

- Two Review Quizzes
- Multiple Regression
- Regression Review
- Answer Questions

24.0 Multiple Regression

distributed).
be independent or other errors, have mean zero, and be normally
data (i.e., the (X^i, Y^i) pairs), and ϵ^i is random error (this should
Here a and b are unknown constants estimated from the observed

$$Y^i = a + b X^i + \epsilon^i.$$

Recall that in simple linear regression one tried to predict Y from X by
assuming:

24.1 Regression: A Review

If one wanted, one could also test whether femur length, $\hat{b} = 1.0941$, is significantly different than chance explanation. This is equivalent to testing whether the estimated coefficient on femur length, $\hat{b} = 1.0941$, is significantly different from 0.

where $e \sim N(0, \sqrt{15.675})$. This relationship explains 96.6% of the variation in humerus length using information on femur length.

$$\text{humerus} = .9207 + 1.0941 * \text{femur} + e$$

In the archaeoptryx example, we predicted the length of the humerus bone from the length of the femur. JMP indicated that the estimated linear relationship was:

and b .)

This is compared to a t -distribution with $n - 2 = 5 - 2 = 3$ degrees of freedom. (Recall: we lose information equivalent to one observation for each estimate we make, and in regression we have had to estimate both a and b .)

$$ts = \frac{\hat{b} - 0}{se} = \frac{1.0941 - 0}{.14998} = 7.3.$$

The test statistic is

$$H^0: b = 0 \text{ vs. } H^A: b \neq 0.$$

To make this test, the null and alternative hypotheses are:

This kind of model is used by wine speculators, but gets more complex.

$$\text{price} = a + b_1(\text{avg. rainfall}) + b_2(\text{avg. temp.}) + b_3(\text{calcium in soil}) + \epsilon.$$

As an example, suppose one wanted to predict the price of wine:

Again, the ϵ_i are independent normal rvs with mean 0.

$$Y_i = a + b_1 X_1^{p_i} + b_2 X_2^{p_i} + \dots + b_d X_d^{p_i} + \epsilon_i.$$

model is

In multiple regression, there is more than one explanatory variable. The

24.2 Multiple Regression

- \log , used when the scatter of the Y values increases as one or more of the X values gets larger.
- arcsine of the square root (useful when Y is a proportion or a percent)
- In practice, one often needs to make some transformation of the response variable (price) to improve the linear relationship. Common transformations include:

We may want to make a transformation of the # of cigarettes per day variable—the difference between 0 and 1 probably has more impact than the difference between 40 and 41.

We expect b_1 to be negative, b_2 to be positive.

$$\text{Lifespan} = a + b_1(\# \text{ cigarettes/day}) + b_2(\text{avg. grandparent age}) + e.$$

As another example,

Age, seniority, and experience were measured in months.

Sex was given a 1 if the person was female, and 0 for a male.

$$b_5(\text{sex}) + \epsilon.$$

$$\text{salary} = a + b_1(\text{sex}) + b_2(\text{seniority}) + b_3(\text{age}) + b_4(\text{educ}) +$$

To assess discrimination, lawyers used multiple regression:

with similar credentials.

In the 1970's, Harris Trust and Savings Bank was accused of gender discrimination in starting salaries. In particular, one main question was whether men in entry-level clerical jobs got higher salaries than women with similar credentials.

The legal question was whether the coefficient b_1 was significantly less than 0. If so, then the effect of gender was to lower the starting salary. The JMP output shows that the estimated model is

$$\text{salary} = 6277.9 - 767.9(\text{sex}) - 22.6(\text{seniority}) + .63(\text{age}) + 92.3(\text{educ}) + 50(\text{exper}) + \epsilon.$$

We observe that the coefficient on sex is negative, which suggests that there may be discrimination. But we still need a significance test. We cannot interpret the size of the effect without one, and it may just be due to chance.

discrimination.

The result is highly significant. Reject the null; there is evidence of degrees of freedom. Since this is off our t -table scale, we use a z -table. This is compared to a t -distribution with $n - p - 1 = 93 - 5 - 1 = 87$

$$ts = \frac{se}{\hat{q}^1 - 0} = \frac{128.9}{767.9 - 0} = 5.95.$$

The test statistic is

$$H_0: q^1 \leq 0 \text{ vs. } H_A: q^1 > 0.$$

The null and alternative hypotheses are: