

EMSE 6765: DATA ANALYSIS

For Engineers and Scientists

Session 12: Comparing Imbedded Models, Forecasting

Version: 4/05/2021



THE GEORGE
WASHINGTON
UNIVERSITY

WASHINGTON, DC

Lecture Notes by: J. René van Dorp¹

www.seas.gwu.edu/~dorpjr

¹ Department of Engineering Management and Systems Engineering, School of Engineering and Applied Science, The George Washington University, 800 22nd Street, N.W., Suite 2800, Washington D.C. 20052. E-mail: dorpjr@gwu.edu.

Regression Analysis: Log(Price) versus Elevation, Sewer, Date, Flood

Analysis of Variance

Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	4	2.2320	0.55800	21.49	0.000
Error	26	0.6753	0.02597		
Total	30	2.9072			

Model Summary

S	R-sq	R-sq(adj)
0.161156	76.77%	73.20%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value
Constant	1.4891	0.0915	16.28	0.000
Elevation	0.01411	0.00816	1.73	0.096
Sewer	-0.000044	0.000014	-3.26	0.003
Date	0.00741	0.00122	6.05	0.000
Flood	-0.3183	0.0887	-3.59	0.001

Regression Equation

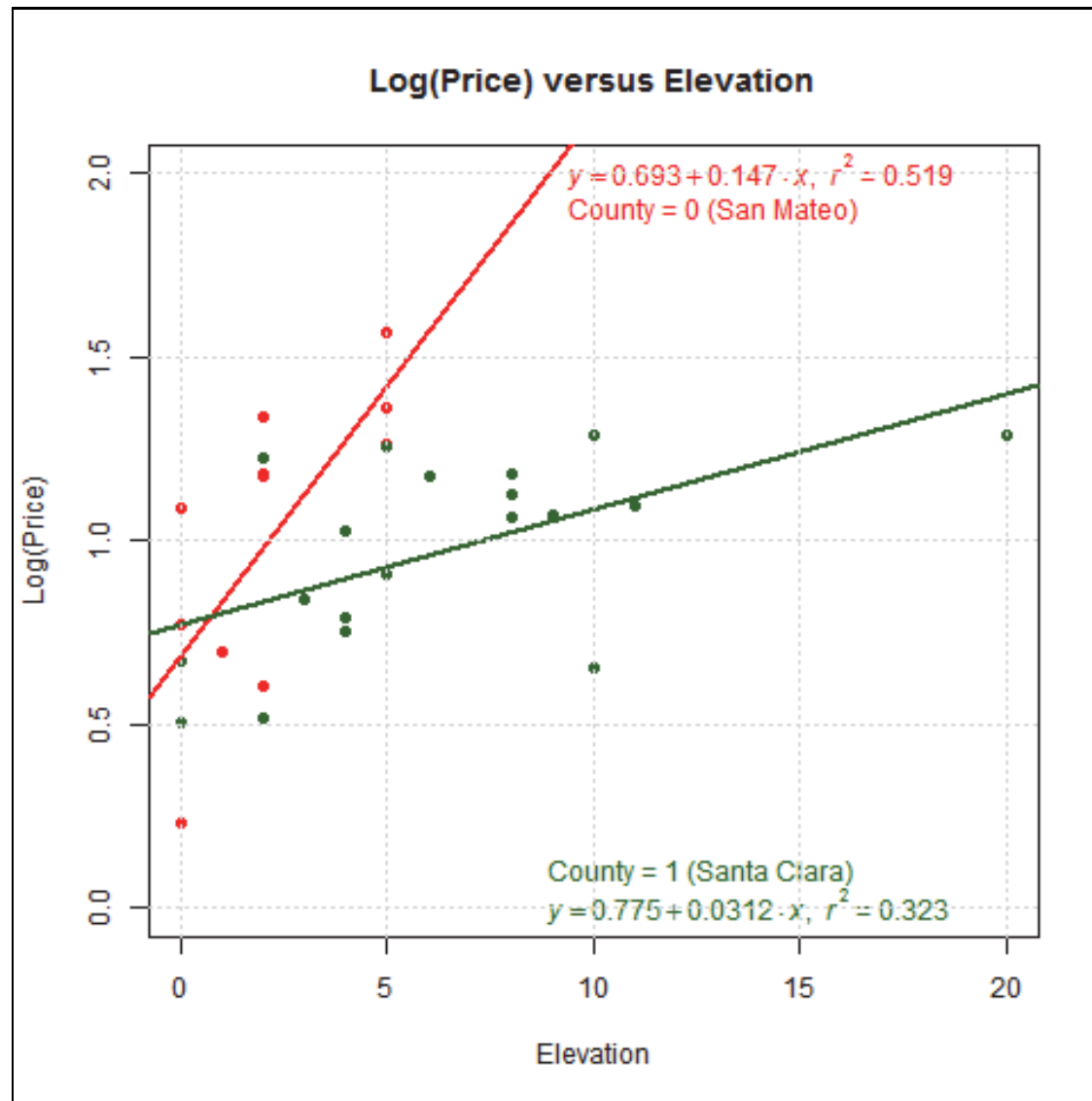
$$\text{Log(Price)} = 1.4891 + 0.01411 \text{ Elevation} - 0.000044 \text{ Sewer} + 0.00741 \text{ Date} - 0.3183 \text{ Flood}$$

- Prediction of Log(Price) using **the smaller model**:

$$Y = \mathbf{x}_0^T \hat{\mathbf{b}} + \epsilon, E[\epsilon] = 0, \epsilon \sim \mathbf{N}(\mathbf{0}, \sigma) \Leftrightarrow E[Y|\mathbf{x}_0] = \mathbf{x}_0^T \hat{\mathbf{b}}$$

