

In an observational study, we would observe the process, with as little interaction or disturbance as possible, in order to obtain the data.

- With adequate planning, an observational study can yield accurate, complete, reliable data
- These studies can lead to ideas on what might be impacting the process
- However, these studies often provide limited information about specific relationships of interest, such as the impact of a variable that is tightly controlled in normal operation

Notes

*LSSV2 student files \ ANOVA linear fit
Worksheet \ Prediction & error 1*

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1																
2																
3																
4																
5																
6																
7																
8																
9																
10																
11																
12																
13																
14																
15																
16																
17																

	X data	Y data	Prediction	Error
	8	6.16	27.90	-21.74
	22	9.88	27.90	-18.02
	35	14.35	27.90	-13.55
	40	24.06	27.90	-3.84
	57	30.34	27.90	2.44
	73	32.17	27.90	4.27
	78	42.18	27.90	14.28
	87	43.23	27.90	15.33
	98	48.76	27.90	20.86
	Sum of squares (SS)	8901.3	= 7007.4	+ 1893.9
	Degrees of freedom (DF)	9	= 1	+ 8
	Root mean square error (RMSE)			15.39
	Average Y	27.90		
	STDEV of Y	15.39		

$$Y = 27.9033 + 0.0000 X$$

Notes

Worksheet \ Prediction & error 2

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1																
2					X data	Y data	Prediction	Error		Y =	0.8387	+	0.4891	X		
3					8	6.16	4.75	1.41								
4					22	9.88	11.60	-1.72								
5					35	14.35	17.96	-3.61								
6					40	24.06	20.40	3.66								
7					57	30.34	= 28.72	+ 1.62								
8					73	32.17	36.54	-4.37								
9					78	42.18	38.99	3.19								
10					87	43.23	43.39	-0.16								
11					98	48.76	48.77	-0.01								
12					Sum of squares (SS)	8901.3	= 8838.0	+ 63.3								
13					Degrees of freedom (DF)	9	= 2	+ 7								
14					Root mean square error (RMSE)				3.007							
15																
16					Average Y	27.90										
17					STDEV of Y	15.39										
18					Adjusted R square	0.962										
19																

Proportion of total Y variation caused by ("explained by") X variation

Notes

1. Run Analyze > Fit Model in JMP to investigate the relationship between y and x
2. Check the p-value for the fit to determine whether the regression is significant. If not, then no need to go further.
3. If the regression is significant, determine the strength of the relationship, using the *Adjusted R²*
4. Check model adequacy by reviewing the residuals plots
 - Residual Normal Quantile Plot
 - Residual by Predicted Plot
 - Studentized Residuals (in run order)

We'll go through these steps and additional analysis details, for simple regression in the following example.

Notes

Simple Regression in JMP

Open: Data sets \ simple regression - generic

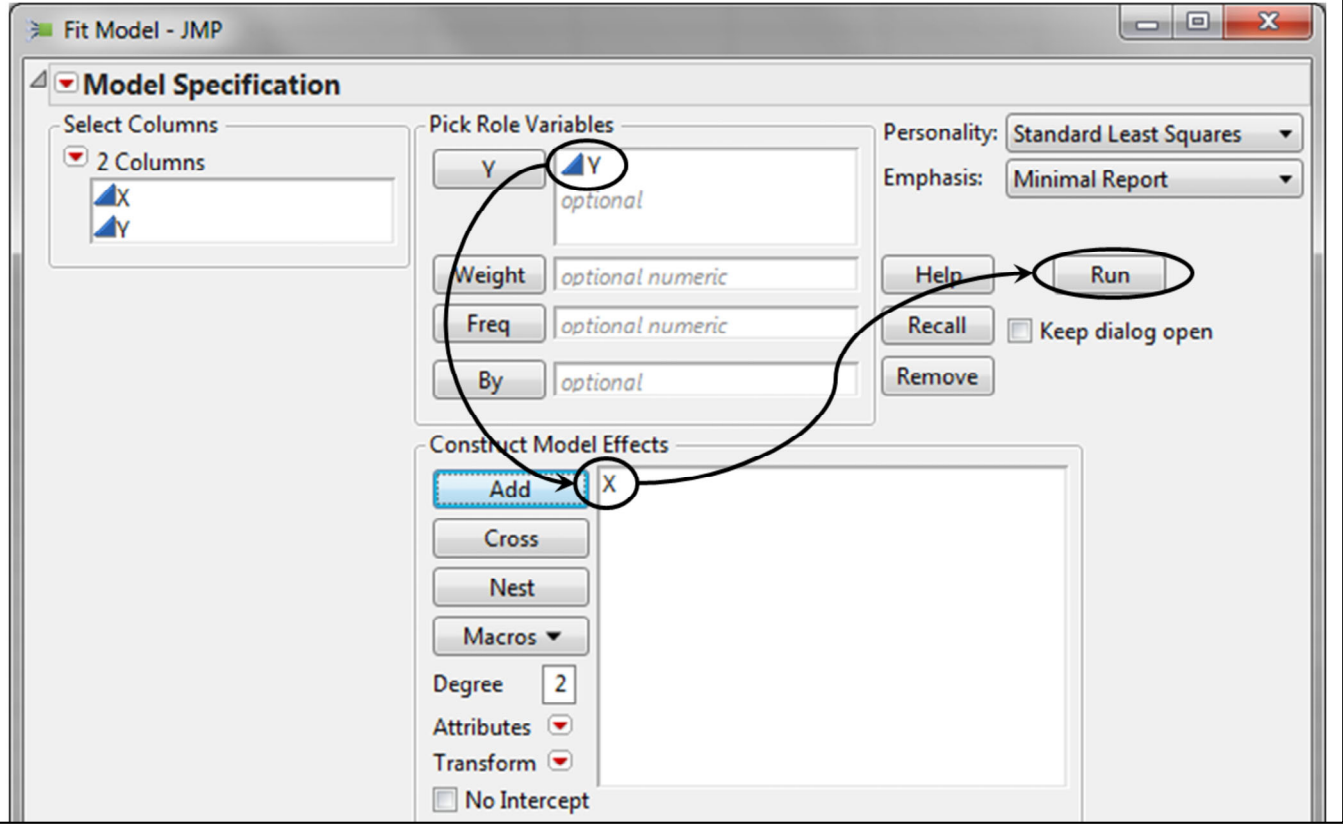
The screenshot shows the JMP software window titled 'simple regression - generic - JMP'. The interface includes a menu bar with options like File, Edit, Tables, Rows, Cols, DOE, Analyze, Graph, Tools, View, Window, and Help. Below the menu bar is a toolbar and a list of columns (X and Y) and rows (All rows, Selected, Excluded, Hidden, Labelled). The main data table contains the following values:

	X	Y
1	8	6.16
2	22	9.88
3	35	14.35
4	40	24.06
5	57	30.34
6	73	32.17
7	78	42.18
8	87	43.23
9	98	48.76

The status bar at the bottom indicates 'evaluations done'.

Notes

Analyze → Fit Model → Set up as shown → Run

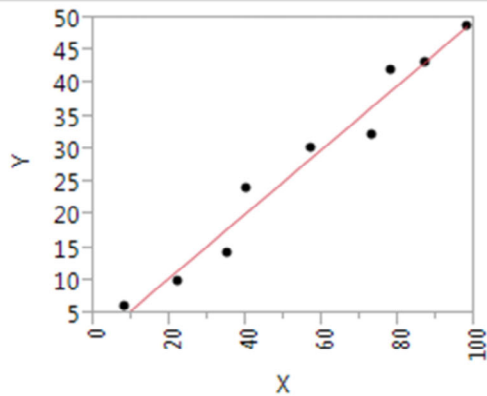


Notes

Analysis details

Response Y

Regression Plot



Summary of Fit

RSquare	0.966581
RSquare Adj	0.961807
Root Mean Square Error	3.006984
Mean of Response	27.90333
Observations (or Sum Wgts)	9

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio	Prob > F
Model	1	1830.6557	1830.66	202.4624	
Error	7	63.2937	9.04		
C. Total	8	1893.9494			

- The **Root Mean Square Error (RMSE)** is the standard deviation of Y caused by factors other than X
- It can be thought of as the **standard deviation** about the fitted line (or model)
- Also known as the “error” or “residual” standard deviation
- Smaller is better

- **P-value** indicates whether the regression is significant
- This low p-value shows that it is significant

Notes

Summary of Fit

RSquare	0.966581
RSquare Adj	0.961807
Root Mean Square Error	3.006984
Mean of Response	27.90333
Observations (or Sum Wgts)	9

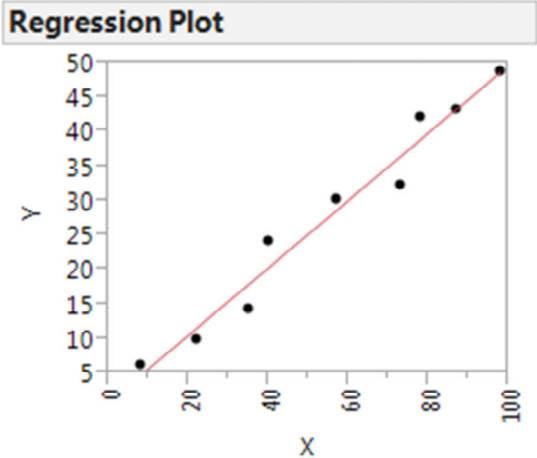
R²
"Coefficient of Determination"

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	1830.6557	1830.66	202.4624
Error	7	63.2937	9.04	Prob > F
C. Total	8	1893.9494		<.0001*

- Proportion of the variation in Y that is "explained by" variation in X.
- Varies from 0 to 1.
- Larger is better
- Unitless

Notes



Summary of Fit	
RSquare	0.966581
RSquare Adj	0.961807
Root Mean Square Error	3.006984
Mean of Response	27.90333
Observations (or Sum Wgts)	9

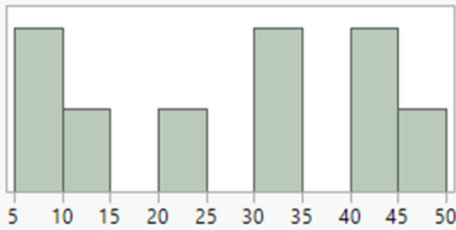
Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	1830.6557	1830.66	202.4624
Error	7	63.2937	9.04	Prob > F
C. Total	8	1893.9494		<.0001*

- **Adjusted R²** also gives us the proportion of Y variation explained by the model (a line in simple regression)
- Varies from 0 to 1
- Larger is better
- **Always use the Adjusted R² value, not R²**
- Adjusted R² takes the number of model terms into account and penalizes for including insignificant terms
- In this example, the simple regression model explains much of the variation in Y.

Notes

Distributions

Y



Summary Statistics

Mean	27.903333
Std Dev	15.386477
N	9
Minimum	6.16
Maximum	48.76
Median	30.34

Standard Deviation (STDEV) of the data set

$$R^2 = 1 - \frac{SS_{Error}}{SS_{Total}}$$

$$R^2_{Adj} = 1 - \frac{SS_{Error}/(n - p)}{SS_{Total}/(n - 1)} = 1 - \left(\frac{RMSE}{STDEV} \right)^2$$

p = number of terms in the model (including the intercept)

n = sample size (number of measurements in the data set)

SS_{Total} is the sum of squares of the data (measurements in the data set)

SS_{Error} is the sum of squares of the Errors or residuals

We saw the sum of squares calculations earlier, in the ANOVA

Notes

There is a potential problem with R^2 :

- R^2 always increases when terms are added to a model, even when the terms are not significant
- **This is particularly a problem in multiple regression**, as it can lead to “overfitting,” giving false confidence in using the model, especially for prediction.
- Adjusted R^2 corrects for this by considering the number of terms in the model
- Adjusted R^2 can actually decrease if non-significant terms are added to a model

Adjusted R^2 is the recommended statistic for determining the proportion of variation in Y explained by the model

Notes

Red triangle next to *Response Y* → *Regression Reports* → *Parameter Estimates*

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	1830.6557	1830.66	202.4624
Error	7	63.2937	9.04	Prob > F
C. Total	8	1893.9494		<.0001*

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	0.8386661	2.150023	0.39	0.7081
X	0.4891205	0.034375	14.23	<.0001*

- In regression of Y on a single X, the Analysis of Variance P-value is the same as the P-value for the slope of the line.
- The P-value for the slope of the line indicates the evidence of a correlation between Y and X.
- Significance of individual model terms are determined by testing whether their regression coefficient is equal to 0, using the t statistic. Hypotheses are:

$$H_0: b_i = 0$$

$$H_1: b_i \neq 0$$

- This is a test of the contribution of the model term, given the other terms in the model.

Notes

Estimates and P-values for the slope and intercept

Parameter Estimates				
Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	0.8386661	2.150023	0.39	0.7081
X	0.4891205	0.034375	14.23	<.0001*

Model: $Y = 0.84 + 0.50X + error$

- In this example, the P-value for the slope of the line indicates very strong evidence of a correlation between Y and X.
- The P-value for the Intercept indicates that it is not significant.
 - **Best practice is to leave the Intercept in the model, whether or not the P-value indicates that it is significant**
 - Regression equations are developed, and are only valid, over the region of the regressor variables (x's) contained in the data set
 - Forcing the model to pass through (0, 0) by removing the intercept, can create problems in the region being modeled

Notes

Both the Adjusted R^2 and the p-values must be considered, in order to understand what has been learned in the analysis.

When the resulting model has:

- **High Adjusted R^2 and significant model term p-values**, this is ideal. Factors driving the response have been identified and the variation is largely explained. A decent model has been created.
- **Low Adjusted R^2 and significant model term p-values**, more work must be done. Some significant factors influencing the response have been identified, but the low Adjusted R^2 indicates that other important factors exist. These need to be found, for the model to be useful.
- **High R^2 and insignificant model terms**, this is usually due to the data violating the assumptions of the regression analysis. There is more information on this scenario in upcoming slides.
- **Low Adjusted R^2 and insignificant model terms**, no relationship between X and Y variables have been found. Usually this means that new ideas about which factors influence Y must be developed, although it can occasionally be due to missing higher order terms.

Notes

This page intentionally left blank

2 Checking Model Adequacy

In least squares fit regression (continuous Y), the analysis methods used to calculate regressor coefficients and their p-values, depend on certain assumptions being met.

Assumptions:

- Errors (residuals) are normally and independently distributed with mean zero and constant variance (σ^2)
- Observations are adequately described by the model

Whether performing regression from “file cabinet” data or analyzing the results of a designed experiment, these assumptions must be validated.

Notes

**To validate that these assumptions have been met,
the *residuals* are examined:**

1. Normal Probability Plot of Residuals

- Validate that the residuals are normally distributed
- In JMP, this is the *Residual Normal Quantile Plot*

2. Residuals vs. Predicted (or Fitted) Values

- Validate constant variance and mean 0
- In JMP, this is the *Residual by Predicted Plot*

3. Residuals vs. Run Order

- Verify independence of errors
- There should be no patterns over the timeframe of the data
- In JMP, the best graph to use is *Studentized Residuals*
- The JMP data table must be in run order for *Studentized Residuals* to graph the residuals in run order

Notes

A fitted model, the equation generated during regression, gives the predicted mean value of the response variable as a function of the predictor variables. These predicted mean values are also called *predicted values*, or just *predicted* for short. The *residual value* is the data (observation) value minus the predicted value. Residual values are called *residuals* for short.

