

Math 583 Exam 3 is on Wednesday, Dec. 6 with emphasis on homeworks 9-11 and quizzes 9-11. You are allowed 9 sheets of notes and a calculator. Any needed tables will be provided. CHECK FORMULAS: YOU ARE RESPONSIBLE FOR ANY ERRORS ON THIS HANDOUT! Final: Wednesday, December 13 2:45-4:45.

113) In *supervised classification*, there are G known groups or populations and m test cases. Each case is assigned to exactly one group based on its measurements \mathbf{w}_i . Assume that for each population there is a probability density function (pdf) $f_j(\mathbf{z})$ where \mathbf{z} is a $p \times 1$ vector and $j = 1, \dots, G$. Hence if the random vector \mathbf{x} comes from population j , then \mathbf{x} has pdf $f_j(\mathbf{z})$. Assume that there is a random sample of n_j cases $\mathbf{x}_{1,j}, \dots, \mathbf{x}_{n_j,j}$ for each group. The $n = \sum_{j=1}^G n_j$ cases make up the training data. Let $(\bar{\mathbf{x}}_j, \mathbf{S}_j)$ denote the sample mean and covariance matrix for each group. Let the i th test case \mathbf{w}_i be a new $p \times 1$ random vector from one of the G groups, but the group is unknown. *Discriminant analysis = classification* attempts to allocate (classify) the \mathbf{w}_i to the correct groups for $i = 1, \dots, m$.

Assume that costs of correct and incorrect allocation are unknown or equal, and assume that the probabilities $\pi_j = \rho_j(\mathbf{w}_i)$ that \mathbf{w}_i is in group j are unknown or equal: $\pi_j = 1/G$ for $j = 1, \dots, G$. Often it is assumed that the G groups have the same covariance matrix $\Sigma_{\mathbf{x}}$. Then the pooled covariance matrix estimator is

$$\mathbf{S}_{pool} = \frac{1}{n - G} \sum_{j=1}^G (n_j - 1) \mathbf{S}_j$$

where $n = \sum_{j=1}^G n_j$. Let $(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j)$ be the estimator of multivariate location and dispersion for the j th group, e.g. the sample mean and sample covariance matrix $(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) = (\bar{\mathbf{x}}_j, \mathbf{S}_j)$.

114) Assume that $G = 2$ and that there is a group 0 and a group 1. Let $\rho(\mathbf{w}) = P(\mathbf{w} \in \text{group 1})$. Let $\hat{\rho}(\mathbf{w})$ be the logistic regression (LR) estimate of $\rho(\mathbf{w})$. Logistic regression produces an estimated sufficient predictor $ESP = \hat{\alpha} + \hat{\boldsymbol{\beta}}^T \mathbf{w}$. Then

$$\hat{\rho}(\mathbf{w}) = \frac{e^{ESP}}{1 + e^{ESP}} = \frac{\exp(\hat{\alpha} + \hat{\boldsymbol{\beta}}^T \mathbf{w})}{1 + \exp(\hat{\alpha} + \hat{\boldsymbol{\beta}}^T \mathbf{w})}$$

The *logistic regression discriminant rule* allocates \mathbf{w} to group 1 if $\hat{\rho}(\mathbf{w}) \geq 0.5$ and allocates \mathbf{w} to group 0 if $\hat{\rho}(\mathbf{w}) < 0.5$. Equivalently, the LR rule allocates \mathbf{w} to group 1 if $ESP \geq 0$ and allocates \mathbf{w} to group 0 if $ESP < 0$. The response plot is as in point 107).

115) Given LR output, as shown below in symbols and for a real data set, and given \mathbf{x} to classify, be able to a) compute ESP, b) classify \mathbf{x} in group 0 or group 1, c) compute $\hat{\rho}(\mathbf{x})$.

Label	Estimate	Std. Error	Est/SE	p-value
Constant	$\hat{\alpha}$	$se(\hat{\alpha})$	$z_{o,0}$	for Ho: $\alpha = 0$
x_1	$\hat{\beta}_1$	$se(\hat{\beta}_1)$	$z_{o,1} = \hat{\beta}_1 / se(\hat{\beta}_1)$	for Ho: $\beta_1 = 0$
\vdots	\vdots	\vdots	\vdots	\vdots
x_p	$\hat{\beta}_p$	$se(\hat{\beta}_p)$	$z_{o,p} = \hat{\beta}_p / se(\hat{\beta}_p)$	for Ho: $\beta_p = 0$

Binomial Regression Kernel mean function = Logistic
 Response = Status, Terms = (Bottom Left), Trials = Ones
 Coefficient Estimates

Label	Estimate	Std. Error	Est/SE	p-value
Constant	-389.806	104.224	-3.740	0.0002
Bottom	2.26423	0.333233	6.795	0.0000
Left	2.83356	0.795601	3.562	0.0004

