



Workshop: Analyzing the Data You Already Have – An Introduction to Model Fitting

Don Lifke

Corporate Lean / Six Sigma Black Belt

Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.



Objectives

- Learn new ways to look at existing data
- Understand limitations with typical data analysis
- Understand what modeling is and why it is useful
- Learn to use statistical software and understand why Excel is not the best choice for efficient statistical analysis



Define “Statistics”

Define “statistics” in two words:



Exercise

You want to sell your house. It has the following features:

- 2000 square feet
- 0.2 acre lot
- 2 years old
- 3 bedrooms
- 3 full bathrooms

What should your asking price be?



Open *House Data for Summit* *Tutorial.xls*



What type of analyses would you do on this data?

-
-
-

Spend 5 minutes determining what you'd list your house for, using the Excel data.



Exercise Time

5 minutes



List Class Prices

- Listing Prices:
 - \$
 - \$
 - \$
 - \$
 - \$
 - \$

Open *House Data for Summit Tutorial.jmp* Data File

JMP (SANDIA NATIONAL LABORATORIES) - [House Data for Summit Tutorial.JMP] - [House Data for Summit Tutorial]

File Edit Tables Rows Cols DOE Analyze Graph Tools View Window Help

House Data for Summit

	SF	Lot	Age	BR	Bath	Price	Price/sf
1	1373	0.13	7	4	3	204962.96	149.281107
2	1377	0.2	1	2	3	279461.24	202.949339
3	2696	0.21	1	2	3.75	432115.58	160.28026
4	2743	0.2	11	3	2.75	291085.68	106.11946
5	1128	0.19	14	5	3.25	163331.78	144.797677
6	3721	0.16	5	3	2	417458.93	112.189984
7	3372	0.05	19	4	2.5	291889.4	86.5626928
8	1342	0.1	20	4	4	91196.94	67.9559911
9	1317	0.23	17	3	3	118951.58	90.3201063
10	2370	0.25	19	3	2.5	186523.51	78.701903
11	1645	0.18	9	5	4	277864.81	168.914778
12	2306	0.08	0	4	2.75	339135.6	147.066609
13	1356	0.23	1	2	3.25	254317.26	187.549602
14	2421	0.08	20	3	3.75	176160.96	72.7637175
15	1801	0.17	11	4	3.75	245049.51	136.063026
16	2195	0.19	17	2	3.5	195129.12	88.8970934
17	2172	0.15	13	4	4	253373.93	116.654664
18	2002	0.17	1	4	2.75	360202.51	179.921334
19	1851	0.2	11	4	3.5	261394.12	141.217785
20	2520	0.1	13	3	4	259948.18	103.15404
21	2102	0.05	14	2	3.75	177637.02	84.5085728
22	2533	0.08	11	3	3.25	285993.76	112.90713
23	2983	0.11	0	4	2	442720.07	148.414371
24	3249	0.23	2	5	3	468637.93	144.240668
25	1585	0.2	19	2	3.75	135501.42	85.4898549
26	1560	0.24	0	5	3	360846.46	231.311833
27	3319	0.14	2	2	2.75	442828.35	133.422221
28	3691	0.21	10	3	3.25	450724.35	122.114427
29	1484	0.1	20	2	2	49746.35	33.5217992
30	3619	0.11	17	5	2	338789.82	93.6142083

Columns (8/0)

- SF
- Lot
- Age
- BR
- Bath
- Price
- Price/sf
- Bogus

Rows

All rows 30

Selected 0

Excluded 0

Hidden 0

Labelled 0



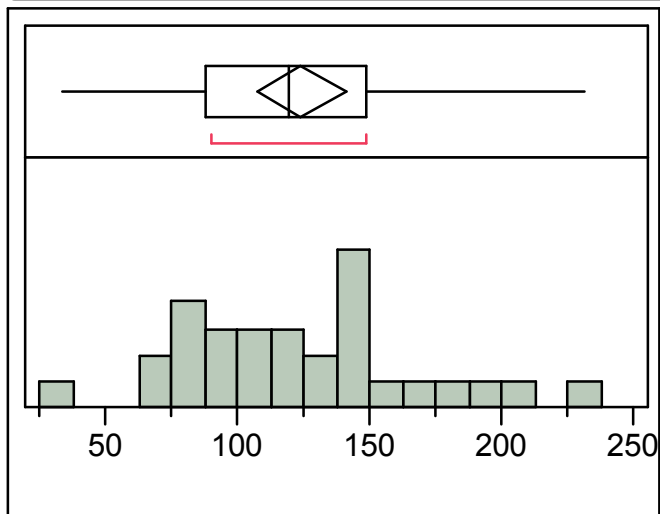
Exercise:

What Will Your Listing Price Be?

Analyze the distribution of Price/sf
Average = \$124.36 per SF

Distributions

Price/sf



Quantiles

100.0%	maximum	231.31
99.5%		231.31
97.5%		231.31
90.0%		186.79
75.0%	quartile	148.63
50.0%	median	119.38
25.0%	quartile	88.31
10.0%		73.36
2.5%		33.52
0.5%		33.52
0.0%	minimum	33.52

Moments

Mean	124.36354
Std Dev	44.020787
Std Err Mean	8.0370594
upper 95% Mean	140.80117
lower 95% Mean	107.92591
N	30

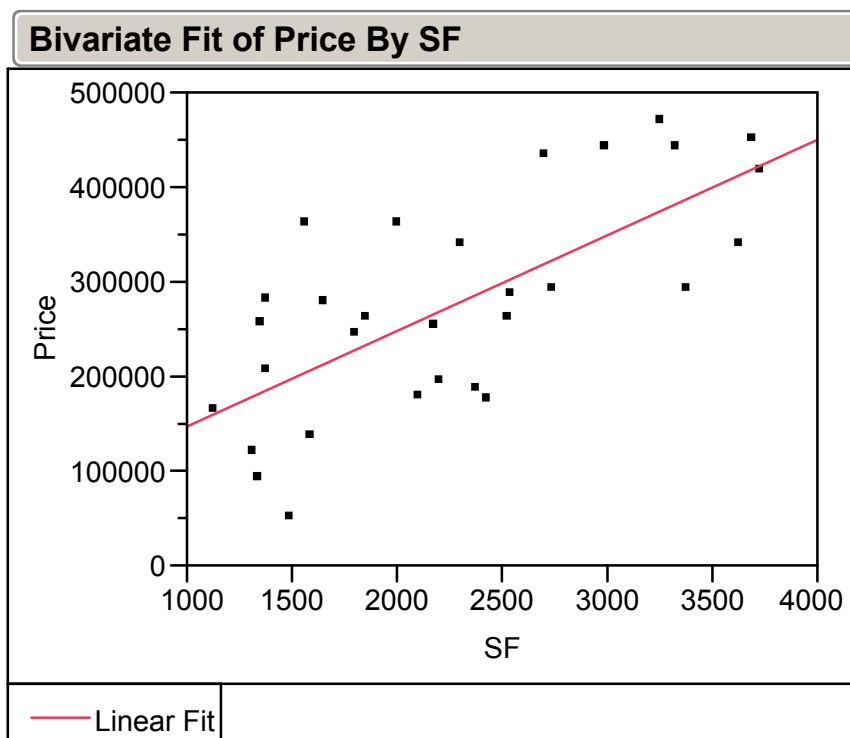


Exercise:

What Will Your Listing Price Be?

- Perform a Fit Y by X for Price vs. SF
- Add a Line Fit

$$\text{Price} = \$45,962 + \$101.34 * \text{SF}$$



Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	1.8791e+11	1.879e+11	28.1954
Error	28	1.8661e+11	6.6647e+9	Prob > F
C. Total	29	3.7452e+11		<.0001*

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	45962.927	45654	1.01	0.3227
SF	101.33845	19.0847	5.31	<.0001*



Exercise Time

Using this new information, Spend 5 minutes determining what you'd list your house for, using the JMP data.

5 minutes

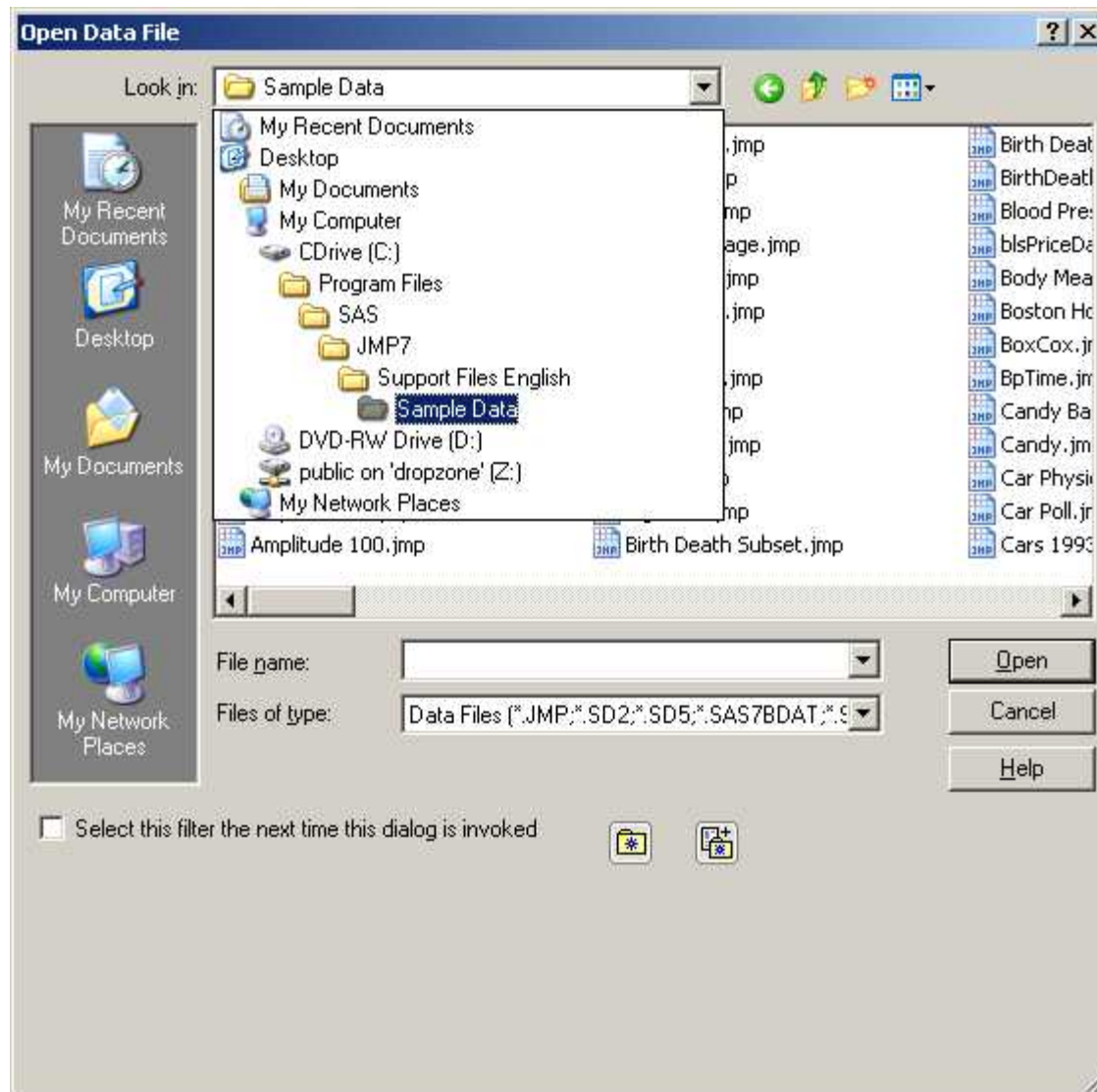


Exercise:

What Will Your Listing Price Be?

- Write Student's Listing Prices on the Board
- Put These in JMP for Future Reference
- We'll Return to this Problem Later

Analyzing the Data You Have: Open Cereal File





Before We Proceed

- Rows → Clear Row States
- File → Preferences
 - Click Reports
 - Change Graph Marker Size to medium.
 - Click Platforms
 - Select Distribution. Under Options, select Stack.
 - Click OK.

Does Calorie Content Depend on the Manufacturer? Analyze → Fit Y by X

Y by X - Contextual

Rows Cols DOE Analyze Graph Tools View Window Help

Distribution of Y for each X. Modeling types determine analysis.

