

1) P^{P} Regression is the study of the conditional distribution $Y | \underline{X}$ of a response variable Y or $Y | B^T \underline{X}$

given a vector of predictors \underline{X} . In a 1D regression, Y is independent of \underline{X} given a single linear combination of the predictors $SP = B^T \underline{X}$, written $Y \perp\!\!\!\perp \underline{X} | SP$.

2) A linear model is $Y = B^T \underline{X} + e$.

Multiple linear regression has at least one continuous predictor. In experimental design, all of the predictors are categorical.

3) P3 iid = independent and identically distributed

4) P3 response variable = dependent variable
predictor variables = carriers, covariates,
explanatory variables or independent variables

5) Skim or skip most of Ch 1.

Ch 2
 P^{P} The response variable Y is the variable that you want to predict, The predictor variables X_1, \dots, X_p are the

Variables used to predict Y. 1.5

ex] predict brain weight from size of head
 y x

2] Let $\underline{x} = (x_1, \dots, x_p)^T$. \underline{x} is a

column vector, \underline{x}^T is a row vector,

3] p17 A quantitative variable takes on numerical values while a qualitative variable takes on categorical values.

ex] brain weight in grams is quantitative
sex = F or M is qualitative
0 or 1

4] p18 know Suppose the response variable y is quantitative and at least one predictor variable x_i is, too. Then the multiple linear regression (MLR) model is

$$y_i = x_{i1} \beta_1 + x_{i2} \beta_2 + \dots + x_{ip} \beta_p + e_i \\ = \underline{x}_i^T \underline{\beta} + e_i \quad \text{for } i=1, \dots, n.$$

Here n is the sample size and e_i is a random variable called the ith error.

5] Notation Often the subscript i will be suppressed especially when describing models. So the MLR model is

$$Y = \underline{x}^T \underline{B} + e,$$

6] * P18 In matrix notation, the n equations become $\begin{matrix} Y \\ \sim \\ n \times 1 \end{matrix} = \begin{matrix} \underline{X} \\ \sim \\ n \times p \end{matrix} \underline{B} + \underline{e}$

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix} \begin{bmatrix} B_1 \\ ; \\ \vdots \\ B_p \end{bmatrix} + \begin{bmatrix} e_1 \\ ; \\ \vdots \\ e_p \end{bmatrix}$$

$\underbrace{\begin{matrix} X_1 & X_2 & \cdots & X_p \end{matrix}}$
columns of \underline{X}

7] Often the 1st column of \underline{X} is $X_1 = \underline{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$, the $n \times 1$ vector of 1's.

8] The i th case (\underline{x}_i^T, Y_i) corresponds to the i th row \underline{x}_i^T of \underline{X} and the i th element Y_i of \underline{Y} .

9] For the MLR model $Y = \underline{x}^T \underline{B} + e$,

Y and e are random variables but we only

have the observed values y_i and x_i . (2.5)

If the e_i are iid, $E(e_i) = 0$, $V(e_i) = \sigma^2$, then the unknown parameters are $\beta = (\beta_1, \dots, \beta_p)^T$ and σ^2 .

[10] p18 The constant Variance MLR model has e_1, \dots, e_n iid with $E(e_i) = 0$, $V(e_i) = \sigma^2$ and the errors are independent of the predictor variables. Then the cases (x_i^T, y_i) are independent for $i = 1, \dots, n$.

[11] p19 know The unimodal MLR model has the same assumptions as [10], but also the iid errors come from a unimodal distribution that is not highly skewed.

[12] p19 The normal (or Gaussian) MLR model adds the assumption that the e_i are iid $N(0, \sigma^2)$ to [10].
[12] is also a special case of [11]

Since the normal distribution is symmetric,

[13] The normality assumption is too strong, but the MLR model works best if the error distribution is roughly symmetric and not too far from normal.
Larger sample sizes n are needed if the error distribution is not roughly "mound shaped."

14) ^{p19} Notation $A \equiv B = f(c)$ means M484 3

A and B are equivalent and equal and $f(c)$ is the formula used to compute A and B .

15) ^{p17} MLR is used to study the conditional distribution of $Y | \underline{X}$ (Y given \underline{X}), to or $Y | B^T \underline{X}$

summarize the relationship between Y and \underline{X} , and for prediction.

