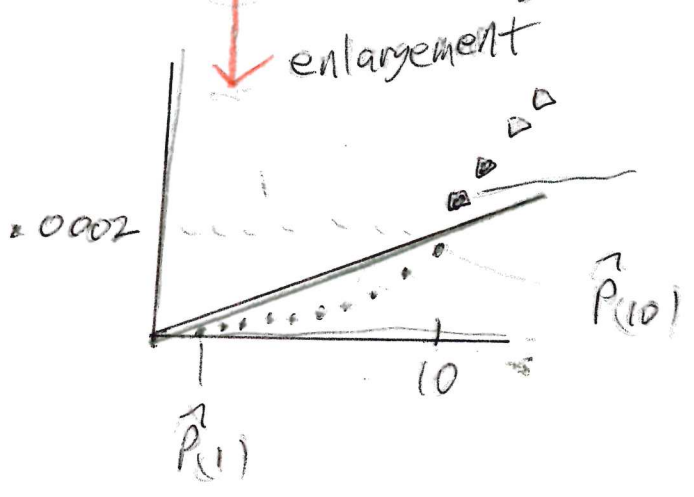


$$\frac{0.05}{2500} = 0.00002$$

Bonferroni cutoff



$\hat{P}(11)$   $\hat{K} = 10$

$$10 \cdot \frac{0.05}{2500} = 0.0002$$

ex)  $n = 102$  men: 52 prostate cancer patients, 50 normal controls.

Each man's gene expression levels were measured on a panel of 6033 genes:  $x_{ij}$  = activity of  $i$ th gene for  $j$ th man.

2 sample  $t$  statistic  $t_i = \frac{\bar{x}_{2i} - \bar{x}_{1i}}{\hat{\sigma}_{\text{pool}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$  was

*red*  $\hat{\sigma}_{\text{pool}}$

2 = cancer, 1 = normal.

computed for gene  $i$   
 then get  $\hat{P}_1, \dots, \hat{P}_m$  = p-values for  $i = 1, \dots, m = 6033$ .

# Multivariate Linear Regression

75.5

ch 10

1] For multivariate linear regression (mreg), there are  $m \geq 2$  response variables  $y_1, \dots, y_m$ . At least one predictor variable takes on many values.  
ex] systolic and diastolic blood pressure

2] The mreg model is

$$\underbrace{y_i}_{m \times 1} = \underbrace{B^T}_{m \times p} \underbrace{x_i}_{p \times 1} + \underbrace{\epsilon_i}_{m \times 1} \quad \text{for } i=1, \dots, n$$

with  $m \geq 2$  response variables

$y_1, \dots, y_m$  and  $p$  predictor variables

$x_1=1, x_2, \dots, x_p$ . The nontrivial predictors.

$i$ th case =  $(\underline{x}_i^T, \underline{y}_i^T)^T$  where

$x_{i1}=1$  is often omitted.

data matrix

$$\begin{pmatrix} \underline{x}_1^T, \underline{y}_1^T \\ \vdots \\ \underline{x}_n^T, \underline{y}_n^T \end{pmatrix}$$

HD 76

but column of 1s  
is usually omitted.

In matrix form,  $\underline{Z}_1 = \underline{X} \underline{B} + \underline{E}$  where

$$\underline{Z}_1 = (\underline{y}_1, \underline{y}_2, \dots, \underline{y}_m) = \begin{pmatrix} \underline{y}_1^T \\ \vdots \\ \underline{y}_n^T \end{pmatrix},$$

$$\underline{B} = [\underline{B}_1 \ \underline{B}_2 \ \dots \ \underline{B}_m], \quad \underline{E} = [\underline{e}_1, \dots, \underline{e}_m] = \begin{pmatrix} \underline{\varepsilon}_1^T \\ \vdots \\ \underline{\varepsilon}_n^T \end{pmatrix}.$$

$$E[\underline{\varepsilon}_k] = \underline{0}, \quad \text{cov}(\underline{\varepsilon}_k) = \underline{\Sigma}_{\varepsilon} = (\sigma_{ij}), \quad k=1, \dots, m,$$

$$E(\underline{e}_i) = \underline{0}, \quad \text{cov}(\underline{e}_i, \underline{e}_j) = \sigma_{ij} \underline{I}_n \quad \text{for} \\ i, j = 1, \dots, m.$$

3] Use  $\hat{\underline{\varepsilon}}$ ,  $\hat{\underline{\Sigma}}$ , and  $\hat{\sigma}_{ij}$  for residuals.