Select Columns

- Name
- Manufacturer
- Mfr
- Hot/Cold
- Calories
- Protein
- Fat
- Sodium
- Fiber
- Complex Carbo
- Tot Carbo
- Sugars
- Calories fr Fat
- Potassium
- Enriched

Cast Selected Columns into Roles

Y, Response: Calories (optional)

X, Factor: Manufacturer (optional)

Block: (optional)

Weight: (optional Numeric)

Freq: (optional Numeric)

By: (optional)

Action

OK

Cancel

Remove

Recall

Help

Oneway

Bivariate

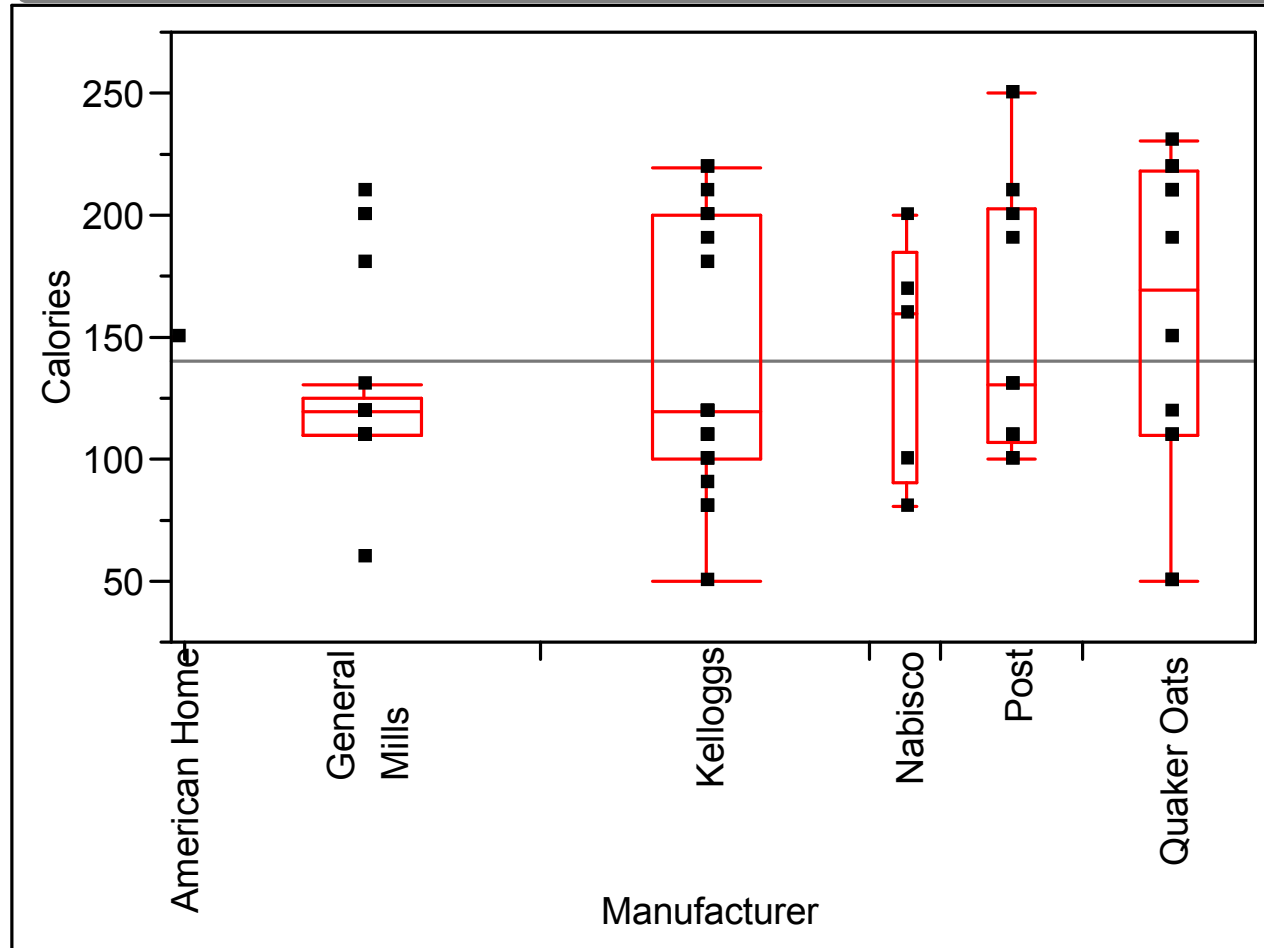
Oneway

Logistic

Contingency

Right-click Title Bar and Select Quantiles

Oneway Analysis of Calories By Manufacturer





Student's t-test

Background:

This test has evolved over the years.

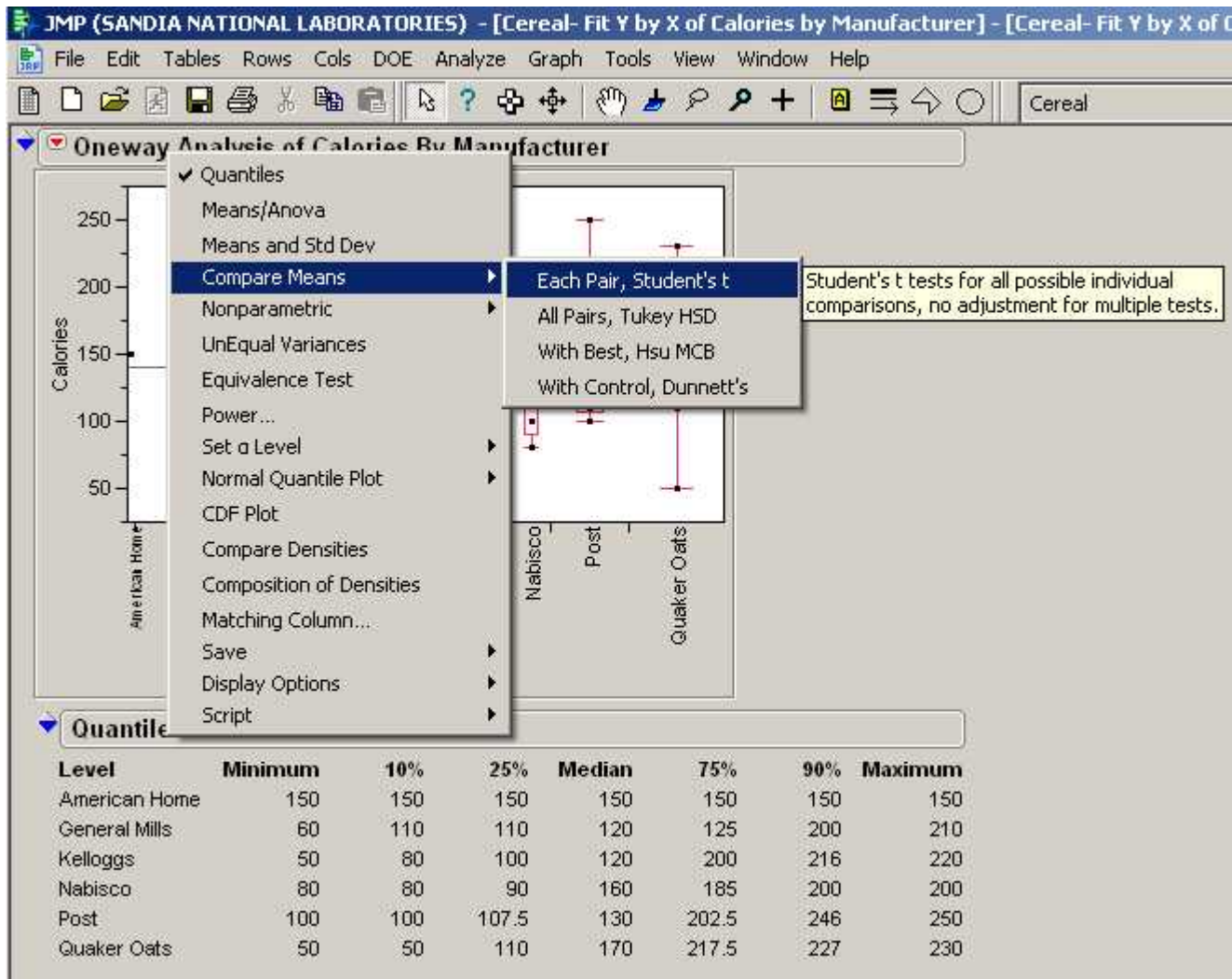
- The original test was called a z-test, which uses a normal distribution as a reference distribution. However, the z-test requires knowing the true population variances - not usually the case.
- A statistician who called himself 'Student' improved the test by basing it on the t distribution, which uses variance estimates from samples. Thus the name Student's t-test.
- The Student's t-test was adapted to work if variances in the two sample groups were different. Sometimes this approach is called the Aspin-Welch Student t-test.
- Then, F.E. Satterthwaite developed a better approximation for degrees of freedom.

Thus the full name of this improved test is the:

Aspin-Welch-Satterthwaite-Student's t-test.

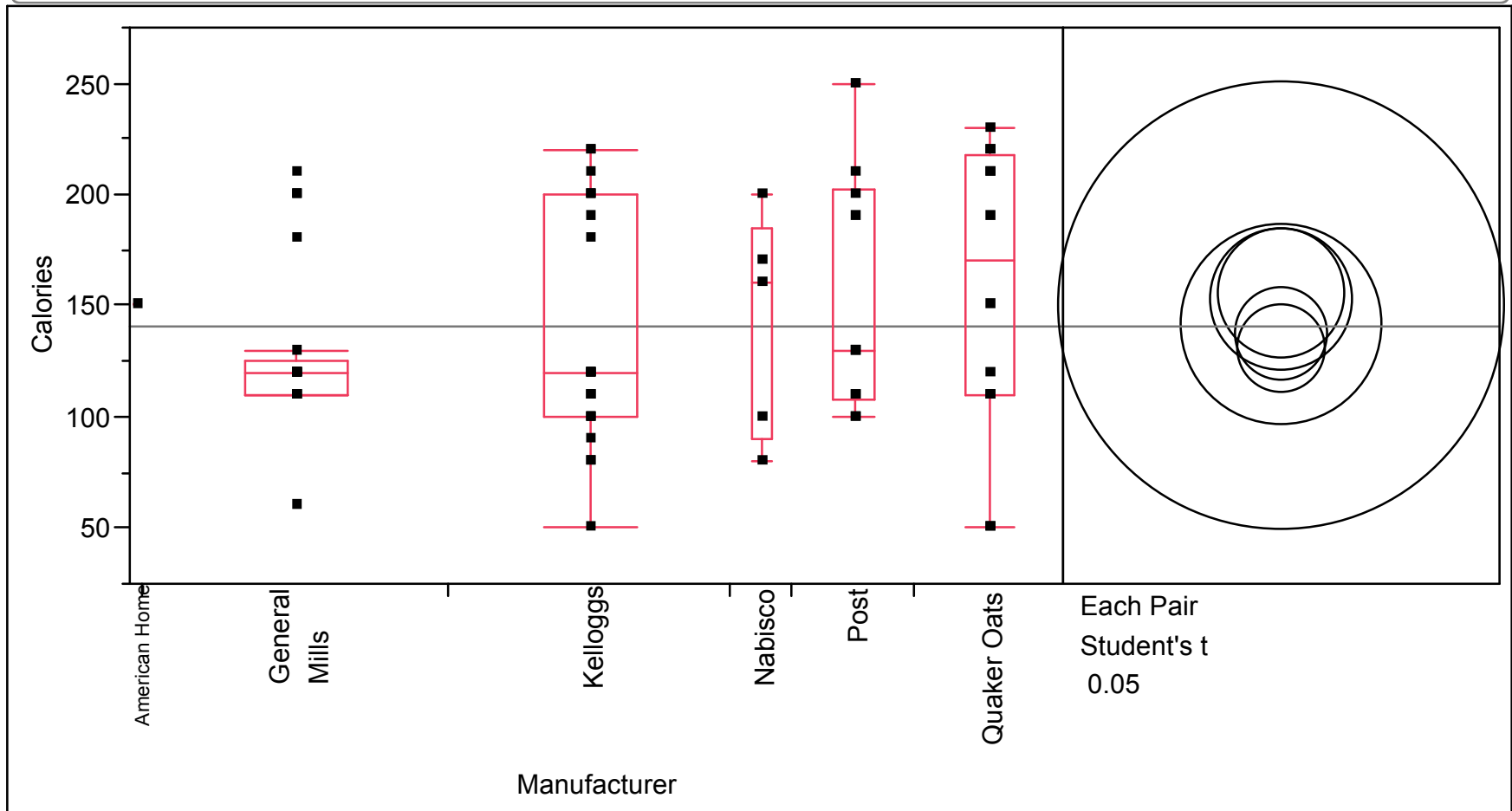
(commonly called the t-test for short).