116) Let $\hat{\Sigma}_{pool}$ be a pooled dispersion estimator such as \mathbf{S}_{pool} . Then the *linear discriminant rule* is allocate \mathbf{w} to the group with the largest value of

$$d_j(\mathbf{w}) = \hat{\boldsymbol{\mu}}_j^T \hat{\Sigma}_{pool}^{-1} \mathbf{w} - \frac{1}{2} \hat{\boldsymbol{\mu}}_j^T \hat{\Sigma}_{pool}^{-1} \hat{\boldsymbol{\mu}}_j = \hat{\alpha}_j + \hat{\boldsymbol{\beta}}_j^T \mathbf{w}$$

where $j = 1, \dots, G$. *Linear discriminant analysis* (LDA) uses $(\hat{\boldsymbol{\mu}}_j, \hat{\Sigma}_{pool}) = (\bar{\mathbf{x}}_j, \mathbf{S}_{pool})$. LDA is robust to nonnormality and somewhat robust to the assumption of equal population covariance matrices. LDA can be useful if the population dispersion matrices are equal: $\Sigma_j \equiv \Sigma$ for $j = 1, \dots, G$. LDA can be useful if some of the $n_j < 10p$. In high dimensions, replace \mathbf{S}_{pool} by a regularized estimator such as $\hat{\Sigma}_{pool} = \mathbf{I}_p$ or $\hat{\Sigma}_{pool} = \text{diag}(\mathbf{S}_{pool})$.

117) Suppose there is training data \mathbf{x}_{ij} for $i = 1, \dots, n_j$ for group j . Hence it is known that \mathbf{x}_{ij} came from group j where there are $G \geq 2$ groups. Use the discriminant analysis method to classify the training data. If m_j of the n_j group j cases are correctly classified, then the *apparent error rate for group j* is $1 - m_j/n_j$. If $m_A = \sum_{j=1}^G m_j$ of the $n = \sum_{j=1}^G n_j$ cases were correctly classified, then the *apparent error rate AER* = $1 - m_A/n$.

118) Get apparent error rates for LDA with the following commands.

```
out2 <- lda(x,group)
1-mean(predict(out2,x)$class==group)
```

Get the AERs for the methods that use variables x_1, x_3 , and x_7 with the following commands.

```
out <- lda(x[,c(1,3,7)],group)
1-mean(predict(out,x[,c(1,3,7)])$class==group)
```

Get the AERs for the methods that leave out variables x_1, x_4 , and x_5 with the following commands.

```
out <- lda(x[, -c(1,4,5)],group)
1-mean(predict(out,x[, -c(1,4,5)])$class==group)
```

119) Expect the apparent error rate to be too low: the method works better on the training data than on the new test data to be classified.

120) leave one out cross validation (CV): for $i = 1, \dots, n$ where the training data has n cases, compute the discriminant rule with case i left out and see if the rule correctly

classifies case i . Let m_C be the number of cases correctly classified. Then the CV error rate is $1 - m_C/n$.

121) leave out a validation set V (data splitting): Suppose the training data has n cases. Randomly select a subset V of n_v cases to be left out when computing the discriminant rule. Hence $n - n_v$ cases are used to compute the discriminant rule. Let m_L be the number of cases from subset V that are correctly classified. Then the “validation set” error rate is $1 - m_L/n_v$. Here n_v should be large enough to get a good rate. Often use n_v between $0.1n$ and $0.5n$.

122) The k -fold CV is similar to that for MLR. Let m_k be the number of cases that are correctly classified. Then the k -fold CV error rate is $1 - m_k/n$.

multiple testing

123) Suppose there are m tests for $g \in I = \{1, \dots, m\}$ with $H_{0g} : \mu_{gA} = \mu_{gB}$ versus $H_{1g} : \mu_{gA} \neq \mu_{gB}$. The m tests result in m estimated p-values (pvalues) $\hat{p}_1, \dots, \hat{p}_m$. Let \hat{R} give a set of indices i for hypotheses H_{0i} , that are rejected. The empty set $\hat{R} = \emptyset$ is possible. Often $\alpha = 0.1, 0.05$, or 0.01 . Let $\hat{p}_{(1)}, \dots, \hat{p}_{(m)}$ be the order statistics of the \hat{p}_i .

124) Let $I_0 = \{i \in I : H_{0i} \text{ is true}\}$. We call *false positive* the indices $i \in \hat{R} \cap I_0$ and *true positive* the indices $i \in \hat{R} \cap I_0^c$. Let FP = number of false positives and TP = number of true positives. The false discovery proportion $FDP = \frac{FP}{FP + TP}$ with $0/0 = 0$. The false discovery rate $FDR = E \left[\frac{FP}{FP + TP} I(FP + TP \geq 1) \right]$.