These terms are easiest to visualize in the simple linear model shown above. A predicted value is the fitted line evaluated at some X value. A residual is the difference between a measured (observed) Y value and the predicted value at the corresponding X.

Residuals contain information about the magnitude and direction of variability in the data relative to the fitted model.

- An unusually large residual might signal a measurement error, data entry error or some other type of outlier.
- A systematic trend or pattern in the residuals might signal an inadequacy in the fitted model.

Notes

If residuals are normally distributed, the plot will be approximately a straight line.

Emphasis should be on the central values of the plot, rather than the ends

It is common for plots to bend upward at the high end and downward at the low end.

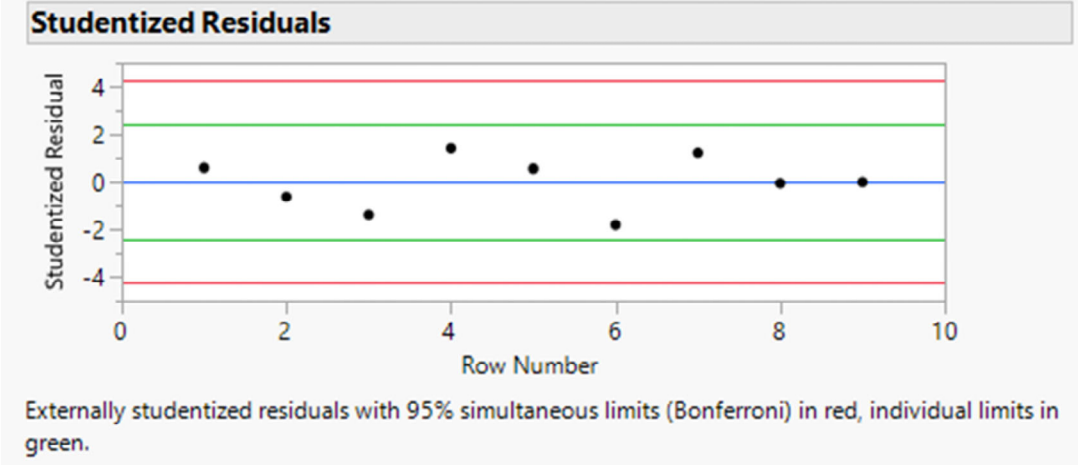
Small sample sizes, such as from experiments, often appear more non-normal

Use the “Fat Pencil” Rule: If a “fat pencil” placed over the central points would cover them on the plot, then the residuals are approximately normal (good enough). Hyperbolic bands displayed in JMP plots give these bounds.

A curve throughout the plot is a strong indication of non-normality. In this case, a transformation would be needed.

The plot above shows an error (residuals) distribution that is approximately normal, so it is not concerning.

Notes



In viewing the Studentized Residuals for the *simple regression-generic*, the best form for checking residuals by run order, we can see whether there are any patterns over the timeframe of the data.

Note that the data table must be in run order for this plot.

Notes

Again, on this graph, healthy residuals look like a random scatter around 0.

Runs (points in a row) of positive-negative-positive-negative residuals indicate correlation between runs. This implies that the assumption of independence has been violated. **In designed experiments, randomization protects against this! Do it every time!**

This plot can also show a change in variance over the time span of the experiment. This could be due to increased skill as the experiment progresses, a process drift, operator fatigue, tool wear, etc. This type of problem would show as an increase or decrease in spread or “scatter” of the residuals across the graph. Increasing or decreasing variance indicates the need for a transformation.

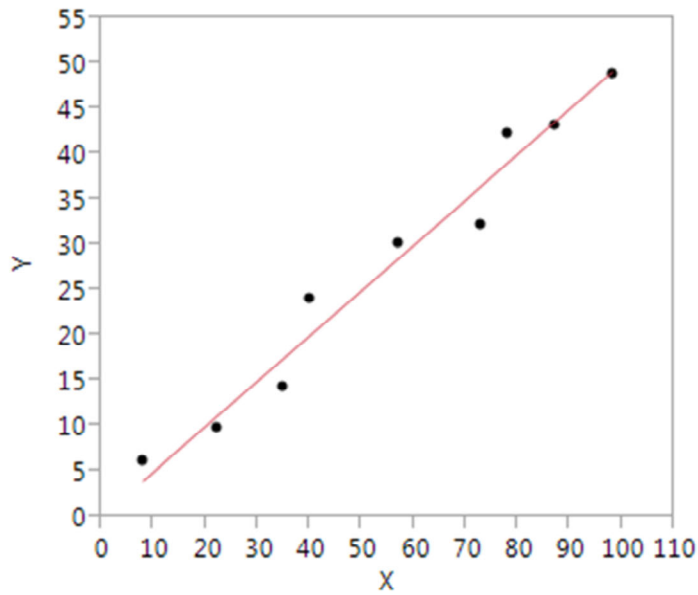
Notes

In this section, we'll see how we can:

- Use the Root Mean Square Error (RMSE) in predicting our future process variation,
- Use JMP's Prediction Profiler to help us optimize our process, and
- Estimate our future % defective, using the t distribution calculator.

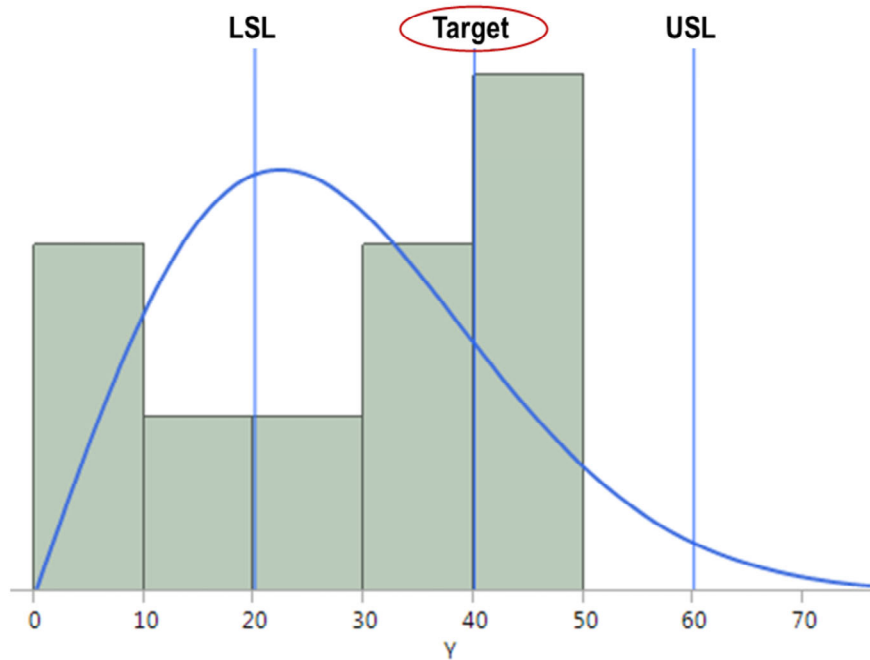
Notes

When Y is correlated with a controllable X variable,



how can we use the regression to improve the Y capability?

Notes



Suppose we are not happy with our current process capability

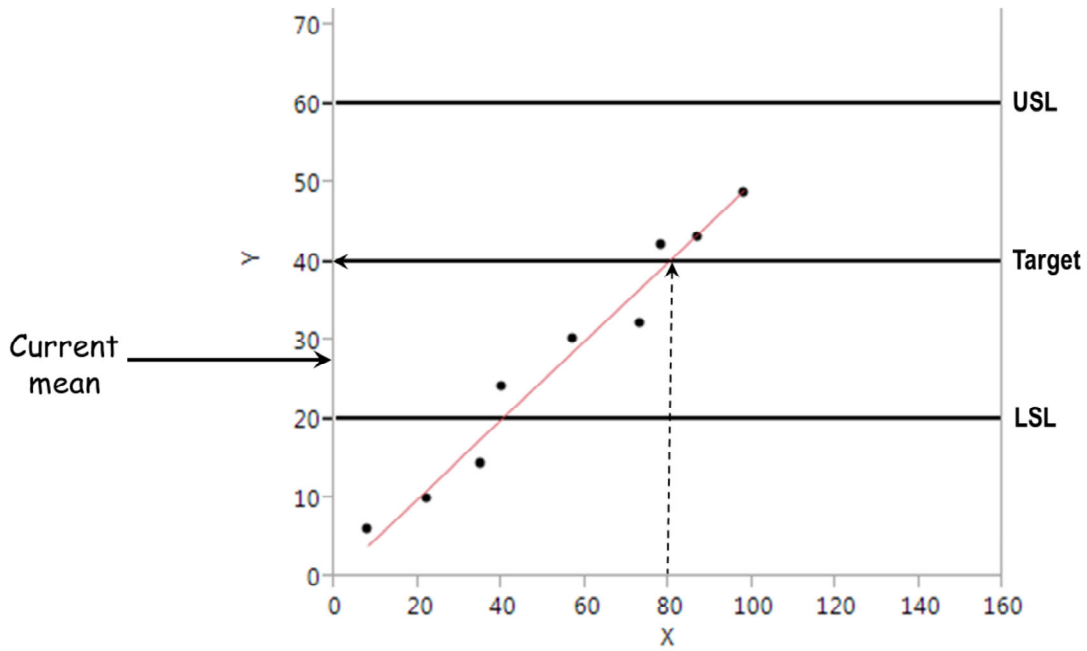
Mean = 27.9, Std dev = 15.4

Defective in the data: 33.3%

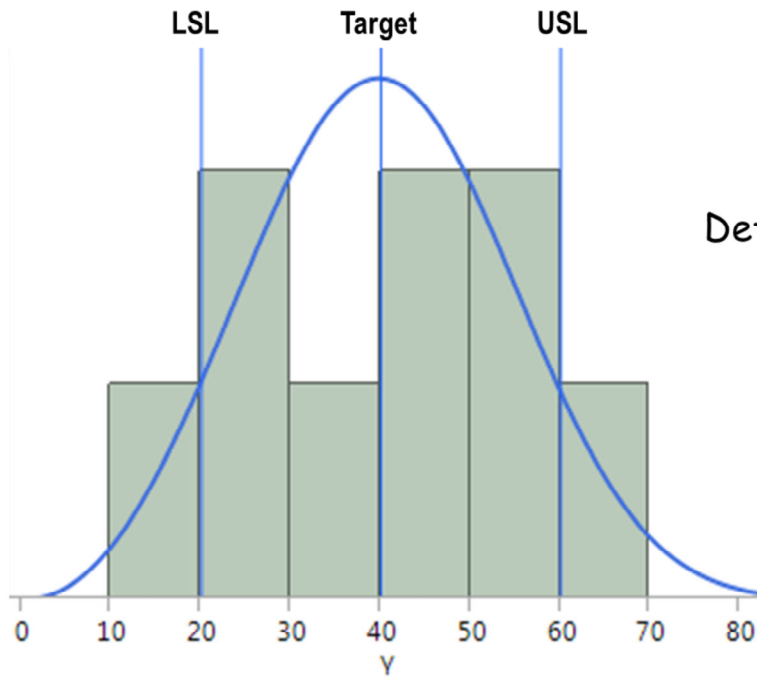
Predicted from distribution curve: 35.8%

Notes

If we control X at 80, the mean will change from 27.9 to 40



Notes

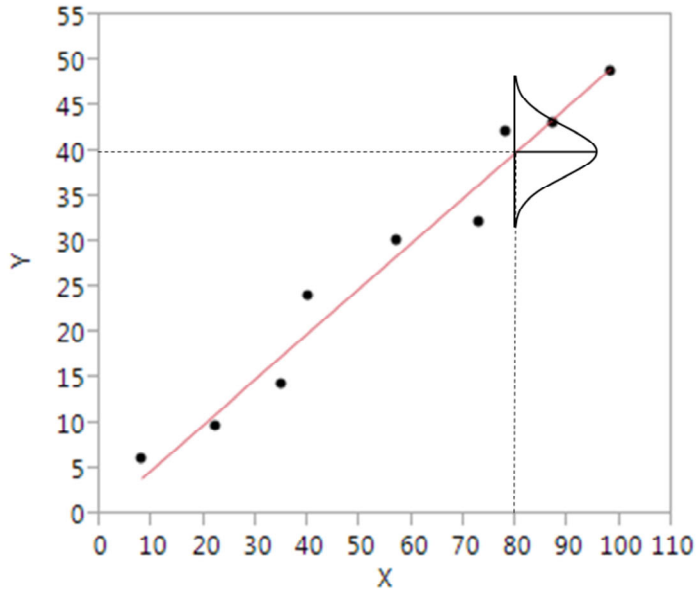


Mean = 40.0
Std dev = 15.4
Defective in the data: 22.2%
Distribution curve: 15.9%

- Moving mean Y to the center of the spec range does reduce % defective
- Is the mean the only thing that changes when we control X at 80?

Notes

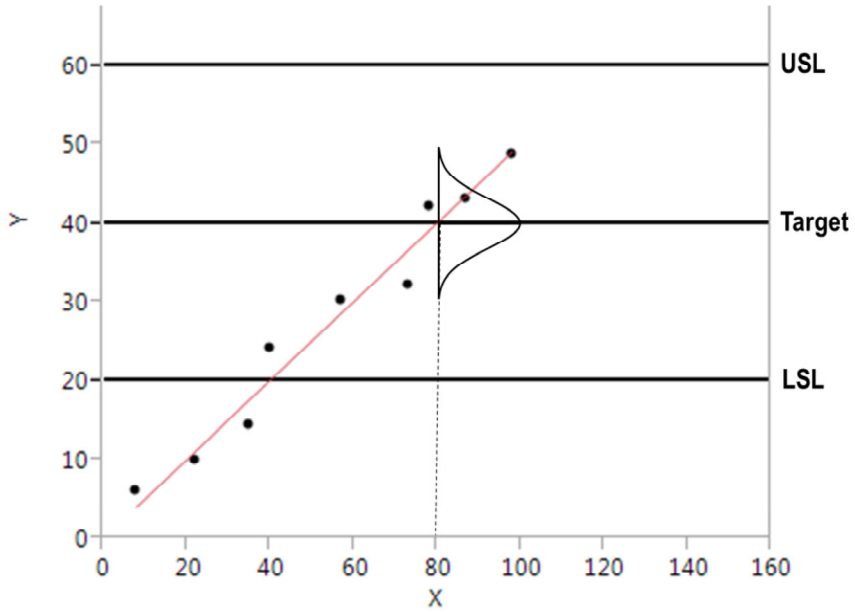
By definition, RMSE is the standard deviation of Y that would result from eliminating the variation in X



$$\sigma = \text{RMSE} = 2.84$$

Notes

When we control X at 80, we don't just move the mean from 27.9 to 40
 — we also reduce the standard deviation from 15.4 to 2.84!

