

Q1.

1. What makes a model linear? Linear in what?

When the independent and dependent terms have a linear relationship (makes a line rather than a curve) in the expression.

2. How do you interpret the coefficient for a dummy/one-hot-encoded variable? (This is a trick question, and the trick involves how you handle the intercept of the model.)

The coefficient for a dummy variable represents the average difference in the dependent variable between the group represented by the dummy variable and the reference group.

3. Can linear regression be used for classification? Explain why, or why not.

Linear regression predicts continuous values, but classification requires discrete values.

4. What are signs that your linear model is over-fitting?

High variance is one sign that the model is over-fitting meaning the model is sensitive to small changes in the training data. Also, if there is a big difference in performance (accuracy) of the training and testing set. For example, if the training set is highly accurate, but then the model does not perform well on new data, this indicates that the model is likely over-fitting.

5. Clearly explain multi-collinearity using the two-stage least squares technique.

Multicollinearity is when the independent variables are highly correlated. The two stage least squares technique can minimize the effects of this problem. In the first stage, you regress the endogenous variable on the exogenous variable. The goal here is to predict an endogenous variable based on the exogenous one to predict a "clean" value. In the second stage, the dependent variable is regressed on the predicted values (from stage one). The goal here is to estimate the relationship between the dependent and predicted variable to lessen the impact of multicollinearity.

6. How can you incorporate nonlinear relationships between your target/response/dependent/outcome variable y and your features/control/response/independent variables x into your analysis?

One way to incorporate nonlinear relationships between x and y variables is to use tree models like random forests, which can capture complex nonlinear relationships.

7. What is the interpretation of the intercept? A slope coefficient for a variable? The coefficient for a dummy/one-hot-encoded variable?

The intercept is the predicted value of y when the x is zero. The slope is the value of a one unit increase in x . The coefficient of a dummy variable is the difference of the mean of the dependent variable between the category and reference category.

$$\textcircled{16} \quad y_i = b_0 + b_1 z_{i1} + b_2 z_{i2}$$

$$\frac{1}{N} \sum_{i=1}^N z_{ij} = 0$$

1. SSE for the model

$$SSE = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

2. partial derivatives w respect to b_0, b_1, b_2

$$\frac{\partial \hat{y}_i}{\partial b_0} = 1 \quad \frac{\partial \hat{y}_i}{\partial b_1} = z_{i1} \quad \frac{\partial \hat{y}_i}{\partial b_2} = z_{i2}$$

3. average error is zero $e \cdot z = 0$

• error term $e_i = y_i - \hat{y}_i$

$$\sum_{i=1}^N e_i = 0 \quad \sum_{i=1}^N e_i z_{ij} = 0 \quad \text{for } j=1,2$$

$$\sum_{i=1}^N (y_i - b_0 - b_1 z_{i1} - b_2 z_{i2}) = 0$$

$$\text{average error} = \frac{1}{N} \sum_{i=1}^N e_i = 0$$

$$\sum_{i=1}^N z_{i1} (y_i - b_0 - b_1 z_{i1} - b_2 z_{i2}) = 0$$

$$\sum_{i=1}^N z_{i2} (y_i - b_0 - b_1 z_{i1} - b_2 z_{i2}) = 0$$

$$\sum_{i=1}^N z_{i1} e_i = 0 \quad \sum_{i=1}^N z_{i2} e_i = 0$$

4. optimal intercept $b_0 = \bar{y}$

$$\frac{\partial SSE}{\partial b_0} = \sum_{i=1}^N 2(y_i - b_0 - b_1 z_{i1} - b_2 z_{i2})(-1)$$

$$= -2 \sum_{i=1}^N (y_i - b_0 - b_1 z_{i1} - b_2 z_{i2})$$

$$= \sum_{i=1}^N y_i - \sum_{i=1}^N (b_0 + b_1 z_{i1} + b_2 z_{i2})$$

$$\frac{1}{N} \sum_{i=1}^N y_i = b_0 + b_1 \frac{1}{N} \sum_{i=1}^N z_{i1} + b_2 \frac{1}{N} \sum_{i=1}^N z_{i2}$$

$$\frac{1}{N} \sum_{i=1}^N z_{i1} = 0, \frac{1}{N} \sum_{i=1}^N z_{i2} = 0$$

$$b_0 = \frac{1}{N} \sum_{i=1}^N y_i = \bar{y} \quad (b_0 = \bar{y})$$

5. results as matrix equation $Ab=c$

$$\begin{bmatrix} \sum z_{i1}^2 & \sum z_{i1}z_{i2} \\ \sum z_{i1}z_{i2} & \sum z_{i2}^2 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} \sum z_{i1}y_i \\ \sum z_{i2}y_i \end{bmatrix}$$

A B = C

6. What is matrix A? What is vector c?

Explain intuition.

$$\frac{1}{N} \begin{bmatrix} \sum z_{i1}^2 & \sum z_{i1}z_{i2} \\ \sum z_{i1}z_{i2} & \sum z_{i2}^2 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \frac{1}{N} \begin{bmatrix} \sum z_{i1}y_i \\ \sum z_{i2}y_i \end{bmatrix}$$

$$A = \frac{1}{N} \begin{bmatrix} \sum_{i=1}^N z_{i1}^2 & \sum_{i=1}^N z_{i1}z_{i2} \\ \sum_{i=1}^N z_{i1}z_{i2} & \sum_{i=1}^N z_{i2}^2 \end{bmatrix} \quad \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \frac{1}{N} \begin{bmatrix} \sum_{i=1}^N z_{i1}y_i \\ \sum_{i=1}^N z_{i2}y_i \end{bmatrix}$$

$$\left[\begin{array}{l} \sum_{i=1}^N z_{i1} z_{i2} - \frac{1}{N} \sum_{i=1}^N z_{i1} \sum_{i=1}^N z_{i2} \\ \sum_{i=1}^N z_{i2}^2 - \frac{1}{N} \left(\sum_{i=1}^N z_{i2} \right)^2 \end{array} \right]$$

$$c = \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N z_{i1} \bar{y}_i \\ \frac{1}{N} \sum_{i=1}^N z_{i2} \bar{y}_i \end{bmatrix}$$

The vector c represents the covariances between the centered predictor & centered response variable.

Q7

$$a^0 = \bar{y}$$

$$b^0 = \frac{\sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^N (x_i - \bar{x})^2}$$

1. when will b^0 be large/small w/ relationship between x, y & variance of x ?

- the numerator is the covariance between y & x . Since they are both in the numerator, this means if y increases, then x does & vice versa. The denominator is the variance of x (amount x varies around the mean (\bar{x})). b^0 will be large when x & y have a strong relationship (covariance) or x has a small variance. b^0 will be small if x & y have a weak relationship or x has a large variance.

2. affect intercept of regression? What happens to b^0 coefficient? How does it affect ability to predict?

- The intercept is independent, so the intercept will not be affected by measurement error in X_i .

- When there is measurement error, b^0 will tend to be biased toward zero.

- The ability to predict ^{accurately} will lessen.

There will be greater uncertainty

3. noise n_i is independent

- Basically, the noise does not introduce bias (bias if zero mean + independent), but it does reduce precision of our estimate of the relationship of X & Y . The added noise of X inflates the variance.

This leads to attenuation of b^0 .

4. Attenuation is important in the cost-benefit analysis of data quality.

When measurement error is large, cleaner data/ accurate measurements are largely beneficial for predictions. Alternatively, if the measurement error is small, more effort to clean data is unlikely not worth the cost.