

→ see exp 2 Notes (slides 77-83)

- 1.) State the geometric reason that for a dummy variable model w/ a single dummy variable (i.e. $y_i = \alpha_0 + \alpha_1 x_i + \epsilon_i$, where $x_i = 1$ if success, 0 if failure) s.t. the first 5 observations are successes and the last 5 are failures ($n=10$), the sum of the first 5 residuals equals zero: $\sum_{i=1}^5 \hat{\epsilon}_i = 0$.

• $\sum_{i=1}^5 \hat{\epsilon}_i = \hat{\epsilon} \cdot \mathbf{1}$. The vector $\mathbf{1} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \in \text{span}(X)$ (where X is our design matrix $X = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}$ in the case given).
The vector $\hat{\epsilon}$ is the distance between the vector y and its projection \hat{y} onto X . Thus $\hat{\epsilon} \perp \text{span}(X)$: The dot product of any two orthogonal vectors is 0.

ask if
this is the answer
prof wants

- 2.) A researcher is interested in how consumption of fat and sugar affects weight gain in rats. Assume for a moment that fat and sugar do not interact; that is, the effect of sugar on weight gain is the same whether or not a rat is consuming a high fat diet, and vice versa. (Q. for me: is it the case that if there are no interactions between variables x_i, x_j in our model then we are assuming these two variables are independent?) In an experiment, each of 6 rats are fed high-fat and/or high-sugar diets as follows: each rat takes 100g fat or sugar every day and then either a much regular rat food as desired. The response variable measured is the amount of weight gain of the rats (wt weight loss measured as negative values) after two weeks on the diet.

R_1 :	1 mg fat
R_2 :	1 mg sugar
R_3 :	1 mg fat + 1 mg sugar
R_4 :	2 mg fat + 2 mg sugar
R_5 :	2 mg fat + 1 mg sugar
R_6 :	1 mg fat + 2 mg sugar

Let β_0 be the average weight gain due to consuming fat; β_1 be the average weight gain due to consuming sugar. Write down a model using a matrix equation to estimate β_0, β_1 , giving your design matrix X .

model: $y = X\beta + \epsilon \Leftrightarrow y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \epsilon_i$

$$X = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} \quad \hat{\beta} = (X'X)^{-1}X'y$$

3.) Instead of using the y-intercept as in the notes & textbook, suppose we wanted to create a linear model using two dummy variables like this one:

$$y_i = \alpha_1 x_1 + \alpha_2 x_2 + \varepsilon_i, \quad i = 1, 2, \dots, n.$$

You might think of the calcium supplement - blood pressure problem from class (see chap 2 notes pg 68-76), but this time, there are m people in the 1st group and $n-m$ people in the 2nd group. Our dummy variables are then defined as follows:

$$x_1 = \begin{cases} 1 & i = 1, 2, \dots, m \\ 0 & \text{o.w.} \end{cases}$$

$$x_2 = \begin{cases} 1 & m+1, m+2, \dots, n \\ 0 & \text{o.w.} \end{cases}$$

(a) Define the parameters α_1 & α_2 in the context of the problem.

- in the context of the example problem from chap 2 (see slides 68-76),
 x_1 is an indicator ^{variable} that equals 1 if the person received a placebo;
 x_2 is an indicator ^{variable} that equals 1 if the person received a calcium supplement.
- Then our model for blood pressure would look like:

$$y_i = \alpha_1 x_1 + \alpha_2 x_2 + \varepsilon_i, \quad i = 1, 2, \dots, n.$$

• For placebo group:

$$y_i = \alpha_1 + \varepsilon_i$$

• For calcium group:

$$y_i = \alpha_2 + \varepsilon_i.$$

- α_1 is the average blood pressure for an individual receiving the placebo
- α_2 is the average blood pressure for an individual receiving the calcium supplement.

(b) Use the formula $\hat{\alpha} = (X'X)^{-1} X'y$ to solve for the parameter estimates of

α_1 & α_2 .

$$\bullet (X'X) = \begin{bmatrix} \overbrace{1 \dots 1}^m & \overbrace{0 \dots 0}^{n-m} \\ \overbrace{0 \dots 0}^m & \overbrace{1 \dots 1}^{n-m} \end{bmatrix} = \begin{bmatrix} m & 0 \\ 0 & n-m \end{bmatrix}$$

$$\bullet (X'X)^{-1} = \frac{1}{m(n-m)} \begin{bmatrix} n-m & 0 \\ 0 & m \end{bmatrix} = \begin{bmatrix} \frac{1}{m} & 0 \\ 0 & \frac{1}{n-m} \end{bmatrix}$$

$$\bullet X'y = \begin{bmatrix} \overbrace{1 \dots 1}^m & \overbrace{0 \dots 0}^{n-m} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m y_i \\ \sum_{i=m+1}^n y_i \end{bmatrix}$$

$$(X'X)^{-1} X'y = \begin{bmatrix} \frac{1}{m} & 0 \\ 0 & \frac{1}{n-m} \end{bmatrix} \begin{bmatrix} \sum_{i=1}^m y_i \\ \sum_{i=m+1}^n y_i \end{bmatrix}$$

$$\begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{m} \sum_{i=1}^m y_i \\ \frac{1}{n-m} \sum_{i=m+1}^n y_i \end{bmatrix} = \begin{bmatrix} \bar{y}_{\text{placebo}} \\ \bar{y}_{\text{calcium}} \end{bmatrix}$$

