

P142-4
6) Let $\underline{z} = \underline{z}_x$ or $\underline{z} = \underline{y}$.

MV 24

The k th pair of pop canonical variables

is $U_k = \underline{a}_k^T \underline{w}$ and $V_k = \underline{b}_k^T \underline{y}$ and the

k th pop canonical correlation $\rho_k = \text{corr}(U_k, V_k) = \sqrt{\lambda_k}$.

Sample analogs use S or S_U or R or R_U .

$$\hat{U}_k = \hat{\underline{a}}_k^T \underline{w} \quad \hat{V}_k = \hat{\underline{b}}_k^T \underline{y}, \quad \hat{\rho}_k = \text{corr}(\hat{U}_k, \hat{V}_k) = \sqrt{\hat{\lambda}_k}$$

7) P144 Let $\underline{z} = S$ or R . (\hat{U}_1, \hat{V}_1) is the pair of linear combinations (\hat{U}, \hat{V}) having unit sample variances that maximizes $\text{corr}(\hat{U}, \hat{V})$, and the max is $\hat{\rho}_1$. The i th pair (\hat{U}_i, \hat{V}_i) is the pair of linear combinations with unit sample variances that maximizes $\text{corr}(\hat{U}, \hat{V})$ among all choices uncorrelated with the previous $i-1$ canonical variable pairs.

Want linearity in scatterplot matrix and DD plot.

8) want $n > 10p$ for classical CCA, $n > 20p$ for RCCA.

9) Kendall 1980 p69 interpretation difficulties mean few convincing CCA applications in the liter.
Interpretation is hard. a) If $\underline{z} = (\underline{w}^T, \underline{y}^T)^T$ and $\underline{w} \perp \underline{y}$, $\hat{\rho}_1$ is not close to 0 until sample size is quite large.

n

b) ^(p455) $\text{corr}(w_i, a_j)$ and $\text{corr}(y_i, b_j)$ 24.5

are not proportional to a_j and b_j .

Computing $\text{corr}(w_i, a_j)$, $\text{corr}(w_i, b_j)$ for $i=1, \dots, m$
and $\text{corr}(y_i, a_j)$, $\text{corr}(y_i, b_j)$ for $i=1, \dots, g$ can help.

c) If some w_i variables are highly correlated and some y_i variables are highly correlated, the canonical variables can be hard to interpret: Multicollinearity.

ex]
$$w \approx \begin{pmatrix} \log \text{ shell mass} = S \\ \log \text{ muscle mass} = M \end{pmatrix}$$

S and M are highly correlated.

$$\hat{a}_1 = \begin{pmatrix} .150 \\ .032 \end{pmatrix}, \quad \text{but } \hat{c} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \quad \text{with } \begin{matrix} c_1 \geq 0 \\ c_2 \geq 0 \end{matrix}$$

may produce $\hat{c}^T w$ that is highly correlated
with $\hat{a}_1^T w$, in particular, $\hat{c} = \begin{pmatrix} .032 \\ .150 \end{pmatrix}$ would work.

d) w_i 's should have similar variances S_i^2
 y_j 's " S_j^2 "
otherwise, components of \hat{a}_i and \hat{b}_i are hard to interpret.

10) 1st pair is "most important" $M \cup 25$
 pairs with low $\hat{\beta}_k$ can be ignored.

11) a) want to know which w_i variables are most important for \hat{a}_i .

b) y_i \hat{b}_i

c) want to know which w_i variable most explains $\hat{b}_i^T y$ see ex;

