

§3.3

HD52}

36) Let $D = \mathbb{1}_x$ or ρ_x ,
correlation matrix

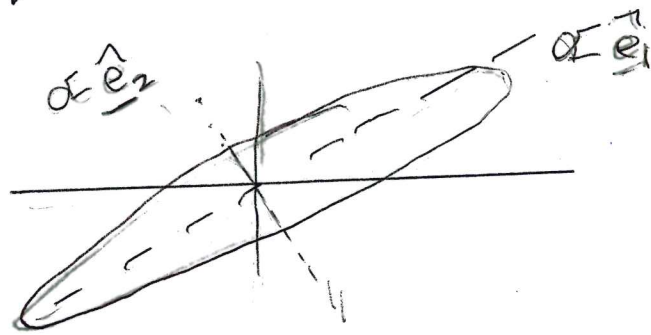
Let $(\hat{\lambda}_1, \hat{e}_1), \dots, (\hat{\lambda}_p, \hat{e}_p)$ be the eigenvalue eigenvector pairs of \hat{D} with $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$. The \hat{e}_i are chosen to be orthonormal:

$$\hat{e}_i^T \hat{e}_i = 1 \quad \text{and} \quad \hat{e}_i^T \hat{e}_j = 0, \quad i \neq j.$$

Then $\hat{m}_j = \hat{e}_j$ and $\hat{v}_i = \begin{pmatrix} \hat{e}_1^T x_i \\ \vdots \\ \hat{e}_p^T x_i \end{pmatrix} = \begin{pmatrix} v_{i1} \\ \vdots \\ v_{ip} \end{pmatrix}$.

This procedure is known as a

principal components analysis, PCA, a big Mult analysis Math 585 topic. PCA gives a change of axes.



37) PCR regresses Y on the $\hat{v}_i, i=1, \dots, p$.

38) Problems: a) no reason why

52.5

the 1st k $\hat{e}_1^T x_i, \dots, \hat{e}_k^T x_i$ should be the most highly correlated with y .

b) $\hat{\beta}_{k, PCR} \neq \hat{\beta}_{OLS}$ unless $\hat{\beta}_{OLS}$ is in the span of $(\hat{e}_1, \dots, \hat{e}_k)$ (which will occur if $k = p$).

39) usually fit $\hat{\beta}_{1, PCR}, \dots, \hat{\beta}_{J, PCR}$

where $J \leq \min(n-1, p)$. Then

use a criterion like 10 fold CV to choose $\hat{\beta}_{MS, PCR} = \hat{\beta}_{k^*, PCR}$.

The resulting model is often linear

and in low dimensions, $k^* = p$ is common.

40) Partial Least Squares (PLS)

can be computed in several equivalent ways. One way uses

$$\hat{m}_j = \frac{1}{\sum_{i=1}^n x_i^2} \sum_{i=1}^n x_i y_i \quad \text{for } j=1, \dots, p.$$

Then $\hat{\beta}_{kPLS}$ regresses Y on

$$\hat{m}_1^T x, \dots, \hat{m}_k^T x. \quad \text{Hence } \hat{\beta}_{kPLS} = \hat{\beta}_{kOPLS}.$$

If $\hat{\beta}_{kPLS} = \hat{\beta}_{kOPLS} = \hat{\beta}_{OLS}$, it can be

shown that $\hat{\beta}_{kPLS} = \hat{\beta}_{OLS}$ for $k=1, \dots, p$.

Also $\hat{\beta}_{pPLS} = \hat{\beta}_{OLS}$.

41] A second equivalent way to compute

$\hat{\beta}_{kPLS}$ is to roughly do a Gram-Schmidt orthogonalization of the \hat{m}_i in 40].

$$\text{Then } \hat{m}_1 = \frac{\sum_{x,y} \hat{x}_i y}{\sqrt{\sum_{x,y} \hat{x}_i^2}}, \quad \hat{m}_i^T \hat{m}_i = 1,$$

$$\hat{m}_i^T \hat{m}_j = 0, \quad i \neq j.$$

42] It can be shown that (if $\hat{v} = \hat{x}$) the m th principal component direction $\hat{m}_m = \hat{e}_m$