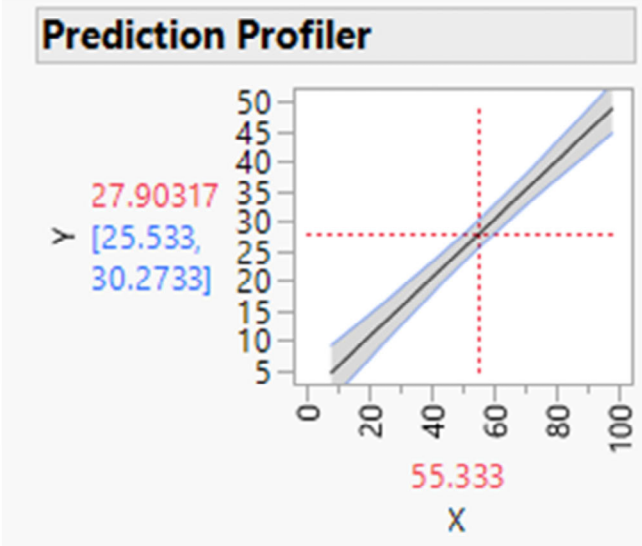


Notes

This page intentionally left blank

4. Introduction to the Prediction Profiler

JMP's Prediction Profiler helps us use our regression model to make predictions and optimize our process.

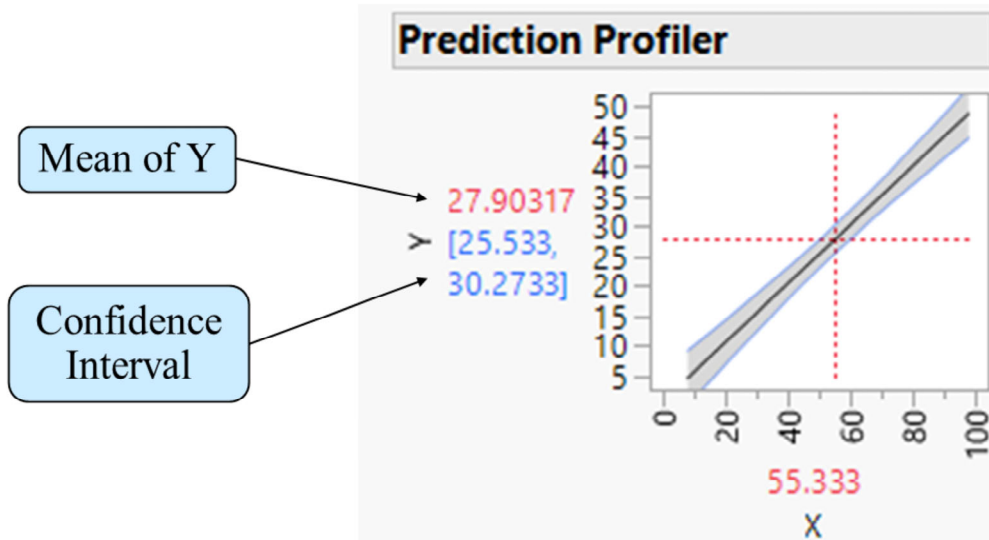


Follow these steps to access the prediction profiler:

- Analyze > Fit Model > Y = Y, Model Effects = X > Run > Red Triangle > Factor Profiling > Profiler

Notes

JMP's Prediction Profiler helps us use our regression model to make predictions and optimize our process.

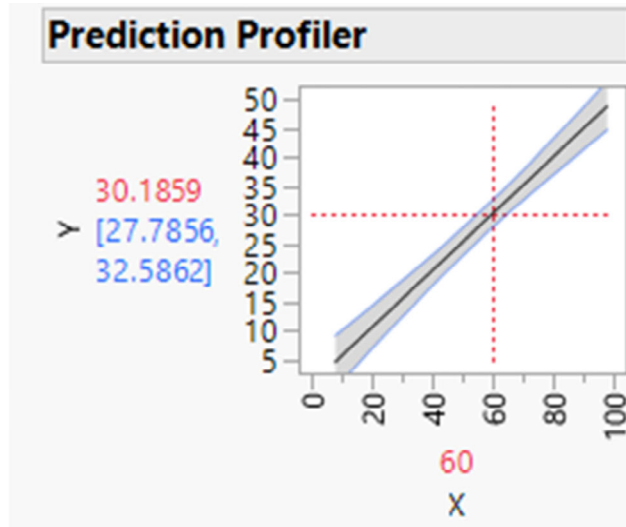


- Calculates predicted *mean Y* as a function of X
- Calculates **confidence intervals** for predicted **means**

Notes

Continuing with the *simple regression-generic* data:

- Suppose we are interested in the predicted mean Y for $X = 60$
- Click on the 55.333, change it to 60

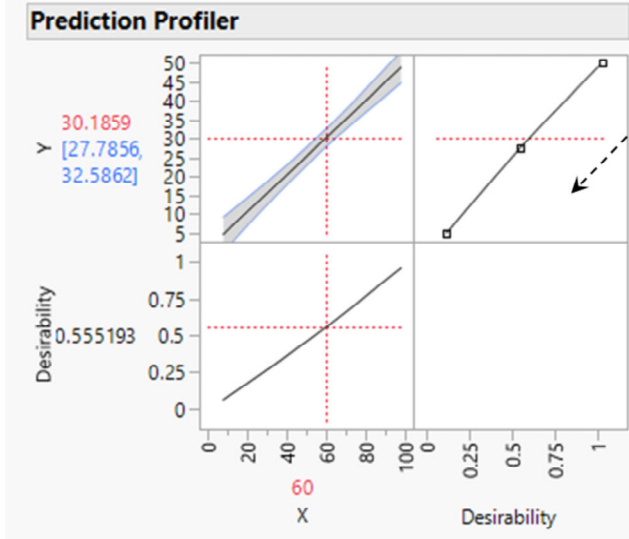


- Predicted mean Y (based on the data) is 30.19
- With 95% confidence, the population mean lies between 27.79 and 32.59

Notes

Simple example of optimization

- Suppose we want to find the X value that predicts a mean Y value of 25
- Red triangle next to *Prediction Profiler* → *Optimization and Desirability* → *Desirability Functions*



- Double click in here (don't touch the line plot)
- Modify the **Response Goal** dialog as shown below
- Click OK

Y	Values	Desirability
High:	30	0.0183
Middle:	25	1
Low:	20	0.0183
Importance:	1	

Notes

- The **95% Confidence Interval on the Mean Response** gives the range which will contain the “true” mean, μ , 95% of the time

- For a sample, the confidence interval is calculated:

$$\bar{Y} - t_{.025, n-1} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{Y} + t_{.025, n-1} \frac{s}{\sqrt{n}}$$

- For a regression, calculation of the confidence interval is similarly structured, but considerably more complicated, involving matrix math.
- A **95% Prediction Interval** gives the range which will contain future individual response observations 95% of the time.
 - The prediction interval is wider than the confidence interval, because it is to contain individual measurements, not averages.
 - Calculation of this interval is complicated, involving matrix math.

Notes

- a) Continuing with *simple regression-generic*, find the X value that predicts a mean Y value of 35. Give the confidence limits for the predicted mean.

- b) The overall standard deviation of Y is 15.39. The RMSE from the regression is 2.84. Which of these would be the standard deviation of Y if we controlled X to a constant value?

- c) Save your script, close and save the data table.

Notes

Data sets \ production vs capacity.

- (a) Fit a regression for *Production qty* as a function of *Capacity utilized (%)* (using *Fit Model*, of course). Is there a correlation? Give the appropriate P-value and strength of evidence.

- (b) For this exercise, we will not review the residuals plots. Use your model to find the capacity utilization level that predicts a mean daily production quantity of 3500. Give the confidence limits.

- (c) The overall standard deviation of *Production qty* is 733.5 (not shown in Fit Model output—calculated in Distribution Platform). The RMSE from the analysis in (a) is 409.732. Which of these would be the standard deviation if capacity utilization was held constant?

- (d) Save your scripts, close and save the data table.

Notes

Once we determine the level at which we want to control our x , we can use the root mean square error (RMSE) and other regression results to estimate the % defective in the improved process.

Remember that by definition, the RMSE is the standard deviation of the improved process, with x 's held at desired levels.

The *t distribution calculator* helps us calculate the future % defective.

Notes

Data sets \ production vs capacity.jmp.

In this process data, on 75% of the days production quantity fell below 3000.

Based on the best fit distribution, the Lognormal, the expected % of days that production quantity will fall below 3000 is 71.8%.

- a) We found earlier that capacity utilization 52.1% gives a mean daily production quantity of 3500. The RMSE was 409.7, the error degrees of freedom was 34. Assuming 52.1% capacity utilization, use the *t distribution calculator* to find the predicted % of days on which production quantity will be less than 3000.

- b) Save your scripts, close and save the data table.

Notes

Open *Data sets \ outgassing process*. *Current* (the Y variable) is the current required to heat a filament to a target temperature. *Resist* (the X variable) is the electrical resistance of the filament. *Machine* is the processing unit. This example shows how to reduce % defective by separate optimization of each machine.

- a) For this process, the % of *Current* data values that fall outside the interval (1.9, 2.1) is 8.87%.
- b) Fit a regression for *Current* as a function of *Resist*, using *Machine* as the *By variable*. For each machine, give the RMSE, the error degrees of freedom, and the resistance that predicts a mean current of 2.

Machine	RMSE	DF	Resistance	% Outside
A				
B				
C				

- c) Assuming we use the indicated resistance values, use the *t distribution calculator* to find for each machine the % of *Current* values predicted to fall outside the interval (1.9, 2.1).
- d) Save your scripts, close and save the data table.

Notes

5 Multiple Regression

- Multiple regression model
- Examples
- Fitting regression models
- Interactive effects
- Predicted values and uncertainty
- Modeling and optimization

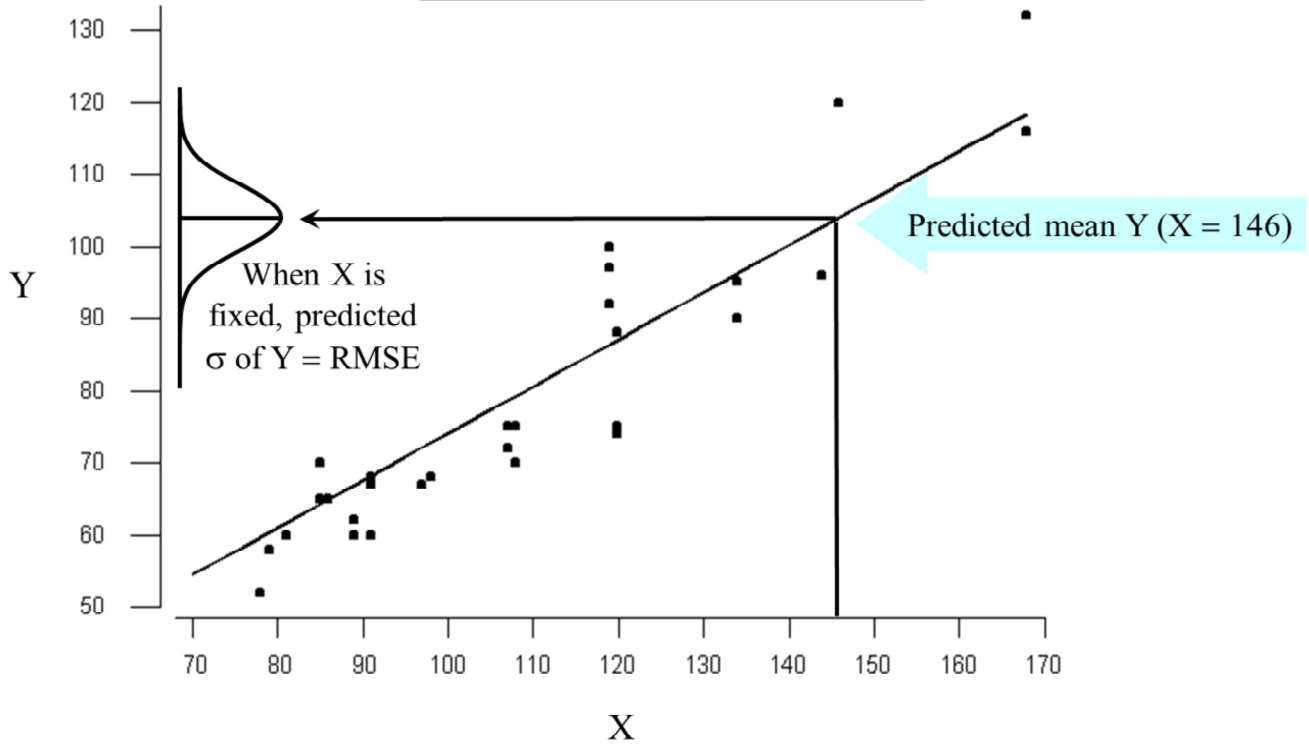
Notes

$$Y = b_0 + b_1X_1 + b_2X_2 + \dots + b_kX_k + \text{“error”}$$

Y	X_1, X_2, \dots, X_k	b_0	b_1, b_2, \dots, b_k	“Error”
Dependent variable	Independent variables	Intercept	Regression coefficients	Residuals
Response variable	Explanatory variables	Parameter	Parameters	Mean = 0
Output	Inputs			Standard deviation = σ (RMSE)
	Predictors			Distribution = Assumed to be Normal
	Regressors			
	Factors (in DOE)			

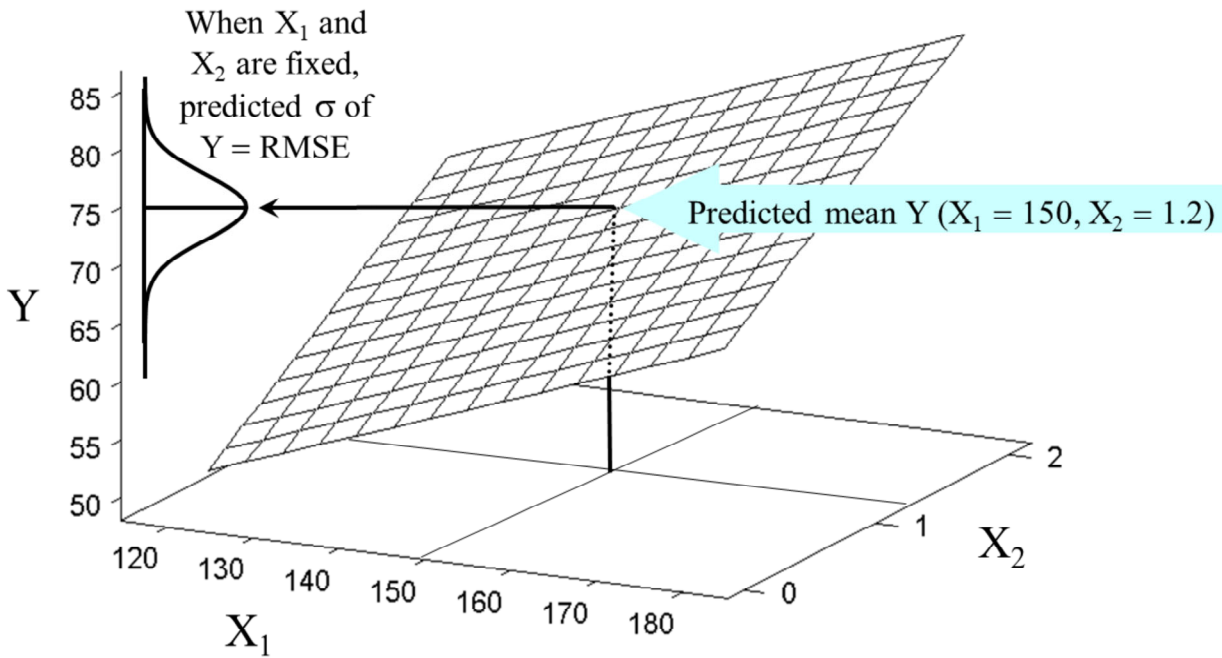
Notes

$Y = b_0 + b_1X + \text{“error”}$



Notes

$$Y = b_0 + b_1X_1 + b_2X_2 + \text{“error”}$$



Notes

Multiple regression examples

Y	X ₁	X ₂	X ₃	X ₄	X ₅
Life of cutting tool	RPM	Tool type	Material	Feed rate	
MPG	Displacement	Horsepower	Weight		
Salary	Education	Experience	Performance	Seniority	Gender
Vending machine service time	Amount of product stocked	Distance from truck to machine			

Fill in examples of interest to you

Notes

Regression model equations

Y	X ₁	X ₂	X ₃	X ₄	X ₅
<i>MPG</i>	Displacement (D)	Horsepower (H)	Weight (W)		

$$MPG = b_0 + b_1D + b_2H + b_3W + \text{error}$$

Y	X ₁	X ₂	X ₃	X ₄	X ₅
<i>Bond strength</i>	Temperature (T)	Dwell time (D)	T × D	T ²	D ²

$$\text{Bond} = b_0 + b_1T + b_2D + b_3TD + b_4T^2 + b_5D^2 + \text{error}$$



Response surface model (RSM) with two continuous Xs.

TD is the interaction term for T and D, T² and D² show curvature.

Notes

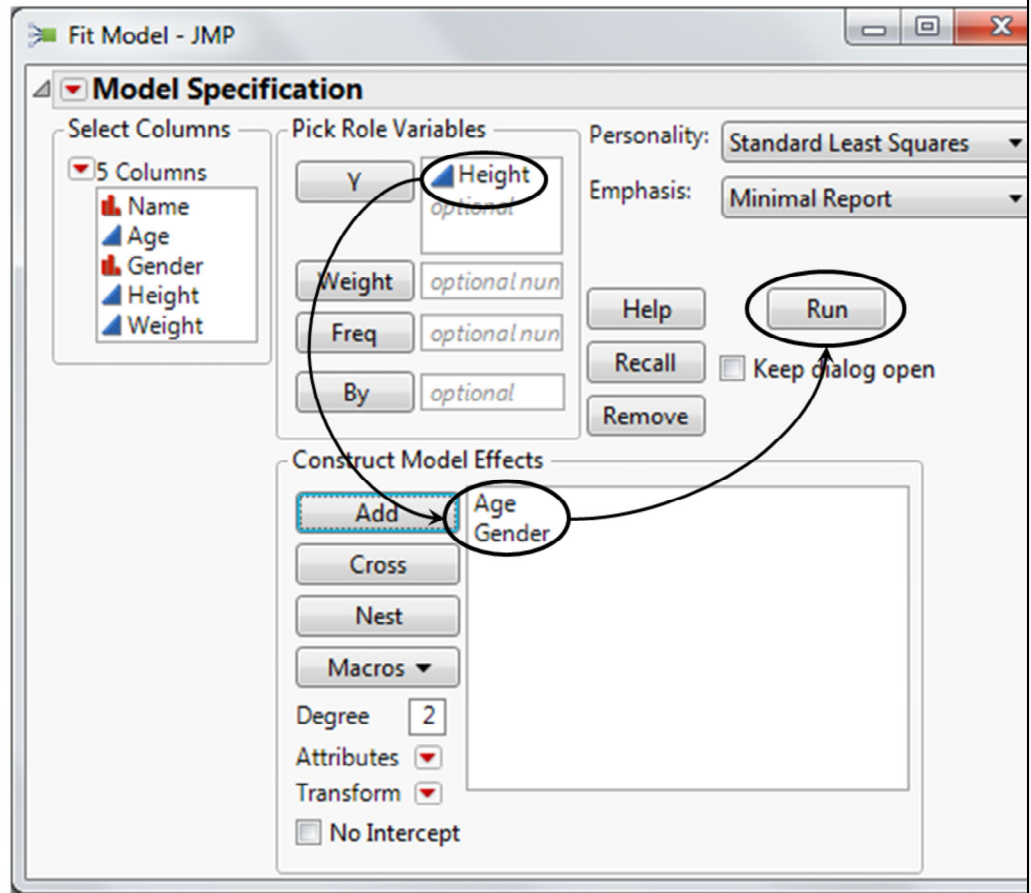
Nonlinear model	Equivalent linear model
$Y = b_0(X_1)^{b_1}(X_2)^{b_2}$	$\log(Y) = \log(b_0) + b_1 \log(X_1) + b_2 \log(X_2)$
$Y = b_0(b_1)^{X_1}(b_2)^{X_2}$	$\log(Y) = \log(b_0) + \log(b_1)X_1 + \log(b_2)X_2$

- In many cases, $\log(Y)$ transformations can successfully linearize nonlinear regression models
- This greatly extends the application of standard multiple regression models

Notes

Say we want to model *Height* as a function of *Age* and *Gender*

Analyze
↓
Fit Model



Notes

Select Options and click OK

Regression Reports

- Summary of Fit
- Analysis of Variance
- Parameter Estimates
- Effect Tests
- Effect Details
- Lack of Fit
- Show All Confidence Intervals
- AICc

Estimates

- Show Prediction Expression
- Sorted Estimates
- Expanded Estimates
- Indicator Parameterization Estimates
- Sequential Tests
- Custom Test
- Multiple Comparisons

- Inverse Prediction
- Parameter Power
- Correlation of Estimates

Effect Screening

- Scaled Estimates
- Normal Plot
- Bayes Plot
- Pareto Plot

Factor Profiling

- Profiler
- Cube Plots
- Box Cox Y Transformation
- Surface Profiler

Row Diagnostics

- Plot Regression
- Plot Actual by Predicted
- Plot Effect Leverage
- Plot Residual by Predicted
- Plot Residual by Row
- Plot Studentized Residuals
- Plot Residual by Normal Quantiles
- Press
- Durbin Watson Test

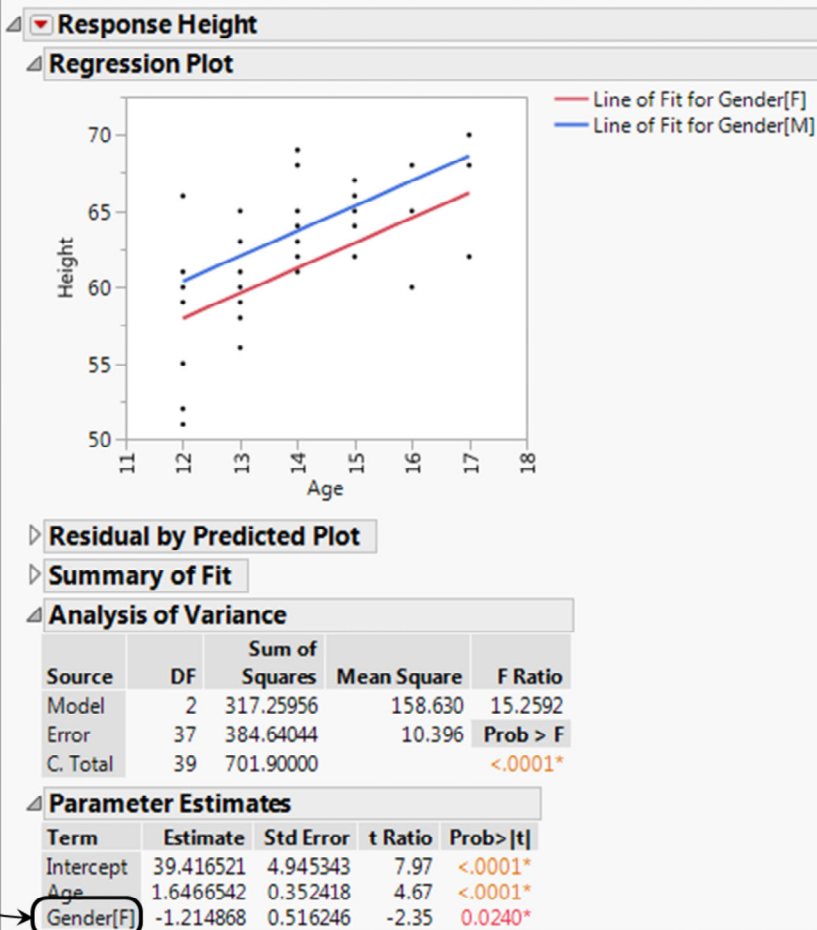
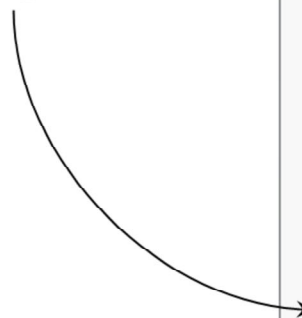
In the last column on the right (not shown), select *Effect Summary*.

Notes

Handling categorical X variables in the model

"Indicator" or "dummy" variables are used to represent categorical variables in regression.

Indicator variable representing the effect of *Gender* in the equation



Notes

In JMP, two-level categorical factors are coded +1 and -1

$$\text{Gender[F]} = \begin{cases} +1 & \text{if Gender is F} \\ -1 & \text{if Gender is M} \end{cases}$$

$$\text{Height} = b_0 + b_1\text{Age} + b_2\text{Gender[F]}$$

$$= \begin{cases} b_0 + b_2 + b_1\text{Age} & \text{if Gender is F} \\ b_0 - b_2 + b_1\text{Age} & \text{if Gender is M} \end{cases}$$

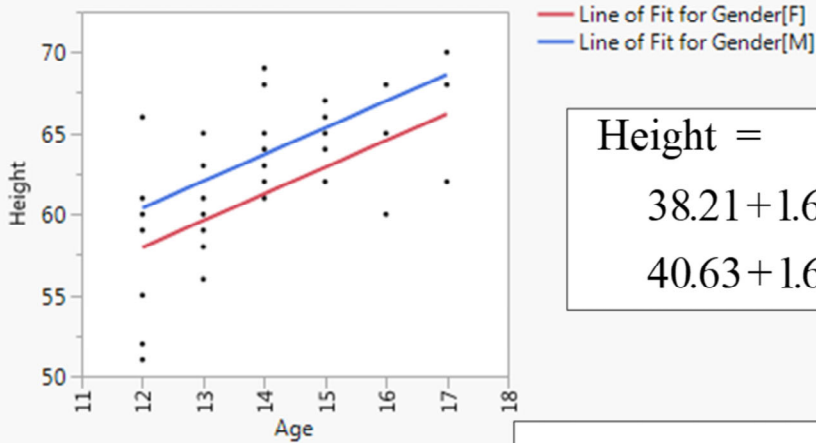
This results in one equation for Females and one equation for Males, with equal slopes (b_1) and different intercepts ($b_0 + b_2$ and $b_0 - b_2$).

An additional indicator variable is added for each additional level of a categorical variable.

Notes

Constructing the model equation

Regression Plot



$$\text{Height} = \begin{cases} 38.21 + 1.65 \text{ Age} & \text{if Gender} = \text{F} \\ 40.63 + 1.65 \text{ Age} & \text{if Gender} = \text{M} \end{cases}$$

Residual by Predicted Plot

Summary of Fit

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	2	317.25956	158.630	15.2592
Error	37	384.64044	10.396	Prob > F
C. Total	39	701.90000		<.0001*

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	39.416521	4.945343	7.97	<.0001*
Age	1.6466542	0.352418	4.67	<.0001*
Gender[F]	-1.214868	0.516246	-2.35	0.0240*

$$\text{Height} = 39.42 + 1.65 \text{ Age} - 1.21 \text{ Gender[F]}$$

If you want to verify the equation:
▼ Response Y → Estimates
→ Show Prediction Expression

Notes

$$\text{Height} = b_0 + b_1 \text{Age} + b_2 \text{Gender}[F] \\ + b_3 \text{Age} * \text{Gender}[F]$$



This product term allows different slopes for M and F

Notes

Adding an interaction effect

The image shows the 'Model Specification' dialog box in SPSS. It is divided into several sections:

- Select Columns:** A list of variables including Name, Age, Gender, Height, and Weight. 'Age' and 'Gender' are highlighted.
- Pick Role Variables:** A section for assigning variables to roles. 'Y' is selected as the dependent variable. Other variables like Height, Weight, Freq, and By are listed as optional.
- Construct Model Effects:** A section for building the model. The 'Cross' button is selected. The list of model effects includes Age, Gender, and Age*Gender.
- Personality and Emphasis:** 'Standard Least Squares' is selected for Personality, and 'Minimal Report' is selected for Emphasis.
- Buttons:** 'Help', 'Run', 'Recall', 'Keep dialog open', and 'Remove' buttons are present.

Three callouts with arrows indicate the steps to add an interaction effect:

- 1. Highlight:** Points to 'Age' and 'Gender' in the 'Select Columns' list.
- 2. Click:** Points to the 'Cross' button in the 'Construct Model Effects' section.
- 3. Interactive effect added to model:** Points to 'Age*Gender' in the list of model effects.

Notes

Regression Plot

— Line of Fit for Gender[F]
— Line of Fit for Gender[M]

The result is one model equation for Females and one for Males, with different slopes and intercepts

$$\text{Height} = \begin{cases} 46.62 + 1.04 \text{ Age} & \text{if Gender=F} \\ 32.30 + 2.24 \text{ Age} & \text{if Gender=M} \end{cases}$$

$$\text{Height} = 39.46 + 1.64 \text{ Age} - 1.23 \text{ Gender[F]} - 0.60 \text{ Gender[F]} * (\text{Age} - 13.98)$$

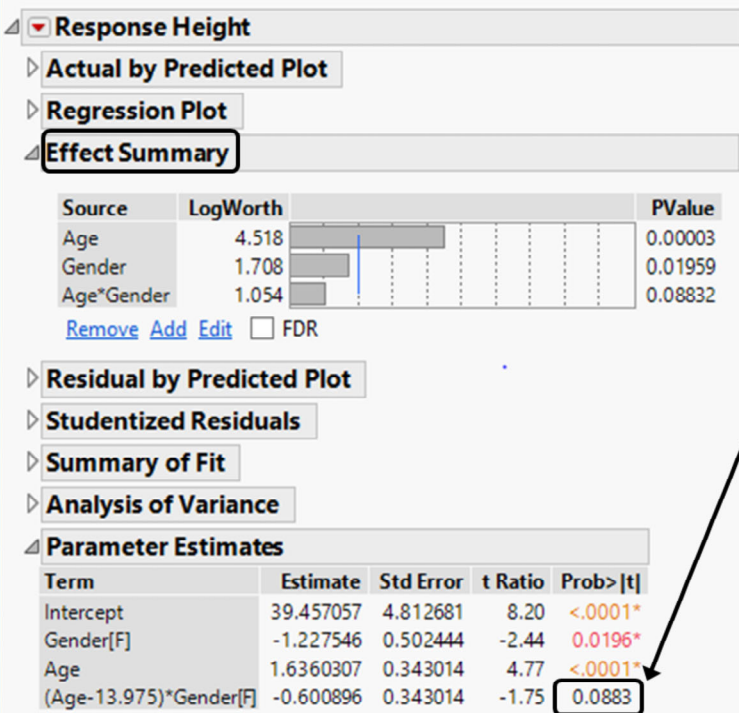
Analysis of Variance

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	39.457057	4.812681	8.20	<.0001*
Age	1.6360307	0.343014	4.77	<.0001*
Gender[F]	-1.227546	0.502444	-2.44	0.0196*
Gender[F]*(Age-13.975)	-0.600896	0.343014	-1.75	0.0883

To verify the equation:
 ▼ Response Y
 → Estimates
 → Show Prediction Expression

Notes



The p-value for Gender*Age indicates some evidence that growth curves for girls and boys have different slopes

- From now on we will use *Effect Summary* to find P-values. It gives the same information and allows model modification.