16) ^{p19} know Let \underline{b} be an estimate of B , then the i th fitted value

$$\hat{Y}_i \equiv \hat{Y}_i(\underline{b}) = \underline{x}_i^T \underline{b} = x_{i1} b_1 + x_{i2} b_2 + \dots + x_{ip} b_p.$$

The i th residual $r_i \equiv r_i(\underline{b}) = Y_i - \hat{Y}_i(\underline{b})$.

17) The (Ordinary) least squares (OLS)

estimator $\hat{B} \equiv \hat{B}_{OLS}$ minimizes

$$Q_{OLS}(\underline{b}) = \sum_{i=1}^n r_i^2(\underline{b}) \quad \text{and} \quad \hat{B}_{OLS} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}.$$

18) know for final The i th fitted value

$$\hat{Y}_i \equiv \hat{Y}_i(\hat{B}) = \underline{x}_i^T \hat{B} = x_{i1} \hat{B}_1 + \dots + x_{ip} \hat{B}_p$$

and the i th residual $r_i = Y_i - \hat{Y}_i$.

19) We almost always use OLS for ch2 & 3 so r_i, \hat{Y}_i, \hat{B} refer to OLS.

20) ^{p19} Let $\hat{Y} = (\hat{Y}_1, \dots, \hat{Y}_n)^T$. Then, (3.5)

$\hat{Y} = \mathbf{X} \hat{\beta} = H Y$ where the hat matrix $H = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ provided the inverse exists.

21) p20 An MLR model is linear in $\underline{\beta}$.
If there is an MLR model in Y and $\underline{\beta}$,
then $Y = \underline{x}^T \underline{\beta} + e$ or $Y = \underline{w}^T \underline{\beta} + e$:

the symbol for the vector \underline{x} or \underline{w} is not important.

22) If $Y = \underline{x}^T \underline{\beta} + e = \underline{x}_1 \underline{\beta}_1 + \dots + \underline{x}_p \underline{\beta}_p + e$,

then $\frac{\partial Y}{\partial \underline{\beta}_i} = \underline{x}_i$. If $Y = \underline{w}^T \underline{\beta} + e$

then $\frac{\partial Y}{\partial \underline{\beta}_i} = w_i$.

ex) $Y_i = \beta_1 + \beta_2 w_i + \beta_3 w_i^2 + e_i$

is an MLR model in Y and $\underline{\beta}$:

take $x_{i1} = 1$, $x_{i2} = w_i$, $x_{i3} = w_i^2$.

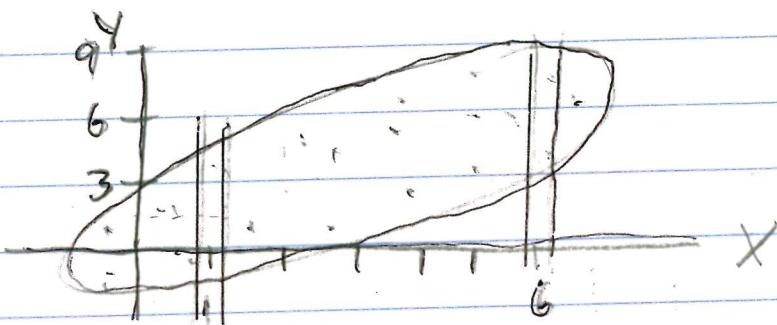
Then $Y_i = \underline{x}_i^T \underline{\beta} + e_i$.

ex) $Y_i = \frac{\beta_1}{1 + \beta_2 e^{\beta_3 x_i}} + e_i$ is not an MLR

model. Note that $\frac{\partial Y}{\partial \beta_1} = \frac{1}{1 + \beta_2 e^{\beta_3 x_i}}$ depends on $\beta_2 \neq \beta_3$.

- 23) P20 A scatterplot is a plot of X vs Y and is used to visualize the conditional distribution $Y|X$

ex



If $X=1$, $E(Y|X=1) \approx 1.5$ and $P(-2 < Y < 5 | X=1)$ is high.

If $X=6$, $E(Y|X=6) \approx 4.5$ and $P(3 < Y < 9 | X=6)$ is high but $P(-2 < Y < 5 | X=6)$ is not high.