Solves
$$\max_{\|\alpha\|=1} \text{Var}(\sum_l \alpha_l x_l)$$

$$\sum_{l=1}^m \hat{m}_l^T x_l \alpha = 0, \quad l=1, \dots, m-1$$

see plot on 36]

and the m th PLS direction $\hat{\underline{n}}_m$

53.5

Solves $\max_{\|\alpha\|=1} [\text{Corr}(Y, \underline{x}_i \alpha)]^2 \text{Var}(\underline{x}_i \alpha)$.

$$\hat{\underline{n}}_l^T \hat{\underline{n}}_l \alpha = 0, l=1, \dots, m-1$$

Hence $[\text{Corr}(Y, \hat{\underline{n}}_1^T \underline{x})]^2 \geq [\text{Corr}(Y, \hat{\underline{n}}_2^T \underline{x})]^2$

$$\geq \dots \geq [\text{Corr}(Y, \hat{\underline{n}}_p^T \underline{x})]^2$$

(Supervised learning)

$$\text{Note that } \underline{x}_i \hat{\underline{n}}_m = \begin{pmatrix} \underline{x}_i^T \hat{\underline{n}}_m \\ \vdots \\ \underline{x}_i^T \hat{\underline{n}}_m \end{pmatrix} = \underline{c}_m =$$

m th component vector - (closely related to the projection of \underline{x}_i on $\hat{\underline{n}}_m$ since the $\hat{\underline{n}}_i$ are orthonormal).

43} There are $p+1$ k -component estimators

eg $\hat{\underline{\beta}}_{k, \text{PLS}}^E$ $k=1, 2, \dots, p$ (or J)

and the model selection estimator

$$\hat{\underline{\beta}}_{\text{MS, PLS}}^E = \hat{\underline{\beta}}_{k^*, \text{PLS}}^E \text{ where } k^* \in \{1, \dots, p\}$$

is picked with model selection eg 10 fold CV.

$$44) \underline{y} = \underline{X} \underline{\beta} + \underline{e}$$

Let the nontrivial predictors

$$\underline{u}_i = (x_{i2}, \dots, x_{ip})^T. \text{ Let the}$$

$n \times (p-1)$ matrix of standardized

nontrivial predictors be $\underline{W} = (w_{ij})$

$$\text{where } \sum_{i=1}^n w_{ij} = 0 \quad \text{and} \quad \sum_{i=1}^n w_{ij}^2 = n$$

j th standardized predictor

$$\text{has } \bar{w}_j = 0$$

$$\text{so } \hat{\sigma}_j^2 = \frac{1}{n} \sum_{i=1}^n w_{ij}^2 = 1$$

method of moments
sample variance
(biased)

See
HW 8.

Then the sample correlation matrix
of the nontrivial predictors is

$$R_{\underline{U}} = \frac{\underline{W}^T \underline{W}}{n}. \text{ Let } \underline{z} = \underline{y} - \bar{y}$$

$$\text{where } \bar{y} = \bar{y} \underline{1}.$$

Many MLR methods fit

$\underline{z} = \underline{w} \underline{n} + \underline{e}$ and then

find $\underline{\hat{\beta}}$ from $\underline{\hat{n}}$,

use $\underline{\hat{y}} = \underline{\hat{z}} + \underline{\bar{y}}$.

45] procedure 44) is used

so $\underline{\hat{\beta}}$ is the same

for anyone who uses the same units of measurement,

ex] $y = ht$ $x_1 = \text{constant}$,

$x_2 = ht$ at shoulder, $x_3 = \text{span}$,

If we fit $\underline{\hat{\beta}}$ in mm and

then convert to ft, tend to get

different answer if we fit $\underline{\hat{\beta}}$ in ft.

point 44] is unit free so the initial estimator uses the "same units!"

ch8 MLR with Heterogeneity

$$1) \text{ MLR1: } \underline{y} = \underline{X}\underline{\beta} + \underline{e}$$

e_i are independent $E(e_i) = 0$

$$V(e_i) = \sigma_i^2, \quad i = 1, \dots, n.$$

$$\text{MLR2: } y = \alpha + \underline{\beta}^T x + e$$

$$y = \alpha \underline{1} + \underline{\delta} \underline{\beta} + e$$

2) weighted least squares (WLS) is often recommended, but there are too many parameters $\underline{\beta}_{p+1}, \sigma_i^2, i = 1, \dots, n.$