Summary of Fit without Interaction

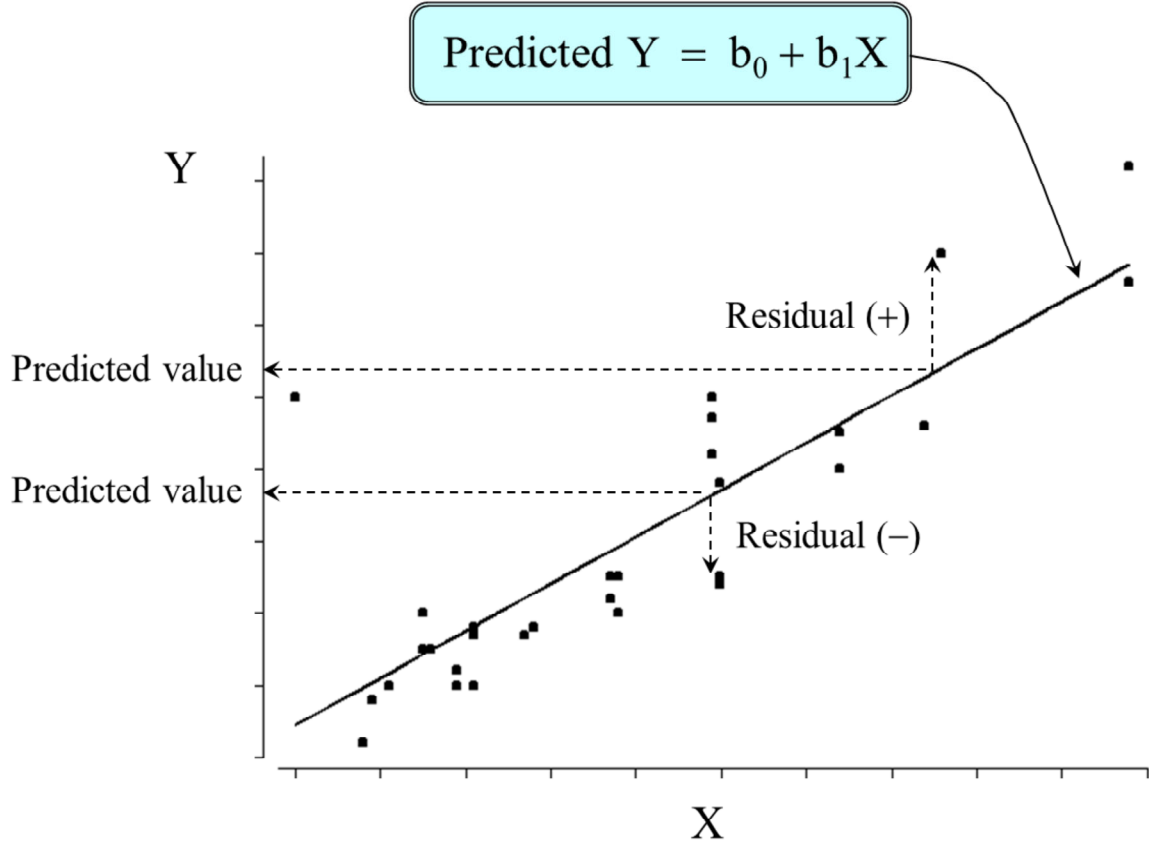
RSquare	0.452001
RSquare Adj	0.42238
Root Mean Square Error	3.224234

- ✓ Adjusted R² went up
- ✓ RMSE went down

Summary of Fit with Interaction

RSquare	0.495046
RSquare Adj	0.452967
Root Mean Square Error	3.137706

Notes



Notes

A fitted model, the equation generated during regression, gives the predicted mean value of the response variable as a function of the predictor variables. These predicted mean values are also called *predicted values*, or just *predicted* for short. The *residual value* is the data (observation) value minus the predicted value. Residual values are called *residuals* for short.

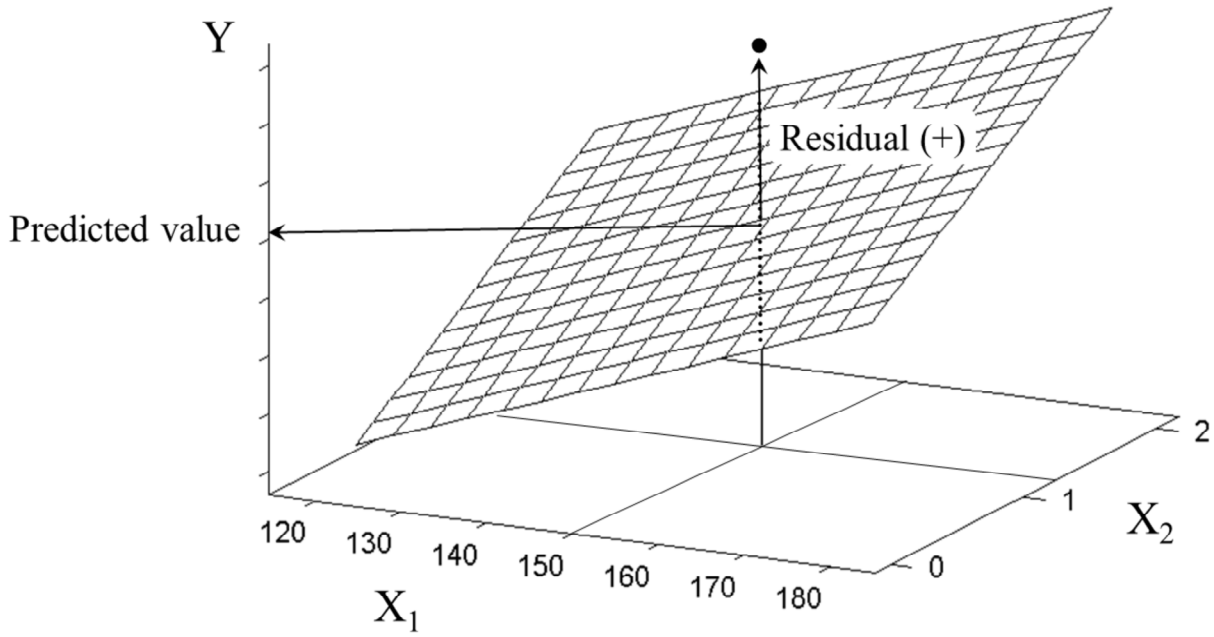
These terms are easiest to visualize in the simple linear model shown above. A predicted value is the fitted line evaluated at some X value. A residual is the difference between a measured (observed) Y value and the predicted value at the corresponding X.

Residuals contain information about the magnitude and direction of variability in the data relative to the fitted model.

- An unusually large residual might signal a measurement error, data entry error or some other type of outlier.
- A systematic trend or pattern in the residuals might signal an inadequacy in the fitted model.

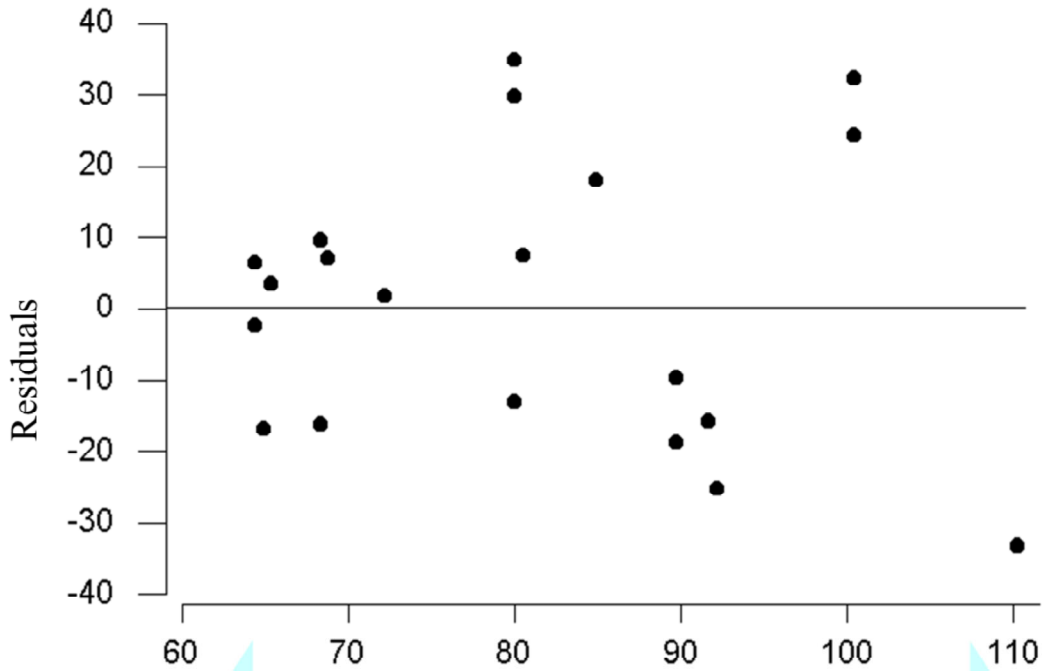
Notes

$$\text{Predicted } Y = b_0 + b_1X_1 + b_2X_2$$



Notes

Plot of residuals by predicted for any number of Xs

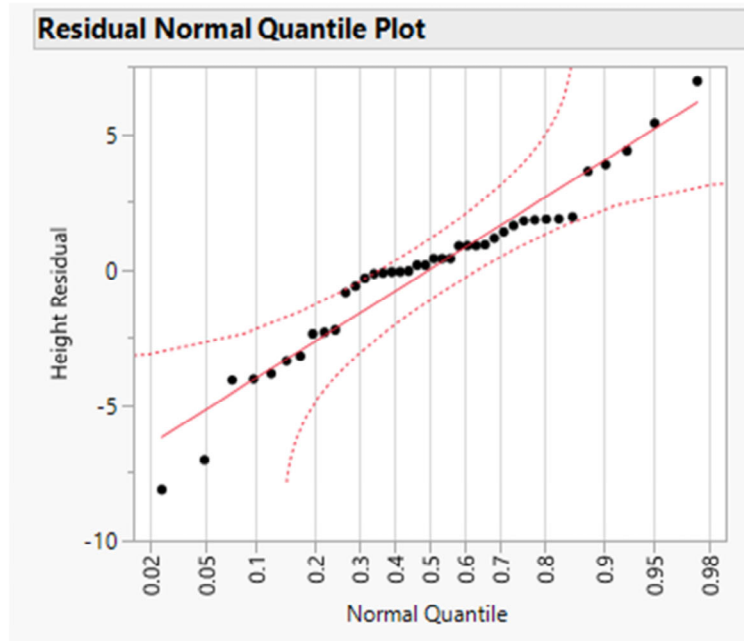


Lower left-hand quadrant of the (X_1, X_2) plane

Predicted Y values

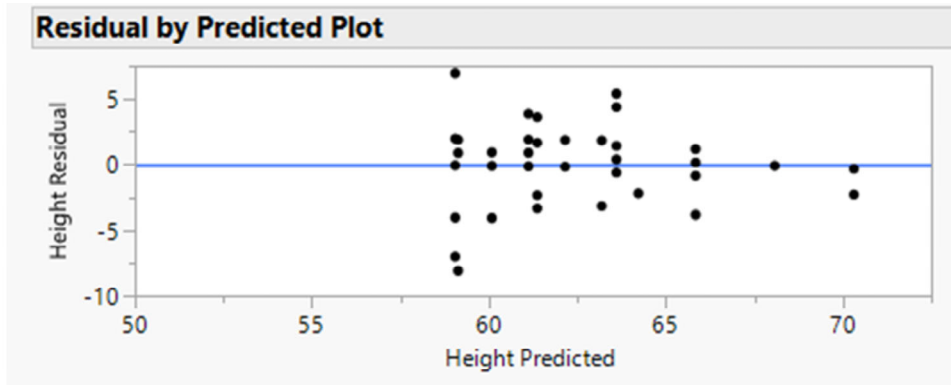
Upper right-hand quadrant of the (X_1, X_2) plane

Notes



We can see points on the hyperbolic bands here, but there is not an obvious curve through the data. Given the small sample size, this is not too concerning.

Notes



In this plot, we can see that the variance in the residuals is decreasing as height increases. This indicates the need for a transformation. We will see how to do this a little later in the course.

Notes

When historical or observational data is used to generate a regression model, an additional test is needed:

- The variance inflation factor (VIF) must be checked
- The VIF indicates whether the regressors (i.e. Xs or predictors) are correlated with each other
 - $VIF = 1$: regressor is independent of all other regressors
 - $1 \geq VIF \geq 5$: regressor is moderately correlated to other regressors
 - $VIF > 5$: regressor is highly correlated with other regressors
- VIFs in the final model need to be less than 5
 - When X variables are correlated (high VIFs), the analysis makes statistical determinations based on the noise between the correlated variables. This will often result in high R^2 values but insignificant p values.
 - VIFs are often lowered when insignificant terms are removed from the model, and terms should be removed one at a time. The first term removed should be the one with the highest p value unless theory implies removing a different one.
 - High VIFs are not an issue in designed experiments, as the designs prevent high correlation between terms/regressors

Notes

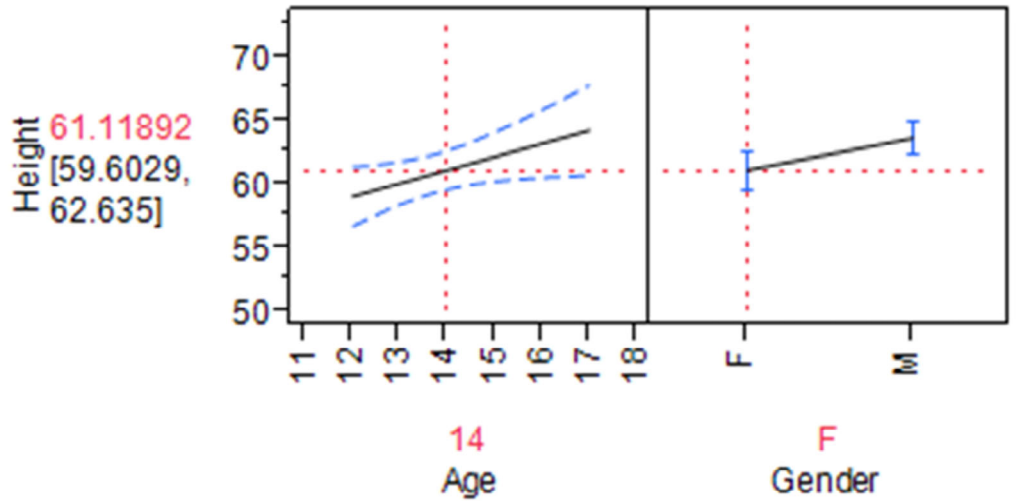
Parameter Estimates					
Term	Estimate	Std Error	t Ratio	Prob> t	VIF
Intercept	39.457057	4.812681	8.20	<.0001*	.
Gender[F]	-1.227546	0.502444	-2.44	0.0196*	1.0154192
Age	1.6360307	0.343014	4.77	<.0001*	1.0155259
(Age-13.975)*Gender[F]	-0.600896	0.343014	-1.75	0.0883	1.0004648

The variance inflation factors for all terms in the model are below 5. There is no concerning level of correlation between model terms.