4)\* Each response variable

76.5

$y_j$  in the `mreg1` model follows a univariate multiple linear regression (MLR)

model  $\underline{y}_j = \underline{X} \underline{\beta}_j + \underline{e}_j$ ,  $j=1, \dots, m$ .

$$y_{ij} = \underline{x}_i^T \underline{\beta}_j + e_{ij} = \beta_{1j} + x_{i2} \beta_{2j} + \dots + x_{ip} \beta_{pj} + e_{ij}$$

with  $x_{i1} \equiv 1$ .  $\underline{X}$  does not depend on  $j$ .

So the same predictors are used for all  $m$  response variables.

5)  $m=1 \rightarrow$  MLR, so for `mreg`,  $m \geq 2$ ,

6) The OLS estimator is

$$\hat{\underline{B}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{z}_1 \quad \text{which gives}$$

$$\hat{\underline{\beta}}_j = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}_j \quad j=1, \dots, m,$$

So  $\hat{\underline{B}}$  can be computed by

computing the OLS MLR model for regressing  $y_j$  on  $X$ ,  $j=1, \dots, m$ .

7) Competing estimators include

A) Regress  $y_j$  on  $X$ ,  $j=1, \dots, m$  to get  $\hat{\beta}_{jE}$ ,  $j=1, \dots, m$  and take  $\hat{B} = [\hat{\beta}_{1E}, \dots, \hat{\beta}_{mE}]$

B) Regress  $y_j$  on  $X$ ,  $j=1, \dots, m$  and get linear combinations

$$\begin{aligned} & \hat{m}_{1|X}^T, \dots, \hat{m}_{1|K_1}^T X \\ & \hat{m}_{2|X}^T, \dots, \hat{m}_{2|K_2}^T X \\ & \vdots \\ & \hat{m}_{m|X}^T, \dots, \hat{m}_{m|K_m}^T X \end{aligned}$$

$$= (w_1, \dots, w_K) \quad \text{with } K = \sum_{i=1}^m K_i$$

Then use  $Z = W B + E$  and

$$\begin{matrix} n \times m & n \times (k+1) & n \times m \end{matrix}$$

$$W = [w_1, \dots, w_K]$$

use OLS mreg to get  $\hat{\beta}$

if  $n > \max\left\{ \frac{1}{\lambda} (m+k)^2, m+k+30, 10k \right\}$ .

ex) RR a) could use

$$\hat{\beta}_R = (\mathbf{X}^T \mathbf{X} + \hat{\lambda} \mathbf{I}_n)^{-1} \mathbf{X}^T \mathbf{Z}$$

$$= [\mathbf{X}^T \mathbf{X} + \hat{\lambda} \mathbf{I}_n]^{-1} [\mathbf{X}^T \mathbf{y}_1 \dots \mathbf{X}^T \mathbf{y}_m],$$

but then the same  $\hat{\lambda}$  is used for all  $m$  regressions.

b) could get  $\hat{\beta}_j = (\mathbf{X}^T \mathbf{X} + \hat{\lambda}_j \mathbf{I}_n)^{-1} \mathbf{X}^T \mathbf{y}_j$

for  $j=1, \dots, m$ .

c) could get  $\hat{\beta}_j$  as in a) or b),

form  $w_j = \hat{\beta}_j^T \mathbf{x}$   $j=1, \dots, m$

and do OLS mreg of  $\mathbf{Z}$  on

$$\mathbf{W} = [\mathbf{1}, \mathbf{w}_1, \dots, \mathbf{w}_m]$$

8) Mreg2 model

← nontrivial predictors

$$\underline{y}_i = \underline{\alpha} + B_s^T \underline{x}_i + \underline{\varepsilon}_i$$

$n \times 1$        $m \times 1$        $m \times p$        $p \times 1$        $m \times 1$

$$\underline{Z} = \begin{bmatrix} \underline{\alpha}^T \\ \vdots \\ \underline{\alpha}^T \end{bmatrix} + \underline{X}_1 B_s + E$$

$n \times m$        $n \times p$        $p \times m$        $n \times m$

$$= \underline{X} \begin{pmatrix} \underline{\alpha}^T \\ B_s \end{pmatrix} + E = \underline{X} B + E$$

$(p+1) \times m$

(\*)

$$\underline{X} = \begin{pmatrix} 1 & \underline{X}_1 \end{pmatrix}$$

$n \times (p+1)$

The  $i$ th row of  $\underline{Z}$  is  $\underline{y}_i^T = \underline{\alpha}^T + \underline{x}_i^T B_s + \underline{\varepsilon}_i^T$ .