Prediction for Log(Price)			
Regression Equation			
Log(Price) = 1.4891 + 0.01411 Elevation - 0.000044 Sewer + 0.00741 Date - 0.3183 Flood			
Settings			
Variable	Setting		
Elevation	0		
Sewer	0		
Date	0		
Flood	0		
Prediction			
Fit	SE Fit	95% CI	95% PI
1.48907	0.0914849	(1.30102, 1.67712)	(1.10816, 1.86999)

$\hat{y} = \mathbf{x}_0^T \hat{\mathbf{b}} \approx 1.48907$ is **both a prediction for r.v. Y and mean $E[Y|\mathbf{x}_0]$**
 (1.301, 1.677) is **conf. interval for true mean $E[Y|\mathbf{x}_0]$** , no prob. interpretation
 (1.108, 1.870) is **pred./cred. interval for r.v. Y** , with prob. interpretation



Suggestion: **Capture interaction effect using an interaction term**

$$\log(\text{PRICE}) = b_0 + b_1 \text{ELEVATION} + b_2 \text{SEWER} + b_3 \text{DATE} \\ + b_4 \text{FLOOD} + b_5 \text{COUNTY} + b_6 (\text{COUNTY} \times \text{ELEVATION})$$

When $\text{COUNTY} = 0$ **the above equation reduces to**

$$\log(\text{PRICE}) = \mathbf{b_0} + \mathbf{b_1} \text{ELEVATION} \\ + b_2 \text{SEWER} + b_3 \text{DATE} + b_4 \text{FLOOD}.$$

When $\text{COUNTY} = 1$ **the above equation reduces to**

$$\log(\text{PRICE}) = (\mathbf{b_0} + \mathbf{b_5}) + (\mathbf{b_1} + \mathbf{b_6}) \text{ELEVATION} \\ + b_2 \text{SEWER} + b_3 \text{DATE} + b_4 \text{FLOOD}$$

Thus, **the interaction effect here allows for different intercepts and slopes by counties.**

Model Summary

S	R-sq	R-sq(adj)	R-sq(pred)
0.146820	82.20%	77.76%	69.05%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	1.424	0.127	11.19	0.000	
Elevation	0.0483	0.0286	1.69	0.104	21.57
Sewer	-0.000048	0.000012	-3.86	0.001	1.32
Date	0.00548	0.00135	4.06	0.000	1.52
Flood	-0.394	0.102	-3.88	0.001	2.00
County	-0.113	0.110	-1.03	0.312	4.11
County*Elevation	-0.0307	0.0297	-1.03	0.312	28.12

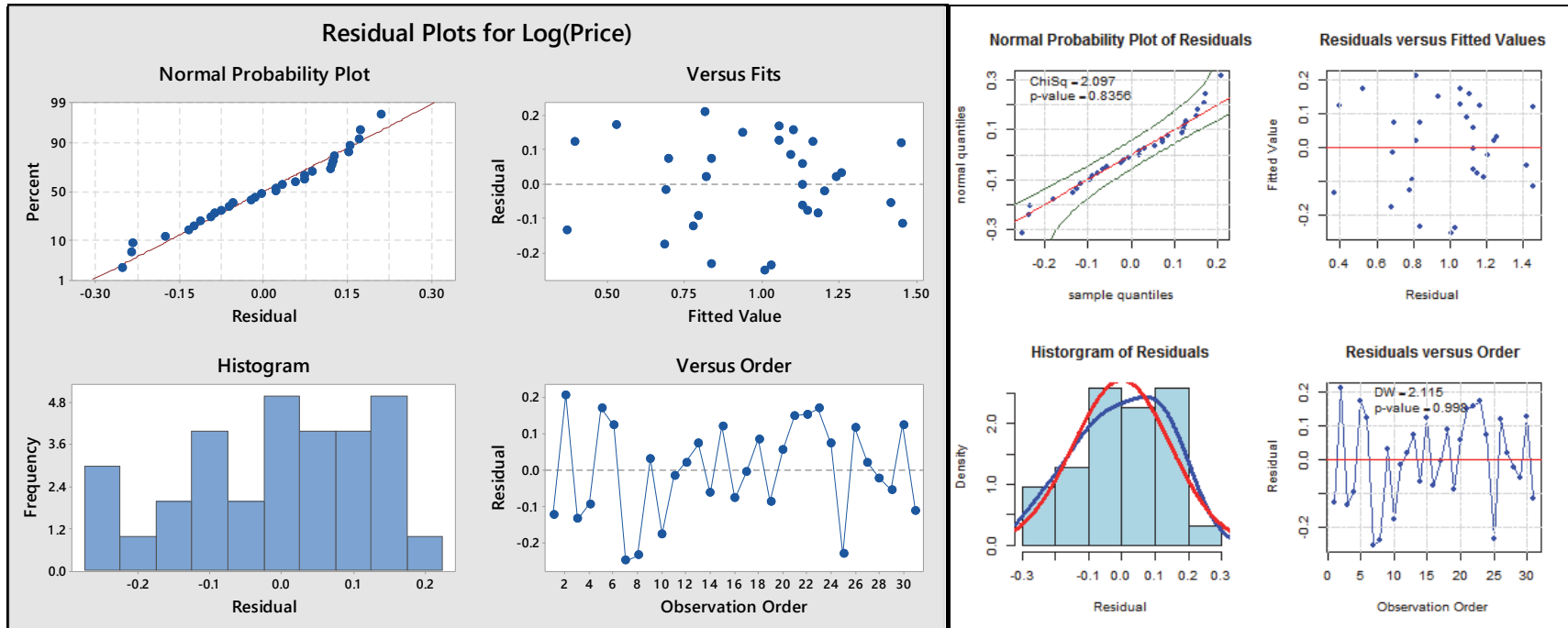
Regression Equation

$$\text{Log(Price)} = 1.424 + 0.0483 \text{ Elevation} - 0.000048 \text{ Sewer} + 0.00548 \text{ Date} - 0.394 \text{ Flood} - 0.113 \text{ County} - 0.0307 \text{ County*Elevation}$$

Durbin-Watson Statistic

$$\text{Durbin-Watson Statistic} = 2.11482$$

$R_{adj}^2 = 77.8\% \uparrow$ (Previously $R_{adj}^2 = 73.2\%$), DW -Statistic ≈ 2.11



Normality and independence assumption of residuals seem reasonable.
(although we can observe one outlier)

$e_2^* \approx 1.72$, $DFIT_2 \approx 1.02 > 2\sqrt{7/31} \approx 0.95 \Rightarrow$ Should be checked.

Model Summary

S	R-sq	R-sq(adj)	R-sq(pred)
0.146820	82.20%	77.76%	69.05%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	1.424	0.127	11.19	0.000	
Elevation	0.0483	0.0286	1.69	0.104	21.57
Sewer	-0.000048	0.000012	-3.86	0.001	1.32
Date	0.00548	0.00135	4.06	0.000	1.52
Flood	-0.394	0.102	-3.88	0.001	2.00
County	-0.113	0.110	-1.03	0.312	4.11
County*Elevation	-0.0307	0.0297	-1.03	0.312	28.12

Regression Equation

$$\text{Log(Price)} = 1.424 + 0.0483 \text{ Elevation} - 0.000048 \text{ Sewer} + 0.00548 \text{ Date} - 0.394 \text{ Flood} - 0.113 \text{ County} - 0.0307 \text{ County*Elevation}$$

Durbin-Watson Statistic

$$\text{Durbin-Watson Statistic} = 2.11482$$

Is the improvement in the R^2 -values from 76.7% to 82.2% (a jump of about 5.4%) statistically significant ?

- When the simpler model description is **completely contained within** the description of the larger model, we can perform **an F -hypothesis test**:

Explanatory Variables in the full model					
County	Elevation	Sewer	Date	Flood	County*Elevation
Explanatory Variables in the restricted/small model					
	Elevation	Sewer	Date	Flood	
Conclusion: All variables of small/restricted model are variables in the full model and "the increase in R^2 test" can be performed					

H_0 : No model improvement , H_1 : Model Improvement

R_f^2 : R^2 -value of **the full model**, R_r^2 : R^2 -value of **restricted model**

df_f : Degrees of Freedom of **Residual/Error Term** in **full model**

df_r : Degrees of Freedom of **Residual/Error Term** in **restricted model**

$$F = \frac{(R_f^2 - R_r^2)/(df_r - df_f)}{(1 - R_f^2)/df_f} \sim F_{(df_r - df_f), df_f}$$

Full Model	
R Square	82.20%
Degrees of Freedom	24
Small Model	
R Square	76.77%
Degrees of Freedom	26

	Value	Df
Numerator	0.0272	2
Denominator	0.0074	24
F-Statistic	3.663	
α	5%	
Critical Value	3.403	
Conclusion	Model Improvement	
p-value	4.09%	
Conclusion	Model Improvement	

- $R_f^2 = 82.2\%$, $df_f = 24$, $R_r^2 = 76.7\%$, $df_r = 26 \Rightarrow F$ -statistic ≈ 3.663 .
 $3.663 > F_{2,24,0.95} \approx 3.403 \Rightarrow F$ -statistic observation in 5% tail of $F \sim F_{2,24}$

Conclusion: Reject H_0 in favor of H_1 and an improvement in the model is detected in terms of the increased R^2 -value.

- This test requires acceptance of **normality assumption of residuals** for both models!
- This test is usefull when small increases in R^2 are observed, but **it requires** that **the smaller model to be nested in the larger model**.

However, the improvement in R^2 and reduction in the Standard Error (SE) here comes at a cost! Note the increase in VIF Factors of the coefficients:

Model Summary					
S	R-sq	R-sq(adj)	R-sq(pred)		
0.146820	82.20%	77.76%	69.05%		

Coefficients					
Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	1.424	0.127	11.19	0.000	
Elevation	0.0483	0.0286	1.69	0.104	21.57
Sewer	-0.000048	0.000012	-3.86	0.001	1.32
Date	0.00548	0.00135	4.06	0.000	1.52
Flood	-0.394	0.102	-3.88	0.001	2.00
County	-0.113	0.110	-1.03	0.312	4.11
County*Elevation	-0.0307	0.0297	-1.03	0.312	28.12