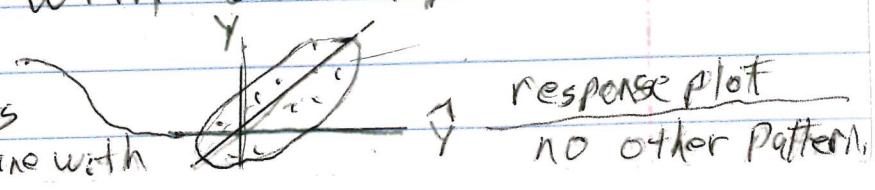
- 24) P21 know for final The response plot is a plot of $\hat{Y} = \underline{x}^T \hat{\beta}$ vs Y .

It is used to visualize the MLR model in the background of the data if $\underline{x}^T \hat{\beta}$ is a good estimator of $\underline{x}' \beta$ ($n \geq 5p$).

$$25) Y = \hat{Y} + r = \hat{Y} + y - \hat{Y}$$

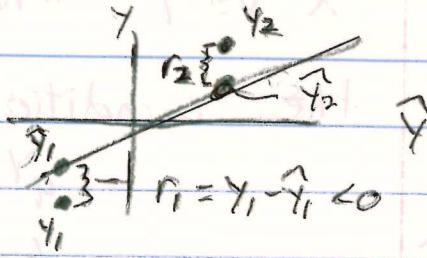
- 26) Ignoring r , $y = \hat{Y}$ is the identity line with unit slope and zero intercept.

Ideal shape: plotted points scatter about the identity line with



27] The vertical deviations from the identity line are the residuals (4.5)

$$Y_i - \hat{Y}_i = r_i$$

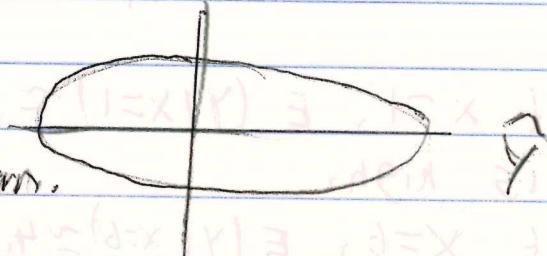


28] p21 know for final A residual plot is a plot of $\hat{Y} = X^T \hat{\beta}$ vs r .

Ideal shape: plotted points

scatter about the

line with no other pattern.



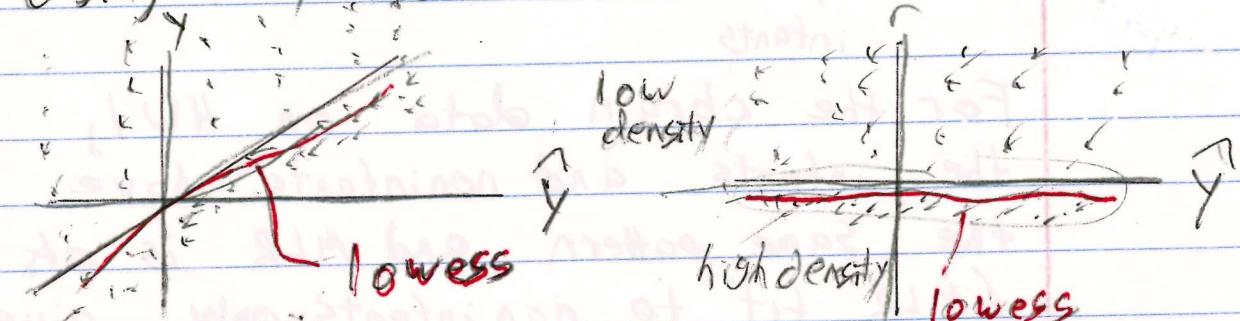
29] p21 know for final For any MLR analysis, always make the response and residual plots.

See Figures 1.2-1.4, 2.1, 2.3, 2.5a, 3.6, 3.7, 3.11

Remark: Plotting \hat{Y} vs Y is a good idea for any model that predicts Y with \hat{Y} (MLR, DOE, time series, neural networks, lasso, Econometric models).

30] The two ideal shapes are for the iid unimodal error model and $n \geq 5p$. If the plots are near the ideal shape and $n \geq 5p$ then expect inference to be approximately correct except for prediction intervals.