9) OLS estimator  $\hat{B} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Z}$ .

2nd way to compute OLS estimator

$$\hat{\underline{\alpha}} = \bar{\underline{y}} - \hat{B}_s^T \bar{\underline{x}}, \quad \hat{B}_s = \hat{\underline{Z}}_x^{-1} \hat{\underline{Z}}_x \underline{z}$$

10) If  $\underline{v}_i = \begin{pmatrix} \underline{y}_i \\ \underline{x}_i \end{pmatrix}$  are iid,  $E(\underline{v}) = \begin{pmatrix} E(\underline{z}) \\ E(\underline{x}) \end{pmatrix} = \begin{pmatrix} \underline{\mu}_y \\ \underline{\mu}_x \end{pmatrix}$ ,

$$\text{cov}(\underline{v}) = \underline{\Sigma}_v = \begin{pmatrix} \underline{\Sigma}_y & \underline{\Sigma}_{yx} \\ \underline{\Sigma}_{xy} & \underline{\Sigma}_x \end{pmatrix}, \text{ then}$$

$$\underline{\alpha}_{OLS} = \underline{\alpha} = \underline{\mu}_y - B_s^T \underline{\mu}_x, \quad B_s = \underline{\Sigma}_x^{-1} \underline{\Sigma}_{xy} = B_{s,OLS}$$

11) Suppose  $\begin{pmatrix} \underline{y} \\ \underline{x} \end{pmatrix} \sim N_{m+p} \left[ \begin{pmatrix} \underline{\mu}_y \\ \underline{\mu}_x \end{pmatrix}, \begin{pmatrix} \underline{\Sigma}_y & \underline{\Sigma}_{yx} \\ \underline{\Sigma}_{xy} & \underline{\Sigma}_x \end{pmatrix} \right]$  (78.9)

where  $\underline{\Sigma}_x^{-1}$  exists.

Then  $\underline{y} | \underline{x} \sim N_m(\underline{\mu}_{y|x}, \underline{\Sigma}_{y|x})$  where

$$\underline{\mu}_{y|x} = E(\underline{y} | \underline{x}) = \underline{\mu}_y + \underline{\Sigma}_{yx} \underline{\Sigma}_x^{-1} (\underline{x} - \underline{\mu}_x)$$

$$= \underline{\mu}_y - \underline{B}_s^T \underline{\mu}_x + \underline{B}_s^T \underline{x} = \underline{\alpha} + \underline{B}_s^T \underline{x},$$

$$\text{and } \underline{\Sigma}_{y|x} = \text{cov}(\underline{y} | \underline{x}) = \underline{\Sigma}_y - \underline{\Sigma}_{yx} \underline{\Sigma}_x^{-1} \underline{\Sigma}_{xy}.$$

Hence  $\underline{y} | \underline{x}$  follows the OLS mreg2 model.

12) Let  $A$  be a full rank  $k \times p$  matrix with  $1 \leq k \leq p$ . Let  $\underline{w} = A\underline{x}$ .

If  $\begin{pmatrix} \underline{y} \\ \underline{x} \end{pmatrix}$  is as in 11], then

$$\begin{pmatrix} \underline{y} \\ \underline{w} \end{pmatrix} \sim N_{m+k} \left[ \begin{pmatrix} \underline{\mu}_y \\ \underline{\mu}_w \end{pmatrix}, \begin{pmatrix} \underline{\Sigma}_y & \underline{\Sigma}_{yw} \\ \underline{\Sigma}_{wy} & \underline{\Sigma}_w \end{pmatrix} \right],$$

and  $\underline{y} | \underline{w}$  follows an OLS mreg2 model as



in 11) with  $\underline{\mu}_w = \underline{\mu}_{Ax} = A \underline{\mu}_x$ ,

$$\underline{\Sigma}_w = \underline{\Sigma}_{Ax} = A \underline{\Sigma}_x A^T, \quad \underline{\Sigma}_{yw} = \underline{\Sigma}_{y, Ax} = \underline{\Sigma}_{yx} A^T,$$

$$\underline{\Sigma}_{wy} = \underline{\Sigma}_{Ax, y} = A \underline{\Sigma}_{xy}.$$