4.) (From Stapleton, 1995). ϕ we have an ordinary household scale such as might be used in a kitchen. When an object is placed on the scale, the reading is a combination of the true weight plus random error. You have two coins of unknown weights β_1 & β_2 . To estimate the weight of the coins you take four observations.

- Put coin 1 on the scale and observe y_1 .
- Put coin 2 on the scale and observe y_2 .
- Put both coins on the scale and observe y_3 .
- Put both coins on the scale again & observe y_4 .

Suppose the random errors are i.i.d.

Write a linear model in matrix form and find the least-squares estimates of the coins weights using the usual formula $(X'X)^{-1}X'y$.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \quad X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

$$(X'X) = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix} \Rightarrow (X'X)^{-1} = \frac{1}{5} \begin{bmatrix} 3 & -2 \\ -2 & 3 \end{bmatrix}$$

$$X'y = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} y_1 + y_3 + y_4 \\ y_2 + y_3 + y_4 \end{bmatrix}$$

$$(X'X)^{-1}X'y = \frac{1}{5} \begin{bmatrix} 3 & -2 \\ -2 & 3 \end{bmatrix} \begin{bmatrix} y_1 + y_3 + y_4 \\ y_2 + y_3 + y_4 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 3(y_1 + y_3 + y_4) - 2(y_2 + y_3 + y_4) \\ -2(y_1 + y_3 + y_4) + 3(y_2 + y_3 + y_4) \end{bmatrix}$$

$$\hat{\beta} = \begin{bmatrix} \frac{1}{5}(3y_1 - 2y_2 + y_3 + y_4) \\ \frac{1}{5}(3y_2 - 2y_1 + y_3 + y_4) \end{bmatrix} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}$$

Q: Why do these answers make sense?

- $3y_1$ - we have 3 obs involving coin 1. y_1 is an estimate of just coin 1's weight.
- $(y_3 + y_4) - 2y_2$ - y_3 & y_4 are both estimates of the combined weight of coin 1 & coin 2, two to get the estimated contribution of coin 1 to both total weights, we subtract off y_2 from each, b/c y_2 is an estimate of just coin 2's weight.
- $1/5$ - we use 5 weights to estimate the weight of coin 1. y_1 involves 1 weight, y_3 & y_4 make 2 weights each.
- Analogous reasoning for $\hat{\beta}_2$.

5.) (Chp 2 Question 5.)

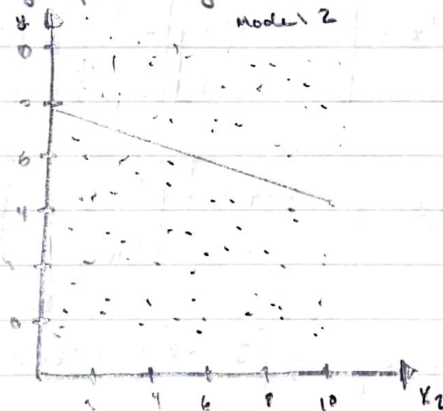
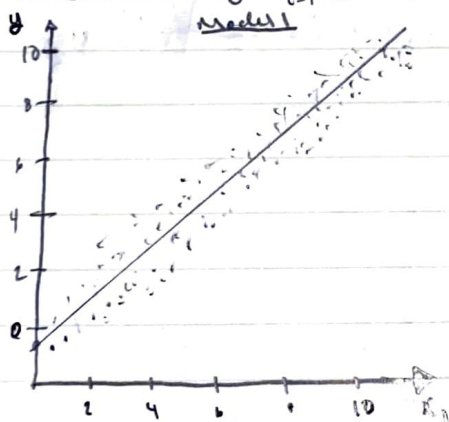
(5) Two alternative straight line regression models have been proposed for Y . In the first model, Y is a linear function of X_1 , while in the second model Y is a linear function of X_2 . The plot in the first column of figure 2.8 is that of Y against X_1 , while the second plot is that of Y against X_2 . These plots also show the least squares regression lines. In the following statements RSS stands for residual sum of squares while SS_{reg} stands for regression sum of squares. Which of the following is true?