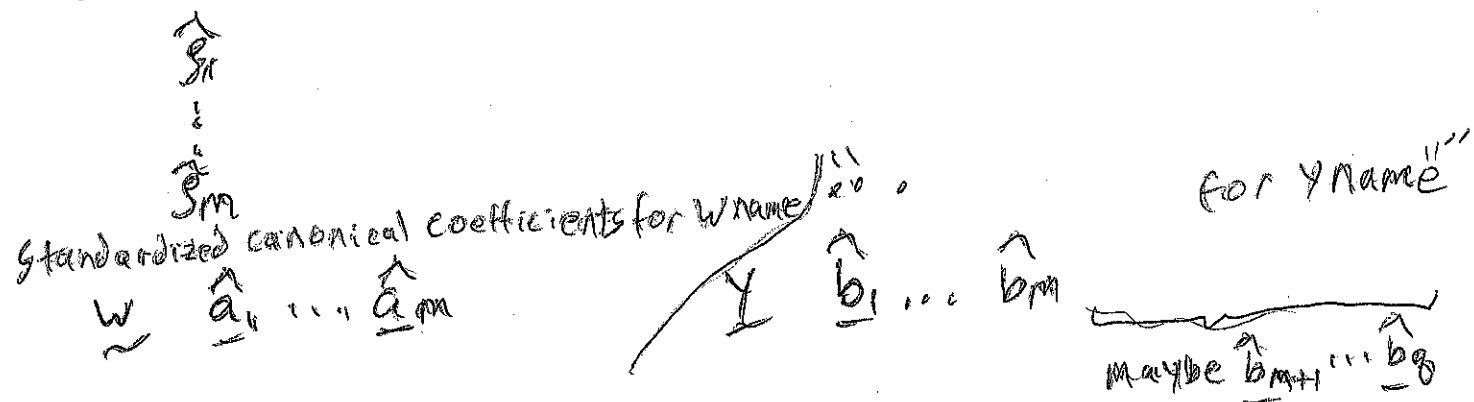
y_i

$\hat{a}_i w$

d) Are $\hat{a}_i^T w$ and $\hat{b}_i^T y$ meaningful?

Often have a "wow factor" \circ
 the client says "wow that makes sense,"
 but often interpretation is difficult.

12) SAS PROC CORR output
 canonical correlation \dots other stuff



ex) weight, waist, pulse

chins situps jumps

\underline{w} = ^{phys} measurements

\underline{y} = exercises

can corr
.7956
.2006
.0726

waist most explains $b_1^T \underline{y}$
chins most explains $a_1^T \underline{w}$

weight -0.0314
waist 0.493
pulse -0.0082

Exer 1
chins -0.0661
situps -0.0168
jumps 0.0140

bigger waist

↔ fewer chinups

13) $\hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{11} & \hat{\Sigma}_{12} \\ \hat{\Sigma}_{21} & \hat{\Sigma}_{22} \end{pmatrix}$

and elements of $\hat{\Sigma}_{12}$
measure pairwise associations

14) $\text{corr}(u, v) = \frac{a^T \hat{\Sigma}_{12} b}{\sqrt{a^T \hat{\Sigma}_{11} a} \sqrt{b^T \hat{\Sigma}_{22} b}}$ if $\hat{\Sigma} = \hat{\Sigma}_x$ or $\hat{\Sigma}_y$

$V(\hat{\Sigma}_w) = V(u) = a^T \hat{\Sigma}_{11} a = 1$

$V(\hat{\Sigma}_y) = V(v) = b^T \hat{\Sigma}_{22} b = 1$

153 Theory for robust CCA is similar to that of robust PCA.

Discriminant Analysis Ch 8 (11)

1) P148 In supervised classification, there are $G \geq 2$ known groups and m cases to be classified. Each case is assigned to exactly one group based on its measurements \underline{w}_i .

ex] patient suffering heart attack symptoms:

3 tests are done $\underline{w} = (w_1, w_2, w_3)^T$

and person is classified as i) had a heart attack or ii) did not have a heart attack.

ex] person applies for credit card based on (salary, credit rating).^T

person is classified as acceptable or not acceptable.

2] Suppose there are $G \geq 2$ groups or populations with pdf $f_j(\underline{z})$ for $j=1, \dots, G$ and \underline{z} is $p \times 1$. so if \underline{x} comes from pop j , then \underline{x} has pdf $f_j(\underline{z})$.