Regression Equation

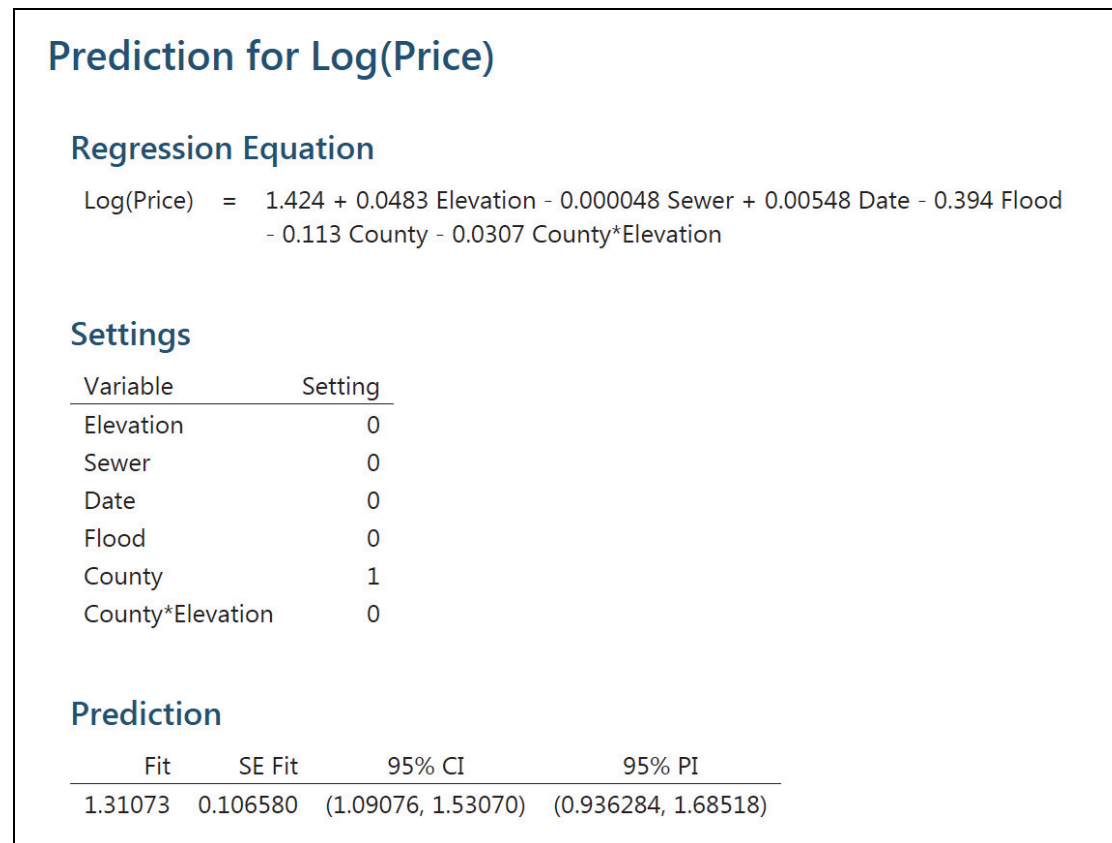
$$\text{Log(Price)} = 1.424 + 0.0483 \text{ Elevation} - 0.000048 \text{ Sewer} + 0.00548 \text{ Date} - 0.394 \text{ Flood} - 0.113 \text{ County} - 0.0307 \text{ County*Elevation}$$

Durbin-Watson Statistic

Durbin-Watson Statistic = 2.11482

Conclusion: While the standard error in the residual error term has reduced from 0.1611 to 0.1468, the uncertainty in standard errors of the coefficient estimators have increased.

- Both the standard errors** in the coefficients and the residuals contribute to the standard error of the prediction, **i.e. its uncertainty**. For the 247 acres property at hand we have:



Prediction interval width is now: $1.68518 - 0.936284 \approx 0.748896$

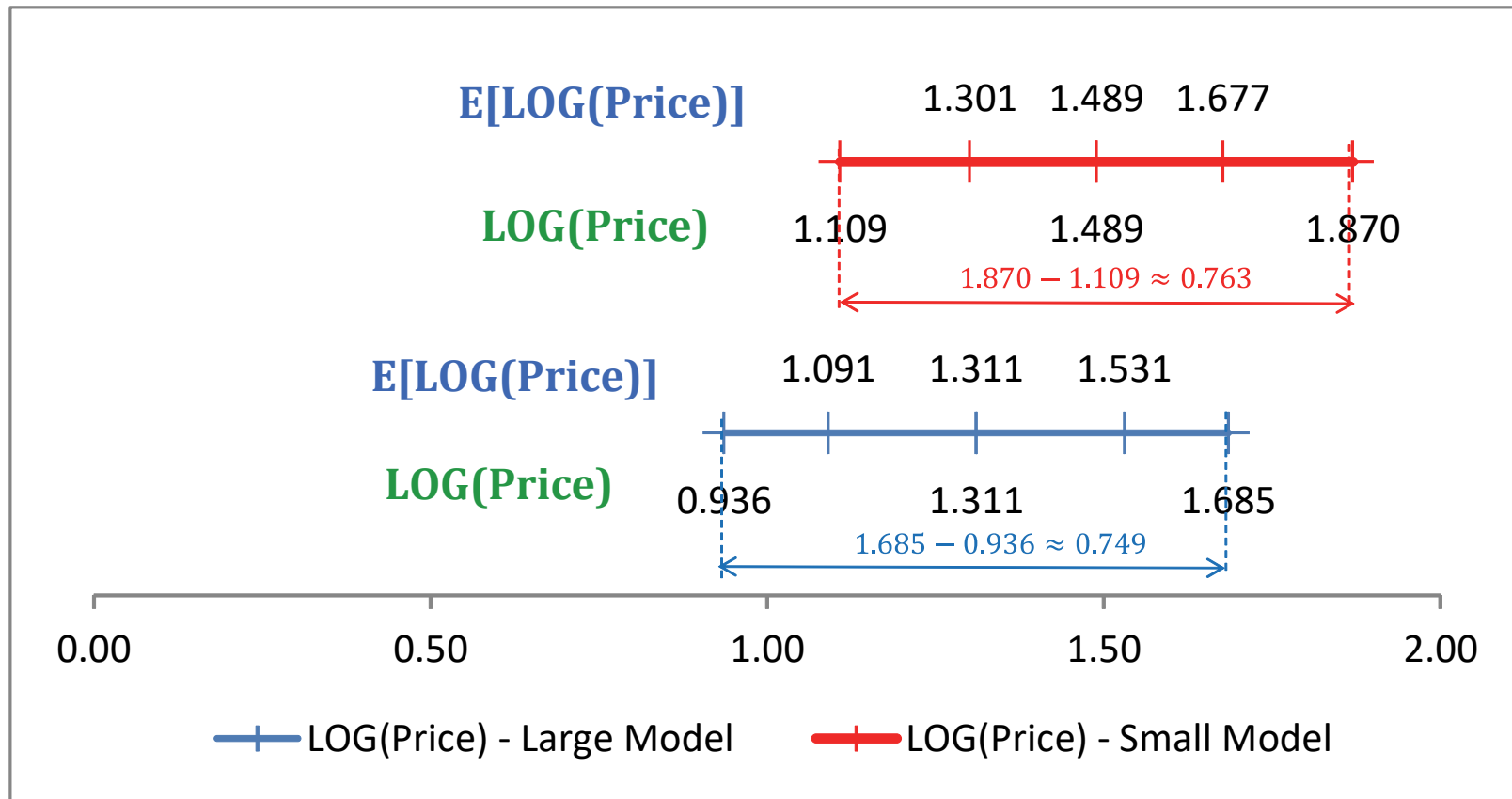
- Prediction of Log(Price) using **the smaller model**:

Prediction for Log(Price)				
Regression Equation				
Log(Price) = 1.4891 + 0.01411 Elevation - 0.000044 Sewer + 0.00741 Date - 0.3183 Flood				
Settings				
Variable	Setting			
Elevation	0			
Sewer	0			
Date	0			
Flood	0			
Prediction				
Fit	SE Fit	95% CI	95% PI	
1.48907	0.0914849	(1.30102, 1.67712)	(1.10816, 1.86999)	

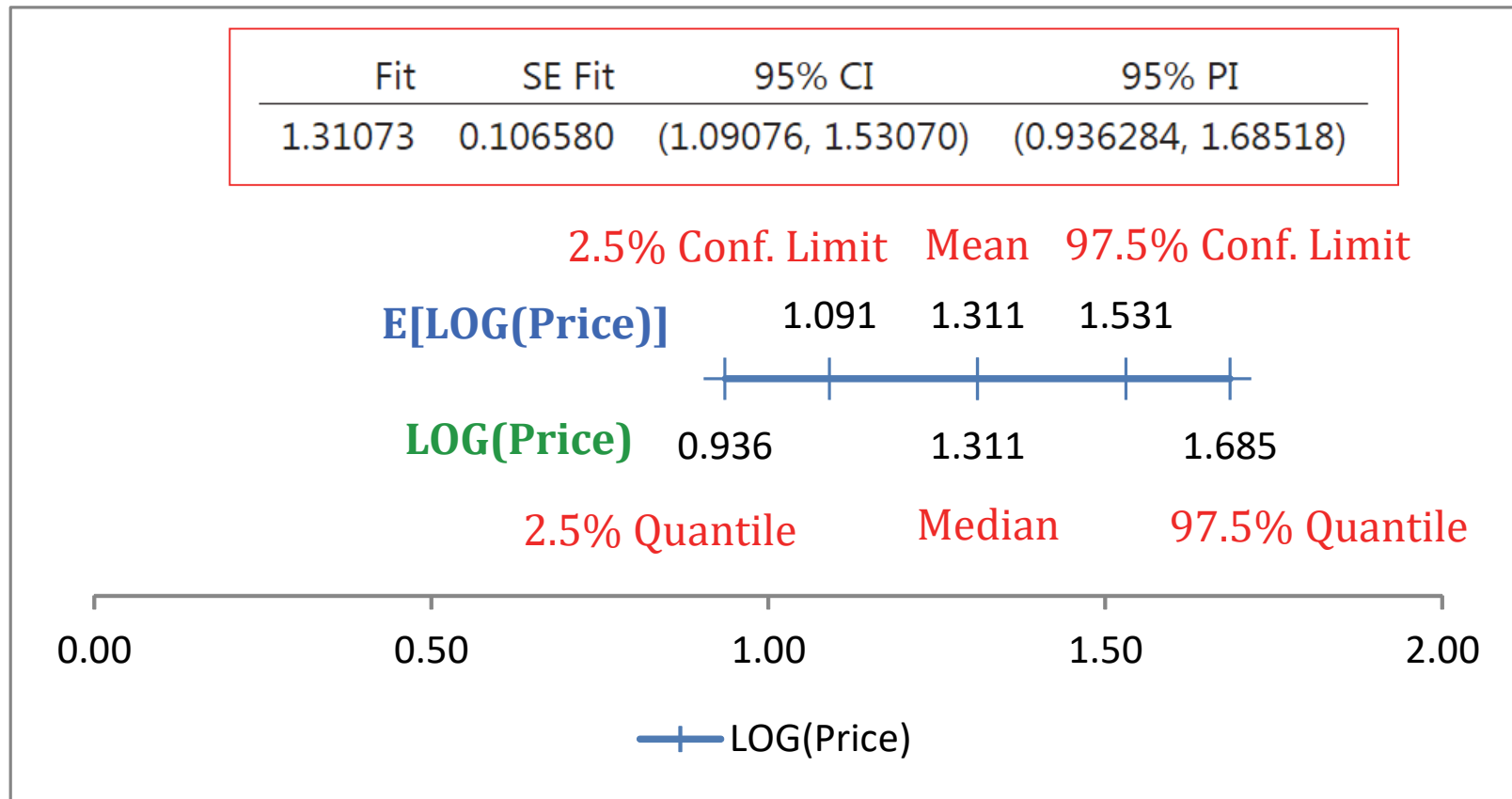
Prediction interval width was: $1.86999 - 1.10816 \approx 0.76183$

Conclusion: Prediction interval width of the smaller model is larger than the prediction interval width of the full model **despite the large VIF factors.**

Thus continue to predict\forecast with the full\larger model!