125) The Bonferroni correction uses $\hat{R}_B = \{i \in I : \hat{p}_i \leq \alpha/m\}$. Often $\alpha = 0.1, 0.05$, or 0.01 .

126) The Benjamani-Hochberg procedure uses $\hat{R}_{BH} = \emptyset$ if $\{k : \hat{p}_{(k)} \leq \alpha k/m\} = \emptyset$. Otherwise, let $\hat{k} = \max\{k \in I : \hat{p}_{(k)} \leq \alpha k/m\}$, and $\hat{R}_{BH} = \{i \in I : \hat{p}_i \leq \alpha \hat{k}/m\} = \{i_1, \dots, i_{\hat{k}} \text{ corresponding to } \hat{p}_{(1)}, \dots, \hat{p}_{(\hat{k})}\}$.

Note that the Bonferroni correction uses α/m while 126) uses $\hat{k}\alpha/m$.

ch. 10

127) The **multivariate linear regression mreg1 model**

$$\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \boldsymbol{\epsilon}_i$$

for $i = 1, \dots, n$ has $m \geq 2$ response variables Y_1, \dots, Y_m and p predictor variables x_1, x_2, \dots, x_p where $x_1 \equiv 1$ is the trivial predictor. The i th case is $(\mathbf{x}_i^T, \mathbf{y}_i^T) = (1, x_{i2}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})$ where the 1 could be omitted. The model is written in matrix form as $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$ where the matrices are defined below. The model has $E(\boldsymbol{\epsilon}_k) = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\epsilon}_k) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} = (\sigma_{ij})$ for $k = 1, \dots, n$. Also $E(\mathbf{e}_i) = \mathbf{0}$ while $\text{Cov}(\mathbf{e}_i, \mathbf{e}_j) = \sigma_{ij} \mathbf{I}_n$ for $i, j = 1, \dots, m$ where \mathbf{I}_n is the $n \times n$ identity matrix and \mathbf{e}_i is defined below. Then the $p \times m$ coefficient matrix $\mathbf{B} = [\boldsymbol{\beta}_1 \ \boldsymbol{\beta}_2 \ \dots \ \boldsymbol{\beta}_m]$ and the $m \times m$ covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$ are to be estimated, and $E(\mathbf{Z}) = \mathbf{X}\mathbf{B}$ while $E(Y_{ij}) = \mathbf{x}_i^T \boldsymbol{\beta}_j$. The $\boldsymbol{\epsilon}_i$ are assumed to be iid. The data matrix $\mathbf{W} = [\mathbf{X} \ \mathbf{Y}]$ except usually the first column $\mathbf{1}$ of \mathbf{X} is omitted. The $n \times m$ matrix

$$\mathbf{Z} = \begin{bmatrix} Y_{1,1} & Y_{1,2} & \dots & Y_{1,m} \\ Y_{2,1} & Y_{2,2} & \dots & Y_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{n,1} & Y_{n,2} & \dots & Y_{n,m} \end{bmatrix} = [\mathbf{Y}_1 \ \mathbf{Y}_2 \ \dots \ \mathbf{Y}_m] = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_n^T \end{bmatrix}.$$

The $n \times p$ design matrix of predictor variables is

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{bmatrix} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_p] = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

where $\mathbf{v}_1 = \mathbf{1}$. The $p \times m$ matrix

$$\mathbf{B} = \begin{bmatrix} \beta_{1,1} & \beta_{1,2} & \cdots & \beta_{1,m} \\ \beta_{2,1} & \beta_{2,2} & \cdots & \beta_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{p,1} & \beta_{p,2} & \cdots & \beta_{p,m} \end{bmatrix} = [\beta_1 \quad \beta_2 \quad \cdots \quad \beta_m].$$

The $n \times m$ matrix

$$\mathbf{E} = \begin{bmatrix} \epsilon_{1,1} & \epsilon_{1,2} & \cdots & \epsilon_{1,m} \\ \epsilon_{2,1} & \epsilon_{2,2} & \cdots & \epsilon_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \epsilon_{n,1} & \epsilon_{n,2} & \cdots & \epsilon_{n,m} \end{bmatrix} = [\mathbf{e}_1 \quad \mathbf{e}_2 \quad \cdots \quad \mathbf{e}_m] = \begin{bmatrix} \boldsymbol{\epsilon}_1^T \\ \vdots \\ \boldsymbol{\epsilon}_n^T \end{bmatrix}.$$

Considering the i th row of \mathbf{Z} , \mathbf{X} and \mathbf{E} shows that $\mathbf{y}_i^T = \mathbf{x}_i^T \mathbf{B} + \boldsymbol{\epsilon}_i^T$.