multitude of models

$$\underline{\tilde{B}} = [\text{diag } \hat{\underline{\Sigma}}_x]^{-1} \hat{\underline{\Sigma}}_{xy}$$

does not fit data well

ex) OPLS type estimators

Let  $w_1 = \hat{\underline{\Sigma}}_{xy_1}^T \underline{x}, \dots, w_m = \hat{\underline{\Sigma}}_{xy_m}^T \underline{x}.$

a) Do SLR of  $y_i$  on  $w_i$  to get

$$\hat{\lambda}_i \hat{\underline{\Sigma}}_{xy_i} = \hat{B}_{OPLS}(i).$$

can be done even if  $m \gg n$  and  $p \gg n$

$$\underline{\hat{B}}_s = [\hat{B}_{OPLS}(1), \dots, \hat{B}_{OPLS}(m)].$$

b) Do m reg of  $\underline{z}_i$  on  $(w_1, \dots, w_m)^T = \underline{w}$

to get  $\hat{B}_s = \hat{\underline{\Sigma}}_w^{-1} \hat{\underline{\Sigma}}_{wz}$

or  $\hat{B} = (\underline{W}^T \underline{W})^{-1} \underline{W}^T \underline{z}$

want

$$\underline{W} = \begin{pmatrix} 1 & \underline{w}_1^T \\ \vdots & \vdots \\ 1 & \underline{w}_m^T \end{pmatrix}$$

$$n \geq \max(4m^2, m^2 + 30, 10m)$$

ex] PCR type estimators

79.5

$$\text{Let } w_1 = \underline{\hat{e}}_1^T x, \dots, w_p = \underline{\hat{e}}_p^T x,$$

$$\underline{w} = (w_1, \dots, w_p)^T,$$

regress  $\bar{y}$  on  $w_1,$

including a constant

then  $w_1, w_2$

$\vdots$

$w_1, \dots, w_k$

want  $n \geq \max((m+k)^2, m+k+30, 10^4)$

---

Back to classification.

1] It has been observed that often

$\hat{\beta}_{OLS} \approx k \hat{\beta}_E$  for several estimators

such as  $\hat{\beta}_E = \hat{\beta}_{LR}$  and  $\hat{\beta}_E = \hat{\beta}_{PR}$ .

2] Let  $\pi_j = P(Y=j)$ ,  $\mu_j = E(X|Y=j)$  and

Let  $N_j = \# Y's = j$  for  $j=0,1$ .

Let  $n = N_0 + N_1$ , Then  $\hat{\pi}_1 = \frac{N_1}{n}$

and  $\hat{\pi}_1 = 1 - \hat{\pi}_0$ .

Let  $\hat{\mu}_i = \bar{x}_i = \frac{1}{N_i} \sum_{j: Y_j=i} x_j$  = sample mean of the  $x_j$  corresponding to  $Y_j = i$  for  $i=0,1$ . HD 80

Claim:  $\hat{\beta}_{xy} = \hat{\pi}_1 \hat{\pi}_0 (\hat{\mu}_1 - \hat{\mu}_0)$ .

Hence for iid cases)

$$\underline{\eta} = \hat{\beta}_{xy} = \pi_1 \pi_0 (\underline{\mu}_1 - \underline{\mu}_0).$$

Proof]  $\hat{\beta}_{xy} = \frac{1}{n} \sum_i x_i y_i - \bar{x} \bar{y}$

$$= \frac{1}{n} \left[ \sum_{j: Y_j=1} x_j (1) + \sum_{j: Y_j=0} x_j (0) \right] - \bar{x} \hat{\pi}_1$$

$$= \frac{1}{n} N_1 \hat{\mu}_1 - \frac{1}{n} \underbrace{(N_1 \hat{\mu}_1 + N_0 \hat{\mu}_0)}_{\sum_{i=1}^n x_i} \hat{\pi}_1$$

$$= \hat{\pi}_1 \hat{\mu}_1 - \hat{\pi}_1^2 \hat{\mu}_1 - \hat{\pi}_1 \hat{\pi}_0 \hat{\mu}_0$$

$$= \hat{\pi}_1 \underbrace{(1 - \hat{\pi}_1)}_{\hat{\pi}_0} \hat{\mu}_1 - \hat{\pi}_1 \hat{\pi}_0 \hat{\mu}_0 = \hat{\pi}_1 \hat{\pi}_0 (\hat{\mu}_1 - \hat{\mu}_0), \quad \square$$

3) If the  $(x_i, y_i)$  are iid,

compute  $\hat{\eta}^T x_i = w_i$  for  $i=1, \dots, n$ .