To display the VIFs, right click in the Parameter Estimates section, click Columns, then VIF.

Notes

Prediction Profiler



Predicted avg. height in the population of 14 year old girls	61.12
95% confidence interval for avg. height of 14 year old girls	[59.60, 62.64] 61.12 ± 1.52

Notes

The model without interaction gave 61.25 ± 1.55 (slightly larger margin of error).

Notes

1. Run Analyze > Fit Model in JMP to investigate the relationship between y and x's. Use the Response Surface Model (all factors, all interactions, quadratic terms for continuous variables/factors)
2. Check model adequacy by reviewing the residuals plots:
 - Residual Normal Quantile Plot
 - Residual by Predicted Plot
 - Studentized Residuals (in run order)
3. Transform the data and resolve other issues, if needed.
4. Verify all VIFs < 5. Address the issue if any are over 5.
5. Remove insignificant terms from the model, that are not needed to maintain model hierarchy (main effects must be included if a higher order term of that variable remains in the model).
6. Use *Adjusted R²* to determine the amount of variation in Y that is explained by the model.

Notes

Your instructor will go through Exercise 5.4 as an example.

Notes

a) In the table below, record the Adjusted R^2 and RMSE from the analysis of *Height* in this section. Also, record the P-values from *Effects Tests*. Run the same analysis for *Weight* and record the corresponding results.

Response	Adj. R^2	RMSE	P-values		
			Age	Gender	Age*Gender
Height					
Weight					

b) Which variable (*Height* or *Weight*) has the greater proportion of variation explained by *Age* and *Gender*?

b) Explain why it wouldn't make sense to compare the two models in terms of RMSE.

Notes

- d) Both *Age* and *Gender* were statistically significant for predicting *Height*. Is this true for *Weight*?

- e) For *Height* we found evidence that the growth curves for girls and boys have different slopes. Is this true for *Weight* as well? Give the P-value that is relevant to this question and explain what it means.

- f) Give the predicted average *Weight* in the population of 15-year-old boys. Give a 95% confidence interval for this average.

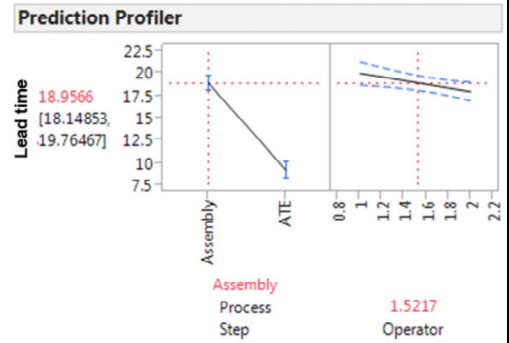
- g) Save your scripts, close and save the data table.

Notes

Exercise 5.2

Data sets \ lead time 2.

- a) Fit a model for *Lead time* including the terms *Process Step*, *Operator*, and their interactive effect. **Be sure you have the correct modeling type for Operator.** (If you got the upper right profiler, the modeling type for Operator is not correct. The lower right profiler is correct.)
- b) Note anything concerning in the residuals plots.
- c) Remove terms under *Effect Summary* with P-values exceeding 0.15 (*Remove* button). Which terms are left? Any issues with VIFs?
- d) Based on the profiler, which factor has the larger effect on lead time (steeper slope)? Does this correlate with the P-values? Please explain.
- e) Save your script, close and save the data table.



Notes

Data sets \ number and size of defects.jmp.

- a) Fit a model for *Max size* including the terms *Welder*, *# Defects*, their interactive effect, and the quadratic effect for *# Defects* (cross it with itself). This is the *Response Surface Model (RSM)* for one categorical factor and one continuous factor.

- b) Do you see anything concerning in the residuals plots?

- c) Using the *Effect Summary*, remove terms with P-values exceeding 0.15 (use the *Remove* button). Which terms are left in the model? Do all remaining terms have VIFs < 5?

- d) Based on the profiler, which factor has the larger effect on *Max size*? Does this correlate with the P-values? Please explain.

- e) Save your script, close and save the data table.

Notes

Exercise 5.4 [Instructor to demonstrate]

In this example you will analyze data from an optimization experiment concerning the removal of excess metal from castings by belt grinding.

The belt supplier had been recommending that belts be discarded when they are “50% used up.” This rule was based on tests conducted by the supplier to define the usage point at which the total of labor and belt costs will be minimized. One of the grinders thought the supplier’s rule caused grinders to discard belts too soon. Aside from being suspicious that the supplier just wanted to sell more belts, he argued that the supplier’s tests did not take into account the time lost to belt changes.

This grinder developed a new standard under which belts would be discarded only after they were “75% used up.” He wanted to do a comparative study to show that his method was cheaper overall. After he explains the study with his fellow grinders, 3 additional factors are added to the experiment.

Each casting in the experiment was weighed before and after the grinding operation. A technician kept track of how many belts were used and how long it took the grinder to complete each casting. From this information the total cost per unit of metal removed was calculated for each casting.

Data sets \ belt grinding.

Notes

Exercise 5.4 (cont'd) [Instructor to demonstrate]

94

- Y variable: *cost per unit of metal removed*
- X variables:
 - Contact wheel land-groove ratio (LGR): Low or High
 - Contact wheel material (MATL): Steel or Rubber
 - Belt usage limit (USAGE): "50%" or "75%"
 - Belt grit size (GRIT): 30 or 50
- **Run the *Fit Model* script provided in the left panel**, by clicking on the green triangle. This is the response surface model for 4 categorical X variables.
- Check the residuals plots. Any problems?
- Using the *Effect Summary*, remove insignificant terms not needed to maintain model hierarchy, starting with the group of terms with $P > 0.20$, then one at a time. Which terms are left in the model?
- Use the *Prediction Profiler* to find the minimum cost factor settings.
- What do you expect the mean and standard deviation of *Cost* to be after implementing the optimal factor settings?
- Save your script, close and save the data table.

Notes

Exercise 5.5

In this example you will analyze data from an optimization experiment concerning the bond strength of potato chip bags.

Chips ‘R’ Us was receiving customer complaints about stale chips, especially from customers on airplanes. They traced the problem to the bag sealing process. The current process involved a temperature of 150°C, a pressure of 100 psi and a dwell time of 1.1 secs. The current average bond strength was about 85 psi.

Process Engineer Chip Kettle ran an experiment to increase the bond strength. Production Manager Justin Thyme reminded Chip that he would very much like to avoid an increase in the dwell time.

Justin is able to free up a bag sealer for only so much time each shift. Chip realizes he will need two shifts to complete the experiment. He decides to include *Shift* as an additional variable in the analysis just in case there is an operator and/or equipment effect.

Data sets \ heat sealing 1.

Notes

- Y variable: *bond strength*
- X variables and feasible ranges:

➤ Temperature (TEMP):	120 to 180
➤ Pressure (PRESS):	50 to 150
➤ Dwell time (DWELL):	0.2 to 2.0
➤ Shift:	1 or 2
- **Run the *Fit Model* script provided in the left panel.** This is the response surface model (RSM) for 3 continuous X's. Is anything concerning in the residuals plots?
- Remove from the model insignificant terms that are not needed to maintain model hierarchy ($P > 0.15$), using the *Effect Summary*. Which terms are left?
- Use the *Prediction Profiler* to maximize the average bond strength. If your solution requires a long dwell time, manually move things around in the profiler to find another solution with a short dwell time.
- What do you expect the mean and standard deviation of *bond* to be after implementing the optimal factor settings?
- Save your script, close and save the data table.

Notes

Data sets \ outgassing process. *Current* (the Y variable) is the electrical current required to heat a filament to a specified temperature. *Resist* (one of the X variables) is the electrical resistance of the filament. *Machine* (the other X variable) identifies which of three processing units was used. We want to develop a model for *Current* as a function of *Resist* and *Machine*.

- a) Fit a response surface model for *Current*. (The terms will be *Resist*, *Machine*, the interaction term *Resist*Machine*, and the quadratic term *Resist*Resist*. To get the quadratic term, highlight *Resist* both under Select Columns and under Construct Model Effects, then click Cross.)
- b) Do you see anything concerning in the residuals plots?
- c) Remove any terms under *Effect Summary* with P value exceeding 0.15. (Use the *Remove* button.) Record the RMSE.
- d) Use the *Prediction Profiler* to find the predicted average *Current* for each machine if we always use filaments with resistance 52.

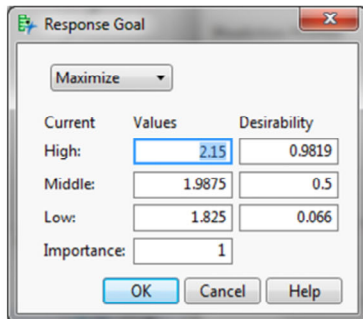
Notes

Exercise 5.6 (cont'd)

98

- e) The target value for *Current* is 2. For each machine, we want to find the resistance for which the average current is 2. On the *Prediction Profiler* red triangle, select *Desirability Functions*. It should look like this:

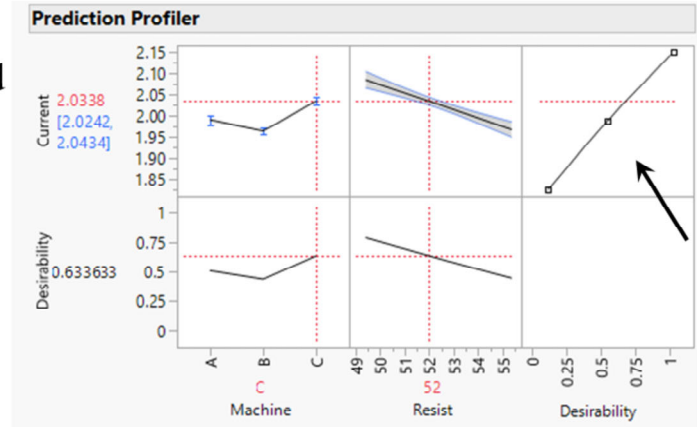
- f) Double click in the upper right hand panel of the profiler. (Try to avoid the plotted line.) You should get the dialog shown below.



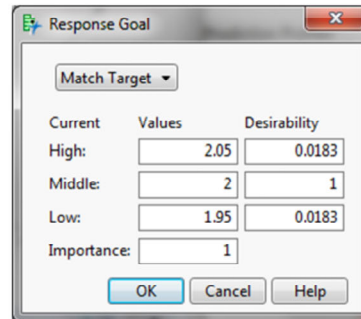
Current	Values	Desirability
High:	2.15	0.9819
Middle:	1.9875	0.5
Low:	1.825	0.066

Importance: 1

OK Cancel Help



- g) Modify the dialog as shown to the right, then select OK. Proceed to the next slide.



Current	Values	Desirability
High:	2.05	0.0183
Middle:	2	1
Low:	1.95	0.0183

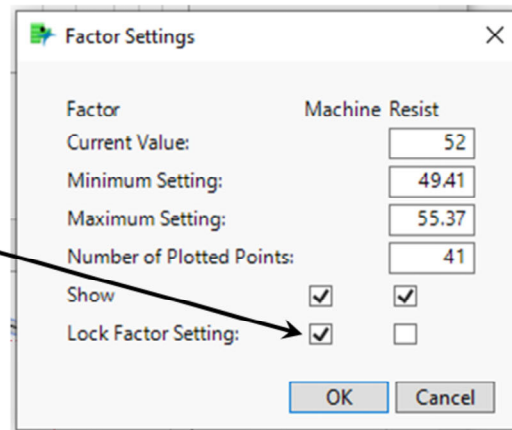
Importance: 1

OK Cancel Help

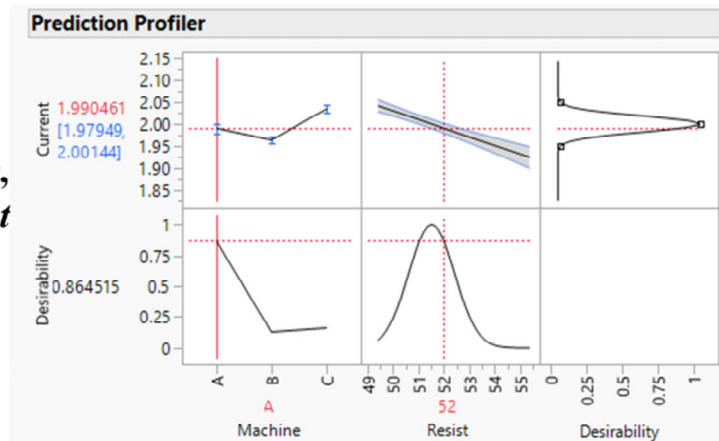
Notes

Exercise 5.6 (cont'd)

- h) On the *Prediction Profiler* red triangle, select *Reset Factor Grid*. We want to lock the factor setting for *Machine*, so check the *Lock Factor Setting* box as shown here.



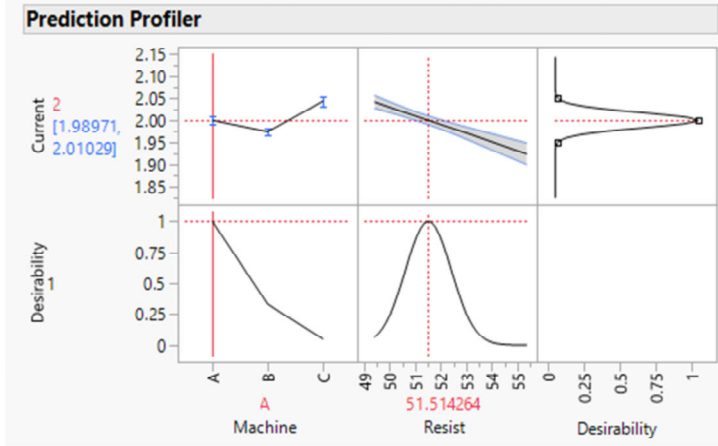
- i) The vertical line for *Machine* should now be solid instead of dotted. **This will hold the machine setting in place during *Maximize Desirability*, which allows you to optimize *Resist* separately for each machine.** On the *Prediction Profiler* red triangle, select *Maximize Desirability*. Proceed to the next slide.



Notes

Exercise 5.6 (cont'd)

j) The optimal resistance value for Machine A is 51.5. Drag the solid vertical line across to B, then click *Maximize Desirability* to find the optimal resistance value for Machine B. Do the same for Machine C.



k) What will the average current be if we always use the optimal resistance values for each machine?

l) What will the standard deviation of current be if we always use the optimal resistance values?

m) Save your scripts, close and save the data table.

