# 9 Inference for curve and surface fitting

Previously, we have discussed how to describe relationships between variables (Ch. 4). We now move into formal inference for these relationships starting with relationships between two variables and moving on to more.

## 9.1 Simple linear regression

Recall, in Ch. 4, we wanted an equation to describe how a dependent (response) variable, $y$, changes in response to a change in one or more independent (experimental) variable(s), $x$.

We used the notation

$$y = \beta_0 + \beta_1 x + \epsilon$$

where $\beta_0$ is the intercept.

$\beta_1$ is the slope.

$\epsilon$ is some error. In fact,

**Goal:** We want to use inference to get interval estimates for our slope and predicted values and significance tests that the slope is not equal to zero.

### 9.1.1 Variance estimation

What are the parameters in our model, and how do we estimate them?

We need an estimate for $\sigma^2$ in a regression, or "line-fitting" context.

**Definition 9.1.** For a set of data pairs $(x_1, y_1), \ldots, (x_n, y_n)$ where least squares fitting of a line produces fitted values $\hat{y}_i = b_0 + b_1 x_i$ and residuals $e_i = y_i - \hat{y}_i$,

$$s_{LF}^2 = \frac{1}{n-2} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \frac{1}{n-2} \sum_{i=1}^{n} e_i^2$$

is the *line-fitting sample variance*. Associated with it are $\nu = n - 2$ degrees of freedom and an estimated standard deviation of response $s_{LF} = \sqrt{s_{LF}^2}$.

$s_{LF}^2$ estimates the level of basic background variation $\sigma^2$, whenever the model is an adequate description of the data.

### 9.1.2 Inference for parameters

We are often interested in testing if $\beta_1 = 0$. This tests whether or not there is a *significant linear relationship* between $x$ and $y$. We can do this using

1.

2.

Both of these require

It can be shown that since $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ and $\epsilon_i \overset{iid}{\sim} N(0, \sigma^2)$, then

$$b_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum(x - \bar{x})^2}\right)$$

So, a $(1 - \alpha)100\%$ CI for $\beta_1$ is

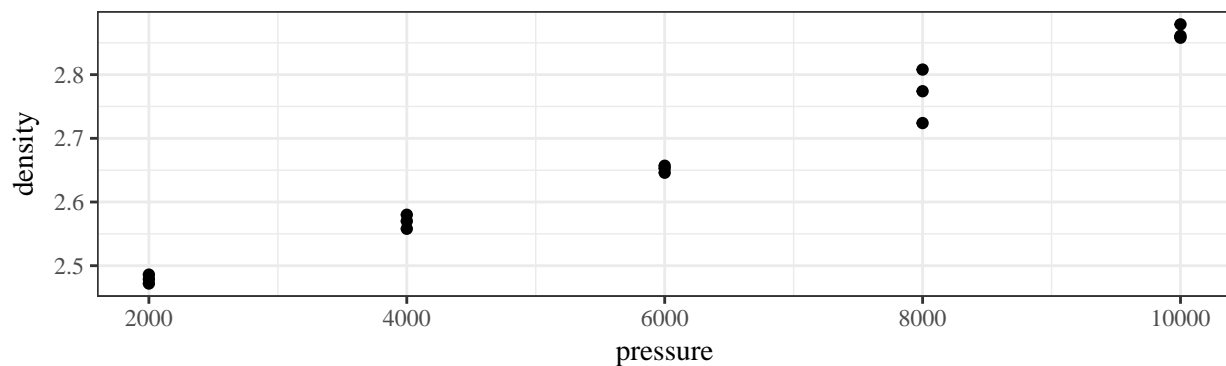and the test statistic for $H_0 : \beta_1 = \#$ is

**Example 9.1** (Ceramic powder pressing). A mixture of $Al_2O_3$, polyvinyl alcohol, and water was prepared, dried overnight, crushed, and sieved to obtain 100 mesh size grains. These were pressed into cylinders at pressures from 2,000 psi to 10,000 psi, and cylinder densities were calculated. Consider a pressure/density study of $n = 15$ data pairs representing

$$x = \text{ the pressure setting used (psi)}$$
$$y = \text{ the density obtained (g/cc)}$$

in the dry pressing of a ceramic compound into cylinders.