Assume there is a random sample of n_j cases $\underline{x}_{1j}, \dots, \underline{x}_{n_j j}$ for each group.

Let $(\bar{\underline{x}}_j, S_j)$ be the sample mean and covariance matrix for each group. Let \underline{w}_i be a new $p \times 1$ random vector from 1 of the G groups, but the group is unknown. Usually there are many \underline{w}_i

and discriminant analysis attempts to ^{26.5}

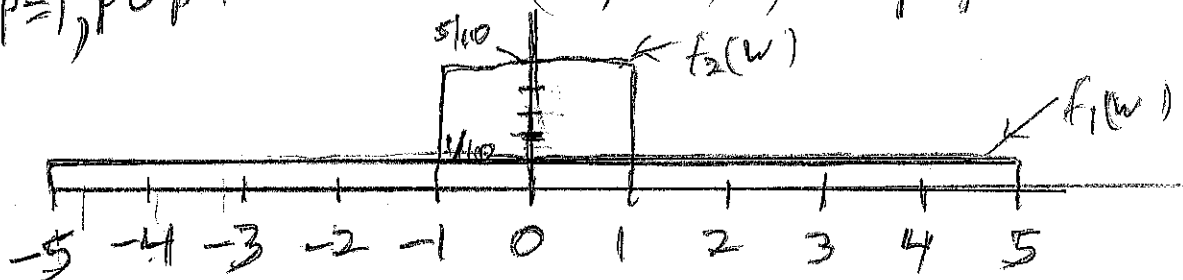
allocate the \underline{w}_i to the correct groups.

3] The maximum likelihood rule allocates case \underline{w} to group a if $\hat{f}_a(\underline{w})$ maximizes $\hat{f}_j(\underline{w})$ for $j=1, \dots, G$.

This rule is robust to nonnormality of the k groups, but it is hard to estimate \hat{f}_j for $p > 1$.

4] E2 question given $p=1$ and graph of f_1, \dots, f_G , give the maximum likelihood rule.

ex] $p=1$, Pop 1 $\sim U(-5, 5)$, Pop 2 $\sim U(-1, 1)$



ML rule: allocate w to pop 2 if $-1 < w < 1$
allocate w to pop 1 if $-5 < w < -1$ or $1 < w < 5$.

5) For now, assume costs of correct (M027) and incorrect allocation are unknown or equal and that probabilities $P_a(\underline{w}_i)$ that \underline{w}_i is in group a are unknown or equal $P_a(\underline{w}_i) = \frac{1}{G}$ for $a=1, \dots, G$.

6) If the G groups are assumed to have the same covariance matrix Σ , then the pooled covariance matrix estimator is

$$S_{\text{pool}} = \frac{1}{N-G} \sum_{j=1}^G (n_j - 1) S_j.$$

by assumption, can be used if n_j are not zero

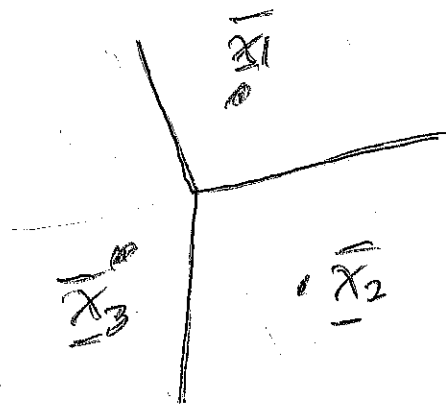
7) Assume the pop dispersion matrices Σ are equal: $\Sigma_j \equiv \Sigma$ for $j=1, \dots, k$. Let $\hat{\Sigma}_{\text{pool}}$ be an estimator of Σ . The linear discriminant rule is to allocate \underline{w} to the group with the largest value of $d_j(\underline{w}) = \hat{\mu}_j^T \hat{\Sigma}_{\text{pool}}^{-1} \underline{w} - \frac{1}{2} \hat{\mu}_j^T \hat{\Sigma}_{\text{pool}}^{-1} \hat{\mu}_j$
 $= \hat{\alpha}_j + \hat{\beta}_j^T \underline{w}$.

