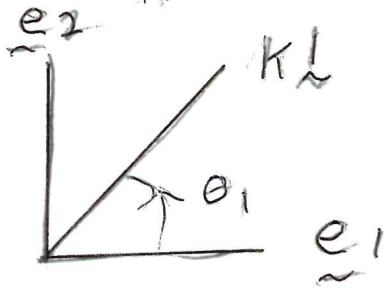


17) Let $\underline{e}_i = (0, \dots, 0, 1, \dots, 0)^T$, $i=1, \dots, p$ HD 16

be a basis for \mathbb{R}^p . Let \underline{k}_\perp be the diagonal of a hypercube or hypersphere.



a rotation does not change the angles, but makes $\cos(\theta_i)$ harder to compute.

Then the angle θ_i between \underline{e}_i and \underline{k}_\perp

$$\text{satisfies } \cos(\theta_i) = \frac{(\underline{e}_i, \underline{k}_\perp)}{\|\underline{e}_i\| \|\underline{k}_\perp\|} = \frac{\sqrt{\underline{k}_\perp^T \underline{e}_i \underline{e}_i^T \underline{k}_\perp}}{\sqrt{1} \sqrt{p}}$$

$$= \frac{1}{\sqrt{p}} \rightarrow 0 \text{ as } p \rightarrow \infty. \text{ So the diagonal}$$

of a hypercube \approx orthogonal to all axes edges of the hypercube in p dimensional space. For large p . For a vector \underline{v} , $\cos(\theta_i) =$

$$\frac{(\underline{e}_i, \underline{v})}{\|\underline{e}_i\| \|\underline{v}\|} = \frac{v_i}{\sqrt{\sum_{j=1}^p v_j^2}} = \frac{v_i}{\sqrt{\sum_{j=1}^p v_j^2}}$$

which often $\rightarrow 0$ as $p \rightarrow \infty$.

19) Suppose $\underline{y} = \underline{X}\underline{\beta} + \underline{e}$ and 16.5

$(\underline{X}^T \underline{X})^{-1}$ exists (eg columns of \underline{X} are orthonormal so $(\underline{X}^T \underline{X})^{-1} = \underline{I}_p$).

Let A be a $p \times n$ matrix.

$$\text{Then } E(\|\underline{A}\underline{e}\|^2) = E(\underline{e}^T \underline{A}^T \underline{A} \underline{e}) = \overset{\text{tr(scalar)}}{\leftarrow} = \text{scalar}$$

$$E[\text{tr}(\underline{e}^T \underline{A}^T \underline{A} \underline{e})] = E[\text{tr}(\underline{e} \underline{e}^T \underline{A} \underline{A}^T)]$$

$$\text{tr}(\underline{AB}) = \text{tr}(\underline{BA})$$

$$= \text{tr}(\underbrace{E(\underline{e} \underline{e}^T)}_{\text{cov}(\underline{e}) = \sigma^2 \underline{I}} \underline{A} \underline{A}^T) = \sigma^2 \text{tr}(\underline{A} \underline{A}^T). \quad (*)$$

Treat \underline{X} as a constant matrix (conditional on $\underline{X}\underline{\beta}$).

$$\text{Then } \hat{\underline{\beta}}_{\text{OLS}} - \underline{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} - \underline{\beta} =$$

$$(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{X} \underline{\beta} + (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{e} - \underline{\beta}. \text{ Thus}$$

$$E[\|\hat{\underline{\beta}} - \underline{\beta}\|^2] = E[\|\underbrace{(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{e}}_A\|^2]$$

$$= \sigma^2 \text{tr} \left[\underbrace{(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{X} (\underline{X}^T \underline{X})^{-1}}_I \right] = \sigma^2 \text{tr}[(\underline{X}^T \underline{X})^{-1}]$$

\uparrow
by (*)

If $(X^T X)^{-1} = I_p$, then $\text{tr}(I_p) = p$

HD 17

and $E[\|\hat{\beta} - \beta\|^2] = p \sigma^2$

If $\Lambda (X^T X)^{-1} \xrightarrow{p} V$, then $\text{tr}(X^T X)^{-1} \approx \frac{\text{tr}(V)}{n}$

$\geq \frac{p \min(V_{ii})}{n}$ which can be

very large if $p \gg n$.