| pressure | density p | ressure d | ensity |
|----------|-----------|-----------|--------|
| 2000 | 2.486 | 6000 | 2.653 |
| 2000 | 2.479 | 8000 | 2.724 |
| 2000 | 2.472 | 8000 | 2.774 |
| 4000 | 2.558 | 8000 | 2.808 |
| 4000 | 2.570 | 10000 | 2.861 |
| 4000 | 2.580 | 10000 | 2.879 |
| 6000 | 2.646 | 10000 | 2.858 |
| 6000 | 2.657 | | |



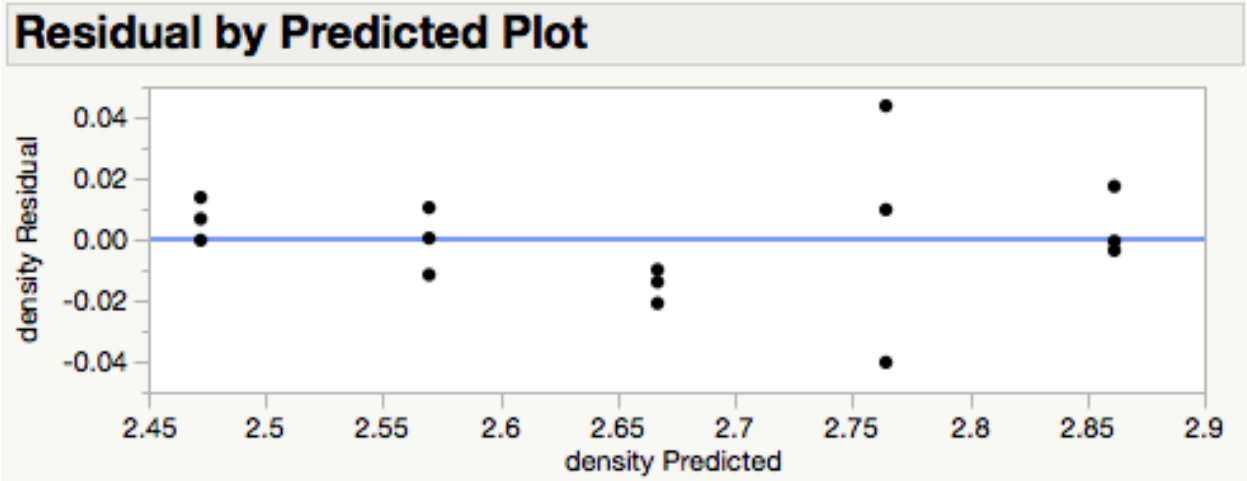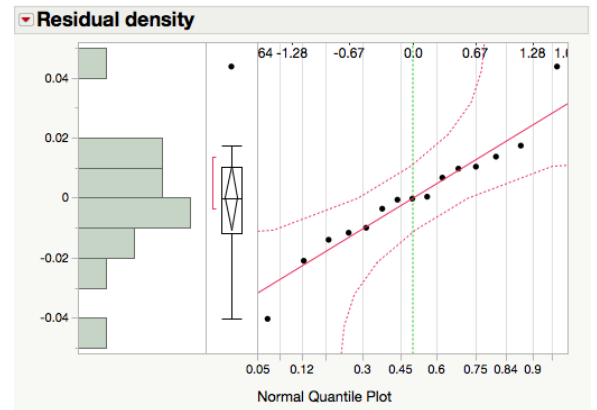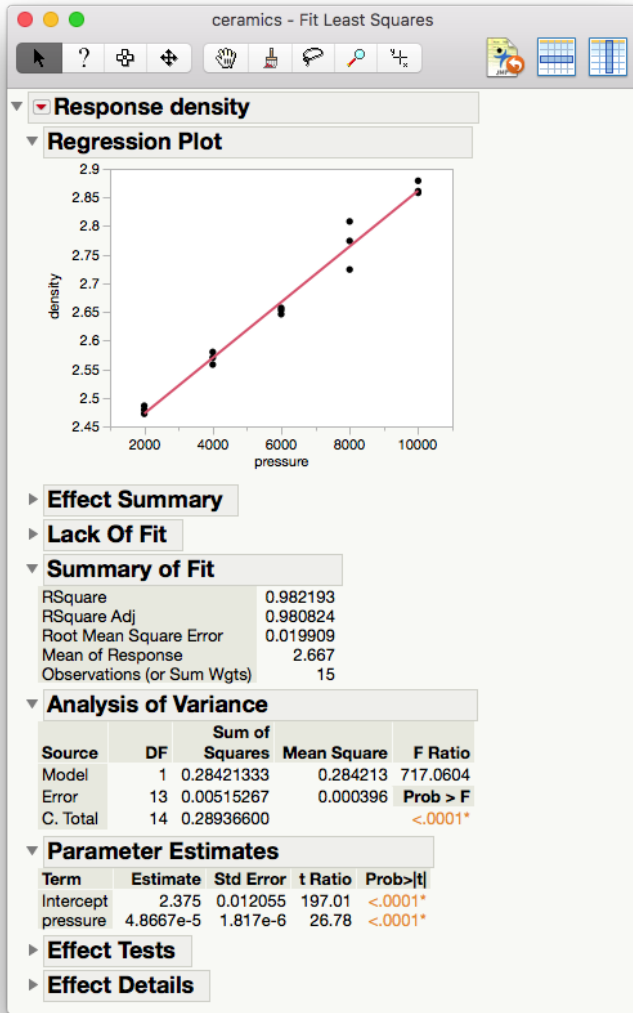A line has been fit in JMP using the method of least squares.