128) Each response variable in a `mreg1` model follows a multiple linear regression model $\mathbf{Y}_j = \mathbf{X}\boldsymbol{\beta}_j + \mathbf{e}_j$ for $j = 1, \dots, m$ where it is assumed that $E(\mathbf{e}_j) = \mathbf{0}$ and $\text{Cov}(\mathbf{e}_j) = \sigma_{jj}\mathbf{I}_n$. Hence the errors corresponding to the j th response are uncorrelated with variance $\sigma_j^2 = \sigma_{jj}$. Notice that the **same design matrix** \mathbf{X} of predictors is used for each of the m models, but the j th response variable vector \mathbf{Y}_j , coefficient vector $\boldsymbol{\beta}_j$ and error vector \mathbf{e}_j change and thus depend on j .

Now consider the i th case $(\mathbf{x}_i^T, \mathbf{y}_i^T)$ which corresponds to the i th row of \mathbf{Z} and the i th row of \mathbf{X} . Then

$$\begin{bmatrix} Y_{i1} = \beta_{11}x_{i1} + \cdots + \beta_{p1}x_{ip} + \epsilon_{i1} = \mathbf{x}_i^T \boldsymbol{\beta}_1 + \epsilon_{i1} \\ Y_{i2} = \beta_{12}x_{i1} + \cdots + \beta_{p2}x_{ip} + \epsilon_{i2} = \mathbf{x}_i^T \boldsymbol{\beta}_2 + \epsilon_{i2} \\ \vdots \\ Y_{im} = \beta_{1m}x_{i1} + \cdots + \beta_{pm}x_{ip} + \epsilon_{im} = \mathbf{x}_i^T \boldsymbol{\beta}_m + \epsilon_{im} \end{bmatrix}.$$

129) The **OLS estimators** for the `mreg1` model are

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z} = [\hat{\boldsymbol{\beta}}_1 \quad \hat{\boldsymbol{\beta}}_2 \quad \cdots \quad \hat{\boldsymbol{\beta}}_m].$$

130) Let $\mathbf{X} = (\mathbf{1} \quad \mathbf{X}_1)$. The **mreg2** model is

$$\mathbf{y}_i = \boldsymbol{\alpha} + \mathbf{B}_S^T \mathbf{x}_i + \epsilon_i$$

for $i = 1, \dots, n$ with

$$\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E} = \mathbf{X} \begin{bmatrix} \boldsymbol{\alpha}^T \\ \mathbf{B}_S \end{bmatrix} + \mathbf{E} = \begin{bmatrix} \boldsymbol{\alpha}^T \\ \vdots \\ \boldsymbol{\alpha}^T \end{bmatrix} + \mathbf{X}_1 \mathbf{B}_S + \mathbf{E}.$$

Now consider the i th case $(\mathbf{x}_i^T, \mathbf{y}_i^T)$ which corresponds to the i th row of \mathbf{Z} and the i th row of \mathbf{X} . The nontrivial predictors are in \mathbf{x}_i for the mreg2 model. Then

$$\begin{bmatrix} Y_{i1} = \alpha_1 + \beta_{11}x_{i1} + \cdots + \beta_{p1}x_{ip} + \epsilon_{i1} = \alpha_1 + \mathbf{x}_i^T \boldsymbol{\beta}_1 + \epsilon_{i1} \\ Y_{i2} = \alpha_2 + \beta_{12}x_{i1} + \cdots + \beta_{p2}x_{ip} + \epsilon_{i2} = \alpha_2 + \mathbf{x}_i^T \boldsymbol{\beta}_2 + \epsilon_{i2} \\ \vdots \\ Y_{im} = \alpha_m + \beta_{1m}x_{i1} + \cdots + \beta_{pm}x_{ip} + \epsilon_{im} = \alpha_m + \mathbf{x}_i^T \boldsymbol{\beta}_m + \epsilon_{im} \end{bmatrix}.$$

131) For the mreg2 model, the OLS estimators are $\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z}$, $\hat{\boldsymbol{\alpha}} = \bar{\mathbf{y}} - \hat{\mathbf{B}}_S^T \bar{\mathbf{x}}$, and $\hat{\mathbf{B}}_S = \hat{\boldsymbol{\Sigma}}_{\mathbf{x}}^{-1} \hat{\boldsymbol{\Sigma}}_{\mathbf{x}\mathbf{y}}$ where $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}\mathbf{y}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{y}_i - \bar{\mathbf{y}})^T$ which has j th column $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}\mathbf{y}_j}$ for $j = 1, \dots, m$. Note that $\hat{\alpha}_j = \bar{Y}_j - \hat{\boldsymbol{\beta}}_j^T \bar{\mathbf{x}}$. Let