P errors $\propto \frac{1}{n}$ or $\frac{1}{n}$ can add up to a big error if $p \gg n$

19] Let $\langle \underline{a}, \underline{b} \rangle = \underline{a}^T \underline{b}$ and

the Euclidean norm $\|\underline{a}\| = \|\underline{a}\|_2 = \sqrt{\underline{a}^T \underline{a}}$
 $= \sqrt{\langle \underline{a}, \underline{a} \rangle}$.

20] The PLS literature incorrectly claims that $\underline{\beta}_{\text{OPLS}} = \underline{\beta}_{\text{OLS}}$ under mild conditions for p fixed and under stronger conditions for HD MLR.

$y|x = \alpha + \underline{\beta}^T x + e \Rightarrow (\alpha, \underline{\beta}) = (\alpha_{\text{OLS}}, \underline{\beta}_{\text{OLS}})$
under reasonable conditions,

OPLS literature assumes

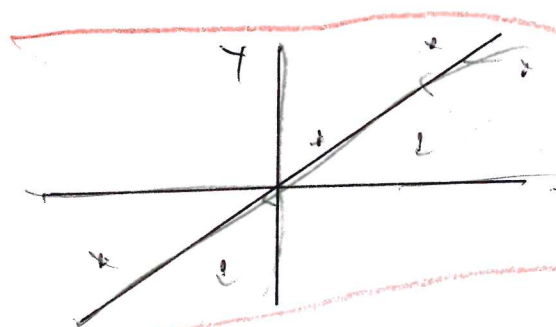
$$y|x = \alpha_{OPLS} + \beta_{OPLS}^T x + e$$

which forces $\beta_{OLS} = \beta_{OPLS}$.

$$y|\beta_{OPLS}^T x = \alpha_{OPLS} + \beta_{OPLS}^T x + e$$

often holds.

21) The conditioning error is made all the time in statistics. For OLS if x is random, statisticians say that OLS results "hold conditional on x " ($y|x$). This is often true for OLS, but not for OPLS. Statisticians also assume that if $y = \alpha + \beta^T x + e$



is reasonable)

$$E y = \alpha_E + \beta_E^T x$$

$$\text{then } y|x = \alpha_E + \beta_E^T x + e.$$

Statisticians also act as if HD18
 there is only one MLR model

$$y = \alpha + \beta^T x + e$$

223 Suppose $\begin{pmatrix} y \\ x \end{pmatrix} \sim N_{p+1} \left[\begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{pmatrix} \Sigma_y & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_x \end{pmatrix} \right]$

with Σ_x nonsingular. Let $\underline{w} = A \underline{x} = \begin{pmatrix} w_1 \\ \vdots \\ w_k \end{pmatrix}$
 where A is a full rank constant matrix. k linear combos

Then $Y | \underline{w} = Y | A \underline{x}$ follows the
 pop OLS model of regressing Y on $A \underline{x} = \underline{w}$.

Proof] Let $B \begin{pmatrix} y \\ x \end{pmatrix} = \begin{pmatrix} 1 & \underline{0}^T \\ \underline{0} & A \end{pmatrix} \begin{pmatrix} y \\ x \end{pmatrix} =$

$$\begin{pmatrix} y \\ A \underline{x} \end{pmatrix} \sim N_{k+1} (B \underline{\mu}, B \Sigma B^T) \quad \leftarrow \text{cov}(y, A \underline{x})$$

$$\sim N_{k+1} \left[\begin{pmatrix} \mu_y \\ A \mu_x \end{pmatrix}, \begin{pmatrix} \Sigma_y & \Sigma_{yx} A^T \\ A \Sigma_{xy} & A \Sigma_x A^T \end{pmatrix} \right]$$