Conclusion: Prediction interval width of the smaller model is larger than the prediction interval width of the full model **despite the large VIF factors.**
Thus continue to predict\forecast with the full\larger model!



Observe the Median and the Mean are of the same value! Why?

$$Y = \mathbf{x}_0^T \hat{\mathbf{b}} + \epsilon, E[\epsilon] = 0, \epsilon \sim \mathbf{N}(0, \sigma) \Leftrightarrow E[Y | \mathbf{x}_0] = \mathbf{x}_0^T \hat{\mathbf{b}}$$

Prediction			
Fit	SE Fit	95% CI	95% PI
1.31073	0.106580	(1.09076, 1.53070)	(0.936284, 1.68518)

$$Pr(\text{Log}(\text{Price}) \leq 1.31073 | x_0) \approx 50\% \Leftrightarrow Pr(10^{\text{Log}(\text{Price})} \leq 10^{1.31073} | x_0) \approx 50\%$$

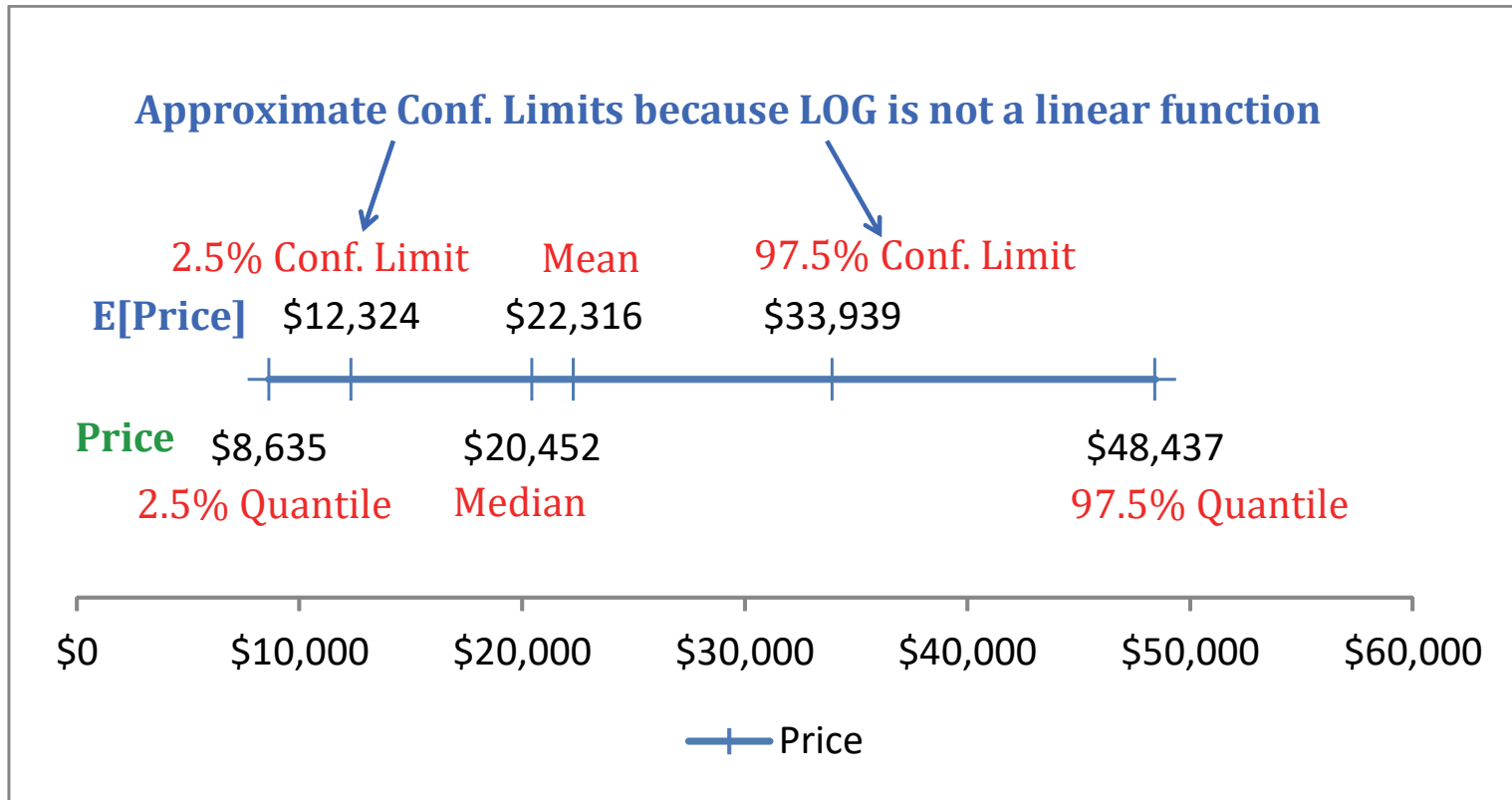
$$\Leftrightarrow Pr(\text{Price} \leq \$20452 | x_0) \approx 50\% \text{ (Recall Price was measured in \$000's)}$$

Hence: \$20452 is **the median estimate for the Price per Acre**

Thus we have here that: $Med[\text{Log}(\text{Price})] = \text{Log}(Med[\text{Price}])$

95% Confidence Interval	
LB E[LOG(PRICE)]	1.09076
UB E[LOG(PRICE)]	1.53070
Approximate 95% Confidence Interval	
LB E[PRICE]	\$12,324.22
UB E[PRICE]	\$33,939.08

95% Prediction Interval (or Credibility Interval)			
LB LOG(PRICE)	0.936284		
UB LOG(PRICE)	1.685175		
95% Prediction Interval (or Credibility Interval)			
PRICE	\$8,635.44		
PRICE	\$48,436.78		