Plug  $w_i$  into a binary regression estimator

Such as LR to get

$$\hat{\beta}_B = \hat{\lambda} \hat{\eta}$$

binary

and do inference using  $\sqrt{n} (\hat{\eta} - \eta) \xrightarrow{D} N_p(0, \Sigma_w)$   
as for MLR.

$n_i$  was  $N_i$

4) Suppose  $x_{11}, \dots, x_{1n_1}$  are iid

$x_{01}, \dots, x_{0n_0}$  are iid

and the 2 samples are ind.

Let  $\eta = \mu_1 - \mu_0$  and  $\hat{\eta} = \bar{x}_1 - \bar{x}_0$

Assume  $n_i = \pi_i n$  and

( $\frac{n_i}{n} \rightarrow \pi_i$  works)

$$\sqrt{n_i} (\bar{x}_i - \mu_i) \xrightarrow{D} N_p(0, \Sigma_i)$$

$$\Sigma_i = \Sigma_{x_i} = \text{cov}(x_i), \quad i = 0, 1.$$

So  $\sqrt{n} (\bar{x}_i - \mu_i) \xrightarrow{D} N_p(0, \frac{\Sigma_i}{\pi_i})$  and

$$\sqrt{n} \begin{pmatrix} \bar{x}_1 - \mu_1 \\ \bar{x}_0 - \mu_0 \end{pmatrix} \xrightarrow{D} N_{2p} \left[ 0, \begin{pmatrix} \frac{\Sigma_1}{\pi_1} & 0 \\ 0 & \frac{\Sigma_0}{\pi_0} \end{pmatrix} \right]$$



Thus  $\sqrt{n} [\bar{X}_1 - \bar{X}_0 - (\mu_1 - \mu_0)] \xrightarrow{D} N_p(0, \Sigma_w)$  (\*)

where  $\Sigma_w = \frac{\Sigma_1}{\pi_1} + \frac{\Sigma_0}{\pi_2}$  and

$$\hat{\Sigma}_w = \frac{n \hat{\Sigma}_1}{n_1} + \frac{n \hat{\Sigma}_0}{n_2}$$

Take  $\hat{\eta}_B = \bar{X}_1 - \bar{X}_0$  and  $\underline{m} = \underline{\mu}_1 - \underline{\mu}_0$ .

Then  $\sqrt{n} (\hat{\eta}_B - \underline{m}) \xrightarrow{D} N_p(0, \Sigma_w)$ .

Let  $w_i = \hat{\eta}_B^T x_i$  for  $i=1, \dots, n$ .

Plug  $w_i$  into a binary reg estimator such as LR to get  $\hat{\beta}_B = \hat{\lambda} \hat{\eta}_B$  and do inference using (\*).

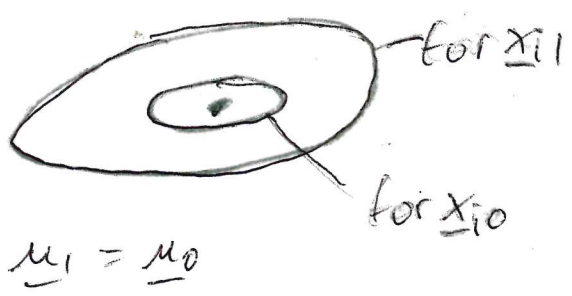
PhD topic

5)



$\bar{X}_1 - \bar{X}_0$  good for classification

6)



$\bar{X}_1 - \bar{X}_0$  bad for classification

7] If the model depends on

$\underline{x}$  only through  $h(\underline{x})$ ,  $h: \mathbb{R}^p \rightarrow \mathbb{R}$ ,

then  $h(\underline{x})$  is a sufficient predictor (SP)

and  $\hat{h}(\underline{x})$  is an estimated sufficient predictor ESP.

ex]  $ESP = \hat{\beta}^T \underline{x}$

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

8] Logistic regression is used a lot in biostatistics and epidemiology where the focus is statistical inference ( $n \geq 10p$ ). Support vector machines (SVMs) are used in machine learning where the goal is classification accuracy.

9] when  $p \gg n$ , there is often a hyperplane

that perfectly separates the two groups. AD 82

The launching point for SVMs was finding the "optimal" separating hyperplane.

