R_M})^\top$ and

$$\hat{\beta}_j = \mathbb{P}_n Y 1_{\mathbf{x} \in R_j} / \mathbb{P}_n 1_{\mathbf{x} \in R_j}$$

- Note that this is a local constant model with

$$\hat{\beta}_n(\mathbf{x}) = \arg \min_{\beta} \mathbb{P}_n (Y - \beta)^2 K_{\sigma}(\mathbf{X}, \mathbf{x})$$

and $K_{\sigma}(\mathbf{x}, \mathbf{x}') = 1_{\mathbf{x}, \mathbf{x}' \in R_j}$, for some j

Potpourri: trees cont'd

- ▶ Estimating optimal partition R_1, \dots, R_M generally combinatorially hard
- ▶ Two approaches:
 - ▶ A1: restrict regions to rectangles and optimize greedily
 - ▶ A2: randomly generate a bunch of (overlapping) regions aka tile coding

Potpourri notes: greedy optimization

Potpourri notes: tile coding

Warm-up (5 minutes)

- ▶ Explain to your stats group
 - ▶ What are random forests?
 - ▶ What is the Gram matrix?
 - ▶ What is a Voronoi partition?
- ▶ True or false
 - ▶ Voronoi partitions are credited to Georgy Voronoy
 - ▶ Voronoi partitions were a key part of a CSI-Cyber episode
 - ▶ Voldemort should have done a better job hiding his horcruxes

Kernel Based Random Forests (KeRFs)

- ▶ Focus on regression setting
 - ▶ Goals for today:
 - ▶ Formally define a class of random forests
 - ▶ Show random forests representable as kernels

Setup

- ▶ Observe $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ w/ $\mathbf{X} \in \mathbb{R}^p$ and $Y \in \mathbb{R}$
- ▶ Creating a forest
 - ▶ Sequence of trees $\mathbf{x} \mapsto \mu_n(\mathbf{x}; \Theta_j), j = 1, \dots, M$, where $\Theta_1, \dots, \Theta_M$ drawn i.i.d. from P_Θ
 - ▶ Θ encodes actions like randomly subsetting/bootstrapping, selecting variables for splits, choosing split points, etc.
 - ▶ Prediction at $\mathbf{X} = \mathbf{x}$

$$\mu_{n,M}(\mathbf{x}; \bar{\Theta}_M) = \frac{1}{M} \sum_{j=1}^M \mu_n(\mathbf{x}; \Theta_j),$$

where $\bar{\Theta}_M = (\Theta_1, \dots, \Theta_M)$

Some forests

- ▶ Breiman: bootstrap data, randomly select set of variables to split on, choose splits to minimize variance
- ▶ Extremely randomized trees (ERTs): randomly select set of variables to split on, small set of candidate splits chosen at random, best chosen in terms of minimizing variance
- ▶ Gump: simple man with a big heart lives an extraordinary life. I'm not crying, you're crying. FINE, I AM CRYING. SORRY FOR HAVING EMOTIONS YOU CYLON.

Formalizing random forests

- ▶ Let $A_n(\mathbf{x}, \Theta_j)$ be the node (subset of \mathbb{R}^p) in the tree indexed by Θ_j to which \mathbf{x} belongs
- ▶ Random forest estimator

$$\begin{aligned}\mu_{M,n}(\mathbf{x}; \bar{\Theta}_M) &= \frac{1}{M} \sum_{j=1}^M \left\{ \sum_{i=1}^n \frac{Y_i 1_{\mathbf{x}_i \in A_n(\mathbf{x}, \Theta_j)}}{\sum_{\ell=1}^n 1_{\mathbf{x}_\ell \in A_n(\mathbf{x}, \Theta_j)}} \right\} \\ &= \frac{1}{M} \sum_{j=1}^M \left\{ \sum_{i=1}^n \frac{Y_i 1_{\mathbf{x}_i \in A_n(\mathbf{x}, \Theta_j)}}{N_n(\mathbf{x}; \Theta_j)} \right\},\end{aligned}$$

where $N_n(\mathbf{x}; \Theta_j) = \sum_{\ell=1}^n 1_{\mathbf{x}_\ell \in A_n(\mathbf{x}, \Theta_j)}$ is the number of cellmates, i.e., cellies, of \mathbf{x} in tree j