These are approximate Confidence Limits since we know:

$$E[\text{Log}(\text{Price})] \neq \text{Log}(E[\text{Price}])$$

How do we get? $\hat{E}[\text{Price}|\mathbf{x}_0] \approx \$22,316$

Prediction			
Fit	SE Fit	95% CI	95% PI
1.31073	0.106580	(1.09076, 1.53070)	(0.936284, 1.68518)

$$Y = \mathbf{x}_0^T \hat{\mathbf{b}} + \epsilon, E[\epsilon] = 0, \epsilon \sim N(0, SE) \Rightarrow E[Y|\mathbf{x}_0] = \mathbf{x}_0^T \hat{\mathbf{b}}$$

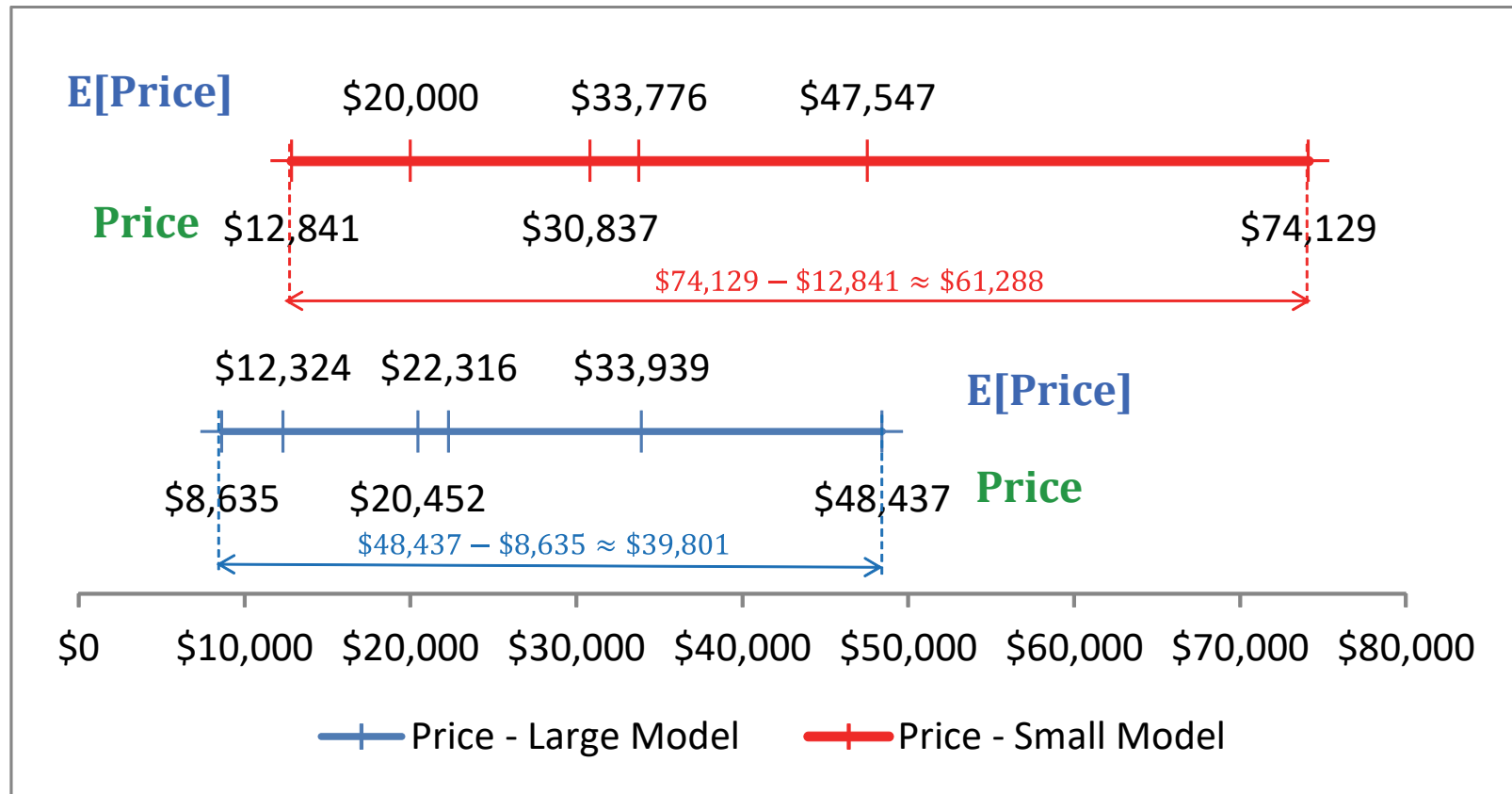
$$\begin{aligned} V[Y|\mathbf{x}_0] &= V[\mathbf{x}_0^T \hat{\mathbf{b}}] + V[\epsilon] \Rightarrow V[Y] = (SE_{Fit})^2 + (SE_{Residuals})^2 \\ &= (\mathbf{0.106580})^2 + (\mathbf{0.14682})^2 = (\mathbf{0.181426})^2 \end{aligned}$$

Summarizing:

$$\begin{cases} E[\text{Log}(\text{Price})|\mathbf{x}_0] = 1.31073 = \mu, \\ V[\text{Log}(\text{Price})|\mathbf{x}_0] = (\mathbf{0.181426})^2 = \sigma^2 \end{cases} \Rightarrow E[\text{Price}] = 10^{\mu + \text{Ln}(10) \times \frac{\sigma^2}{2}}$$

$$E[\text{Price}] \approx 10^{1.311 + 2.302 \times \frac{(0.181)^2}{2}} \approx 22.316 (\times \$1000 \approx \$22,316).$$

Note the following formula in the book is wrong: $E[\text{Price}] = 10^{\mu + \frac{\sigma^2}{2}}$



Conclusion: Prediction interval width of the smaller model is larger than the prediction interval width of the full model **despite the large VIF factors.**
Thus predict\forecast with the full\larger model!