Compare Means → Each Pair, Student's t

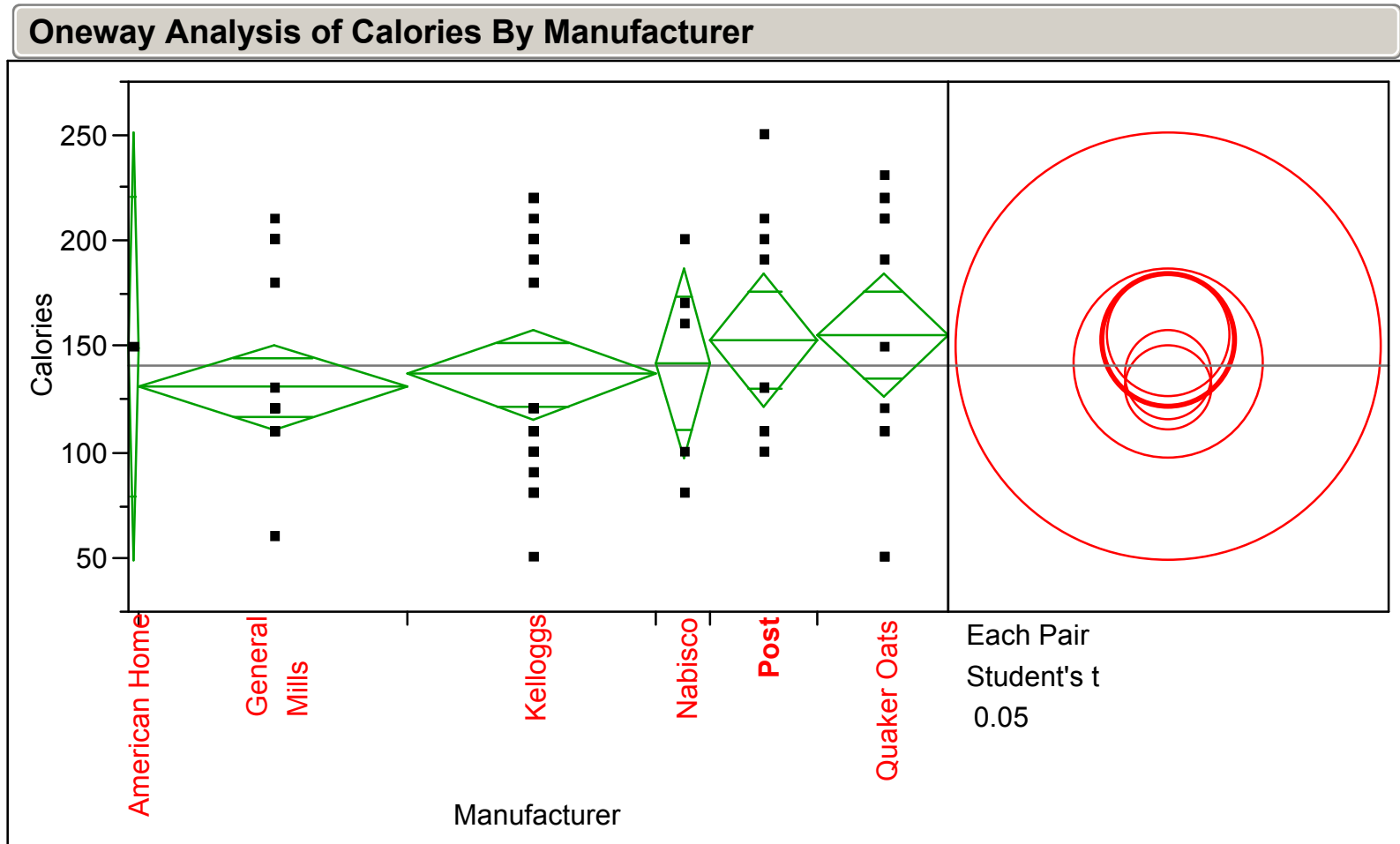


Click Any Circle – All Others Turn Red

Oneway Analysis of Calories By Manufacturer



Right-click Title Bar – Click Means/ANOVA, Uncheck Quantiles

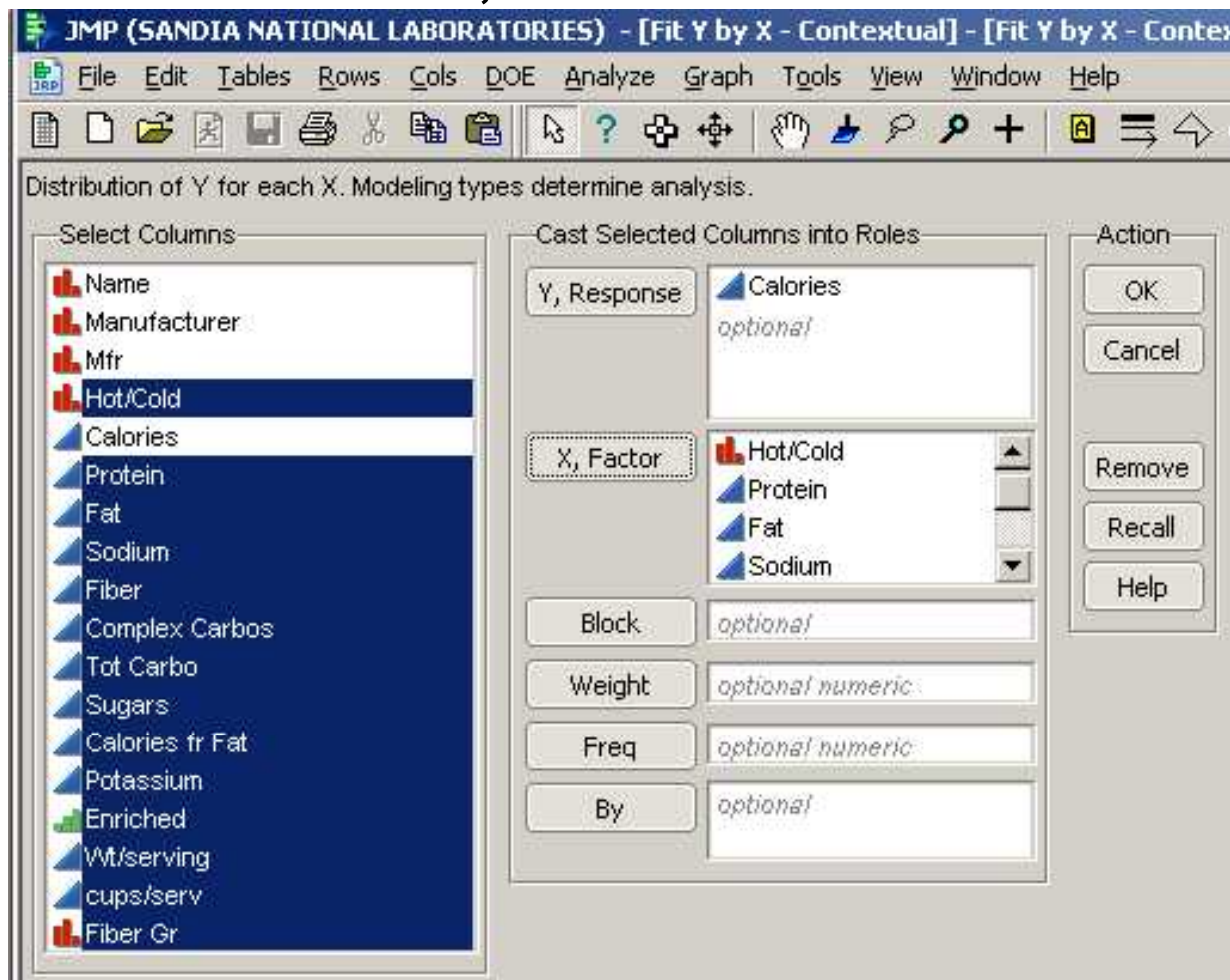


None of these are different, statistically

What Else Affects Calories?

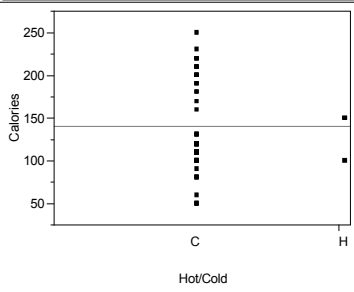
Analyze → Fit Y by X

Select Calories for Y, those below for X

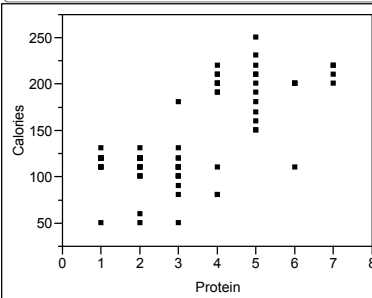


Fit Y by X Plots

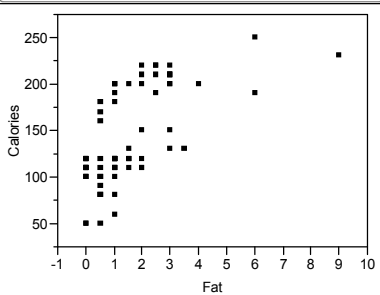
Oneway Analysis of Calories By Hot/Cold



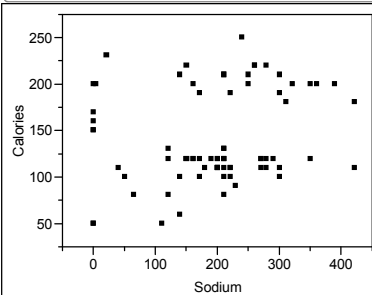
Bivariate Fit of Calories By Protein



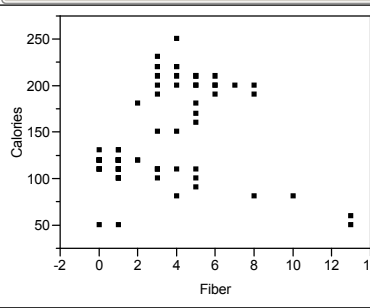
Bivariate Fit of Calories By Fat



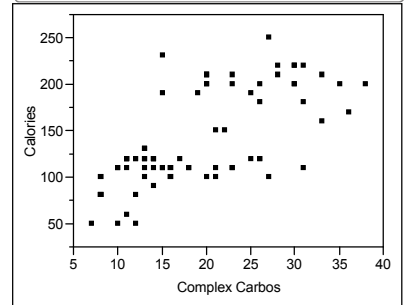
Bivariate Fit of Calories By Sodium



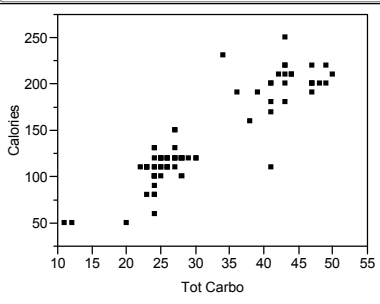
Bivariate Fit of Calories By Fiber



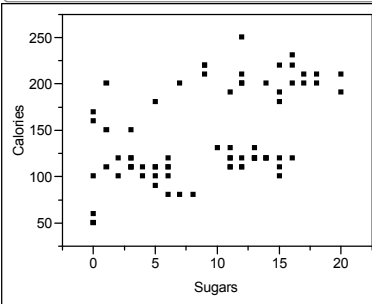
Bivariate Fit of Calories By Complex Carbos



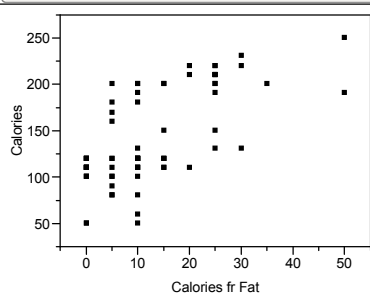
Bivariate Fit of Calories By Tot Carbo



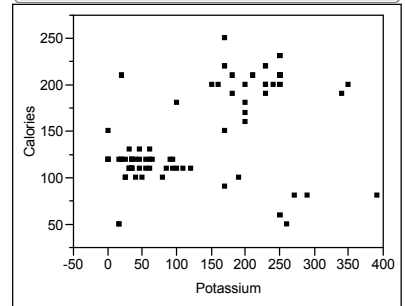
Bivariate Fit of Calories By Sugars



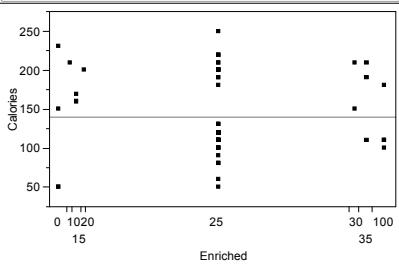
Bivariate Fit of Calories By Calories fr Fat



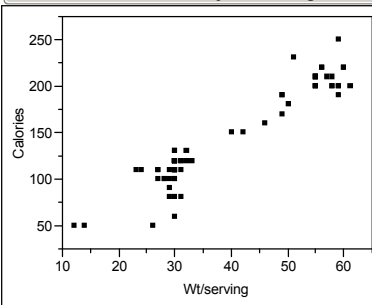
Bivariate Fit of Calories By Potassium



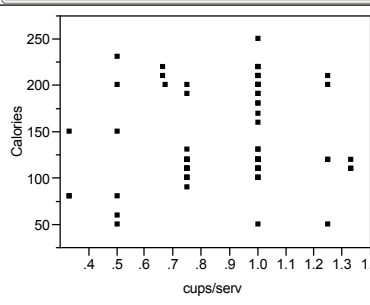
Oneway Analysis of Calories By Enriched



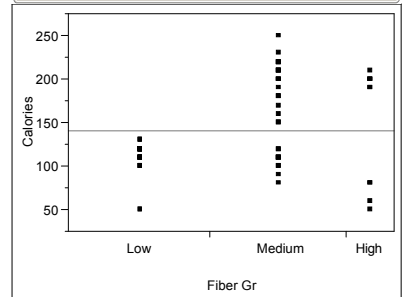
Bivariate Fit of Calories By Wt/serving



Bivariate Fit of Calories By cups/serv



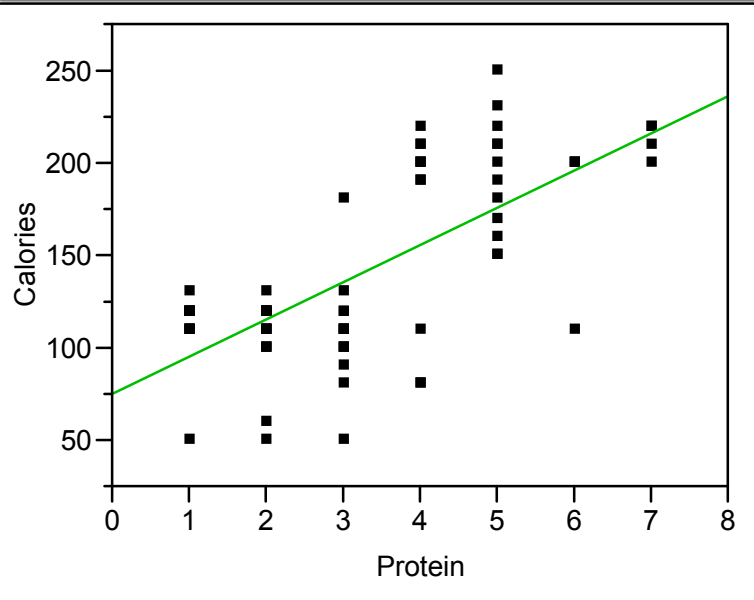
Oneway Analysis of Calories By Fiber Gr



Look Closer at Protein

Right-click title bar – Fit Line

Bivariate Fit of Calories By Protein



Linear Fit

$$\text{Calories} = 74.874142 + 20.200669 \text{ Protein}$$

It appears that each g of protein adds 20.2 calories

Summary of Fit

RSquare	0.495772
RSquare Adj	0.488958
Root Mean Square Error	35.46409
Mean of Response	140.5263
Observations (or Sum Wgts)	76

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	91509.03	91509.0	72.7589
Error	74	93069.92	1257.7	Prob > F
C. Total	75	184578.95		<.0001*

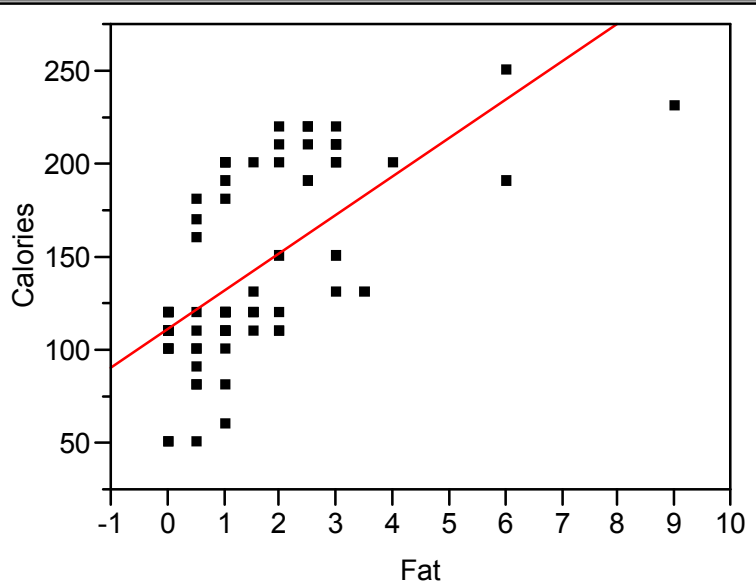
Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	74.874142	8.705646	8.60	<.0001*
Protein	20.200669	2.368223	8.53	<.0001*

Look Closer at Fat

Right-click title bar – Fit Line

Bivariate Fit of Calories By Fat



— Linear Fit

Linear Fit

$$\text{Calories} = 110.77451 + 20.555796 \text{ Fat}$$

It appears that each g of fat adds 20.6 calories

Summary of Fit

RSquare	0.4173
RSquare Adj	0.409425
Root Mean Square Error	38.12395
Mean of Response	140.5263
Observations (or Sum Wgts)	76

Analysis of Variance

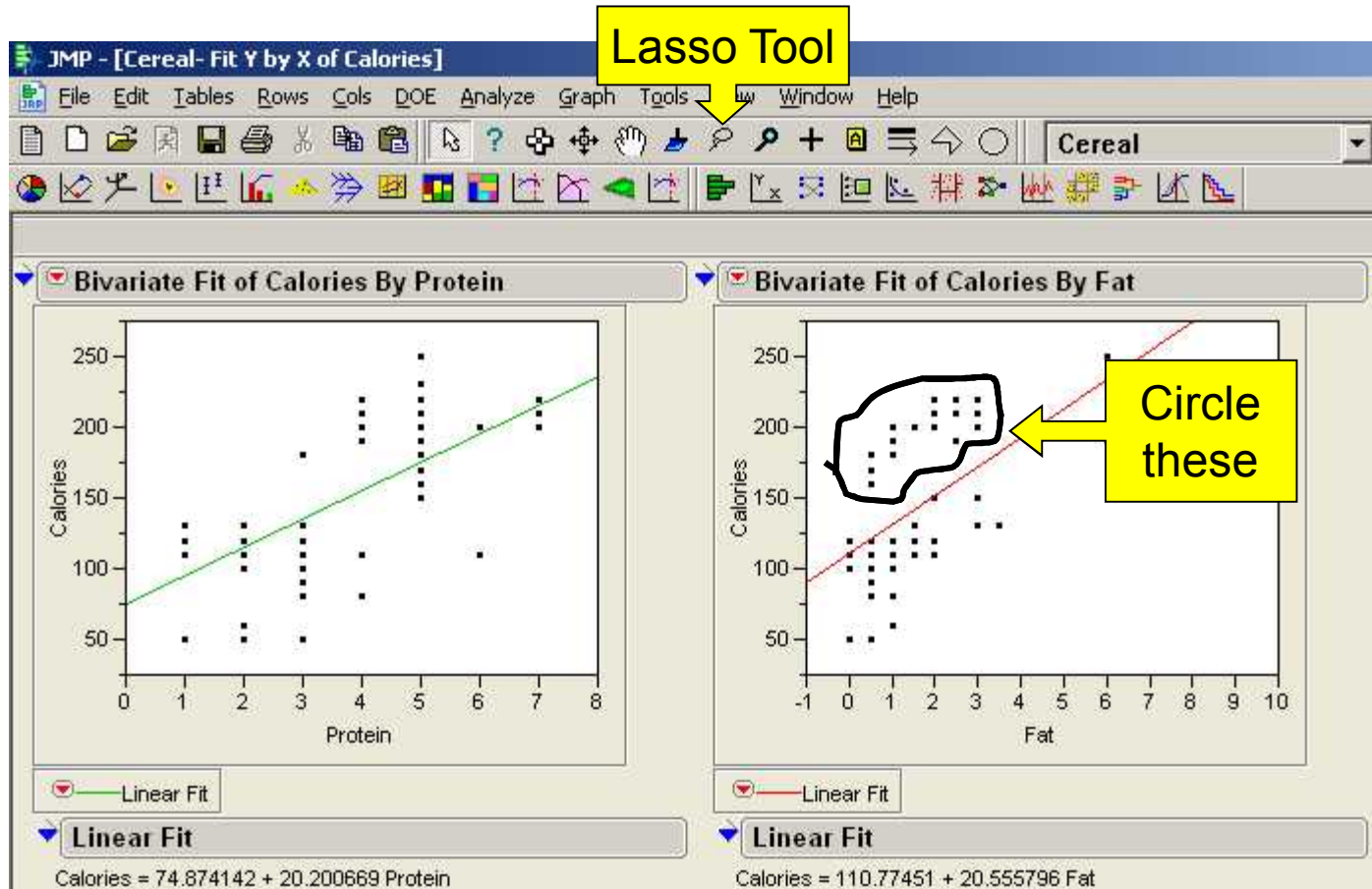
Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	77024.73	77024.7	52.9949
Error	74	107554.22	1453.4	Prob > F
C. Total	75	184578.95		<.0001*

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	110.77451	5.985571	18.51	<.0001*
Fat	20.555796	2.82369	7.28	<.0001*

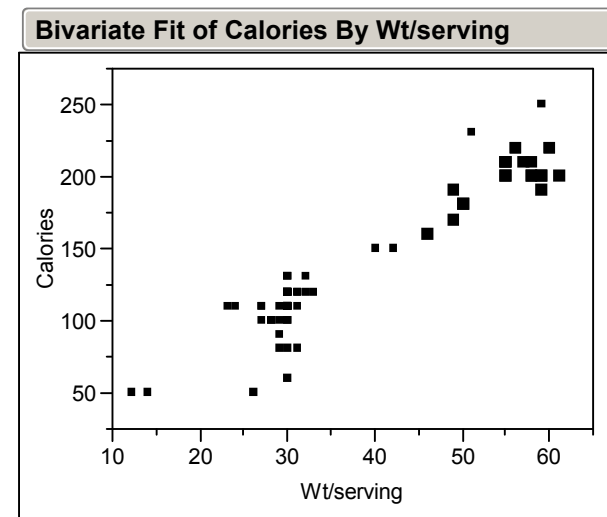
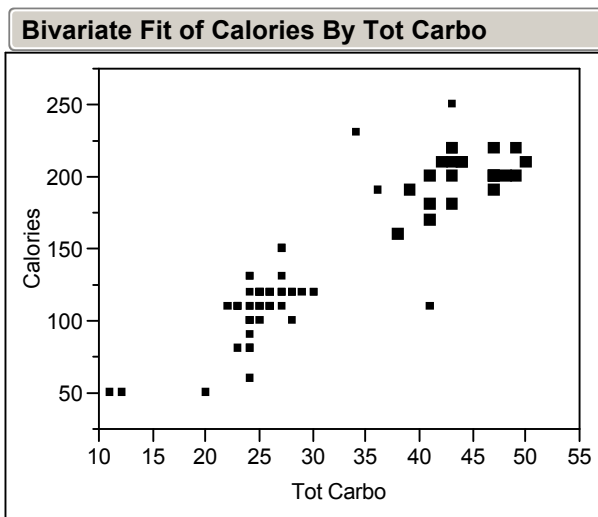
Fat Looks Bimodal

Let's investigate a little further.
Click the Lasso tool and circle the points shown.



Now Look at the Fit Y by X Plots

What do you observe?



It appears these are the cereals with higher Total Carbs or higher Wt/serving.

Let's Take a Closer Look

Analyze → Fit Y by X

Distribution of Y for each X. Modeling types determine analysis.

Select Columns

- Hot/Cold
- Calories
- Protein
- Fat
- Sodium
- Fiber
- Complex Carbos
- Tot Carbo
- Sugars
- Calories fr Fat
- Potassium
- Enriched
- Wt/serving
- cups/serv
- Fiber Gr

Cast Selected Columns into Roles

Y, Response: Tot Carbo (optional)

X, Factor: Wt/serving (optional)

Block: (optional)

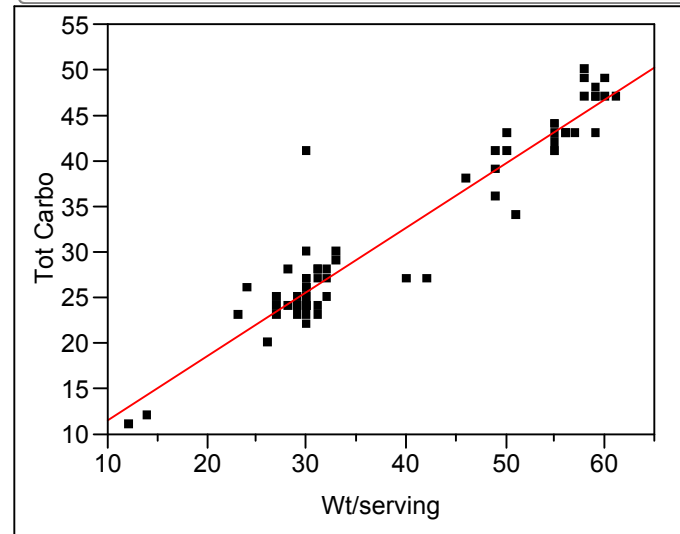
Weight: (optional Numeric)

Freq: (optional Numeric)

By: (optional)

Action: OK, Cancel, Remove, Recall, Help

Bivariate Fit of Tot Carbo By Wt/serving



Linear Fit

$$\text{Tot Carbo} = 4.3400284 + 0.7073036 \text{ Wt/serving}$$

Summary of Fit

RSquare	0.909686
RSquare Adj	0.908449
Root Mean Square Error	2.937065
Mean of Response	31.45333
Observations (or Sum Wgts)	75