$$\mathbf{v} = \begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix}, \quad E(\mathbf{v}) = \boldsymbol{\mu}_{\mathbf{v}} = \begin{pmatrix} E(\mathbf{y}) \\ E(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_{\mathbf{y}} \\ \boldsymbol{\mu}_{\mathbf{x}} \end{pmatrix}, \quad \hat{\boldsymbol{\mu}}_{\mathbf{v}} = \bar{\mathbf{v}} = \begin{pmatrix} \bar{\mathbf{y}} \\ \bar{\mathbf{x}} \end{pmatrix},$$

$$\text{Cov}(\mathbf{v}) = \boldsymbol{\Sigma}_{\mathbf{v}} = \begin{pmatrix} \boldsymbol{\Sigma}_{\mathbf{y}} & \boldsymbol{\Sigma}_{\mathbf{y}\mathbf{x}} \\ \boldsymbol{\Sigma}_{\mathbf{x}\mathbf{y}} & \boldsymbol{\Sigma}_{\mathbf{x}} \end{pmatrix}, \quad \text{and} \quad \hat{\boldsymbol{\Sigma}}_{\mathbf{v}} = \mathbf{S}_{\mathbf{v}} = \begin{pmatrix} \hat{\boldsymbol{\Sigma}}_{\mathbf{y}} & \hat{\boldsymbol{\Sigma}}_{\mathbf{y}\mathbf{x}} \\ \hat{\boldsymbol{\Sigma}}_{\mathbf{x}\mathbf{y}} & \hat{\boldsymbol{\Sigma}}_{\mathbf{x}} \end{pmatrix}.$$

Let the vector of constants be $\boldsymbol{\alpha}^T = (\alpha_1, \dots, \alpha_m)$ and the matrix of slope vectors $\mathbf{B}_S = \begin{bmatrix} \boldsymbol{\beta}_1 & \boldsymbol{\beta}_2 & \dots & \boldsymbol{\beta}_m \end{bmatrix}$. For iid case \mathbf{v}_i , the population least squares coefficient matrix is

$$\mathbf{B} = \begin{pmatrix} \boldsymbol{\alpha}^T \\ \mathbf{B}_S \end{pmatrix}$$

where $\boldsymbol{\alpha} = \boldsymbol{\mu}_{\mathbf{y}} - \mathbf{B}_S^T \boldsymbol{\mu}_{\mathbf{x}}$ and $\mathbf{B}_S = \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \boldsymbol{\Sigma}_{\mathbf{x}\mathbf{y}}$ provided the relevant matrices exist.

The OLS mreg2 estimator can be calculated by computing the sample mean and sample covariance matrix $(\bar{\mathbf{v}}, \mathbf{S}_{\mathbf{v}}) = (\bar{\mathbf{v}}, \hat{\boldsymbol{\Sigma}}_{\mathbf{v}})$ of the \mathbf{v}_i , and then plug in the results into the formulas for $\hat{\boldsymbol{\alpha}}$ and $\hat{\mathbf{B}}_S$.

More on classification

132) For two groups with $Z = 1$ or $Z = -1$, let $SP = \beta_0 + \boldsymbol{\beta}^T \mathbf{x}$. Classify \mathbf{x} in group 1 if $ESP > 0$ and classify \mathbf{x} in group -1 if $ESP < 0$. So the classifier $\hat{C}(\mathbf{x}) = \text{sign}(ESP)$.

133) Suppose the two groups of training data in 132) are separable by a hyperplane. The estimated *optimal separating hyperplane* ESP has the largest margin on the training data. The hyperplanes parallel to the ESP that form the boundaries of the margin are called fences. The fences pass through as least 2 training data set cases forming the support set S of support vectors. The margin M is the distance between the fences. A separating hyperplane has $SP > 0$ if $\mathbf{x} \in$ group 1 and $SP < 0$ if $\mathbf{x} \in$ group -1 . Hence $Z_i SP_i = Z_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) > 0$ for $i = 1, \dots, n$. Think of the hyperplane $\beta_0 + \boldsymbol{\beta}^T \mathbf{x}$ as dividing \mathbb{R}^P into two halves. The SVM split tries to make the halves homogeneous.

134) *Wide data* = ultra-high dimensional data has $p \gg n$. If $n \leq p + 1$ then there is a separating hyperplane unless there are “exact predictor ties across the class barrier.”

135) The optimal margin classifier $(\hat{\beta}_{0M}, \hat{\boldsymbol{\beta}}_M)$ solves $\max_{\beta_0 \in \mathbb{R}, \boldsymbol{\beta} \in \mathbb{R}^P} M$ subject to (*):

$Z_i SP_i = Z_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) \geq M$ for all $i = 1, \dots, n$. Equivalently, solve $\min_{\beta_0, \boldsymbol{\beta}} \|\boldsymbol{\beta}\|_2$ subject to (*). This classifier is called a *hard margin classifier* since no training data cases from either group can pass the fences of the classifier. It turns out that $\hat{\boldsymbol{\beta}}_M = \sum_{i \in S} \hat{\alpha}_i \mathbf{x}_i$.