\uparrow $E(A \underline{x})$ \uparrow $\text{cov}(A \underline{x}, y)$ \uparrow $\text{cov}(A \underline{x})$

MVN determined
 by mean and
 cov matrix

Then $Y | A\underline{x} = Y | \underline{w} \sim N(\mu_{Y|\underline{w}}, \Sigma_{Y|\underline{w}})$.

$$\mu_{Y|\underline{w}} = \mu_Y + \Sigma_{Y\underline{x}} A^T (A \Sigma_{\underline{x}} A^T)^{-1} (A\underline{x} - A \underline{\mu}_x)$$

$$= \mu_Y + \Sigma_{Y\underline{w}} \Sigma_{\underline{w}}^{-1} (A\underline{x} - A \underline{\mu}_x)$$

$$= \underbrace{\mu_Y - \Sigma_{Y\underline{w}} \Sigma_{\underline{w}}^{-1} E(\underline{w})}_{E(Y) - \underbrace{\beta_{\underline{w}, OLS}^T}_{\alpha_{\underline{w}, OLS}} E(\underline{w})} + \underbrace{\Sigma_{Y\underline{w}} \Sigma_{\underline{w}}^{-1}}_{\beta_{\underline{w}, OLS}^T} \underline{w}$$

$$\Sigma_{Y|\underline{w}} = \sigma_{Y|\underline{w}}^2 = \Sigma_Y - \Sigma_{Y\underline{x}} A^T (A \Sigma_{\underline{x}} A^T)^{-1} A \Sigma_{\underline{x}Y}$$

$$= \Sigma_Y - \Sigma_{Y\underline{w}} \Sigma_{\underline{w}}^{-1} \Sigma_{\underline{w}Y}$$

$$\text{So } Y | \underline{w} = \alpha_{\underline{w}, OLS} + \underline{w}^T \beta_{\underline{w}, OLS} + e$$

$$\text{where } e \sim N(0, \sigma_{Y|\underline{w}}^2). \quad \square$$

23} In particular, let $w = \underset{\substack{\uparrow \\ RV}}{n^T x}$ with $A = \underset{\substack{\uparrow \\ 1 \times p}}{n^T}$. HD 19

Then under the conditions of 22},

$$Y | \underset{\sim}{n^T x} \sim N(\alpha_{\sim} + \underset{\sim}{\beta}^T x, \sigma_{\sim}^2) \text{ where}$$

$$\alpha_{\sim} = \mu_y - \underset{\sim}{\beta}^T \mu_x, \quad \underset{\sim}{\beta} = \lambda \underset{\sim}{n}, \quad \sigma_{\sim}^2 = \sigma_y^2 - \underset{\sim}{\beta}^T \underset{\sim}{\Sigma}_{xy}$$

$$\text{and } \lambda = \frac{\underset{\sim}{\Sigma}_{xy}^T \underset{\sim}{n}}{\underset{\sim}{n^T} \underset{\sim}{\Sigma}_x \underset{\sim}{n}}$$

$$\text{Proof} \left\{ \begin{pmatrix} 1 & 0^T \\ 0 & \underset{\sim}{n^T} \end{pmatrix} \begin{pmatrix} Y \\ x \end{pmatrix} \sim \begin{pmatrix} Y \\ \underset{\sim}{n^T x} \\ \underset{\sim}{w} \end{pmatrix} \sim N_2 \left[\begin{pmatrix} \mu_y \\ \underset{\sim}{n^T} \mu_x \end{pmatrix}, \begin{pmatrix} \underset{\sim}{\Sigma}_y & \underset{\sim}{\Sigma}_{xy}^T \underset{\sim}{n} \\ \underset{\sim}{n^T} \underset{\sim}{\Sigma}_{xy} & \underset{\sim}{n^T} \underset{\sim}{\Sigma}_x \underset{\sim}{n} \end{pmatrix} \right]$$