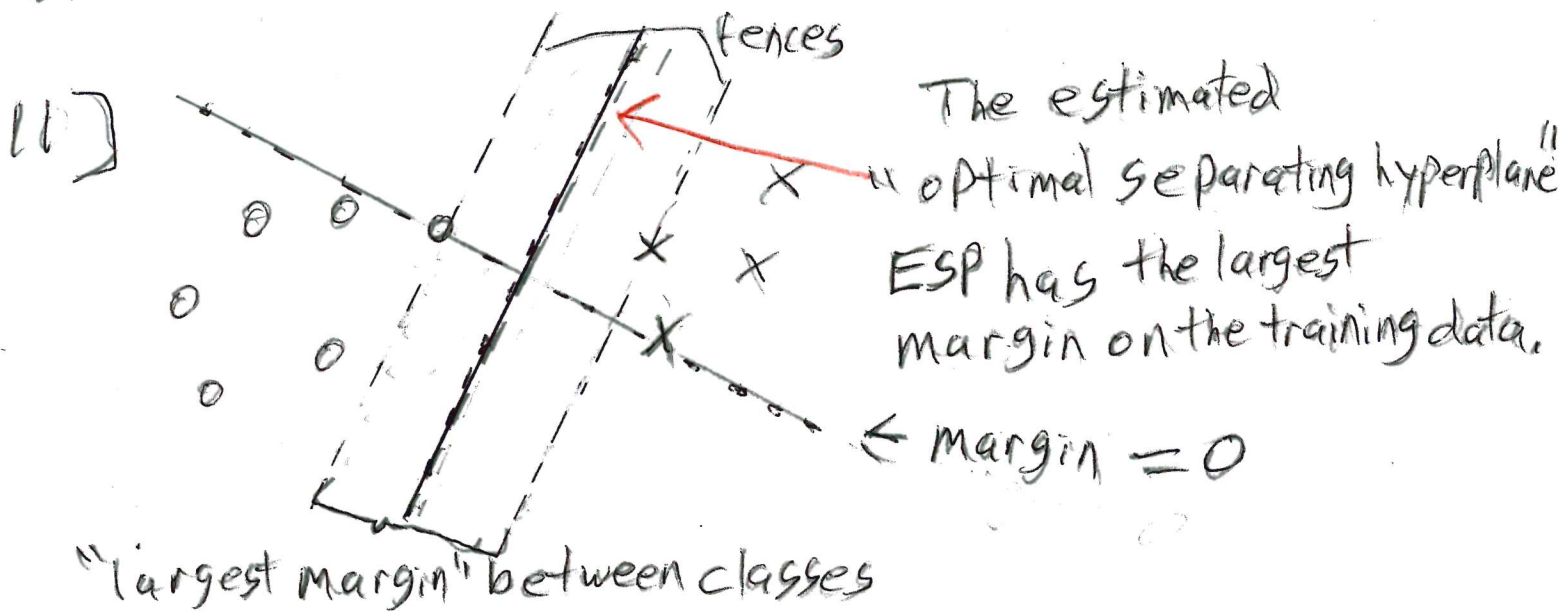
10) For 2 groups, use  $z \in \{-1, 1\}$ .

$$\text{Let } SP = \underbrace{B_0}_{\text{was } \alpha} + \underline{B^T} x$$

Classify  $x$  in group 1 if  $ESP > 0$   
-1 if  $ESP < 0$

(Just like LR but code group 0 as -1).

So classifier  $\hat{c}(x) = \text{sign}(ESP)$ .



12) The SVM split tries to make the 2 "halves" or partitions as homogeneous as possible. || 82.5

13) The hyperplanes parallel to the ESP that form the boundaries of the margin are called fences. The fences pass through at least 2 training data cases. These cases form the Support set S of support vectors.

It turns out that

$$\hat{\beta}_M = \sum_{i \in S} \hat{\alpha}_i \tilde{x}_i = \text{optimal marginal classifier.}$$

14) Wide data = ultra high dimensional data has  $p \gg n$ . If  $n \leq p+1$  there is a separating hyperplane unless there are exact predictor ties across the class barrier (whatever that means).