3) OLS can be shown to be \sqrt{n} consistent and asymptotically normal.

$$\text{MLR 1: } \sqrt{n}(\hat{\beta}_{\text{OLS}} - \beta) \xrightarrow{D} N_p(0, V \Omega V)$$

$$\text{where } \frac{\underline{X}^T \underline{X}}{n} \xrightarrow{P} V^{-1} \text{ and } \frac{1}{n} \sum_{i=1}^n E[e_i^2 \underline{x}_i \underline{x}_i^T] \xrightarrow{P} \Omega.$$

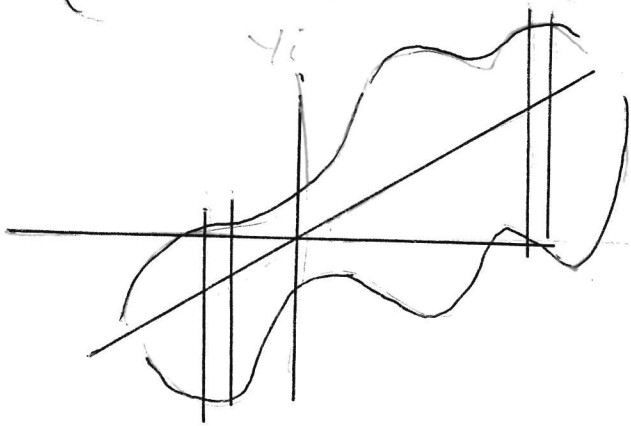
$$\text{Under iid cases, } \Omega = E[e_i^2 \underline{x}_i \underline{x}_i^T].$$

$$4) E(Y | \underline{w}) = m(\underline{w})$$

$$V(Y | \underline{w}) = \mathcal{N}(\underline{w}).$$

$$\text{Hence } E(Y | \hat{\beta}_{\text{OLS}}^T \underline{x}) = \hat{\beta}_{\text{OLS}}^T \underline{x}$$

$$V(Y | \hat{\beta}_{\text{OLS}}^T \underline{x}) = \mathcal{N}(\hat{\beta}_{\text{OLS}}^T \underline{x}) \quad \left. \vphantom{V(Y | \hat{\beta}_{\text{OLS}}^T \underline{x})} \right\} \begin{array}{l} \text{can estimate} \\ \text{from plot} \end{array}$$



$$\hat{\beta}_{\text{OLS}}^T \underline{x}_i = \underline{x}_i^T \hat{\beta}_{\text{OLS}} = w_i$$

5) Visualize $Y | \underset{\substack{\uparrow \\ RV}}{w}$ with a plot of w vs Y and narrow vertical slices.

6) Suppose $(\underline{x}_i^T, \varepsilon_i)^T$ are iid where $E \varepsilon_i = 0$ and $V(\varepsilon_i) = 1$.

$$\text{Let } y_i = \underline{x}_i^T \underline{\beta} + \sigma_i \quad \varepsilon_i = \underline{x}_i^T \underline{\beta} + e_i$$

HD 56

where σ_i^2 is a function of x_i .

Then (y_i, x_i) are iid and earlier MLR OPLS theory still holds.

PHD TOPIC

§4.2 Single Index regression model.

1) we can always write *if (y_i) have a joint dist*

$$y_i = E(y_i | \underline{m}^T \underline{x}_i) + e_i$$

$$= m(\underline{m}^T \underline{x}_i) + e_i$$

$$\text{where } e_i = y_i - m(\underline{m}^T \underline{x}_i).$$

$y_i = m(\underline{m}^T \underline{x}_i) + e_i$ is known as a single index model.