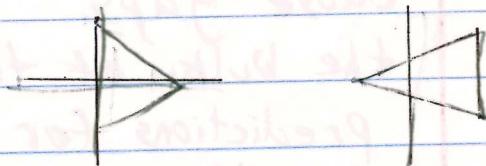
31] If the response, and residual plots suggest iid skewed errors and lowess added to the plot is close to the identity or $r=0$ lines, then expect that much larger values of n ($n >> 5p$) may be needed before inference is approx correct. The scatterplot is smoother. Lowess tries to estimate $E(Y|X)$ without using any model where $V=Y$ or $V=R$.



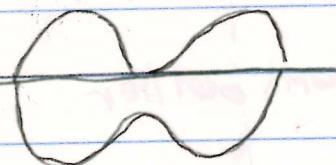
see HW2 D)

32] Do not look at a plot of r_i vs y_i because y_i and r_i are correlated and you could see a pattern.

33] Residual plots



and



all suggest that the

constant variance assumption should be checked.

34] A parabolic residual plot

suggests that one or more squared predictors X_i^2 needs to be added to the model.



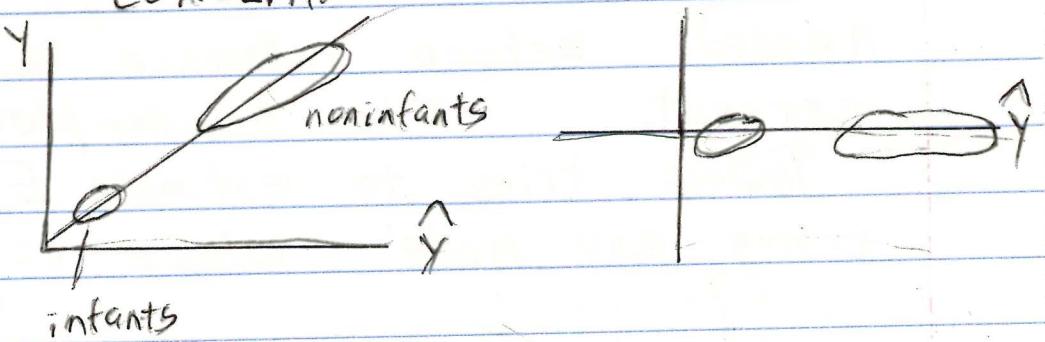
skip
to point 4)
now

Read ex 2.3 carefully.

(5.5)

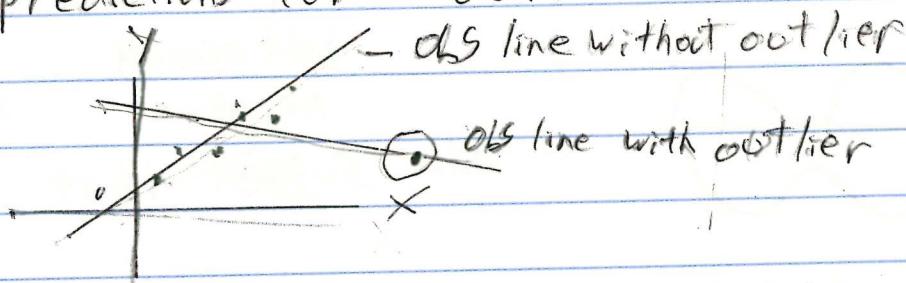
- 35) Gaps in the plots are cause to concern.

get
 $Q = X^T \beta$
from output

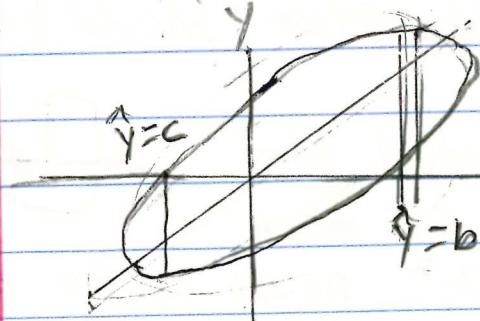


For the cbrain data on HW1, the infants and noninfants have the same pattern and MLR is off (MLR fit to noninfants only gives good predictions for infants).