136) A *soft margin classifier* allows training data cases from either group to pass the fences or to be misclassified. Let the $\epsilon_i \geq 0$ be slack variables. This classifier solves $\min_{\beta_0, \boldsymbol{\beta}} \|\boldsymbol{\beta}\|_2$ subject to $Z_i SP_i = Z_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) \geq 1 - \epsilon_i$ for $i = 1, \dots, n$ and $\sum_{i=1}^n \epsilon_i \leq \beta_0$.

B. Equivalently $\min_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^n [1 - Z_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i)]_+ + \lambda \|\boldsymbol{\beta}\|_2^2$, a criterion similar to that of ridge regression. Here $[w]_+ = w$ if $w \geq 0$ and $[w]_+ = 0$ if $w \leq 0$. The *hinge loss* $[1 - Z_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i)]_+$ is the cost of \mathbf{x}_i being on the wrong side of the margin (which is 0 if \mathbf{x}_i is on the correct side of the margin).

137) A *support vector machine* (SVM) that uses \mathbf{x}_i minimizes the above loss criterion. For separable training data, $(\hat{\beta}_{0,SVM}, \hat{\boldsymbol{\beta}}_{SVM}) \rightarrow (\hat{\beta}_{0,M}, \hat{\boldsymbol{\beta}}_M)$ as $\lambda \rightarrow 0$. The SVM also has fences and a support set S of support vectors with $\hat{\boldsymbol{\beta}}_{SVM} = \sum_{i \in S} \hat{\gamma}_i \mathbf{x}_i$. The *ESP* $= \hat{\beta}_{0,SVM} + \hat{\boldsymbol{\beta}}_{SVM}^T \mathbf{x} = \hat{\beta}_{0,SVM} + \sum_{i \in S} \hat{\gamma}_i \mathbf{x}_i^T \mathbf{x}$. The SVM can be computed with $O(n^2p)$ complexity using the Gram matrix $\mathbf{X}\mathbf{X}^T$ or with $O(np^2)$ complexity using $\mathbf{X}^T\mathbf{X}$. Ridge regression could also be computed this way.

138) A lasso-SVM solves $\min_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^n [1 - Z_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i)]_+ + \lambda \|\boldsymbol{\beta}\|_1$ and does variable selection. For $Z \in \{-1, 1\}$, a “ridged logistic regression” solves

$\min_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^n \log[1 + \exp(-Z_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i))] + \lambda \|\boldsymbol{\beta}\|_2^2$. A “lasso logistic regression” would change the squared norm $\|\boldsymbol{\beta}\|_2^2$ to $\|\boldsymbol{\beta}\|_1$.

139) A ROC curve is used to evaluate binary classifiers, and the overall performance is summarized by the area under the ROC curve (AUC). An ideal ROC curve is close to the top left corner (left and top sides of the rectangle) of the plot. The larger the AUC, the better the classifier. The ROC curve plots the false positive rate versus the true positive rate, so $0 \leq AUC \leq 1$. A classifier with $AUC = 0.5$ does no better than chance. A ROC from test data or validation data is better than a ROC from training data.

140) A *truth table = confusion matrix*.

predict	truth		total
	-1	1	
-1	true negative (TN)	false negative (FN)	N^*
1	false positive (FP)	true positive (TP)	P^*
total	N	P	

The **error rate** is $(FP + FN)/(FP + FN + TN + TP)$. This rate is the AER if training data was used and VER if a validation set was used.

The *false positive rate* $= FP/N = 1 - \text{specificity} \approx$ type I error.

The *true positive rate* $= TP/P = 1 - \text{sensitivity} \approx 1 -$ type II error \approx power \approx recall.

141) A SVM uses a kernel function $K(\mathbf{x}_i, \mathbf{x}_j)$. The SVM in 137) uses a *linear kernel* $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$. A *polynomial kernel of degree d* is $K(\mathbf{x}_i, \mathbf{x}_j) = [1 + \mathbf{x}_i^T \mathbf{x}_j]^d$. A *radial kernel* is $K(\mathbf{x}_i, \mathbf{x}_j) = \exp[-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|_2^2]$. The SVM with a linear kernel is a competitor

of LDA and logistic regression. The SVM with a nonlinear kernel is a competitor of QDA and KNN. The SVM uses $f(\mathbf{x}) = \hat{\beta}_0 + \sum_{i=1}^n \hat{\alpha}_i K(\mathbf{x}, \mathbf{x}_i) = \hat{\beta}_0 + \sum_{i \in S} \hat{\alpha}_i K(\mathbf{x}, \mathbf{x}_i) = ESP$ with $\hat{C}(\mathbf{x}) = \text{sign}(ESP)$. The $n(n-1)/2$ distinct pairs $(\mathbf{x}_i, \mathbf{x}_j)$ are needed to estimate $\hat{\beta}_0$ and the $\hat{\alpha}_i$.