Linear discriminant analysis LDA

uses $\hat{\mu}_j = \bar{\underline{x}}_j$ and $\hat{\Sigma}_{\text{pool}} = S_{\text{pool}}$.

8) LDA is the most used rule, is robust to nonnormality and somewhat robust to the assumption of equal pop covariance matrices.

LDA basically separates groups with G hyperplanes.



9) The quadratic discriminant rule allocates \underline{w} to the group with the largest value of $Q_j(\underline{w}) = -\frac{1}{2} \log(|\hat{\Sigma}_j|) - \frac{1}{2} (\underline{w} - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (\underline{w} - \hat{\mu}_j)$.

Quadratic discriminant analysis QDA uses

$(\hat{\mu}_j, \hat{\Sigma}_j) = (\bar{x}_j, S_j)$ and is somewhat robust to nonnormality.

10) The distance discriminant rule allocates \underline{w} to the group with the smallest squared distance $D_{\underline{w}}^2(\hat{\mu}_j, \hat{\Sigma}_j) =$

$(\underline{w} - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (\underline{w} - \hat{\mu}_j)$ and is robust to

nonnormality and robust to the assumption of equal $\hat{\Sigma}_j$, but needs $n_j > 10p$ $j=1, \dots, K$. QDA and LDA are usually better.

11) p149 Assume $G=2$ and that there is a group 0 and a group 1. Let $S(\underline{w}) = P(\underline{w} \in \text{group } 1)$. Let $\hat{S}(\underline{w})$ be the logistic regression estimate of $S(\underline{w})$. Then the estimated sufficient predictor $ESP = \hat{\alpha} + \hat{\beta}^T \underline{w}$ and $\hat{S}(\underline{w}) = \frac{e^{ESP}}{1 + e^{ESP}} =$

$$\frac{\exp(\hat{\alpha} + \hat{\beta}^T \underline{w})}{1 + \exp(\hat{\alpha} + \hat{\beta}^T \underline{w})}$$

The logistic

regression discriminant rule allocates \underline{w} to group 1 if $\hat{S}(\underline{w}) \geq 0.5$ and allocates \underline{w} to group 0 if $\hat{S}(\underline{w}) < 0.5$.

This rule is robust to nonnormality and to the assumption of equal π_j ; $j=0,1$.

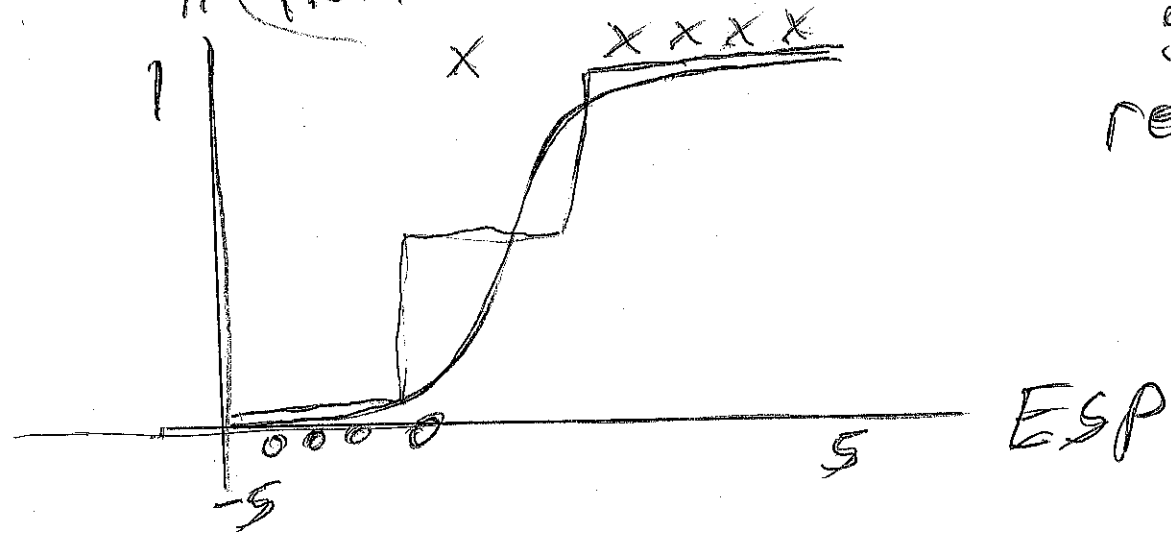
12) A response plot is a plot of ESP vs Y_i (Eq 13)