4

Figure 1: Least squares regression of density on pressure of ceramic cylinders.

1. Write out the model with the appropriate estimates.

2. Are the assumptions for the model met?

3. What is the fraction of raw variation in $y$ accounted for by the fitted equation?

4. What is the correlation between $x$ and $y$?

5. Estimate $\sigma^2$.

6. Estimate $\text{Var}(b_1)$.

7. Calculate and interpret the 95% CI for $\beta_1$

8. Conduct a formal hypothesis test at the $\alpha = .05$ significance level to determine if the relationship between density and pressure is significant.

### 9.1.3  Inference for mean response

Recall our model

$$y_1 = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \overset{\text{iid}}{\sim} N(0, \sigma^2).$$

Under the model, the true mean response at some observed covariate value $x_i$ is

Now, if some new covariate value $x$ is within the range of the $x_i$'s, we can estimate the true mean response at this new $x$

But how good is the estimate?

Under the model,

So we can construct a $N(0, 1)$ random variable by standardizing.

And when $\sigma$ is unknown (i.e. basically always),

To test $H_0 : \mu_{y|x} = \#$, we can use the test statistics

$$K =$$

which has a $t_{n-2}$ distribution if $H_0$ is true and the model is correct.

A 2-sided $(1-\alpha)100\%$ CI for $\mu_{y|x}$ is

**Example 9.2** (Ceramic powder pressing). Return to the ceramic density problem. We will make a 2-sided 95% confidence interval for the true mean density of ceramics at 4000 psi and interpret it.

Now calculate and interpret a 2-sided 95% confidence interval for the true mean density at 5000 psi.

## 9.2 Multiple regression

Recall the summarization the effects of several different quantitative variables $x_1, \ldots, x_{p-1}$ on a response $y$.

$$y_i \approx \beta_0 + \beta_1 x_{1i} + \cdots \beta_{p-1} x_{p-1,i}$$

Where we estimate $\beta_0, \ldots, \beta_{p-1}$ using the *least squares principle* by minimizing the function

$$S(b_0, \ldots, b_{p-1}) = \sum_{i=1}^{n}(y_i - \hat{y})^2 = \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_{1,i} - \cdots - \beta_{p-1} x_{p-1,i})^2$$

to find the estimates $b_0, \ldots, b_{p-1}$.

We can formalize this now as

$$Y_i = \beta_0 + \beta_1 x_{1i} + \cdots \beta_{p-1} x_{p-1,i} + \epsilon_i$$

where we assume $\epsilon_i \overset{\text{iid}}{\sim} N(0, \sigma^2)$.