Formalizing random forests cont'd

- ▶ Weights for Y_i are $1_{\mathbf{x}_i \in A_n(\mathbf{x}; \boldsymbol{\Theta}_j)}$ which are large when $N_n(\mathbf{x}; \boldsymbol{\Theta}_j)$ is small, i.e., cells with little data influential
- ▶ Scornet et al., proposed KeRF (Kernels based on Random Forests)

$$\begin{aligned}\tilde{\mu}_{M,n}(\mathbf{x}; \overline{\boldsymbol{\Theta}_j}) &\triangleq \frac{\sum_{j=1}^M \sum_{i=1}^n Y_i 1_{\mathbf{x}_i \in A_n(\mathbf{x}; \boldsymbol{\Theta}_j)}}{\sum_{j=1}^M N_n(\mathbf{x}; \boldsymbol{\Theta}_j)} \\ &= \frac{\sum_{j=1}^M \sum_{i=1}^n Y_i 1_{\mathbf{x}_i \in A_n(\mathbf{x}; \boldsymbol{\Theta}_j)}}{\sum_{j=1}^M \sum_{i=1}^n 1_{\mathbf{x}_i \in A_n(\mathbf{x}; \boldsymbol{\Theta}_j)}}\end{aligned}$$

Random Forests as kernel estimators

► Note that

$$\tilde{\mu}_{M,n}(\mathbf{x}; \overline{\Theta_j}) = \frac{\sum_{i=1}^n Y_i K_{M,n}(\mathbf{x}, \mathbf{X}_i)}{\sum_{i=1}^n K_{M,n}(\mathbf{x}, \mathbf{X}_i)},$$

where $K_{M,n}(\mathbf{x}, \mathbf{z}) = M^{-1} \sum_{j=1}^M 1_{\mathbf{z} \in A_n(\mathbf{x}; \Theta_j)}$ is the fraction of trees in which \mathbf{x}, \mathbf{z} are in the same cell

Comments on Kernels based on Random Forests

- ▶ Kernel selection is a major problem esp. in high-dim settings
- ▶ Selection of a high-quality kernel can drastically improve predictive performance
- ▶ KeRF is a data-adaptive kernel
 - ▶ Variants that use outcome are especially appealing in that they may identify important variables for prediction
 - ▶ Interesting benefit: use variables that exist in training data to form kernel (laber discusses search engine example)

Blank page for notes



Blank page for notes



Exercise (2 min)

- ▶ How might one use the preceding idea for evil?!

Relating KeRF and RF

(A1) Suppose there exists sequences $\{a_n\}_{n \geq 1}$ and $\{b_n\}_{n \geq 1}$ such that $a_n \leq N_n(\mathbf{x}; \Theta) \leq b_n$ almost surely for all n

Theorem

Let $\mu_{M,n}(\mathbf{x}; \bar{\Theta}_M)$ be a random forest algorithm that satisfies (A1).
Then for each \mathbf{x}

$$\left| \frac{\mu_{M,n}(\mathbf{x}; \bar{\Theta}_j)}{\tilde{\mu}_{M,n}(\mathbf{x}; \bar{\Theta}_j)} - 1 \right| \leq \frac{b_n - a_n}{a_n},$$

almost surely.

Sparsity of the RF Kernel

- ▶ RF Kernel often induces a sparse Gram matrix (why?)
- ▶ Useful in very large prediction problems esp as an emulator for more complex (and expensive) models at runtime

Thank you.

`eric.laber@duke.edu`

`laber-labs.com`