$$E(Y|w) = \mu_y + \frac{\underset{\sim}{\Sigma}_{xy}^T \underset{\sim}{n}}{\underset{\sim}{n^T} \underset{\sim}{\Sigma}_x \underset{\sim}{n}} (\underset{\sim}{n^T x} - \underset{\sim}{n^T} \mu_x) =$$

$$\underbrace{\mu_y - \lambda \underset{\sim}{n^T} \mu_x}_{\alpha_{\sim}} + \underbrace{\lambda \underset{\sim}{n^T x}}_{\underset{\sim}{\beta}^T x}$$

$$\sigma_{\sim}^2 = \sigma_y^2 - \frac{\underset{\sim}{\Sigma}_{xy}^T \underset{\sim}{n} \underset{\sim}{n^T} \underset{\sim}{\Sigma}_{xy}}{\underset{\sim}{n^T} \underset{\sim}{\Sigma}_x \underset{\sim}{n}} = \sigma_y^2 - \frac{(\underset{\sim}{\Sigma}_{xy}^T \underset{\sim}{n})^2}{\underset{\sim}{n^T} \underset{\sim}{\Sigma}_x \underset{\sim}{n}} = \sigma_y^2 - \lambda \underset{\sim}{n^T} \underset{\sim}{\Sigma}_{xy}$$



Note that $\sigma_m^2 < \sigma_y^2$ unless

19.5

$\underline{m}^T \underline{\hat{x}}_{xy} = 0$. OPLS uses $\underline{m} = \underline{\hat{x}}_{xy}$.

24} If $\begin{pmatrix} y \\ \underline{x} \end{pmatrix} \sim N_{k+1}(\underline{\mu}, \underline{\Sigma})$, then there are a multitude of MLR models.
Every linear combination $\underline{m}^T \underline{x}$ gives an MLR model.

$\begin{pmatrix} y \\ \underline{m}^T \underline{x} \end{pmatrix} \sim N_2(\underline{\mu}, \underline{\Sigma})$, so highest density regions are ellipsoids, which are linear.
 cases: $(x_i^T, y_i)^T \text{ iid } N_{k+1}(\underline{\mu}, \underline{\Sigma})$ is a very common lit. assumption

25} Under 22} $y | \underline{x}$ and $y | \underline{x}_T$ follow OLS MLR models, which was known.
 (Hence $y | \underline{w}$ follows an OLS MLR model.)

26} A response plot is a plot of ESP vs y . For MLR add the identity line with unit slope and 0 intercept.

For MLR, $ESP = \hat{\alpha}_E + \underline{x}^T \hat{\beta}_E = \hat{y}$.

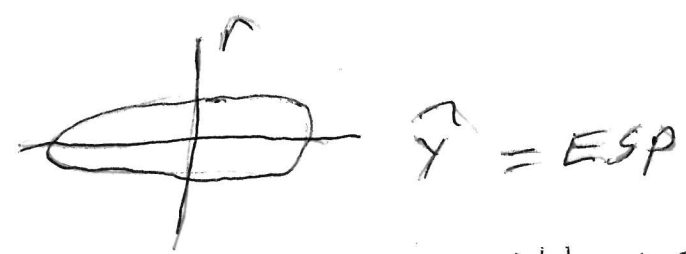
want scatter about the identity line with no other pattern.



$E[(\underline{x}^T \hat{\beta}_E) \pm 0 + 1] (\hat{\alpha}_E + \underline{x}^T \hat{\beta}_E)$

27] For MLR, a residual plot HD 20

is a plot of $ESP = \hat{y}$ vs r
 $= y - \hat{y}$.



want scatter about the $r = 0$ line with no other pattern.

28] Arrange a data set y_1, \dots, y_n
 z_1, \dots, z_n in ascending order.

$$y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n-1)} \leq \underbrace{y_{(n)}}_{\text{max}}$$

min

are the order statistics.