Prediction			
Fit	SE Fit	95% CI	95% PI
1.31073	0.106580	(1.09076, 1.53070)	(0.936284, 1.68518)

- The variance of the term $\mathbf{x}_0^T \mathbf{b}$, where \mathbf{b} is the estimator vector of the coefficients, is estimated by $\sqrt{SE^2 \times \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}$ and is referred to as **the sampling error**. SE is the standard error of the residuals.
- A $100(1 - \alpha)\%$ confidence interval for the mean $E[y|x_0]$ is :

$$\mathbf{x}_0^T \hat{\mathbf{b}} \pm t_{n-p-1, 1-\alpha/2} \times SE \times \sqrt{[\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0]}$$

$$\mathbf{x}_0^T \hat{\mathbf{b}} \pm t_{n-p-1, 1-\alpha/2} \times (0.106580) = (1.09076, 1.53070)$$

Standard Error : sample standard deviation of the residuals.

\mathbf{x}_0 : values of the explanatory variables for which you would to forecast the dependent variable y .

$\hat{\mathbf{b}}$: The estimates of the regression coefficients.

Prediction			
Fit	SE Fit	95% CI	95% PI
1.31073	0.106580	(1.09076, 1.53070)	(0.936284, 1.68518)

- The variance of the residuals SE^2 is referred to as **the model error**. The variance of the prediction $Y = \mathbf{x}_0^T \mathbf{b} + \epsilon$ is **the sum of the sampling error and the model error (there is independence between the two terms)**.

$$\begin{aligned} Var(Y|x_0) &= (\text{Standard Error})^2 \times [\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0] + (\text{Standard Error})^2 \\ &= (0.106580)^2 + (0.14682)^2 = (0.181426)^2 \end{aligned}$$

- A **100(1 - α)% prediction interval for the random variable ($y|x_0$) is :**

$$\mathbf{x}_0^T \hat{\mathbf{b}} \pm t_{n-p-1, 1-\alpha/2} \times \sqrt{Var(Y|x_0)} = (0.936284, 1.68518)$$

$$\mathbf{x}_0^T \hat{\mathbf{b}} \pm t_{n-p-1, 1-\alpha/2} \times (0.181426) = (0.936284, 1.68518)$$

\mathbf{x}_0 : values of the explanatory variables to forecast the dependent variable y .

$\hat{\mathbf{b}}$: The estimates of the regression coefficient.

- **Summarizing**, the value $\hat{y} = \mathbf{x}_0^T \hat{\mathbf{b}}$ is both an estimate of **the random variable** $(Y|\mathbf{x}_0)$, but also of its **expected value** $E[Y|\mathbf{x}_0]$.
- When describing the uncertainty in the estimate for $E[Y|\mathbf{x}_0]$, one only has to account for **the uncertainty in the regression coefficients**, leading to **a confidence interval** for $E[Y|\mathbf{x}_0]$.
- When describing the uncertainty in the random variable $(Y|x_0)$ one has to account for both **the uncertainty in the regression coefficients** and **in the residuals**, leading to **a prediction/credibility interval** for $(y|x_0)$.
- The vector \mathbf{b} is an **estimator -vector** for the regression coefficients, where $\mathbf{b} \sim \mathbf{MVN}(\hat{\mathbf{b}}, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$. For its variance-covariance matrix estimate we have $\hat{\Sigma}(\mathbf{b}) = SE^2(\mathbf{X}^T \mathbf{X})^{-1}$, where SE is **the residual standard error**.
- From the above it follows that:

$$Var(\hat{Y}|\mathbf{x}_0) = SE^2 \times \left[\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0 + 1 \right] = (\mathbf{0.181426})^2$$

- It is known **for the natural logarithm $\text{Ln}(\cdot)$** that:

$$X \sim \text{LN}(\mu, \sigma) \Leftrightarrow \text{Ln}(X) \sim N(\mu_x, \sigma_x) \Rightarrow E[X|\mu_x, \sigma_x] = e^{\mu_x + \sigma_x^2/2} \quad (1)$$

- It is known for **the logarithm $\text{Log}(\cdot)$ with base 10** that:

$$\text{Log}(Y) = \frac{\text{Ln}(Y)}{\text{Ln}(10)} \Leftrightarrow \text{Ln}(Y) = \text{Ln}(10) \times \text{Log}(Y) \quad (2)$$

- Hence from (2), **$\text{Ln}(Y)$ is linear transformation of $\text{Log}(Y)$** :

$$\text{Log}(Y) \sim N(\mu, \sigma) \Rightarrow \text{Ln}(Y) \sim N\{\text{Ln}(10) \times \mu, \text{Ln}(10) \times \sigma\} \quad (3)$$

- With **the expected value expression** in (1) and (3) it now follows that:

$$\begin{cases} \mu_y = \text{Ln}(10) \times \mu \\ \sigma_y = \text{Ln}(10) \times \sigma \end{cases} \Rightarrow \begin{cases} E[Y|\mu, \sigma] = e^{\mu_y + \sigma_y^2/2} = e^{\text{Ln}(10) \times \mu + \{\text{Ln}(10) \times \sigma\}^2/2} \\ = \left[e^{\text{Ln}(10)} \right]^{\mu + \text{Ln}(10) \times \sigma^2/2} \\ = 10^{\mu + \text{Ln}(10) \times \sigma^2/2} \end{cases} \quad (4)$$