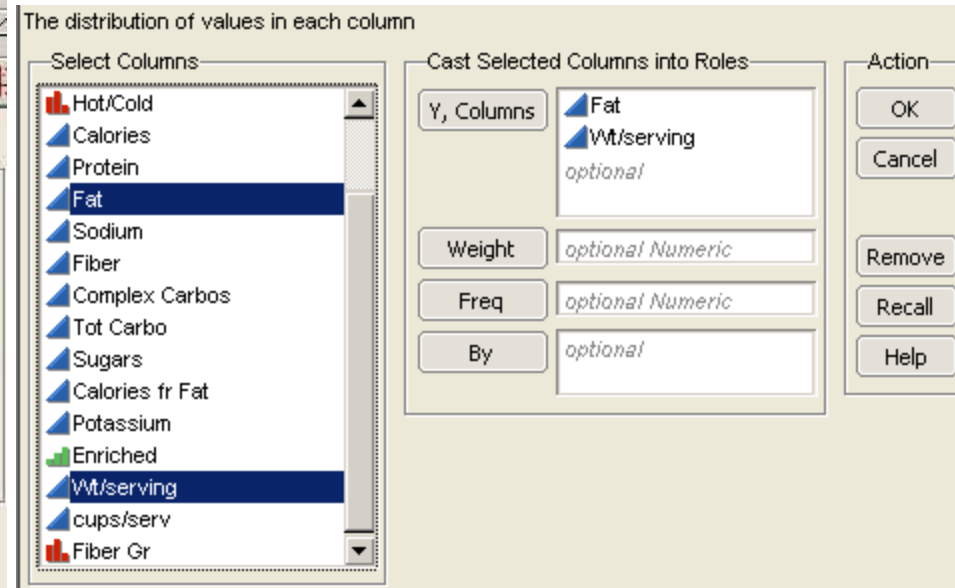
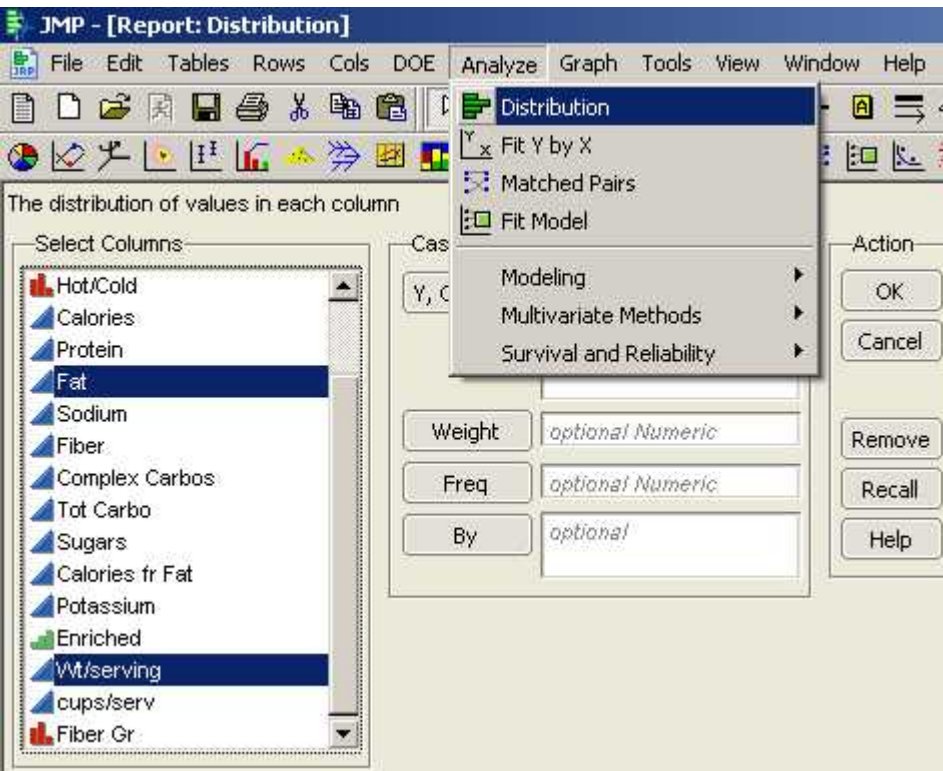
Total Carbs and Wt/serving are closely correlated

Let's Look at the Distributions of Fat and Wt/serving

Analyze → Distribution

Click Fat, Hold [Ctrl] and click Wt/serving (you may need to scroll down)

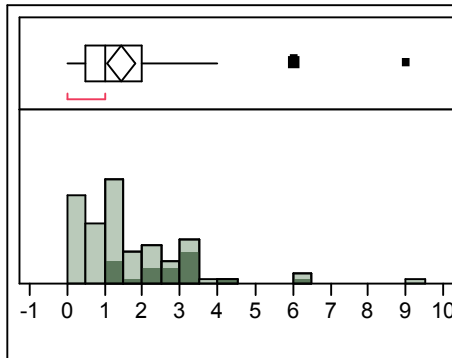
Click [Y, Columns], then [OK]



Distributions

Distributions

Fat



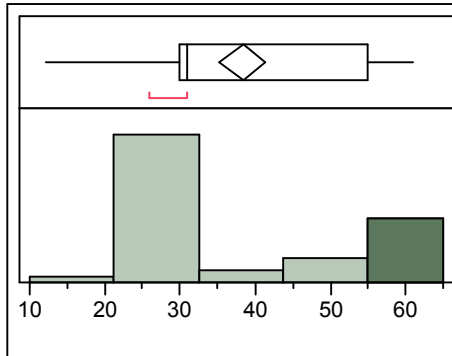
Quantiles

100.0%	maximum	9.0000
99.5%		9.0000
97.5%		6.2250
90.0%		3.0000
75.0%	quartile	2.0000
50.0%	median	1.0000
25.0%	quartile	0.5000
10.0%		0.0000
2.5%		0.0000
0.5%		0.0000
0.0%	minimum	0.0000

Moments

Mean	1.4473684
Std Dev	1.5590145
Std Err Mean	0.1788312
upper 95% Mean	1.8036185
lower 95% Mean	1.0911183
N	76

Wt/serving



Quantiles

100.0%	maximum	61.000
99.5%		61.000
97.5%		60.100
90.0%		58.400
75.0%	quartile	55.000
50.0%	median	31.000
25.0%	quartile	30.000
10.0%		27.000
2.5%		13.800
0.5%		12.000
0.0%	minimum	12.000

Moments

Mean	38.333333
Std Dev	13.089436
Std Err Mean	1.5114379
upper 95% Mean	41.344939
lower 95% Mean	35.321728
N	75

- Click the Grabber (hand) and grab the top of a bar and drag it downward.
- Click the Arrow and then click the right bar on the Wt/serving histogram. Look at the Fat histogram.



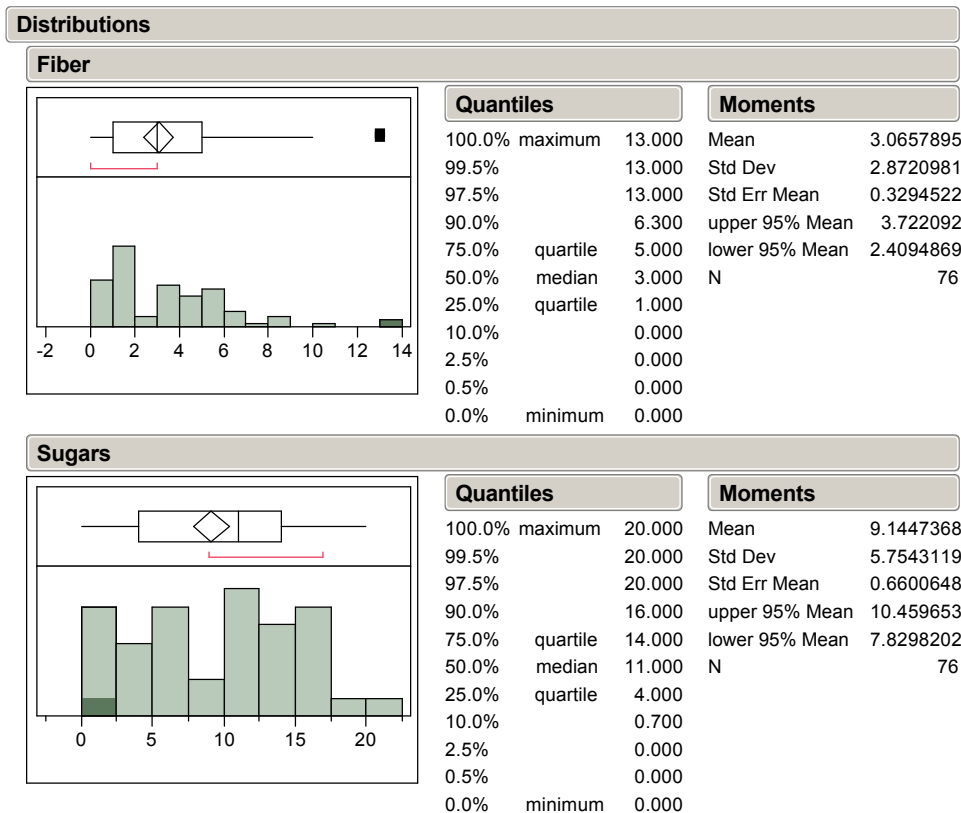
Exercise

Your doctor has told you that your triglycerides are high, and that you need to eat foods high in fiber and low in sugar. Which cereal should you choose?
Hint: Start by looking at the distributions.
(It's probably not Cocoa Puffs.)

5 minutes

Exercise – One Possible Solution

- Click Analyze → Distribution
 - Choose Fiber and Sugars.
 - Click on bars until you find a good choice. The obvious choice is the highest fiber bar. Those happen to appear in the low sugar bar.



Exercise – One Possible Solution

Go back to the data (Window → Cereal) to see that you have two choices:

- All-Bran with Extra Fiber
- Fiber One

Both have 13g fiber and 0g sugar

Third option:

Eat the box. 48g fiber, 0g sugar



JMP - [Cereal]

	Name	
1	100% Bran	Nat
2	100% Nat. Bran Oats & Honey	Qu
3	100% Nat. Low Fat Granola w raisins	Qu
4	All-Bran	Kel
5	All-Bran with Extra Fiber	Kel
6	Almond Crunch w Raisins	Kel
7	Apple Cinnamon Cheerios	Ger
8	Apple Jacks	Kel
9	Banana Nut Crunch	Pos
10	Basic 4	Ger
11	Bran Buds	Kel
12	Bran Flakes	Pos
13	Cap'n Crunch	Qu
14	Cheerios	Ger
15	Cinnamon Toast Crunch	Ger
16	Cocoa Puffs	Ger
17	Complete Oat Bran	Kel
18	Complete Wheat Bran	Kel
19	Corn Chex	Ger
20	Corn Flakes	Kel
21	Corn Pops	Kel
22	Cracklin' Oat Bran	Kel
23	Cream of Wheat (Instant)	Nat
24	Crispix	Kel
25	Fiber One	Ger
26	Franken Berry	Ger
27	French Toast Crisp	Ger

Columns (18/1)

- Name
- Manufacturer +
- Mfr
- Hot/Cold
- Calories
- Protein
- Fat
- Sodium
- Fiber
- Complex Carbo +
- Tot Carbo
- Sugars
- Calories fr Fat
- Potassium
- Enriched
- Wt/serving
- cups/serv
- Fiber Gr +

Rows

All rows 76

Which Would Be Your Worst Choice?

- Clear Row States
- Go back to distributions
- Select the lowest Fiber bar in the histogram
- Hold [Ctrl] and click the sugar bars less than 15.
- Go back to the Data and see that your worst choice is ...
- Golden Crisp (0g fiber, 15g sugar)

Analyze → Fit Model

Calories in Y, all below that in X

The screenshot displays the JMP - [Fit Model] window. The left sidebar shows the project structure with 'Cereal' and 'Fit Model' selected. The main window is titled 'Model Specification' and contains the following sections:

- Select Columns:** A list of variables including Name, Manufacturer, Mfr, Hot/Cold, Calories, Protein, Fat, Sodium, Fiber, Complex Carbos, Tot Carbo, Sugars, Calories fr Fat, Potassium, Enriched, Wt/serving, cups/serv, and Fiber Gr.
- Pick Role Variables:** A section for assigning roles to variables. 'Y' is assigned to 'Calories'. Other roles like Weight, Freq, and By are currently empty.
- Personality:** A dropdown menu set to 'Stepwise'.
- Buttons:** 'Help', 'Run Model', and 'Remove' buttons are present.
- Construct Model Effects:** A section for building the model. It includes 'Add', 'Cross', and 'Nest' buttons, a 'Macros' dropdown, and a list of selected effects: Protein, Fat, Sodium, Fiber, Complex Carbos, Tot Carbo, Sugars, Calories fr Fat, Potassium, and Enriched.
- Options:** 'Degree' is set to 2, 'Attributes' and 'Transform' are set to red, and 'No Intercept' is unchecked.