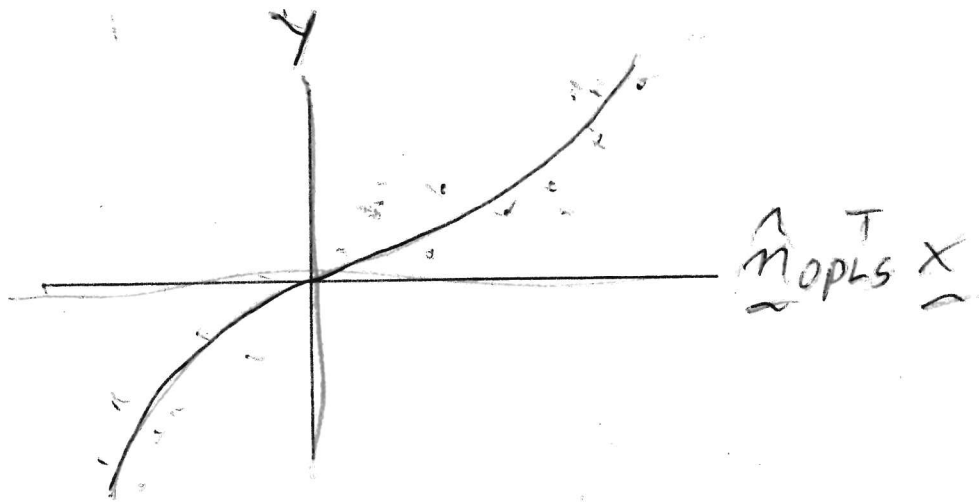
(A multiple index model uses $m(\underline{m}_1^T \underline{x}_i, \dots, \underline{m}_K^T \underline{x}_i)$.)

2] If the cases (y_i, \underline{x}_i) are iid (56.5)

$$\text{let } y_i = m(\underline{m}_{\text{OPLS}}^T \underline{x}_i) + e_i \quad i=1, \dots, n$$

$$E(e_i) = 0, \quad V(e_i) = \sigma_i^2$$

3] visualize m with a response plot



4] The $\hat{\Sigma}_{xy}^T = \hat{m}_{\text{OPLS}}$ theory

still holds $\sqrt{n}(\hat{\Sigma}_{xy} - \Sigma_{xy}) \xrightarrow{D} N_p(0, \Sigma_w)$

PhD topic

§4.4 Poisson Regression

$$1] y_i | \underline{m}^T \underline{x}_i \sim \text{Poisson}[\mu(\alpha + \underline{m}^T \underline{x}_i)]$$

where $E(Y|\underline{n}^T \underline{x}) = \mu(\alpha + \underline{n}^T \underline{x}) = m(\alpha + \underline{n}^T \underline{x})$ HD 57

and $V(Y|\underline{n}^T \underline{x}) = \mu(\alpha + \underline{n}^T \underline{x}) = \nu(\alpha + \underline{n}^T \underline{x})$.

2] A Poisson regression (PR) generalized linear model (GLM)

takes $\mu(\alpha + \underline{\beta}_{PR}^T \underline{x}) = \exp(\alpha + \underline{\beta}_{PR}^T \underline{x})$.

3] MMLE

Fit $Y \sim \text{Poisson}(e^{\alpha + \beta x})$

or $Y_k \sim \text{Poisson}(e^{\alpha + \beta_j x_{jk}})$ $k=1, \dots, n$

to get $\hat{\beta}_j$ for $j=1, \dots, P$.

Then $\hat{\underline{\beta}}_{MMLE} = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_P \end{pmatrix}$.

4] MMLE variable selection:

standardize the predictors (then do $\hat{\beta}_j$ on the w_j)
and take the J x_i corresponding
to the largest $|\hat{\beta}_i|$.

5) Lasso variable selection for 57.5
PR can be used in high dimensions.
No theory but if $\lambda = 0$, $\hat{\beta}_L = \hat{\beta}_{PR}$ if $n > p - s$.