29] The sample median $MED(n) = \begin{cases} y_{(\frac{n+1}{2})} & n \text{ odd} \\ \frac{y_{(\frac{n}{2})} + y_{(\frac{n}{2}+1)}}{2} & n \text{ even} \end{cases}$

(ave of middle order stat or stats)

30] The sample median absolute deviation $MAD(z_1, \dots, z_n) = MAD(n) = MED(|z_i - MED(n)|, i=1, \dots, n)$ is a measure of spread like the standard

deviation.

20.5

ex} 1, 2, 3, 4, 5, 6, 7, 8, 9

↑↑
MED(n) = 5

the median is the middle ordered value or the average of 2 middle values.

$|Y_i - MED(n)|$: +4, +3, +2, -1, 0, 1, 2, 3, 4

ordered 0, 1, 1, 2, 2, 3, 3, 4, 4

↑
MAD(n) = 2

pth variable

31) data $W = \sum_{n \times p} = \begin{pmatrix} x_{11} \\ \vdots \\ x_{n1} \\ \vdots \\ x_{n1} \end{pmatrix} = (v_1, \dots, v_p)$
(x_1 for MLR) ↑
nth "case"

$= \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \leftarrow \text{1st case}$
2nd variable.

32) The coordinatewise median $MED(W)$
 $= (MED(x_1), \dots, MED(x_p))^T$ where

$MED(x_i)$ is the sample median of the i th column.

HD 21

33} Know how to get \bar{x} and the coordinatewise median for a small data set.

ex) 4 measurements on 5 trees ordered (columns \rightarrow rows) \overrightarrow{w} or \overrightarrow{x}

N	E	S	W
72	66	76	77
60	53	66	63
56	57	64	58
41	29	36	38
32	32	35	36

32	41	56	60	72	N
29	32	53	57	66	E
35	36	64	66	76	S
36	38	58	63	77	W

$\underbrace{\hspace{10em}}_{MED(\overrightarrow{w})}$

261	237	277	272	Σ
				sum

show work \rightarrow

$$\frac{1}{5} \begin{pmatrix} 261 \\ 237 \\ 277 \\ 272 \end{pmatrix} = \begin{pmatrix} 52.2 \\ 47.4 \\ 55.4 \\ 54.4 \end{pmatrix} = \bar{x}$$

see HW 3

34] The i th Mahalanobis distance (2/5)

$$D_i = \sqrt{D_i^2} \quad \text{where} \quad D_i^2 = D_i^2(\underline{T}(\omega), G(\omega))$$

\uparrow
vector

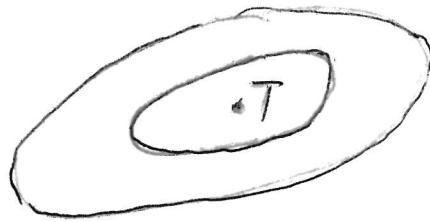
$$= (\underline{x}_i - T(\omega))^T (G(\omega))^{-1} (\underline{x}_i - T(\omega))$$

$$= (\underline{x}_i - T)^T G^{-1} (\underline{x}_i - T) = D_{\underline{x}_i}^2(T, G)$$

$$D_{\underline{x}}^2(T, G) = (\underline{x} - T)^T G^{-1} (\underline{x} - T)$$

\uparrow
center of hyperellipsoid

\leftarrow eigenvectors determine the axes of the hyperellipsoid



contours of constant distance of \underline{x} from T .

35] $D_T^2(\underline{x}, G) = D_{\underline{x}}^2(T, G)$

(distance of \underline{x} from T = distance of T from \underline{x})

36] The POP squared Mahalanobis distance

$$D_{\underline{x}}^2(\underline{\mu}, \Sigma) = (\underline{x} - \underline{\mu}) \Sigma^{-1} (\underline{x} - \underline{\mu}) \quad \text{where}$$

$\underline{\mu}$ = pop. location vector and Σ =

pop dispersion matrix, often the

pop. cov matrix in low dimensions.