Select Personality → Stepwise, Run Model Change Direction to Mixed

Stepwise]

Rows Cols DOE Analyze Graph Tools View Window Help

Cereal

Stepwise Fit

Response: Calories

Stepwise Regression Control

Prob to Enter: 0.250 Enter All
Prob to Leave: 0.100 Remove All

Direction: Forward
Rules: Backward
Mixed

Go [S] e Model

1 rows not used due to excluded rows or missing values.

Current Estimates

	SSE	DFE	MSE	RSquare	RSquare Adj	Cp	AIC
	184152	74	2488.5405	0.0000	0.0000	5920.3953	587.4522

Lock	Entered	Parameter	Estimate	nDF	SS	"F Ratio"	"Prob>F"
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Intercept	140.8	1	0	0.000	1.0000
<input type="checkbox"/>	<input type="checkbox"/>	Protein	0	1	91104.86	71.476	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	Fat	0	1	77629.82	53.200	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	Sodium	0	1	7229.601	2.983	0.0884
<input type="checkbox"/>	<input type="checkbox"/>	Fiber	0	1	6803.278	2.800	0.0985
<input type="checkbox"/>	<input type="checkbox"/>	Complex Carbo	0	1	82335.46	59.033	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	Tot Carbo	0	1	151698.7	341.229	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	Sugars	0	1	48260.7	25.925	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	Calories fr Fat	0	1	83468.88	60.519	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	Potassium	0	1	36211.09	17.868	0.0001
<input type="checkbox"/>	<input type="checkbox"/>	Enriched{0&10&15&20&25&30&35-100}	0	1	1893.429	0.758	0.3867
<input type="checkbox"/>	<input type="checkbox"/>	Enriched{0&10&15&20&25-30&35}	0	2	7358.154	1.498	0.2304
<input type="checkbox"/>	<input type="checkbox"/>	Enriched{0-10&15&20&25}	0	3	9011.016	1.218	0.3097
<input type="checkbox"/>	<input type="checkbox"/>	Enriched{10&15&20-25}	0	4	17304.63	1.815	0.1356
<input type="checkbox"/>	<input type="checkbox"/>	Enriched{10-15&20}	0	5	18137.96	1.508	0.1988
<input type="checkbox"/>	<input type="checkbox"/>	Enriched{15-20}	0	6	18954.63	1.300	0.2689
<input type="checkbox"/>	<input type="checkbox"/>	Enriched{30-35}	0	3	7478.154	1.002	0.3972
<input type="checkbox"/>	<input type="checkbox"/>	Wt/serving	0	1	166749.4	699.477	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	cups/serv	0	1	2503.25	1.006	0.3192
<input type="checkbox"/>	<input type="checkbox"/>	Fiber Gr{ Low-High& Medium}	0	1	48736.88	26.273	0.0000
<input type="checkbox"/>	<input type="checkbox"/>	Fiber Gr{High- Medium}	0	2	50823.59	13.723	0.0000

Step History



Stepwise Regression Control Panel

The Stepwise Regression Control Panel (Control Panel for short) has editable areas, buttons and a popup menu. You use these dialog features to limit regressor effect probabilities, determine the method of selecting effects, begin or stop the selection process, and create a model.

You use the Control Panel as follows:

Prob to Enter

is the significance probability that must be attributed to a regressor term for it to be considered as a forward step and entered into the model. Click the field to enter a value.

Prob to Leave

is the significance probability that must be attributed to a regressor term in order for it to be considered as a backward step and removed from the model. Click the field to enter a value.

Direction

accesses the popup menu shown here, which lets you choose how you want variables to enter the regression equation.



Stepwise Regression Control Panel

Forward brings in the regressor that most improves the fit, given that term is significant at the level specified by **Prob to Enter**.

Backward removes the regressor that affects the fit the least, given that term is not significant at the level specified in **Prob to Leave**.

Mixed alternates the forward and backward steps. It includes the most significant term that satisfies **Prob to Enter** and removes the least significant term satisfying **Prob to Leave**. It continues removing terms until the remaining terms are significant and then it changes to the forward direction.



Stepwise Regression Control Panel

Buttons on the controls panel let you control the stepwise processing:

Go

starts the selection process. The process continues to run in the background until the model is finished.

Stop

stops the background selection process.

Step

stops after each step of the stepwise process

Enter All

enters all unlocked terms into the model.

Remove All

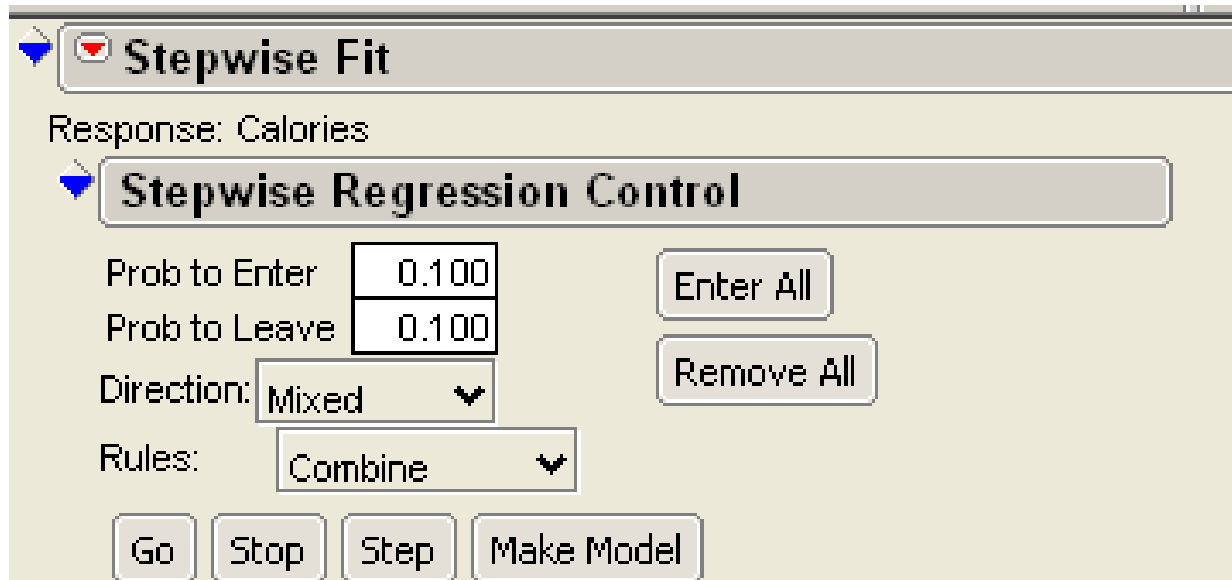
removes all terms from the model.

Make Model

forms a model for the Model Specification Dialog from the model currently showing in the Current Estimates table. In cases where there are nominal or ordinal terms, Make Model can create new data table columns to contain terms that are needed for the model.

A Few Changes First

- You already changed direction to Mixed.
- Change “Prob to Enter” and “Prob to Leave” to 0.100



Stepwise Fit

Response: Calories

Stepwise Regression Control

Prob to Enter: 0.100

Prob to Leave: 0.100

Direction: Mixed

Rules: Combine

Enter All

Remove All

Go Stop Step Make Model

Click Step and Watch Factors Get Added to the Model

JMP - [Cereal- Fit Stepwise]

File Edit Tables Rows Cols DOE Analyze Graph Tools View Window Help

Windows: JMP Starter Cereal Fit Model Cereal- Fit

Stepwise Fit

Response: Calories

Stepwise Regression Control

Prob to Enter: 0.050 Enter All
Prob to Leave: 0.050 Remove All
Direction: Mixed
Rules: Combine
Go Stop **Step** Make Model

1 rows not used due to excluded rows or missing values.

Current Estimates

	SSE	DFE	MSE	RSquare	RSquare Adj	Cp	AIC
	2361.5019	69	34.224666	0.9872	0.9862	13.857241	270.7174
Lock	Entered	Parameter	Estimate	nDF	SS	"F Ratio"	"Prob>F"
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Intercept	-5.2762836	1	0	0.000	1.0000
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Protein	0	1	108.0134	3.259	0.0754
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Fat	7.29376175	1	5262.903	153.775	0.0000
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Sodium	0	1	7.077795	0.204	0.6526
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Fiber	-3.6657725	1	5463.813	159.645	0.0000
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Complex Carbo	0	1	0.061301	0.002	0.9666
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Tot Carbo	0.69015019	1	211.2171	6.171	0.0154
<input type="checkbox"/>	<input type="checkbox"/>	Sugars	0	1	0.061301	0.002	0.9666
<input type="checkbox"/>	<input type="checkbox"/>	Calories fr Fat	0	1	47.93983	1.409	0.2393
<input type="checkbox"/>	<input type="checkbox"/>	Potassium	0	1	44.55929	1.308	0.2568
<input type="checkbox"/>	<input type="checkbox"/>	Enriched(0&10&15&20&25&30&35-100)	0	1	47.83582	1.406	0.2399
<input type="checkbox"/>	<input type="checkbox"/>	Enriched(0&10&15&20&25-30&35)	0	2	48.87546	0.708	0.4963
<input type="checkbox"/>	<input type="checkbox"/>	Enriched(0-10&15&20&25)	0	3	53.50218	0.510	0.6768
<input type="checkbox"/>	<input type="checkbox"/>	Enriched(10&15&20-25)	0	4	97.28987	0.698	0.5959
<input type="checkbox"/>	<input type="checkbox"/>	Enriched(10-15&20)	0	5	102.9973	0.584	0.7123
<input type="checkbox"/>	<input type="checkbox"/>	Enriched(15-20)	0	6	103.0703	0.479	0.8213
<input type="checkbox"/>	<input type="checkbox"/>	Enriched(30-35)	0	3	107.5512	1.050	0.3765
<input type="checkbox"/>	<input checked="" type="checkbox"/>	Wt/serving	3.00809054	1	6234.103	182.152	0.0000
<input type="checkbox"/>	<input checked="" type="checkbox"/>	cups/serv	11.0815235	1	335.1154	9.792	0.0026
<input type="checkbox"/>	<input type="checkbox"/>	Fiber Gr(Low-High& Medium)	0	1	6.129258	0.177	0.6753
<input type="checkbox"/>	<input type="checkbox"/>	Fiber Gr(High- Medium)	0	2	164.2459	2.504	0.0894

Click Make Model

Model]

Rows Cols DOE Analyze Graph Tools View Window Help

Model Specification

Select Columns

- Name
- Manufacturer
- Mfr
- Hot/Cold
- Calories
- Protein
- Fat
- Sodium
- Fiber
- Complex Carbos
- Tot Carbo
- Sugars
- Calories fr Fat
- Potassium
- Enriched
- Wt/serving
- cups/serv
- Fiber Gr

Pick Role Variables

y: ☐ Calories *optional*

Weight:

Freq:

By:

Personality: Standard Least Squares

Emphasis: Effect Leverage

Help Run Model Remove

Construct Model Effects

Add Cross Nest Macros

Degree: 2

Attributes: ☒

Transform: ☒

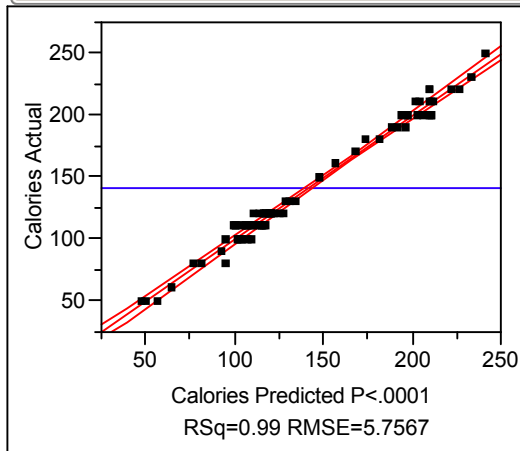
☐ No Intercept

Fat
Fiber
Tot Carbo
Wt/serving
cups/serv

Click Run Model

Whole Model

Actual by Predicted Plot



Summary of Fit

RSquare	0.987763
RSquare Adj	0.986683
Root Mean Square Error	5.756695
Mean of Response	140.8
Observations (or Sum Wgts)	75

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	6	181898.51	30316.4	914.8112
Error	68	2253.49	33.1	Prob > F
C. Total	74	184152.00		<.0001*

Lack Of Fit

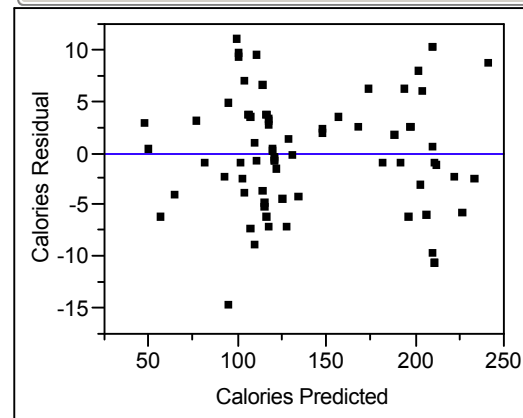
Source	DF	Sum of Squares	Mean Square	F Ratio
Lack Of Fit	65	2253.4886	34.6691	.
Pure Error	3	0.0000	0.0000	Prob > F
Total Error	68	2253.4886		.