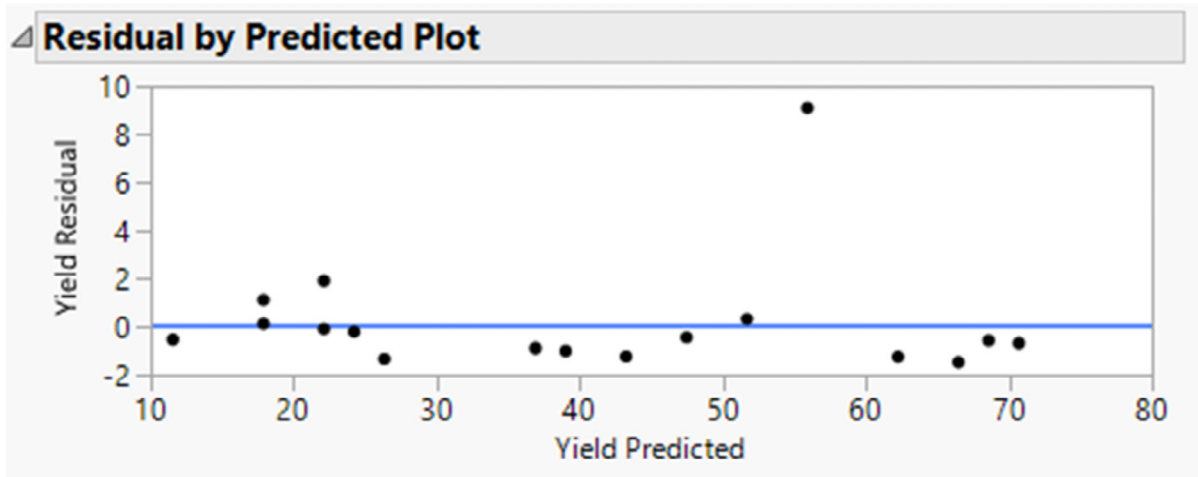
Notes

In this section, we will cover the most common model adequacy issues:

- Outliers
- Pattern in run order plot of residuals
- Multicollinearity (VIFs over 5)
- Unequal variance and non-normal residuals

Notes

Outliers can easily be seen on the Residual by Predicted and Studentized Residuals (residuals by run order) plots



Remember, healthy residuals look like random scatter about zero.

Here, it looks like there might be a suspicious data point.

Notes

- Investigate the data point.
 - If it turns out to be just a data entry error, we simply enter the correct value, then all is well. Most of the time it's not that simple.
- If you have an outlier of unknown origin:
 - Run the analysis with and without the questionable data point.
 - If you're lucky, the results will be pretty much the same both ways, hence no worries. Leave the data point in.
- If excluding the outlier does make a significant difference in the results, then you have a hard decision to make.
 - The official rule is: leave the data point in unless you can identify the cause. The idea is to throw it out only if you can demonstrate that it does not come from the population you want to study. This is the "pure" approach.
 - This should be tempered with the following practical consideration: you don't want your results to be unduly influenced by one extreme outlier, even if you can't explain it.

Notes

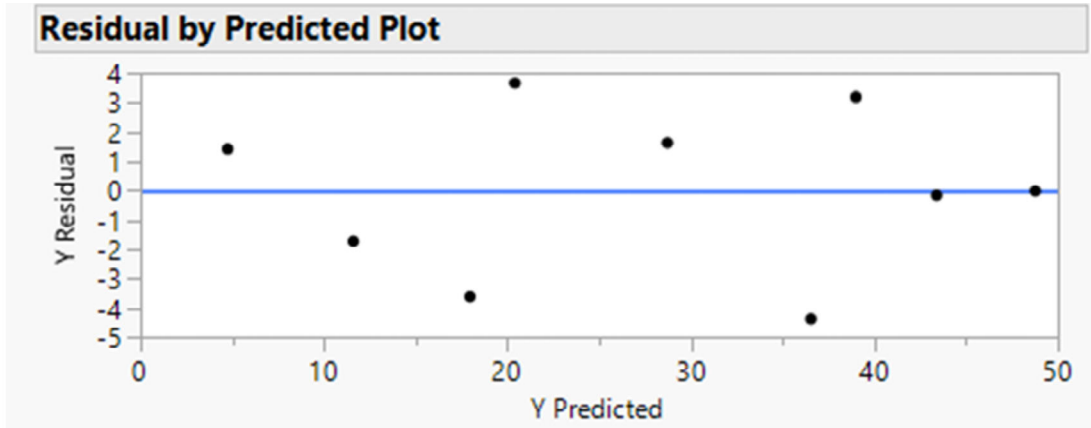
- Runs (points in a row) of positive-negative-positive-negative residuals indicate correlation between runs in an experiment.
 - This implies that the assumption of independence has been violated.
 - **Randomization of an experiment protects against this! Do it every time!**
- This plot can show changes in variance over the time span of the experiment or data collection.
 - This could be due to increased skill as the experiment progresses, a process drift, operator fatigue, tool wear, etc.
 - This type of problem would show as an increase or decrease in spread or “scatter” of the residuals across the graph.
 - If there is x data available to support it, one remedy is to add a factor (time since tool change, number of hours of operator work, etc.)
 - Increasing or decreasing variance can also indicate the need for a transformation.

Notes

Several strategies can be tried for resolving multicollinearity, but they may not be satisfactory, especially if the model will be used for prediction.

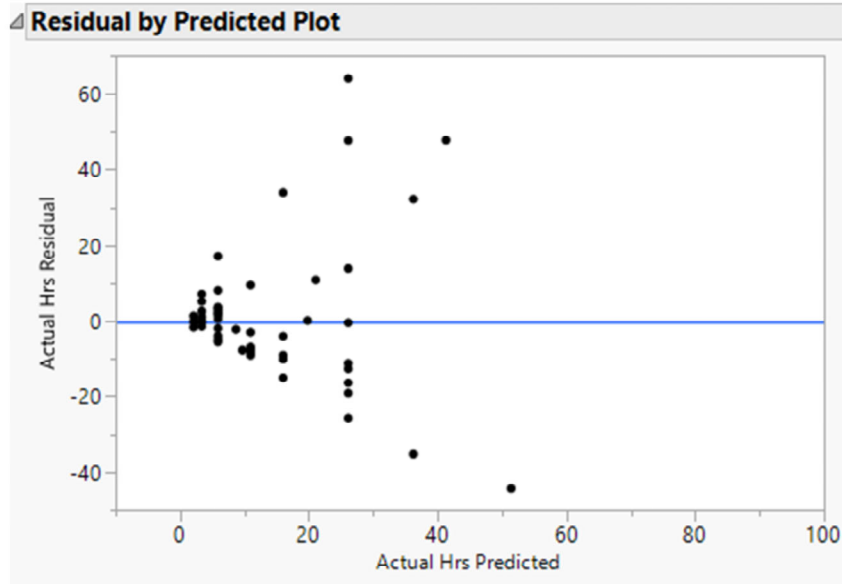
- Collect additional data in a way that breaks up the multicollinearity.
 - Historical data may contain only certain combinations of x-variables, for example, only low levels of x_1 when x_2 is at a low level and only high levels of x_1 when x_2 is at a high level
 - Note: it may not be feasible or possible to collect this additional data.
 - In some cases, the factors (x's) are inherently correlated, for example as may be the case with household income and house size.
- Respecifying the model, can help.
 - If x_1 and x_2 are nearly linearly dependent, use one term, $x = x_1 + x_2$, which preserves the information content of the original variables
 - Try removing the term with the highest p-value, and look at that model. Then, replace it and remove the term with the highest VIF. See which gives the better model.
- Use ridge or principal-component regression (way beyond the scope of this course)

Notes



Remember, the variation in the residuals should be fairly constant across the Residual by Predicted Plot. There is no issue here.

Notes

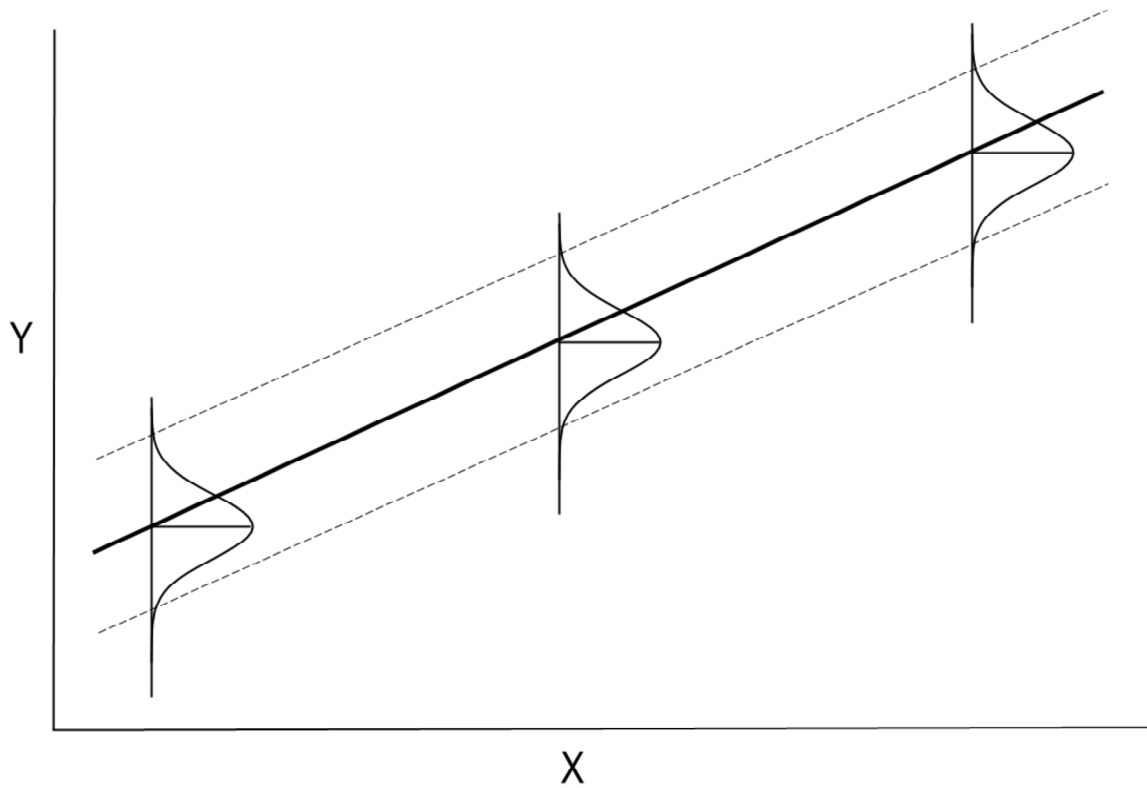


In this plot, we can see an issue.

σ_Y^2 proportional to mean Y \rightarrow "sideways V"

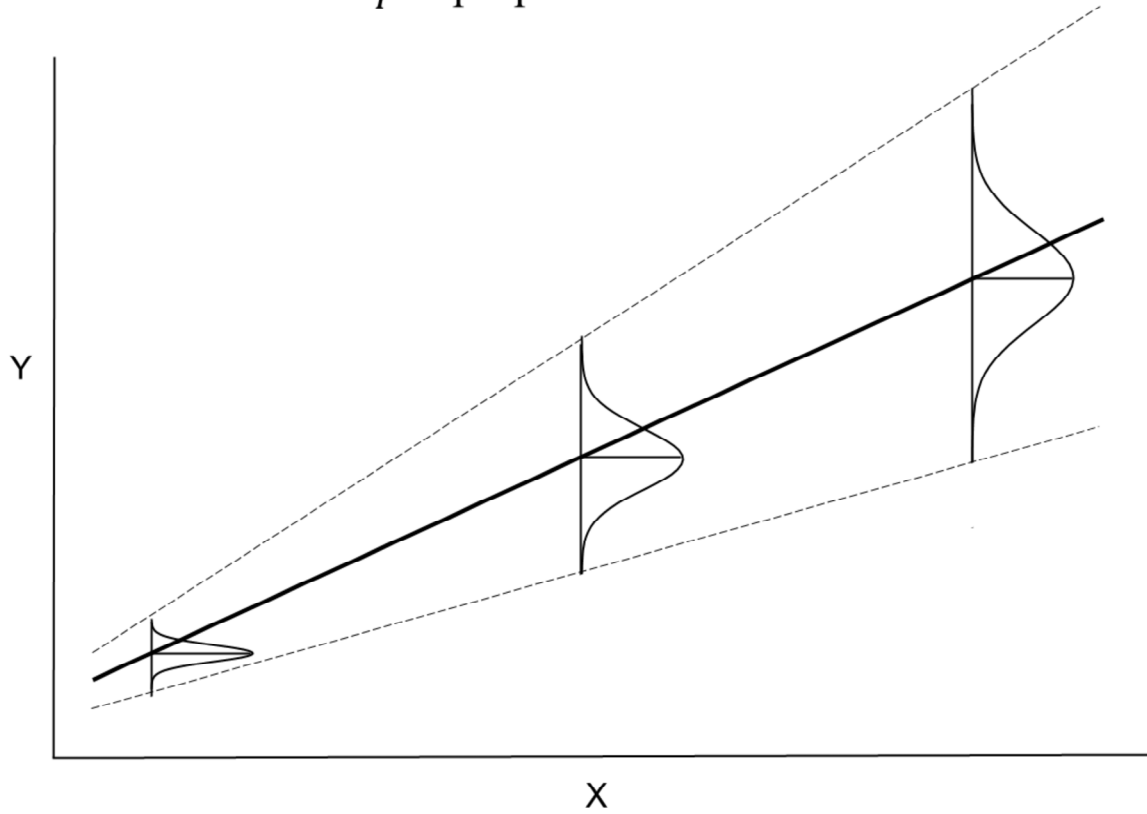
Notes

σ_Y^2 is constant (does not depend on the X's)



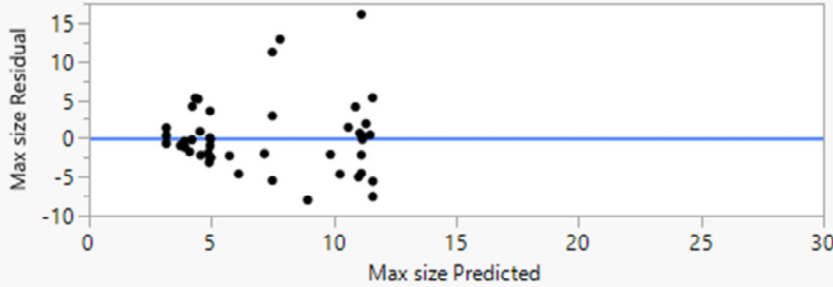
Notes

σ_Y^2 is proportional to mean Y

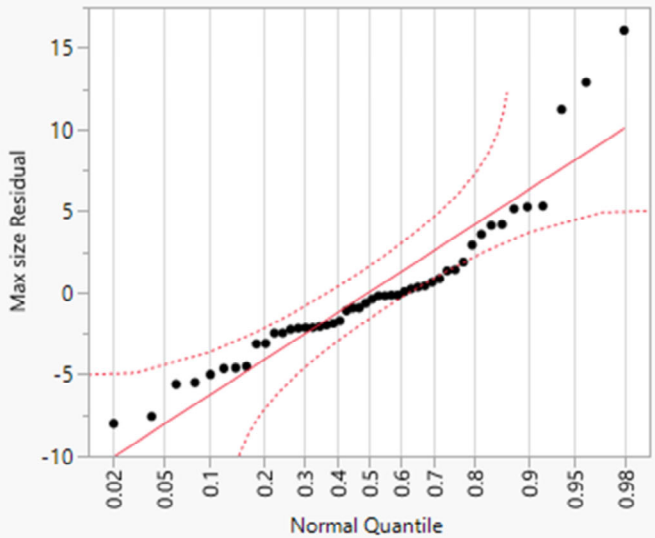


Notes

Residual by Predicted Plot



Residual Normal Quantile Plot



- Often, when there is an issue with constant variance, there is also the issue of non-normal residuals.
- This can be seen in these two plots
- Fortunately, they usually both resolve with the same treatment—a transformation.