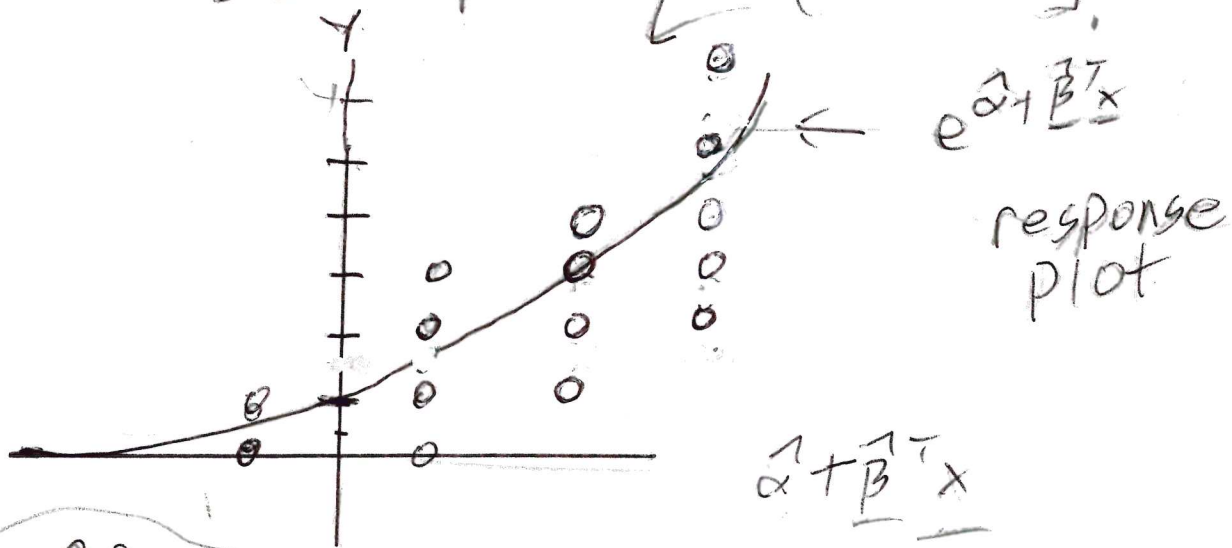
6) Fit the lasso PR model $\hat{\beta}_L$,
then fit the PR (GLM) to the
variables with $\hat{\beta}_i \neq 0$. see HW 8.

6) Let $z_i = \begin{cases} y_i & \text{if } y_i > 0 \\ 0.5 & \text{if } y_i = 0. \end{cases}$

The minimum chi-square estimator
 $(\hat{\alpha}_M, \hat{\beta}_M)$ is found from the WLS
regression of $\log(z_i)$ on \underline{x}_i with
weights $w_i = z_i$. Equivalently,
use the OLS regression (without
intercept) of $\sqrt{z_i} \log(z_i)$ on $\sqrt{z_i} (1, \underline{x}_i^T)^T$.

7) $(\hat{\alpha}_M, \hat{\beta}_M)$ and $(\hat{\alpha}_{PR}, \hat{\beta}_{PR})$ tend to
be consistent estimators of (α, β) .

if $Y | \underline{\beta}^T x \sim \text{Poisson}[\exp(\alpha + \underline{\beta}^T x)]$ HD 58



8] PR

$$\log(z) \approx \alpha + \underline{\beta}^T x + e$$

$$\log(Y+1) \approx \alpha + \underline{\beta}^T x + e$$

MLR with heterogeneity

If $(\log(z_i), x_i)$ are iid ($\log(Y_i, x_i)$ iid)
 or $(\log(Y_i+1), x_i)$ are iid, then

$$\sqrt{n} \left(\begin{matrix} \hat{\underline{\beta}} \\ \hat{\alpha} \end{matrix} - \begin{matrix} \underline{\beta} \\ \alpha \end{matrix} \right) \xrightarrow{D} N_p(\underline{0}, \underline{\Sigma}_w)$$

where $U = z$ or $Y+1$.

PHD topic (MMLE, ...)

9] Let $w_i = \hat{\underline{\eta}}^T x_i$, $\hat{\underline{\eta}} = \hat{\underline{\beta}}_{\text{MMLE}}$ or $\hat{\underline{\beta}}_{x_i, \log(U)}$.

Do PR of Y_i on w_i , get response plot.