### 9.2.1 Variance estimation

Based on our multiple regression model, the residuals are of the form

$$e_i = y_i - \hat{y}_i$$

And we can estimate the variance similarly to the SLR case.

**Definition 9.2.** For a set of $n$ data vectors $(x_{11}, x_{21}, \ldots, x_{p-11}, y), \ldots, (x_{1n}, x_{2n}, \ldots, x_{p-1n}, y)$ where least squares fitting is used to fit a surface,

$$s_{SF}^2 = \frac{1}{n-p} \sum (y - \hat{y})^2 = \frac{1}{n-p} \sum e_i^2$$

is the *surface-fitting sample variance.* Associated with it are $\nu = n - p$ degrees of freedom and an estimated standard deviation of response $s_{SF} = \sqrt{s_{SF}^2}$.

**Note:** the SLR fitting sample variance $s_{LF}^2$ is the special case of $s_{SF}^2$ for $p = 2$.

**Example 9.3** (Stack loss). Consider a chemical plant that makes nitric acid from ammonia. We want to predict stack loss ($y$, 10 times the % of ammonia lost) using

- $x_1$: air flow into the plant

- $x_2$: inlet temperature of the cooling water

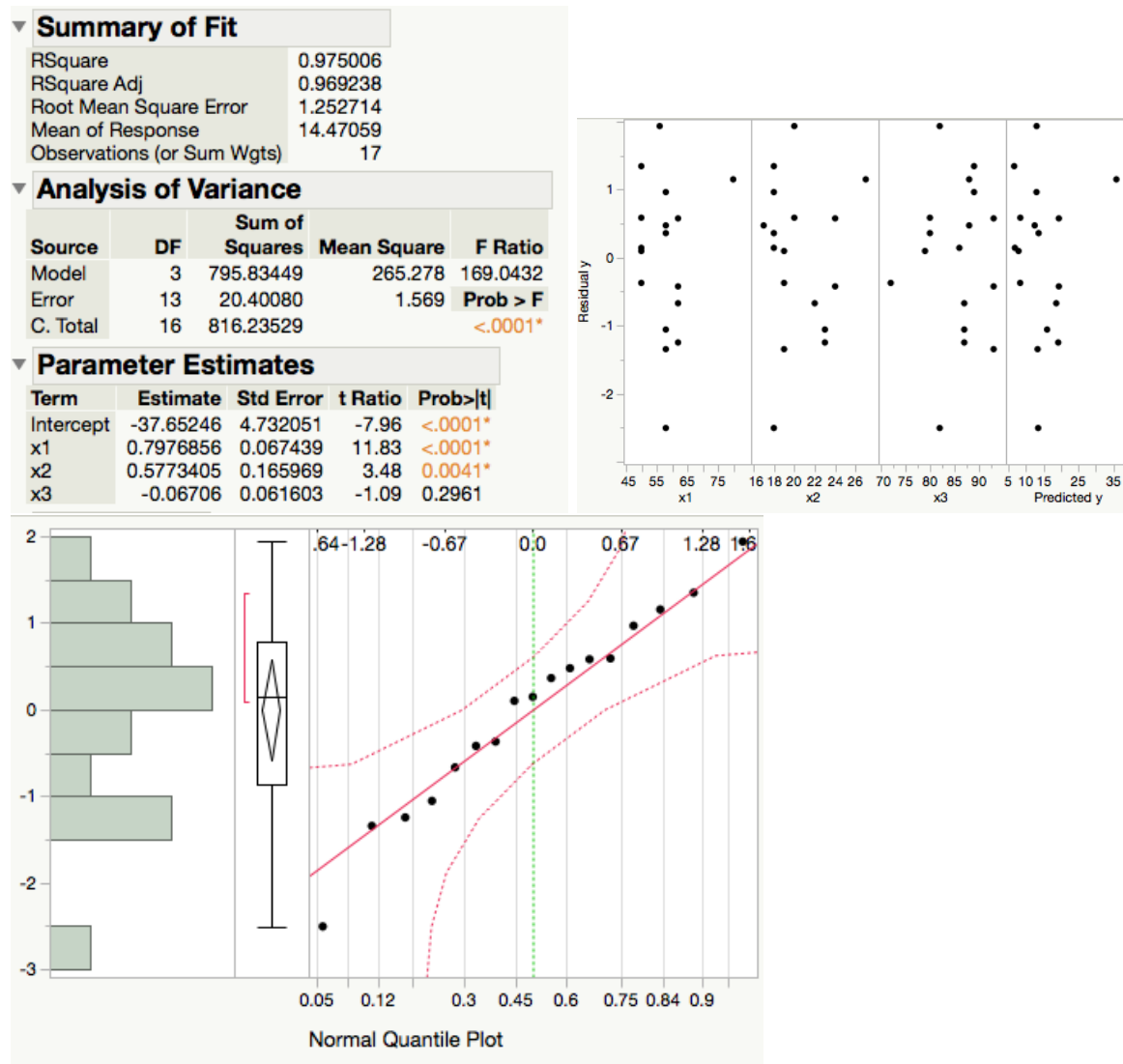- $x_3$: modified acid concentration (% circulating acid -50% ) $\times$ 10