with $\hat{S}(x_i) \equiv \hat{S}(ESP)$ added as a visual aid, ↑
vertical axis
Also divide ESP into J slices with approx the same number of cases in each slice. Compute the

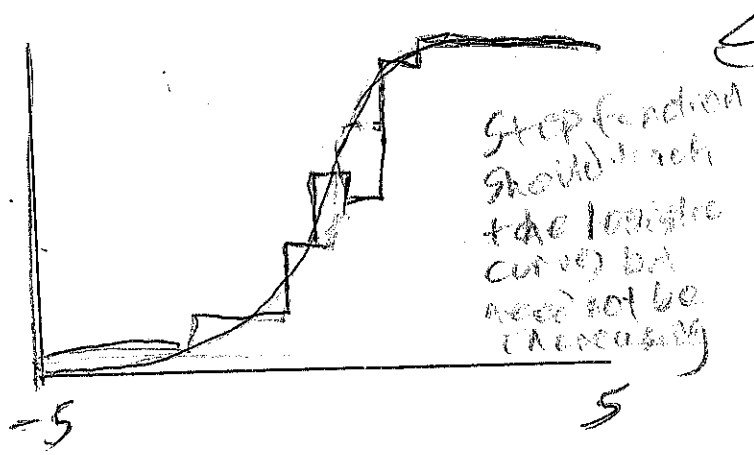
sample mean = sample prop in each slice; (28.5)

$\hat{p}_s = \bar{y}_s = \frac{\sum y_i}{m_s}$ where m_s is the number of cases in slice s . Add the step function to the response plot. If n_0 and n_1 are the sample sizes of groups 0 and 1 and $n_i > 5P$, the logistic regression model is useful if the step function of slice proportions scatter fairly closely about the logistic curve $\hat{p}(ESP)$.

13] know for E_2 y_i from Be able to tell a good LR response plot a bad response plot.