But often outliers, observations far from the bulk of the data, cause gaps. Here the MLR fit to the bulk of the data gives poor predictions for outliers.



36)



$$\hat{y} = \tilde{x}^T \hat{\beta} \approx \tilde{x}^T \beta$$

For \underline{x} such that $\hat{Y} = \underline{x}^T \hat{\beta} = b \approx \underline{x}^T \beta$,

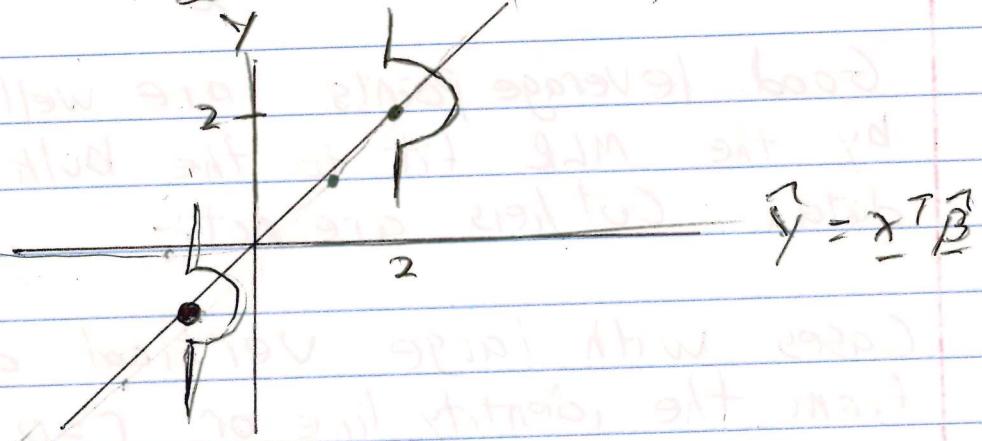
$$E(Y|\underline{x}^T \beta) \approx E(Y|\underline{x}^T \hat{\beta}) \approx \underline{x}^T \hat{\beta} = \hat{Y} = b.$$

$$\text{So } \underline{x} \text{ with } \hat{Y} = 10 \Rightarrow E(Y|\underline{x}^T \beta) \approx E(Y|\underline{x}^T \hat{\beta}) = 10$$

$$\hat{Y} = -5 \Rightarrow E(Y|\underline{x}^T \beta) \approx E(Y|\underline{x}^T \hat{\beta}) = -5$$

37) The iid normal MLR model is especially simple: $Y|\underline{x}^T \beta \sim N(\underline{x}^T \beta, \sigma^2)$

$$\text{and } Y|\underline{x}^T \beta \approx N(\underline{x}^T \hat{\beta}, \hat{\sigma}^2) = N(\hat{Y}, \text{MSE})$$



$$Y|\underline{x}^T \hat{\beta} = 2 \approx N(2, \hat{\sigma}^2)$$

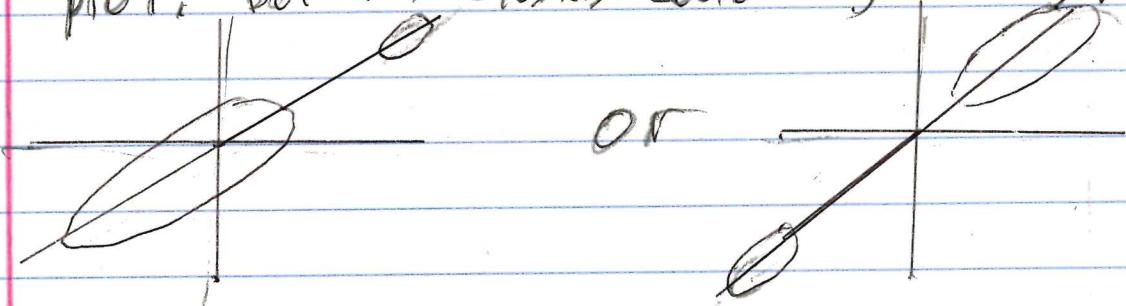
$$Y|\underline{x}^T \hat{\beta} = -1 \approx N(-1, \hat{\sigma}^2)$$

$Y|\underline{x}^T \hat{\beta} = a \approx N(a, \hat{\sigma}^2)$ for any a in the range of \hat{Y} .