Max RSq
1.0000

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	-4.230658	3.478264	-1.22	0.2281
Protein	1.3019557	0.721158	1.81	0.0754
Fat	7.414865	0.582652	12.73	<.0001*
Fiber	-3.874463	0.308004	-12.58	<.0001*
Tot Carbo	0.8724868	0.291431	2.99	0.0038*
Wt/serving	2.755463	0.260157	10.59	<.0001*
cups/serv	10.097625	3.527139	2.86	0.0056*

Residual by Predicted Plot



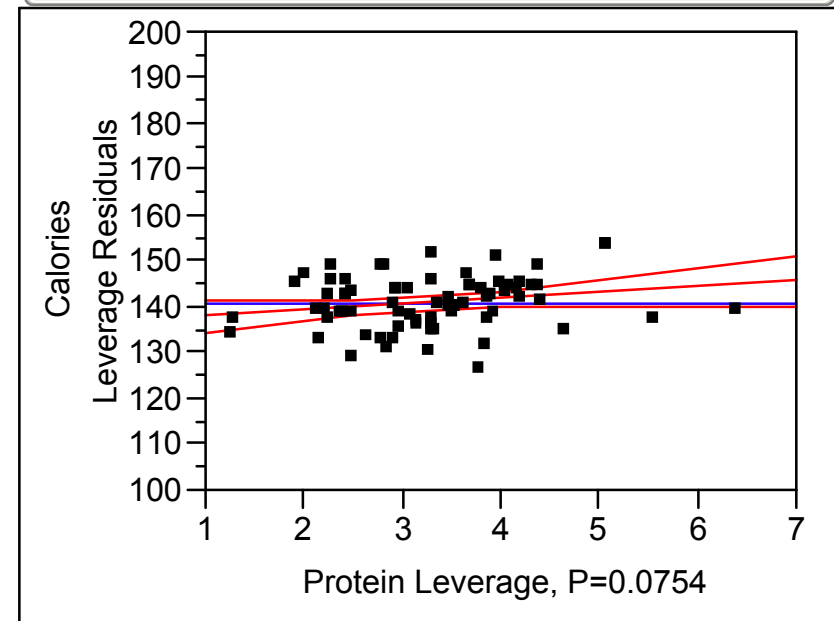
What is the Effect of Protein Now?

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	-4.230658	3.478264	-1.22	0.2281
Protein	1.3019557	0.721158	1.81	0.0754
Fat	7.414865	0.582652	12.73	<.0001*
Fiber	-3.874463	0.308004	-12.58	<.0001*
Tot Carbo	0.8724868	0.291431	2.99	0.0038*
Wt/serving	2.755463	0.260157	10.59	<.0001*
cups/serv	10.097625	3.527139	2.86	0.0056*

Protein

Leverage Plot



Recall that it was 20.2 calories per g when we just looked at calories vs. protein. Nutritionists tell us that the real number is 4 calories per g of fat.

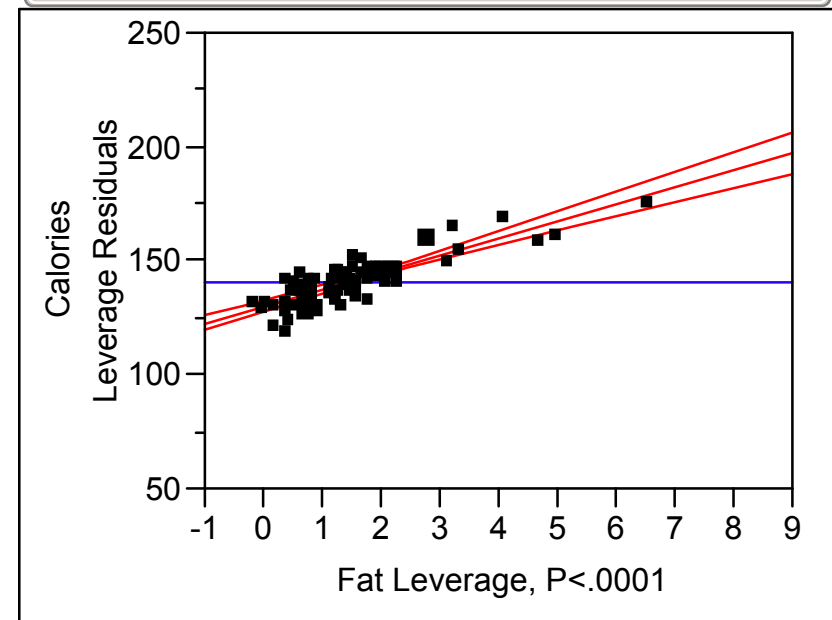
What is the Effect of Fat Now?

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	-4.230658	3.478264	-1.22	0.2281
Protein	1.3019557	0.721158	1.81	0.0754
Fat	7.414865	0.582652	12.73	<.0001*
Fiber	-3.874463	0.308004	-12.58	<.0001*
Tot Carbo	0.8724868	0.291431	2.99	0.0038*
Wt/serving	2.755463	0.260157	10.59	<.0001*
cups/serv	10.097625	3.527139	2.86	0.0056*

Fat

Leverage Plot

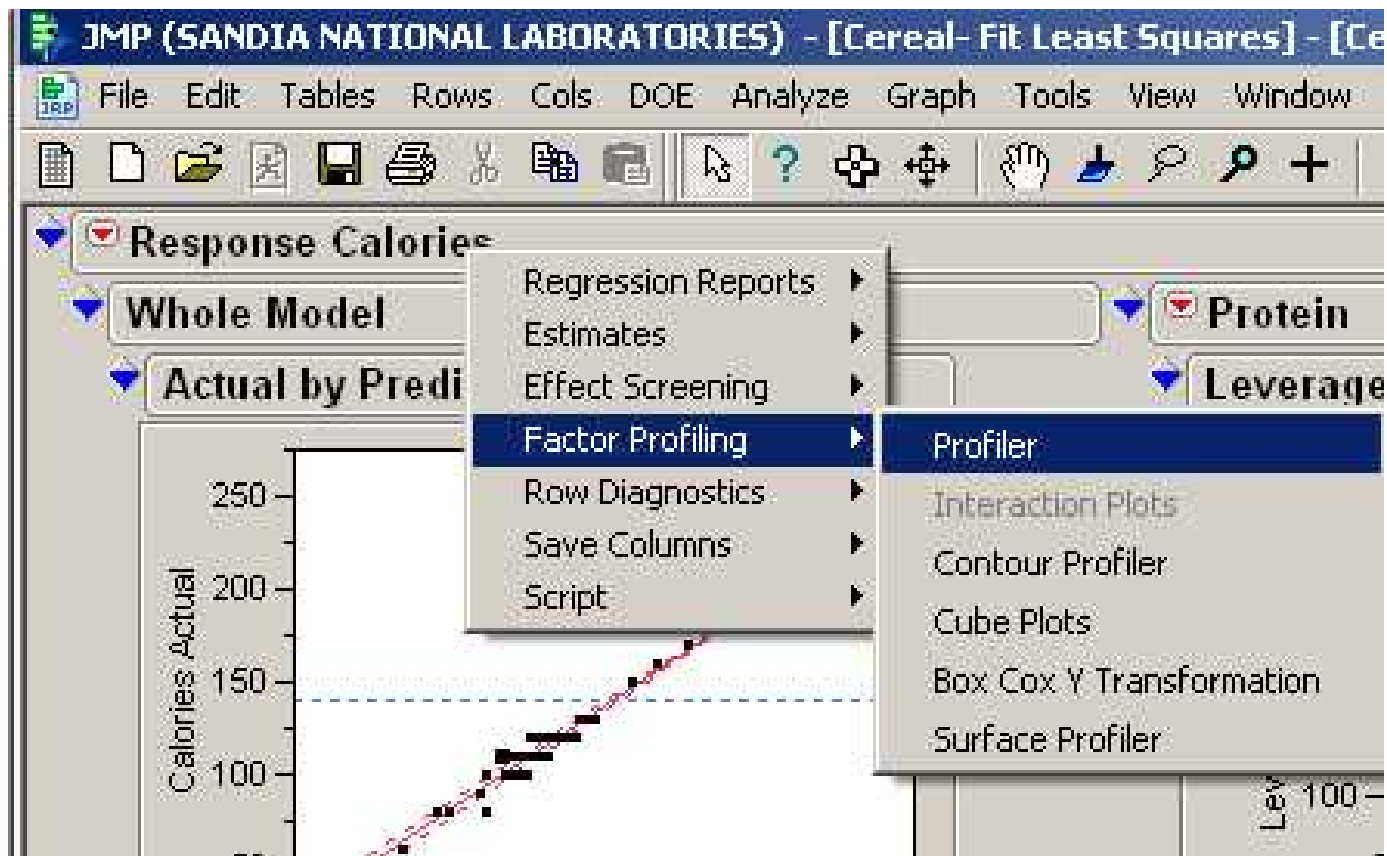


Recall that it was **20.6 calories per g** when we just looked at calories vs. fat.
Nutritionists tell us that the real number is 9 calories per g of fat.

You Need to Look at the Model!

- When we just looked at calories vs. protein, we concluded that each gram of protein adds 20.2 calories.
- When we looked at the entire model, we discovered that each gram of protein really only adds 1.3 calories!
- When we just looked at calories vs. protein, we concluded that each gram of fat adds 20.6 calories.
- When we looked at the entire model, we discovered that each gram of fat really only adds 7.4 calories!
- Looking at the entire model captured all contributing factors, and gave us coefficients closer to what we've been told by nutritionists.

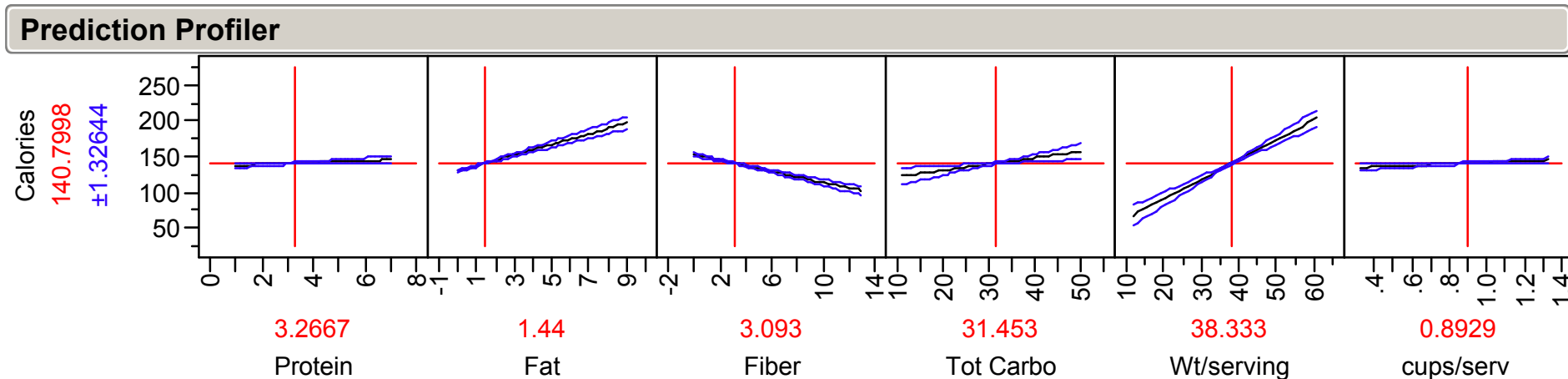
Right-click Response Calories title bar Select Factor Profiling → Profiler



Scroll Down to Prediction Profiler

Try moving the vertical lines.

Which factors could you change in order to reduce calories?