good response plot

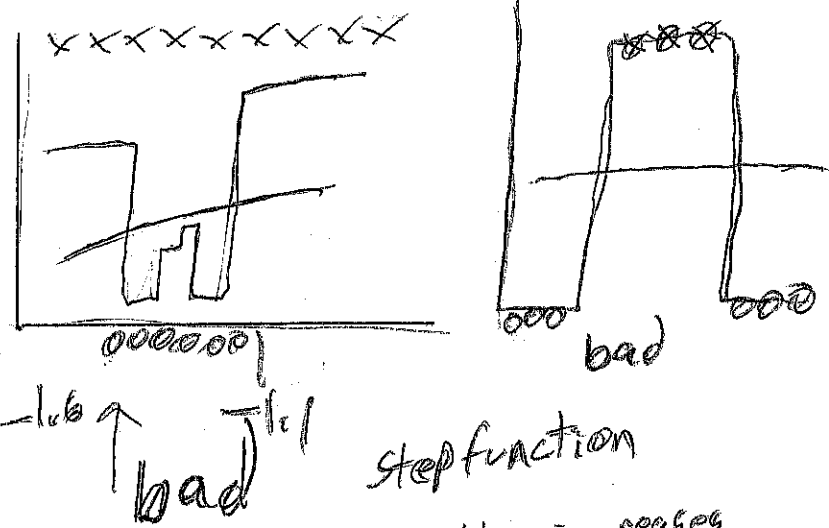
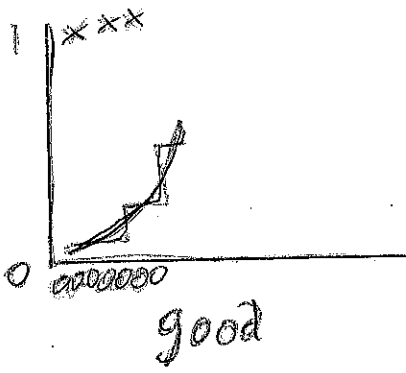
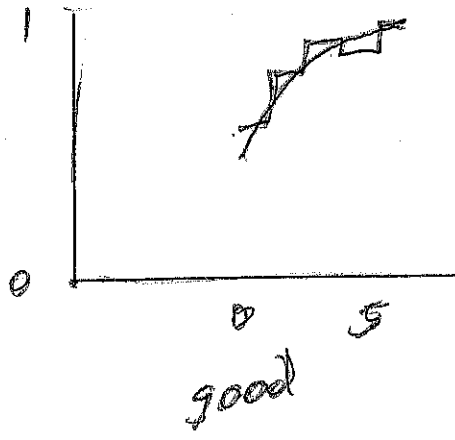
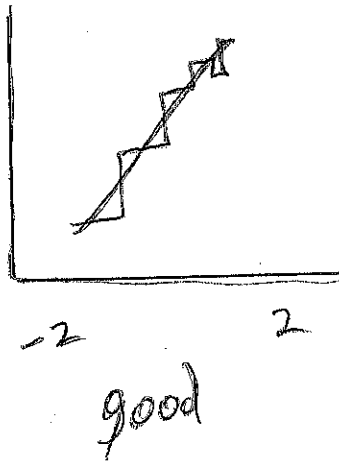


← good response plot



← good response plot

perfect classification of training data



step function
decreases dramatically, then increases dramatically

14] $ESP = 0 \rightarrow \hat{\beta}$ misclassification rate

$ESP = 0 \rightarrow 0.5$

$.25 \text{ or } .75$

$\hat{\beta}$

0.5

$.25$

$\min(\hat{\beta}, 1 - \hat{\beta})$

$ESP > 2$ or $ESP < -2$ will have low error (misclassification) rate if LR model is good

15]

output label constant	estimate	std/Error	Est/SE	pvalue
x_1	$\hat{\beta}_1$			
\vdots	\vdots			
x_p	$\hat{\beta}_p$			

$\frac{e^{\hat{\alpha}}}{1 + e^{\hat{\alpha}}} = \hat{\beta}(x=0)$ makes sense if $x=0$ does

$\alpha=0 \rightarrow \hat{\beta}(x=0) = \frac{1}{2}$ which rarely makes sense

ex] know for final

29.5

sex 0 F 1 M = group

label Estimate
constant -19.7762
Circum 0.0244688
length 0.0371472
head measurements.

Let Circum = $X_1 = 550$, length = $X_2 = 200$

a) Find ESP for X .

$$ESP = \alpha + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 = -19.7762 + 0.0244688(550) + 0.0371472(200) = 1.11108$$

b) Is X classified in group 0 or group 1?

group 1 since $ESP > 0$

c) Find $\hat{p}(X)$.

$$= \frac{e^{ESP}}{1 + e^{ESP}} = \frac{3.0376}{4.0376} = \boxed{0.7523}$$

16) Training data: \underline{X}_{ij} , n_j , group j .

Use a 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 $1 - \frac{m_j}{n_j}$. Suppose $M = \sum_{j=1}^G m_j$