Notes

The standard assumption in all comparison and correlation analyses involving a quantitative Y variable is that the noise (unexplained/error/residual) variation follows a Normal distribution with mean 0 and a standard deviation that does not depend on the X variables.

This simple model has served us well. However, when Normality or constant σ is grossly violated, something must be done. The most common remedy is to use $\log(Y)$ or \sqrt{Y} as the dependent variable instead of Y. This is a transformation. This “trick of the trade” is simple and, in most cases, effective.

Notes

JMPs notation regarding Logs requires some clarification:

- Although JMP expresses the logarithm as “Log”, it is actually base e, or the natural log, which is usually written as Ln. It is not a base 10 logarithm.
- However, the plots that use a log transformed X-axis display use base 10 log for the X-axis. This does not change the interpretation of the chart.

The impact of transformation on R^2 and p-values:

- In the previous example, a transformation was required because the residuals variance wasn't constant over the range of the predicted values.
- After the transformation, the R^2 value went down. This can lead to a belief that the non-transformed model was “better”. However,
- Residuals showing this condition (heteroscedasticity) can cause p-values and R^2 to be over or under stated.
- When this condition occurs, the problem must be corrected. The resulting model, even if R^2 is lower or p-values are higher, is the more “real” model.

Notes

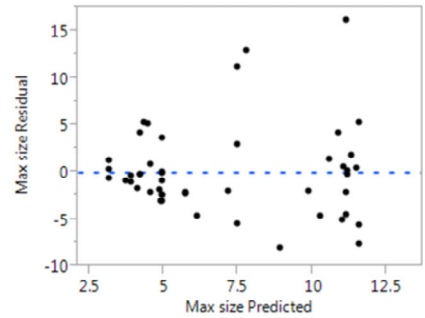
Data sets \ number and size of defects.jmp

- a) Fit a model for *Max size* including the terms *Welder*, *# Defects*, their interactive effect, and the quadratic effect for *# Defects* (*response surface model* for one continuous factor and one categorical factor). You should see a distinct sideways V. Do you see issues in any other residuals plots?

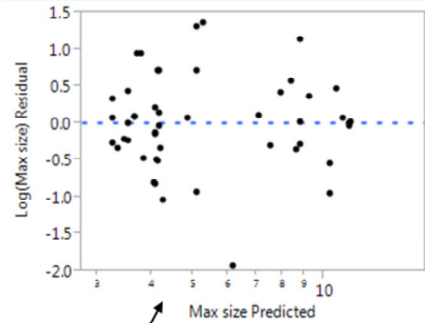
- b) Select *Model Dialog* on the *Response* red triangle menu, apply a Log transformation to *Max size*, re-run the model. The sideways V isn't completely gone, but close enough. Did other residuals plots improve?

- c) Use *Effect Summary* to remove terms with $P > 0.15$.

Residual by Predicted Plot



Residual by Predicted Plot



Remember to change the x-axis on the plot, as well.

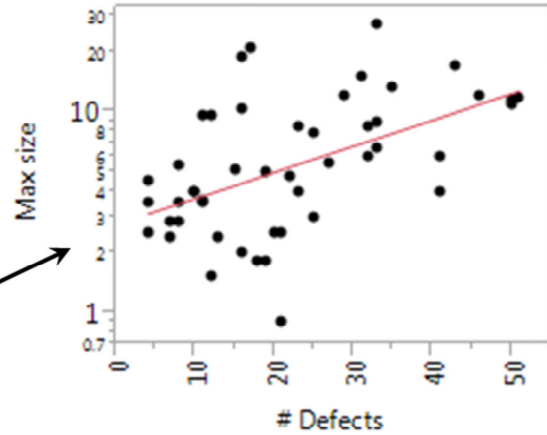
Notes

d) Which terms are left in the model?

e) Now we have a log-linear simple regression.

When you use a Log or square root transformation on Y, it is helpful to use same scale for the Y axes of the plots

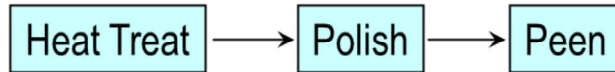
Regression Plot



f) Save your script, close and save the data table.

Notes

An aerospace manufacturer uses integral castings as structural components of jet engines. Integral castings give design engineers more flexibility and simplify the assembly process. Defect-free castings are known to have long cycle fatigue life, but defects often arise in the casting process and must be weld repaired. The engine manufacturer’s metallurgical team has proposed a finishing process of the following type to ensure adequate cycle fatigue life of weld-repaired castings:



The team wants to optimize the first two steps in this process to achieve maximum cycle fatigue life. Also, though other applications of similar processes have included peening, they would like to see if it can be omitted to reduce processing time and cost.

Due to project time constraints and limited availability of test fixtures, the team can perform at most 12 cycle fatigue tests for their experiment.

Notes

- Y variable: *Cycles* (to failure)
- X variables:
 - Heat treat: Anneal or Solution/age
 - Polish: Chemical or Mechanical
 - Peen: Yes or No
- *Data sets \ weldment fatigue.jmp*.
- Run the *Model* script provided in the left panel, run the model.
- Notice the extreme sideways V on the *Residual by Predicted Plot*. Are there issues in any of the other residuals plots? If yes, what are they?
- Rerun the model using a Log transformation on *Cycles*. Did residuals plots improve?
- Remove insignificant terms from the model ($P > 0.15$) that are not needed to maintain model heirarchy.
- Use the *Prediction Profiler* to maximize the cycle fatigue life.

Notes

Exercise 6.3

124

A Black Belt wants to minimize the *leak rate* in plastic containers ultrasonically welded together. The X variables and ranges are:

- Force: 70 to 150
- Energy: 275 to 325
- Amplitude: 70 to 90

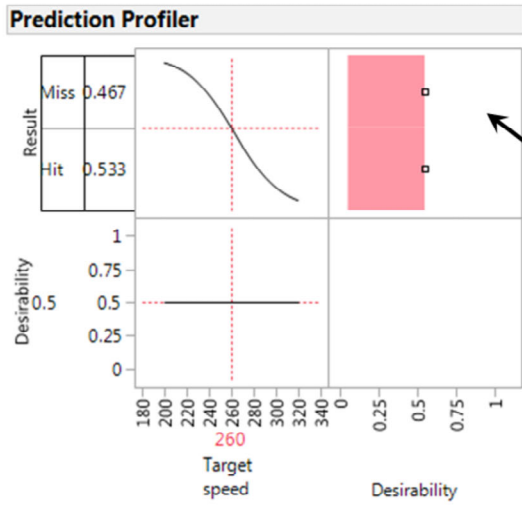
- *Data sets \ ultrasonic welding 1.jmp.*
- Run the *Model* script provided in the left panel.
- What problems do you see in the residuals plots?

Notes

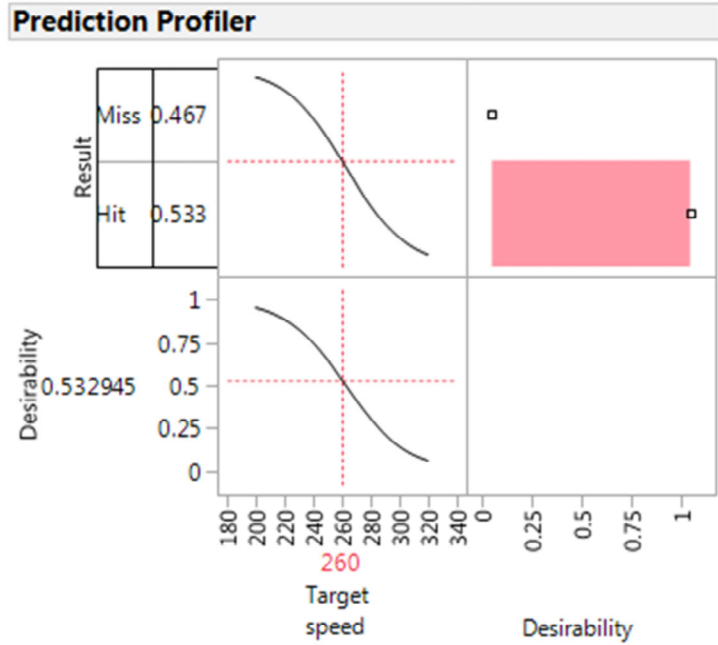
- Rerun the model using the Log transformation on *leak rate*. (Be sure to change the x-scale to Log on the Residual by Predicted Plot.)
- Rerun the model using the Sqrt transformation on *leak rate*. (Be sure to change the x-scale to Sqrt on the Residual by Predicted Plot.)
- Which set of residuals plots looks better? Use whichever transformation looks like it worked better, going forward.
- Remove insignificant term(s) from the model ($P > 0.15$), while maintaining model hierarchy.
- Use the *Prediction Profiler* to minimize the leak rate.

Notes

This page intentionally left blank

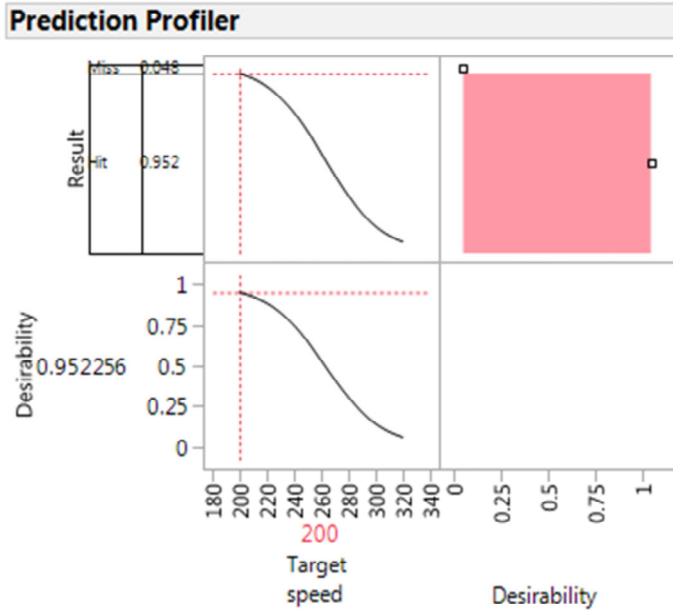


- Red Triangle → Profiler → Prediction Profiler
red triangle → Optimization and Desirability → Desirability Functions
- Double-click in the blank area, enter 1 for Hit and 0 for Miss → OK → OK → next slide



Notes

Prediction Profiler red triangle → Optimization and Desirability → Maximize Desirability



- The target speed of 200 produces the maximum hit probability of 0.952
- The corresponding miss probability is 0.048
- The target speed of 320 produces the minimum hit probability of 0.061
- The corresponding miss probability is 0.939

Notes

	Mins at temp	Result	Freq
1	2	Cracked	0
2	2	Not cracked	100
3	4	Cracked	1
4	4	Not cracked	99
5	6	Cracked	2
6	6	Not cracked	98
7	8	Cracked	3
8	8	Not cracked	97
9	10	Cracked	7
10	10	Not cracked	93
11	12	Cracked	9
12	12	Not cracked	91
13	14	Cracked	12
14	14	Not cracked	88
15	16	Cracked	13
16	16	Not cracked	87
17	18	Cracked	15
18	18	Not cracked	85

Analyze
↓
Fit Model
↓
See next slide
↓
Set up as shown

Notes

Model Specification

Select Columns: 3 Columns
 Mins at temp
 Result
 Freq

Pick Role Variables:
 y: Result (optional)
 Weight: optional numeric
 Freq: Freq
 By: optional

Personality: Nominal Logistic
 Target Level: Cracked

Buttons: Help, Run, Recall, Remove

Construct Model Effects:
 Add: Mins at temp
 Cross, Nest, Macros
 Degree: 2
 Attributes, Transform
 No Intercept

Keep dialog open

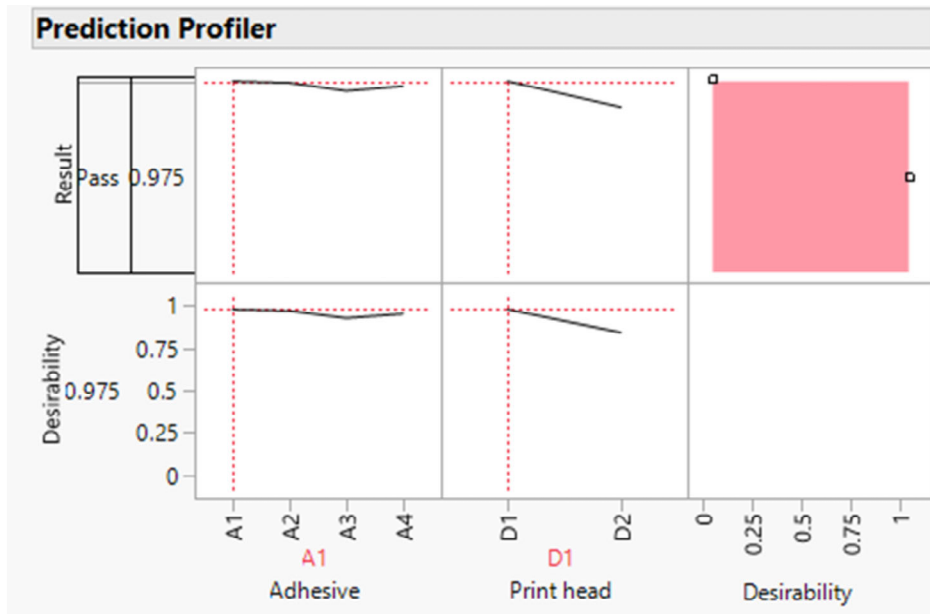
In this data set, instead of a row for each observation, the results are tabulated—there is a count of outcomes for each level of the X variable.

Using the Freq values tells JMP how many times to count each row.

Notes

This page intentionally left blank

- *Prediction Profiler* red triangle → *Optimization and Desirability* → *Maximize Desirability*
- The failure rate predicted from the optimization was 0.025 or 2.5% (current state failure rate was 20% or more)
- Best combination was D1 with A1



Notes

A Black Belt wants to minimize the occurrence of bubbles and ripples in the urethane coating on truck nameplates. The X variables and ranges are:

- Badge temp: 20 to 40
- Mixing ratio: 92.6 to 94.6
- Curing temp: 30 to 55

- *Data sets \ urethane coating pass-fail*
- Run the *Model* script in the left panel. In the *Model Specification*, switch the *Target Level* from *Fail* to *Pass*, then run the model.
- Remove insignificant terms from the *Effect Summary* ($P > 0.15$).
- Use the *Prediction Profiler* to find a factor combination that maximizes the yield.
- The current state yield was about 95%. What is the predicted yield for the improved process?

Notes