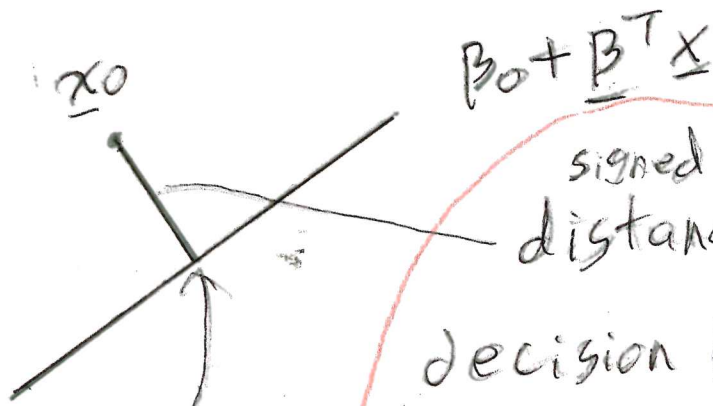
15] Let  $ESP = \hat{\beta}_0 + \hat{\beta}^T x$  and

HD83

$$f(x) = SP = \beta_0 + \beta^T x = SP.$$

So  $\hat{f}(x) = ESP$ ,  $z_i \in \{-1, 1\}$ .

16]



Projection of  $x_0$  onto hyperplane must be  $\beta_0 + \beta^T x_0$

signed distance of  $x_0$  from decision boundary =

$$\frac{\beta_0 + \beta^T x_0}{\|\beta\|_2} = \frac{f(x_0)}{\|\beta\|_2}$$

17] Think of the hyperplane  $\beta_0 + \beta^T x_i$

$= \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$  dividing

$\mathbb{R}^p$  into 2 halves. A separating hyperplane

has  $SP > 0$  if  $x \in$  group 1  
 $< 0$  group -1.

So  $z_i s p_i = z_i (\beta_0 + \beta^T x_i) > 0$  for  $i=1, \dots, n$ .

18] Think of a binary classifier

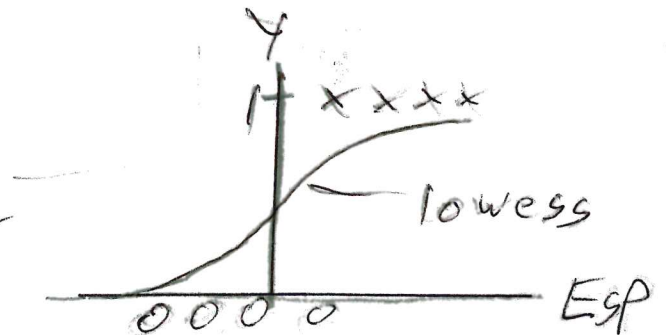
$Y | X \sim \text{bin}[m=1, s(x)]$  where

$s(x) = s(sp) = P(Y=1 | X)$ , where

$z = -1$  is recorded as  $Y=0$ , but

$s(sp)$  is unknown.

response plot



Use lowess or a step function

← (see 11)

19] Let  $M$  be the margin. The optimal margin classifier  $(\hat{\beta}_0, \hat{\beta}_M)$

maximizes  $M$  subject to

$\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p$

$$z_i s p_i = z_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M \quad (*)$$

$\forall i=1, \dots, n$ .

This is called a hard margin

AD 84

classifier since no cases from either group can pass the fences of the classifier.

Equivalently,  $\min_{\beta_0, \beta} \|\beta\|_2$  subject to (\*).

20] A soft margin classifier allows cases from either group to pass the fences and thus be misclassified. This classifier

solves  $\min_{\beta_0, \beta} \|\beta\|_2$  subject to

$$z_i (\beta_0 + x_i^T \beta) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0 \text{ for } i=1, \dots, n$$

$$\text{and } \sum_{i=1}^n \epsilon_i \leq D \quad \text{some constant}$$

Slack variables are used in linear programming,

21] This minimization is equivalent to

$$\min_{\beta_0, \beta} \sum_{i=1}^n \left[ 1 - z_i (\beta_0 + x_i^T \beta) \right]_+ + \lambda \|\beta\|_2^2$$

hinge loss =  $\begin{cases} 0 & \text{if } x_i \text{ is on the correct side} \\ \text{cost of } x_i \text{ being on the wrong side} & \text{of the margin} \end{cases}$

84.5

where  $[w]_+ = \begin{cases} w & w \geq 0 \\ 0 & w \leq 0 \end{cases}$ .