10) PIS: Parametric bootstrap

58.5

$$Y_i^* \stackrel{\text{bootstrap}}{\sim} \text{Poisson}(\hat{\alpha} + \hat{\beta}^T x_i), \quad i=1, \dots, B.$$

100(1- δ)% PI \approx 100(1- δ)% Shorth PI

applied to Y_1^*, \dots, Y_B^* eg $B=100$ or 1000 ,

If $Y_i \sim D(\alpha + \beta^T x_i, \underline{\gamma})$

and the regression method gives

$\hat{\alpha}, \hat{\beta}$, and $\hat{\underline{\gamma}}$, then use

$$Y_i^* \sim D(\hat{\alpha} + \hat{\beta}^T x_i, \hat{\underline{\gamma}}) \text{ for } i=1, \dots, B.$$

(Poisson regression, binomial regression, Weibull regression, etc.)

11) Weibull regression

(Y_i, x_i) follow a Weibull proportional hazards (PH) regression model iff

$$\log(Y_i) = \alpha + \beta^T x_i + \varepsilon_i \text{ follows an}$$

MLR model = Weibull accelerated failure time

(AFT) model.

AD 59

12) Do lasso for PH, then

fit Weibull PH model to the x_j with $\hat{\beta}_j \neq 0$ (lasso vs).

13) MMLE fit the p marginal models.

14) Usually can't fit OLS to the AFT due to the y_i being censored. Need an estimator of $\$xy$ when y are censored.

PhD topic

Binary regression and classification of 2 groups:

1) For binary regression, $y_i \in \{0, 1\}$.

We may use $z_i = 2y_i - 1 \in \{-1, 1\}$.

$y_i | \underline{m}^T x_i \sim \text{binomial}(m=1, \mathcal{S}(\underline{m}^T x_i))$

where $\mathcal{S}(\underline{m}^T x_i) = P(y=1 | \underline{m}^T x_i)$.

2) Multitude of models:

59.5

Since Y is binary, if $\begin{pmatrix} Y \\ \underline{w} \end{pmatrix}$ has a joint dist, then $Y | \underline{w} \sim \text{bin}(m=1, p(\underline{w}))$

and $Y | \underline{m}^T \underline{x} \sim \text{bin}[m=1, p(\underline{m}^T \underline{x})]$.

Hence if $\begin{pmatrix} Y \\ \underline{x} \end{pmatrix}$ has a joint dist, then

$Y | \underline{m}^T \underline{x}$ follows a binary regression model for every $\underline{m} \in \mathbb{R}^p$.

3) Note that $E(Y | \underline{m}^T \underline{x}) = m p(\underline{m}^T \underline{x}) = p(\underline{m}^T \underline{x})$.

4) One binary regression model is the logistic regression LR model with

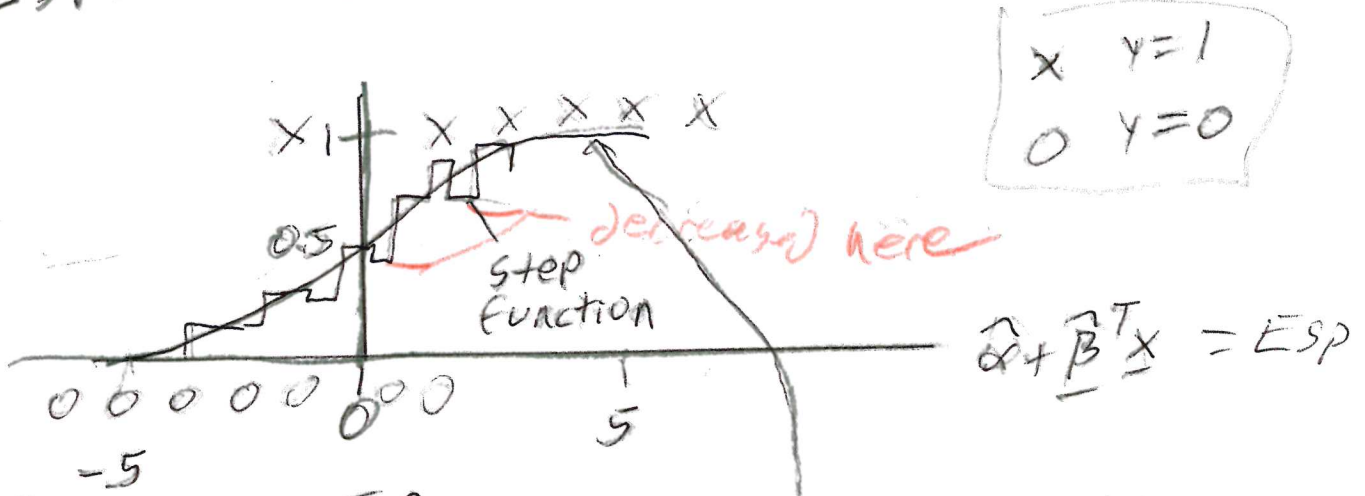
$$p(x) = p(\alpha + \underline{\beta}^T \underline{x}) = \frac{\exp(\alpha + \underline{\beta}^T \underline{x})}{1 + \exp(\alpha + \underline{\beta}^T \underline{x})} \in [0, 1].$$