Think of the picture as a 3 dimensional funnel: the road is $E(Y|\underline{x}^T \beta) = \underline{x}^T \beta = \hat{Y}$ and the different normal curves are $\approx N(\underline{x}^T \beta, \hat{\sigma}^2)$

(Replace a $\hat{\beta}$ by β , $\underline{x}^T \hat{\beta}$ by $\underline{x}^T \beta$, $\hat{\sigma}^2$ by σ^2)
 gives a similar picture

38} p22-23 Outliers lie far from the bulk of the data. If there are gaps in the plots, outliers tend to be in the upper right or lower left corner of the response plot. But these clusters could be "good leverage points".



Good leverage points are well predicted by the MLR fit to the bulk of the data. Outliers are not.

Cases with large vertical deviations from the identity line or $r=0$ line may also be outliers. But beginners tend to label too many cases as outliers. Mentally put a box about the bulk of the data ignoring any potential outliers. Then double the width (about the identity line or $r=0$ line). Cases outside the doubled box are outliers. Or mentally estimate the standard deviation of the residuals in both plots. Look for cases with residuals more than 5 standard deviations from the $r=0$ line.

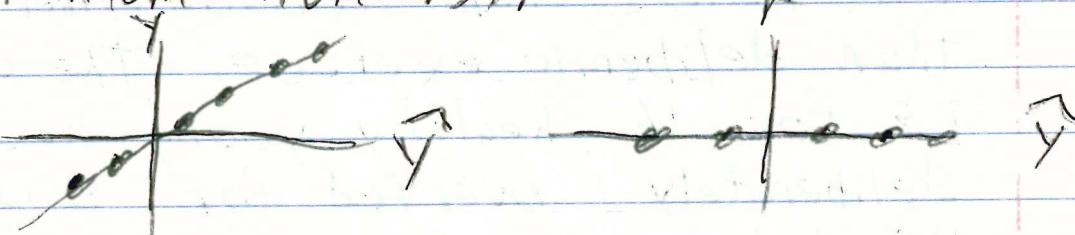
39] Want $n \geq 5p$. If $n < 5p$, try to get rid of unnecessary predictors. If

ML84 7

this is not possible, do not expect your model to work well for inferences, including prediction.

extreme case: $n = p$ then $\hat{Y} = H Y = I_n Y = Y$. So plotted points fall

exactly on the identity line and $r=0$ line regardless of the value of the predictors (the x_2, \dots, x_p could be random numbers).



40) A common mistake is to use an MLR model where $n \approx p$ or $n \approx 2p$ or $n < 5p$. These models "fit noise" and do not work well.

Data - ex's before 35)

41) computer output p 42

	label or predictor estimate or coef	std err	tval	pval
intercept or constant	$\hat{\beta}_0$			
x_2	$\hat{\beta}_2$			
x_3	$\hat{\beta}_3$			
\vdots	\vdots			
x_D	$\hat{\beta}_D$			

model contains constant unless told otherwise (7.5)

42) Know for final Assume $Y = \beta_1 + \beta_2 X_2 + \dots + \beta_p X_p$ etc.

Given X , find $\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X$

or given X_2, \dots, X_p

find $\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p$

ex) label coef stderr tvalue pvalue
constant 3.5051 30.36 11.54 0.0

NEA -0.00344 0.000741 -4.64 0.0

Suppose it is desired to predict weight gain (in kg) from NEA change, the change in energy use (in calories) from activity other than deliberate exercise. The output is for 16 healthy young adults deliberately overfed for 8 weeks.

a) Find regn, b) Predict weight gain if NEA = 485.

SOLN a) $\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X = 3.5051 - 0.00344 X$

b) $\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X = 3.5051 - 0.00344(485)$

$$= \boxed{1.8367}$$

ex) n=36 gestation time^{days}, body wt kg, brain wt g for different animals $Y = \log(\text{brain wt})$

label coef stderr tvalue pvalue
constant -0.457

log(body) 0.551

log(gest) 0.668

Predict $Y = \log(\text{brain wt})$ if $\log(\text{body}) = 2$ & $\log(\text{gest}) = 5$

$$\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 = -0.457 + .551(2) + .668(5) = \boxed{3.985}$$