This technique is similar to ridge reg,

22) A Support Vector Machine that uses  $\underline{x}_i$  minimizes the above criterion. For separable data,  $(\hat{\beta}_{0,svm}, \underline{\beta}_{svm}) \xrightarrow{p} (\hat{\beta}_{0,m}, \underline{\beta}_m)$  as  $\lambda \rightarrow 0$ .

23) It turns out that  $\underline{\beta}_{svm} = \sum_{i \in S} \hat{\gamma}_i \underline{x}_i$

and  $\hat{\beta}_{0,svm} + \underline{x}^T \underline{\beta}_{svm} = \hat{\beta}_{0,svm} + \sum_{i \in S} \hat{\gamma}_i \langle \underline{x}, \underline{x}_i \rangle$

where  $\langle \underline{x}, \underline{x}_i \rangle = \underline{x}^T \underline{x}_i$ . This quantity can be computed using the  $n \times n$  Gram matrix  $\underline{X} \underline{X}^T$  with  $O(n^2 p)$  complexity or using  $\underline{X}^T \underline{X}$  with  $O(np^2)$  complexity.

(RR used similar computations.)

24] A lasso-SVM solves

HD 85

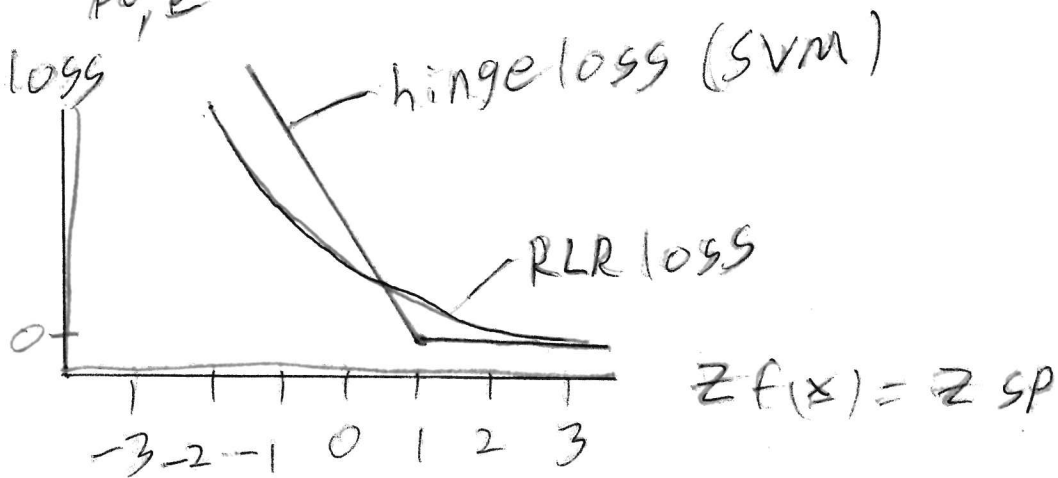
$$\min_{\beta_0, \underline{B}} \sum_{i=1}^n \left[ 1 - z_i (\beta_0 + \underline{x}_i^T \underline{B}) \right]_+ + \lambda \|\underline{B}\|_1$$

like lasso

and does variable selection.

25] "Ridged logistic regression," with  $z_i \in \{-1, 1\}$

$$\text{solves } \min_{\beta_0, \underline{B}} \sum_{i=1}^n \log \left[ 1 + e^{-z_i (\beta_0 + \underline{x}_i^T \underline{B})} \right] + \lambda \|\underline{B}\|_2^2$$



26] \* Truth table = confusion matrix

	truth		
	-1	1	
Predict	-1	3	misclassified $\text{error rate} = 1 - \frac{18+7}{18+7+3}$ $= 1 - \frac{25}{28} = \frac{3}{28} = \boxed{0.1071}$
	1	7	
			correct

85.5

ex)

	truth			
	a	b	c	d
predict a	10	0	0	12
b	0	100	11	0
c	0	5	50	0
d	6	0	0	30

correct

diagonal:  
correctly classified  
off diagonal:  
incorrectly classified

$$\text{error rate} = \frac{6 + 5 + 11 + 12}{6 + 5 + 11 + 12 + 10 + 100 + 50 + 30}$$

$$= \frac{\text{total} - \text{diagonal total}}{\text{total}} = \frac{224 - 190}{224}$$

$$= \frac{34}{224} = \boxed{0.1518}$$

See HW10

27) sometimes 1 or a few observations shift the maximal margin hyperplane.  
The SVM classifier is a soft margin classifier and can do better.