abuse of notation

5) LR response plot

$$ESP = \hat{\alpha} + \hat{\beta}^T x,$$

$$\hat{\beta} = \hat{\beta}_{LR}$$



$$E(Y|SP) = \frac{e^{ESP}}{1 + e^{ESP}} \quad \rightarrow \text{logistic curve}$$

Add a step function with step height $= \bar{y} - \hat{\beta}$ for the cases in the interval.
 want the step function to track the logistic curve closely, but the step function often is not an increasing function.

6) Suppose there are 2 groups 0 and 1.
 Let $f(\underline{w}) = P(\underline{w} \in \text{group } 1) = P(Y=1 | \underline{w})$.

Let $\beta(\beta) = \frac{e^{\beta}}{1+e^{\beta}}$ and let

60.5

$\hat{\beta}(\text{ESP}) = \frac{e^{\text{ESP}}}{1+e^{\text{ESP}}}$ be the LR

estimator of $\beta(\beta)$. The LR
discriminant rule or LR classification
rule allocates \underline{w} to group

$$\begin{cases} 1 & \text{if } \hat{\beta}(\underline{w}) \geq 0.5 \quad (\text{ESP} \geq 0) \\ 0 & \text{if } \hat{\beta}(\underline{w}) < 0.5 \quad (\text{ESP} < 0) \end{cases}$$

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

Given \underline{w} , the training data misclassification

rate $\approx \min(\hat{\beta}, 1 - \hat{\beta})$. $\text{ESP} > 2$ or

$\text{ESP} < -2$ will have low error (misclassification)

rate if the LR model is good. $\text{ESP} = 0$

≈ 0.5 error rate (coin flip).

7] Lasso variable selection: HD 61

Fit lasso LR model $\hat{\beta}_L$.

Fit LR model to the J predictors with nonzero lasso $\hat{\beta}_i$.

8] MMLE:

Fit the LR model that regresses y on x_j to

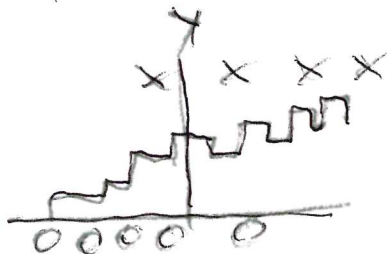
get $\hat{\beta}_j$, Then $\hat{\beta}_{MMLE} = (\hat{\beta}_1, \dots, \hat{\beta}_p)^T$.

9] MMLE variable selection:

Standardize the predictors (then do 8] on the x_j)

and take the J x_j corresponding to the largest $|\hat{\beta}_j|$.

end exam 2 material
10] To visualize $\hat{\beta}$, use response plot



$$\hat{\alpha} + \hat{\beta}^T x = ESP$$

11) * Training data:

x_{1j}, \dots, x_{nj} group j . ($Y=j$)

use a discriminant or classification method to classify the training data

(Optimistic compared to test data: estimated error rate is too low).

If m_j of the n_j group j cases are correctly classified,

the apparent error rate for group j

is $1 - \frac{m_j}{n_j}$. Suppose there are

G groups ($G=2$ for binary), Let

$m_A = \sum_{j=1}^G m_j$. Then the

apparent error rate AER = $1 - \frac{m_A}{n}$

where $n = \sum_{j=1}^G n_j$,