142) Let $Y = 1$ if $Z = 1$ and $Y = 0$ if $Z = -1$. Then $Y|\mathbf{x} \sim \text{binomial}(m = 1, \rho(\mathbf{x}))$ where $\rho(\mathbf{x}) = \rho(SP) = P(Y = 1|\mathbf{x})$ and $\rho(0) = 0.5$. This is a binary regression with ρ unspecified. A response plot is ESP versus Z with lowess added as a visual aid. If $ESP = \hat{\beta}_0 + \hat{\beta}^T \mathbf{x}$ and $n_i \geq 20p$, then the bootstrap with n_i cases selected with replacement from each group is likely useful. Use the prediction region method.

143) If there are $G \geq 2$ classes, the *one versus one* or *all pairs* classifier constructs $\binom{G}{2}$ binary classifiers, one for each distinct pair of groups. Classify \mathbf{x} with $f_{ij}(\mathbf{x})$, and let $m_i =$ number of times \mathbf{x} is predicted to be in class i . Then $\hat{Y}(\mathbf{x}) = \hat{C}(\mathbf{x}) = d$ where $m_d = \max(m_1, \dots, m_G)$. The *one versus all* classifier fits G binary classifiers: group $i = 1$ versus the $G - 1$ other classes with -1 with $f_i(\mathbf{x})$. Then $\hat{Y}(\mathbf{x}) = \hat{C}(\mathbf{x}) = d$ where $\hat{f}_d(\mathbf{x}) = \max(\hat{f}_1(\mathbf{x}), \dots, \hat{f}_G(\mathbf{x}))$. (These are ESPs.)

144) The two classifiers in 143) can be applied to other binary classifiers, and the labels $Y \in \{a, b\}$ can be used. For example $a = 0$ and $b = 1$.

145) A *regression tree* is a flexible method for $Y = m(\mathbf{x}) + e$ or for $Y_i = m(\mathbf{x}_i) + \sigma_i e_i$. A *classification tree* is a flexible method for classification. Both methods produce graphs called *trees* that look like dendrograms. Each branch has a label like $X_i > 7.56$ or $X_i < 3.45$ where X_i is quantitative or a label like $X_j = a, c$ or $X_j \neq d, g$ if X_j is quantitative with levels a, b, \dots, g, h . **Unless told otherwise**, go to the left of the branch if the condition is true, and go to the right of the branch if the condition is false. A *split* is a rule for creating new branches. The bottom of the tree has leaves = *terminal nodes* that give $\hat{Y} = \hat{Y}|\mathbf{x}$ where \hat{Y} is a number for a regression tree and \hat{Y} is a label for the classification group for a classification tree. The tree is binary so a tree with $d \geq 1$ splits (rules) has $d + 1$ terminal nodes.

146) Know how to find \hat{Y} given a tree and \mathbf{x} values. If $\mathbf{x} = (X_1, \dots, X_p)^T$, often not all of the X_i values are needed to find $\hat{Y} = \hat{Y}|\mathbf{x}$.

147) Trees that use recursive partitioning for classification and regression trees use the CART algorithm. In growing the tree the binary CART algorithm recursively splits the data in each node until either the terminal node is homogeneous (the region R_m corresponding to the node has all cases from the same group for classification and $Y \approx$ constant for regression), or until the terminal node has ≤ 5 observations. The region R_m corresponding to the m th terminal node is a hyper-box: a p -dimensional set if \mathbf{x} is $p \times 1$. Hence trees suffer from the curse of dimensionality.

148) The tree divides the p -dimensional predictor space into J distinct and nonoverlapping regions (p -dimensional hyper-boxes) R_1, \dots, R_J . For each observation that falls in region R_j make the same prediction \hat{Y}_{R_j} . For example, $\hat{Y}_{R_j} = \bar{Y}_j$, the sample mean of the training data Y in R_j for regression, and $\hat{Y}_{R_j} = \text{mode}_j$ for classification where mode_j is the training data group that occurred most often for the training data Y in R_j (a lot like KNN where the region is a p -dimensional hypersphere that contains K training data Y 's). For a regression tree, the *response plot* is a plot of \hat{Y} versus Y . Then there are J dot plots, one for each value of \hat{Y}_{R_j} , with n_j values $Y_{1,1}, \dots, Y_{1,n_j}$ where the $Y_{i,j}$ are the

training data in R_j . These J dot plots scatter about the identity line. The residuals corresponding to R_j are $r_{i,j} = Y_{i,j} - \bar{Y}_j$. The *residual plot* of \hat{Y} versus r consists of J dot plots of the residuals, one for each value of $\hat{Y}_{R_j} = \bar{Y}_j$. These dot plots scatter about the $r = 0$ line.