NOTE (see book pg 28)*

• $SST = SY = \sum_{i=1}^n (y_i - \bar{y})^2 = RSS + SS_{reg}$ (Total sample variability)

• $RSS = \sum_{i=1}^n (\hat{y}_i - y_i)^2$ (variability of errors)

• $SS_{reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ (variability explained by model)



• We can see clearly from the above plots that

(i) $RSS(M_1) < RSS(M_2)$. M_1 seems to be a much better fit than M_2 .

we can see that the y values are much closer to the fitted line of M_1

than they are for M_2 . Thus, the $RSS(M_1) < RSS(M_2)$.

(ii) $SS_{reg}(M_1) > SS_{reg}(M_2)$. The line fitted by M_2 is almost a straight line at $y = \bar{y}$, thus the $SS_{reg}(M_2) = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ will be relatively small. This variable y vs X_2 seem to be uncorrelated, that is, $E[Y|X] \approx E[Y] = \bar{y}$.

$\Rightarrow \boxed{D}$

6.) (Chp 2, Question 6)

Show that $SST = SSR + RSS$. To do this, show that

$$\sum (y_i - \hat{y}_i)(y_i - \bar{y}) = 0$$

(a) Show that $(y_i - \hat{y}_i) = (y_i - \bar{y}) - \hat{\beta}_1(x_i - \bar{x})$

$$\begin{aligned} y_i - \hat{y}_i &= y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \quad * \text{We know from previous derivations of } \hat{\beta}_0 \text{ \& } \hat{\beta}_1 \text{ that} \\ &= y_i - \bar{y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \\ &= (y_i - \bar{y}) - \hat{\beta}_1(x_i - \bar{x}) \end{aligned}$$

(b) Show that $(\hat{y}_i - \bar{y}) = \hat{\beta}_1(x_i - \bar{x})$

$$\hat{y}_i - \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y} = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i - \bar{y} = \hat{\beta}_1(x_i - \bar{x})$$

(c) Utilizing the fact that $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$, show that $\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 0$

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n ((y_i - \bar{y}) - \hat{\beta}_1(x_i - \bar{x}))(\hat{\beta}_1(x_i - \bar{x})) \quad * \text{from (a) \& (b)} \\ &= \sum_{i=1}^n (y_i - \bar{y} - \frac{S_{xy}}{S_{xx}}(x_i - \bar{x}))(\frac{S_{xy}}{S_{xx}}(x_i - \bar{x})) \\ &= \frac{S_{xy}}{S_{xx}} \left[\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - \frac{S_{xy}}{S_{xx}} \sum_{i=1}^n (x_i - \bar{x})^2 \right] \\ &= \frac{S_{xy}}{S_{xx}} [S_{xy} - S_{xy}] \\ &= 0 \end{aligned}$$

See Chp 3 notes:
Pg 9, 35

7.) For the simple linear regression model $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, we find the t -statistic for testing $H_0: \beta_1 = 0$ to be $t = \frac{\hat{\beta}_1 - 0}{\text{se}(\hat{\beta}_1)}$.

(a) Which of the usual assumptions of the model must be met in order for the t -statistic to have the t -distribution? Why?

the prob about this on exam.

• Normality of errors: if the ϵ_i are not $\sim N$, then β_1 is not a linear combination of normally distributed random variables $\Rightarrow \beta_1$ is not normally distributed $\Rightarrow \hat{\beta}_1$ doesn't have a t -distribution.

(b) Does having a larger sample size change your answer? Why or why not?

Yes, if we have a larger sample size then the CLT applies.

8.) Suppose \mathbf{x} is a random n dimensional vector and that $E[\mathbf{x}] = \underline{\mu}$. Show that the covariance matrix $\Sigma = E[(\mathbf{x} - \underline{\mu})(\mathbf{x} - \underline{\mu})'] = E[\mathbf{x}\mathbf{x}'] - \underline{\mu}\underline{\mu}'$

$$\begin{aligned} \Sigma &= E[(\mathbf{x} - \underline{\mu})(\mathbf{x} - \underline{\mu})'] = E[\mathbf{x}^2 - 2\underline{\mu}'\mathbf{x} + \underline{\mu}'\underline{\mu}] \\ &= E[\mathbf{x}'\mathbf{x} - 2\underline{\mu}'\mathbf{x} + \underline{\mu}'\underline{\mu}] \\ &= E[\mathbf{x}'\mathbf{x}] - 2E[\underline{\mu}'\mathbf{x}] + E[\underline{\mu}'\underline{\mu}] \\ &= E[\mathbf{x}'\mathbf{x}] - 2\underline{\mu}'\underline{\mu} + \underline{\mu}'\underline{\mu} \end{aligned}$$

$$\Sigma = E[(\mathbf{x} - \underline{\mu})(\mathbf{x} - \underline{\mu})'] = E[\mathbf{x}'\mathbf{x}] - \underline{\mu}'\underline{\mu}$$