**Summary of Fit**

| | |
|---|---|
| RSquare | 0.975006 |
| RSquare Adj | 0.969238 |
| Root Mean Square Error | 1.252714 |
| Mean of Response | 14.47059 |
| Observations (or Sum Wgts) | 17 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Ratio |
|---|---|---|---|---|
| Model | 3 | 795.83449 | 265.278 | 169.0432 |
| Error | 13 | 20.40080 | 1.569 | Prob > F |
| C. Total | 16 | 816.23529 | | <.0001* |

**Parameter Estimates**

| Term | Estimate | Std Error | t Ratio | Prob>|t| |
|---|---|---|---|---|
| Intercept | -37.65246 | 4.732051 | -7.96 | <.0001* |
| x1 | 0.7976856 | 0.067439 | 11.83 | <.0001* |
| x2 | 0.5773405 | 0.165969 | 3.48 | 0.0041* |
| x3 | -0.06706 | 0.061603 | -1.09 | 0.2961 |

Figure 2: Least squares regression of stack loss on air flow, inlet temperature, and modified acid concentration.

### 9.2.2 Inference for parameters

We are often interested in answering questions (doing formal inference) for $\beta_0, \ldots, \beta_{p-1}$ individually. For example, we may want to know if there is a significant relationship between $y$ and $x_2$ (holding all else constant).

Under our model assumptions,
$$b_i \sim N(\beta_i, d_i \sigma^2)$$
for some positive constant $d_i, i = 0, 1, \ldots, p - 1$.

That means

So, a test statistic for $H_0 : \beta_i = \#$ is

and a 2-sided $(1 - \alpha)100\%$ CI for $\beta_i$ is

**Example 9.4** (Stack loss, cont'd). Using the model fit on page 15, answer the following questions:

1. Is the average change in stack loss $(y)$ for a one unit change in air flow into the plant $(x_1)$ less than 1 (holding all else constant)? Use a significance testing framework with $\alpha = .1$.

2. Is the there a significant relationship between stack loss $(y)$ and modified acid concentation $(x_3)$ (holding all else constant)? Use a significance testing framework with $\alpha = .05$.

3. Construct and interpret a 99% confidence interval for $\beta_3$.

4. Construct and interpret a 90% confidence interval for $\beta_2$.

### 9.2.3 Inference for mean response

We can also estimate the mean response at the set of covariate values, $(x_1, x_2, \ldots, x_{p-1})$. Under the model assumptions, the estimated mean response, $\mu_{y|x}$, at $\boldsymbol{x} = (x_1, x_2, \ldots, x_{p-1})$ is

with:

Then, under the model assumptions

And a test statistic for testing $H_0 : \mu_{y|\boldsymbol{x}} = \#$ is

A 2-sided $(1 - \alpha)100\%$ CI for $\mu_{y|\boldsymbol{x}}$ is

**Example 9.5** (Stack loss, cont'd)**.** We can use JMP to compute a 2-sided 95% CI around the mean response at point 3:

$$x_1 = 62, x_2 = 23, x_3 = 87, y = 18$$

Figure 3: How to get predicted values and standard errors.

Figure 4: Predicted values and standard errors.