149) Let $\hat{Y}|\mathbf{x} = \hat{f}(\mathbf{x})$. For a regression tree, draw a sample of size n with replacement from $\mathbf{x}_1, \dots, \mathbf{x}_n$. Fit a tree and find $\hat{f}_1^*(\mathbf{x})$. Repeat B times to get $\hat{f}_1^*(\mathbf{x}), \dots, \hat{f}_B^*(\mathbf{x})$. then the *bagging estimator* $\hat{f}_{bag}^*(\mathbf{x}) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i^*(\mathbf{x})$. For classification take samples of size n_i

with replacement from each group. Let $\hat{f}_i^*(\mathbf{x}) = j_i(\mathbf{x}) \in \{1, \dots, G\}$ be the estimated level (group) of Y given \mathbf{x} . Let $m_k =$ number of $j_i(\mathbf{x}) = k$ for $k = 1, \dots, G$. Take $\hat{f}_{bag}^*(\mathbf{x}) = d$ where $m_d = \max(m_1, \dots, m_G)$, so d is the “mode level group.” For both regression and classification, the trees are not pruned, so terminate when each terminal node has 5 or fewer observations. Bagging a tree typically gives more accuracy than a single tree.

150) The probability of a case not being selected for the i th bootstrap sample is about $e^{-1} \approx 1/3$. These are called out of bag (OOB) observations. Predict \hat{Y} for each OOB observation. Doing this for all B bootstrap samples produces about $B/3$ predictors \hat{Y}_i for each \mathbf{x}_i . Let the OOB predictor $\hat{Y}_{io} =$ average \hat{Y}_i for regression and mode level for classification. Then the OOB MSE = $\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_{io})^2$ for regression and

$\frac{1}{n} \sum_{i=1}^n I(Y_i \neq \hat{Y}_{io})$ for classification. The OOB MSE is “virtually” equivalent to the leave one out CV estimate for sufficiently large B .

151) Random forests use the bootstrap, but at each split, a random sample of $m \approx \sqrt{p}$ predictors is used as split candidates. Random forests produce trees that are less correlated than bagged trees, and tend to have better test error than bagging.

152) Boosting has $\hat{f}(\mathbf{x}) = \sum_{b=1}^B \lambda \hat{f}_b(\mathbf{x})$. First set $\hat{f}(\mathbf{x}) \equiv 0$ and $r_i = Y_i$. For $b = 1, \dots, B$ fit a tree \hat{f}_b with d splits (often $d = 1$ where the tree is a *stump* or $d = 2$) to the training data (\mathbf{X}, \mathbf{r}) . Update the tree and the residuals $\hat{f}(\mathbf{x}) \leftarrow \hat{f}(\mathbf{x}) + \lambda \hat{f}_b(\mathbf{x})$ and $r_i \leftarrow r_i - \lambda \hat{f}_b(\mathbf{x})$. Using stumps ($d = 1$) leads to an additive model: $\hat{f}(\mathbf{x}) = \sum_{j=1}^p \hat{f}_j(X_j)$ where $\mathbf{x} = (X_1, \dots, X_p)$. So boosting with $d = 1$ is a competitor of the additive error GAM $\hat{Y} = \hat{\alpha} + \sum_{j=1}^p \hat{S}_j(X_j)$. Typically $\lambda = 0.01$ or 0.001 .

153) For a binary classification tree with $Y = 0$ or 1 , for a fixed value of \mathbf{x} , the bootstrap produces B estimates $\hat{P}_i^*(Y = 1|\mathbf{x})$ of $P(Y = 1|\mathbf{x})$. Let $\hat{Y}_i^* = 1$ if $\hat{P}_i^*(Y = 1|\mathbf{x}) \geq 0.5$ and $\hat{Y}_i^* = 0$ if $\hat{P}_i^*(Y = 1|\mathbf{x}) < 0.5$. Two common methods to get $\hat{Y}|\mathbf{x}$ are a) $\hat{Y}|\mathbf{x} =$ mode class of 0 or 1, or b) $\hat{Y}|\mathbf{x} = 1$ if $\frac{1}{B} \sum_{i=1}^B \hat{P}_i^*(Y = 1|\mathbf{x}) \geq 0.5$ and

$\hat{Y}|\mathbf{x} = 0$ if $\frac{1}{B} \sum_{i=1}^B \hat{P}_i^*(Y = 1|\mathbf{x}) < 0.5$.