Back to the Exercise

You want to sell your house. It has the following features:

- 2000 square feet
- 0.2 acre lot
- 2 years old
- 3 bedrooms
- 3 full bathrooms

Load *House Data for Summit Tutorial.jmp* Data File

JMP (SANDIA NATIONAL LABORATORIES) - [House Data for Summit Tutorial.JMP] - [House Data for Summit Tutorial]

File Edit Tables Rows Cols DOE Analyze Graph Tools View Window Help

House Data for Summit

	SF	Lot	Age	BR	Bath	Price	Price/sf
1	1373	0.13	7	4	3	204962.96	149.281107
2	1377	0.2	1	2	3	279461.24	202.949339
3	2696	0.21	1	2	3.75	432115.58	160.28026
4	2743	0.2	11	3	2.75	291085.68	106.11946
5	1128	0.19	14	5	3.25	163331.78	144.797677
6	3721	0.16	5	3	2	417458.93	112.189984
7	3372	0.05	19	4	2.5	291889.4	86.5626928
8	1342	0.1	20	4	4	91196.94	67.9559911
9	1317	0.23	17	3	3	118951.58	90.3201063
10	2370	0.25	19	3	2.5	186523.51	78.701903
11	1645	0.18	9	5	4	277864.81	168.914778
12	2306	0.08	0	4	2.75	339135.6	147.066609
13	1356	0.23	1	2	3.25	254317.26	187.549602
14	2421	0.08	20	3	3.75	176160.96	72.7637175
15	1801	0.17	11	4	3.75	245049.51	136.063026
16	2195	0.19	17	2	3.5	195129.12	88.8970934
17	2172	0.15	13	4	4	253373.93	116.654664
18	2002	0.17	1	4	2.75	360202.51	179.921334
19	1851	0.2	11	4	3.5	261394.12	141.217785
20	2520	0.1	13	3	4	259948.18	103.15404
21	2102	0.05	14	2	3.75	177637.02	84.5085728
22	2533	0.08	11	3	3.25	285993.76	112.90713
23	2983	0.11	0	4	2	442720.07	148.414371
24	3249	0.23	2	5	3	468637.93	144.240668
25	1585	0.2	19	2	3.75	135501.42	85.4898549
26	1560	0.24	0	5	3	360846.46	231.311833
27	3319	0.14	2	2	2.75	442828.35	133.422221
28	3691	0.21	10	3	3.25	450724.35	122.114427
29	1484	0.1	20	2	2	49746.35	33.5217992
30	3619	0.11	17	5	2	338789.82	93.6142083

Columns (8/0)

- SF
- Lot
- Age
- BR
- Bath
- Price
- Price/sf
- Bogus

Rows

All rows	30
Selected	0
Excluded	0
Hidden	0
Labelled	0



Exercise:

What Will Your Listing Price Be?

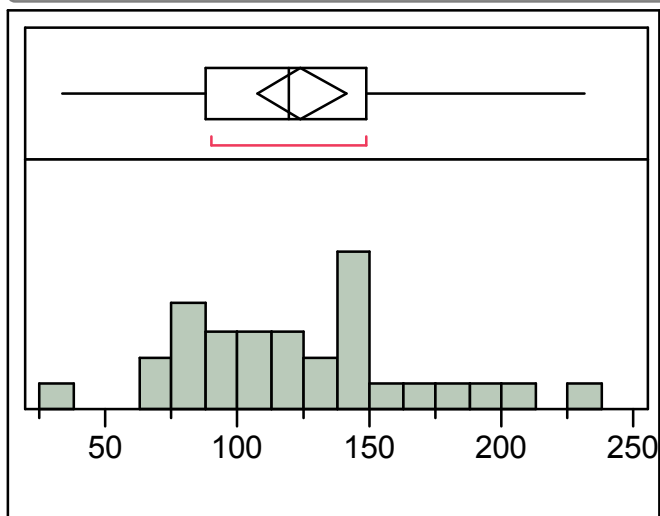
You Analyzed the distribution of Price/sf

Average = \$124.36 per SF

Therefore, $\$124.36 / \text{sf} \times 2,000 \text{ sf} = \$248,720$

Distributions

Price/sf



Quantiles

100.0%	maximum	231.31
99.5%		231.31
97.5%		231.31
90.0%		186.79
75.0%	quartile	148.63
50.0%	median	119.38
25.0%	quartile	88.31
10.0%		73.36
2.5%		33.52
0.5%		33.52
0.0%	minimum	33.52

Moments

Mean	124.36354
Std Dev	44.020787
Std Err Mean	8.0370594
upper 95% Mean	140.80117
lower 95% Mean	107.92591
N	30



Exercise:

What Will Your Listing Price Be?

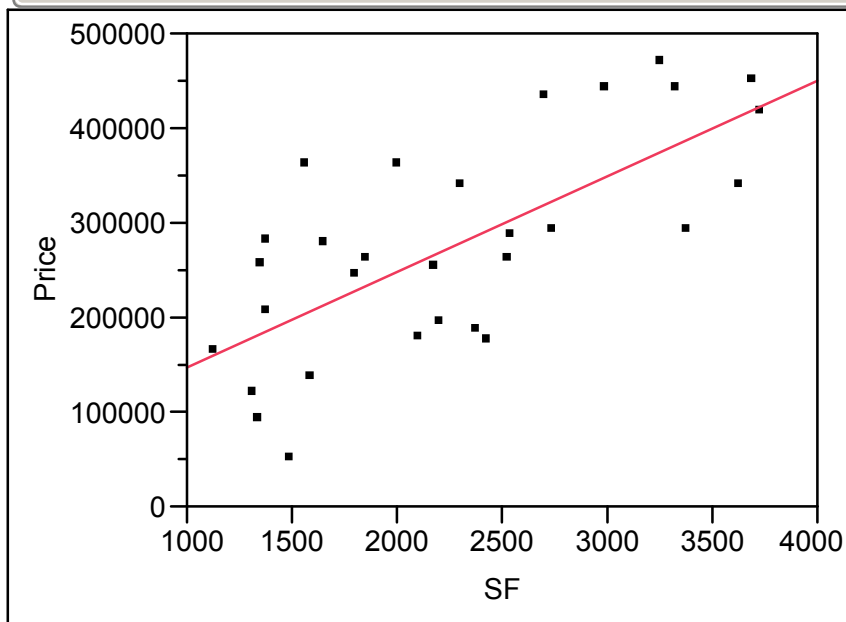
You performed a Fit Y by X for Price vs. SF

You also added a Line Fit

$$\text{Price} = \$45,962 + \$101.34 * \text{SF}$$

$$\text{Therefore, Price} = \$248,642$$

Bivariate Fit of Price By SF



— Linear Fit

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	1.8791e+11	1.879e+11	28.1954
Error	28	1.8661e+11	6.6647e+9	Prob > F
C. Total	29	3.7452e+11		<.0001*

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	45962.927	45654	1.01	0.3227
SF	101.33845	19.0847	5.31	<.0001*




Exercise:

What Will Your Listing Price Be?

- Review Student's Listing Prices previously captured on the Board

You Want to Sell Your House

- Your realtor pulls up the set of data for recent home sales in your zip code, and tells you the average selling price was \$124.36 per square foot. (Data file provided.)
- Your realtor breaks out the calculator and tells you your home is worth $\$124.36/\text{ft}^2 \times 2,000 \text{ ft}^2 = \$248,720$.
- Your realtor tells you to list your house for \$260,000. “That leaves a little room for negotiating,” they explain.
- You’re just about to sign the listing paperwork, but you remember the Workshop from the Black Belt Summit.



Should You Listen to Your Realtor?

Exercise

- Create a model for home price, including only significant factors.
- Determine the value of your home based on the model.
- Capture the students' listing prices on the board.
- Are these much different than what your Realtor recommended?



Exercise Time

15 minutes

Solutions

Create a model for home price, including only significant factors.

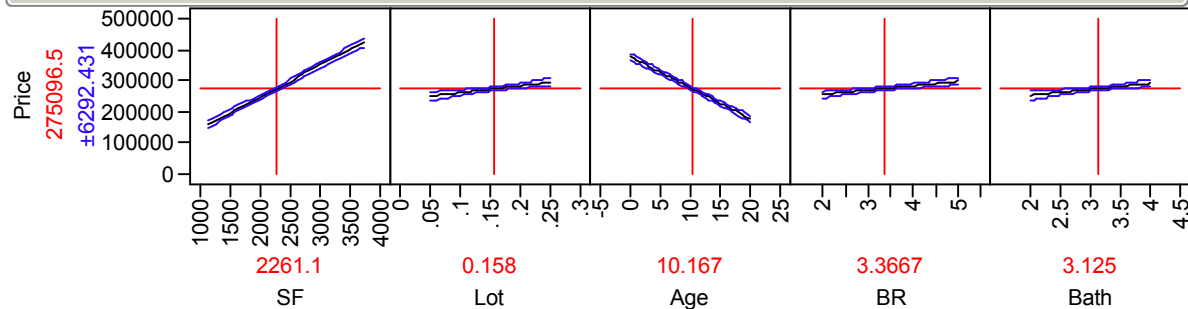
Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	3242.1849	28037.38	0.12	0.9089
SF	100.26837	4.377378	22.91	<.0001*
Lot	228519.2	55577.05	4.11	0.0004*
Age	-9954.605	456.565	-21.80	<.0001*
BR	14362.019	2925.965	4.91	<.0001*
Bath	19803.935	5364.676	3.69	0.0011*

Prediction Expression

+14362.01919530077*BR
+19803.9349350072*Bath

Prediction Profiler






Solutions

Determine the value of your home based on the model.

$$3242.18 + 100.27 * (2000) + 228519.20 * (0.2) - 9954.60 * (2) + 14362.02 * (3) + 19803.93 * (3) =$$
$$\$332,075$$

Should you listen to your Realtor and list your house for \$260,000?



What Does the Model Tell You?

- Which factors are statistically significant?
- What are the coefficients for these factors?
- In particular, what is the coefficient for \$/square foot?



A Word of Caution

- Three types of variables
 - Continuous
 - Time
 - Distance
 - Ordinal
 - Character data with an order (poor, fair, good, better, best)
 - Numerical data with unequal spacing (4 = strongly agree, 3 = agree, 2 = disagree, 1 = strongly disagree)
 - Nominal
 - Character data with no specific order (green, blue, yellow)
 - Numerical data with no specific order (NASCAR car #)
- Should BR and Bath be treated as continuous variables?
- What if we had treated them as Ordinal Variables?

Treating BR and Bath as Ordinal

- If Time Permits, change BR and Bath to Ordinal and redo the analysis

House Data for Sum

	SF	Lot	Age	BR
1	1373	0.13	7	4
2	1377	0.2	1	2
3	2696	0.21	1	2
4	2743	0.2	11	3
5	1128	0.19	14	5
6	3721	0.16	5	3
7	3372	0.05	19	4
8	1342	0.1	20	4
9	1317	0.23	17	3
10	2370	0.25	19	3
11	1645	0.18	9	5
12	2306	0.08	0	4
13	1356	0.23	1	2
14	2421	0.08	20	3
15	1801	0.17	11	4
16	2195	0.19	17	2
17	2470	0.15	10	4

Columns (7/1)

- SF
- Lot
- Age
- BR

Continuous

Ordinal

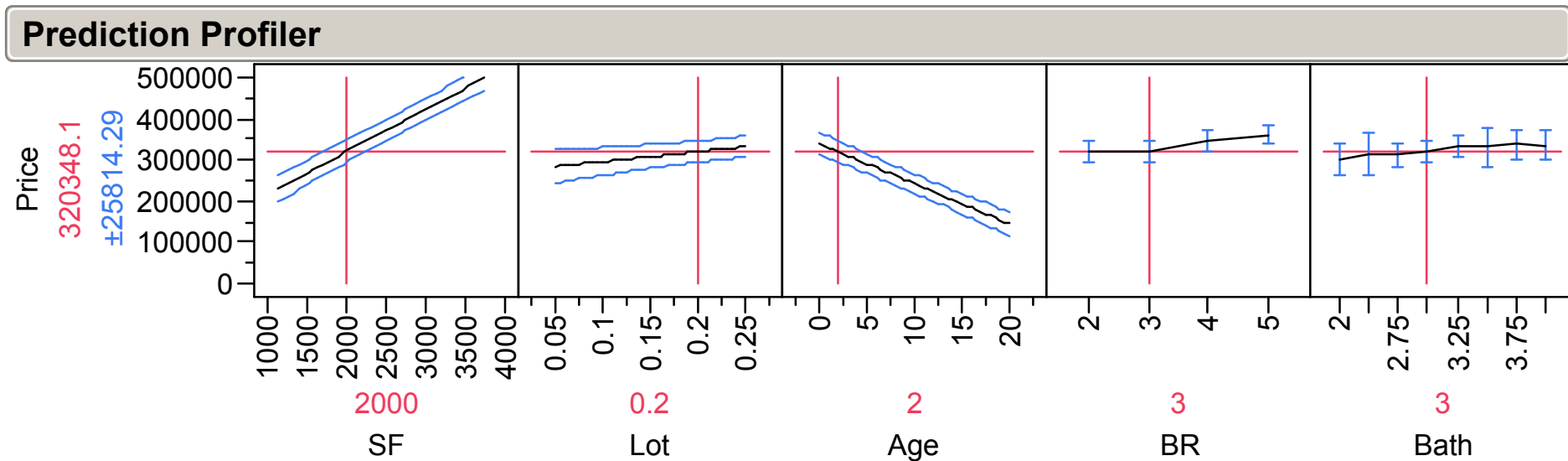
Nominal

Columns (7/10)

- SF
- Lot
- Age
- BR
- Bath
- Price
- Price/sf +



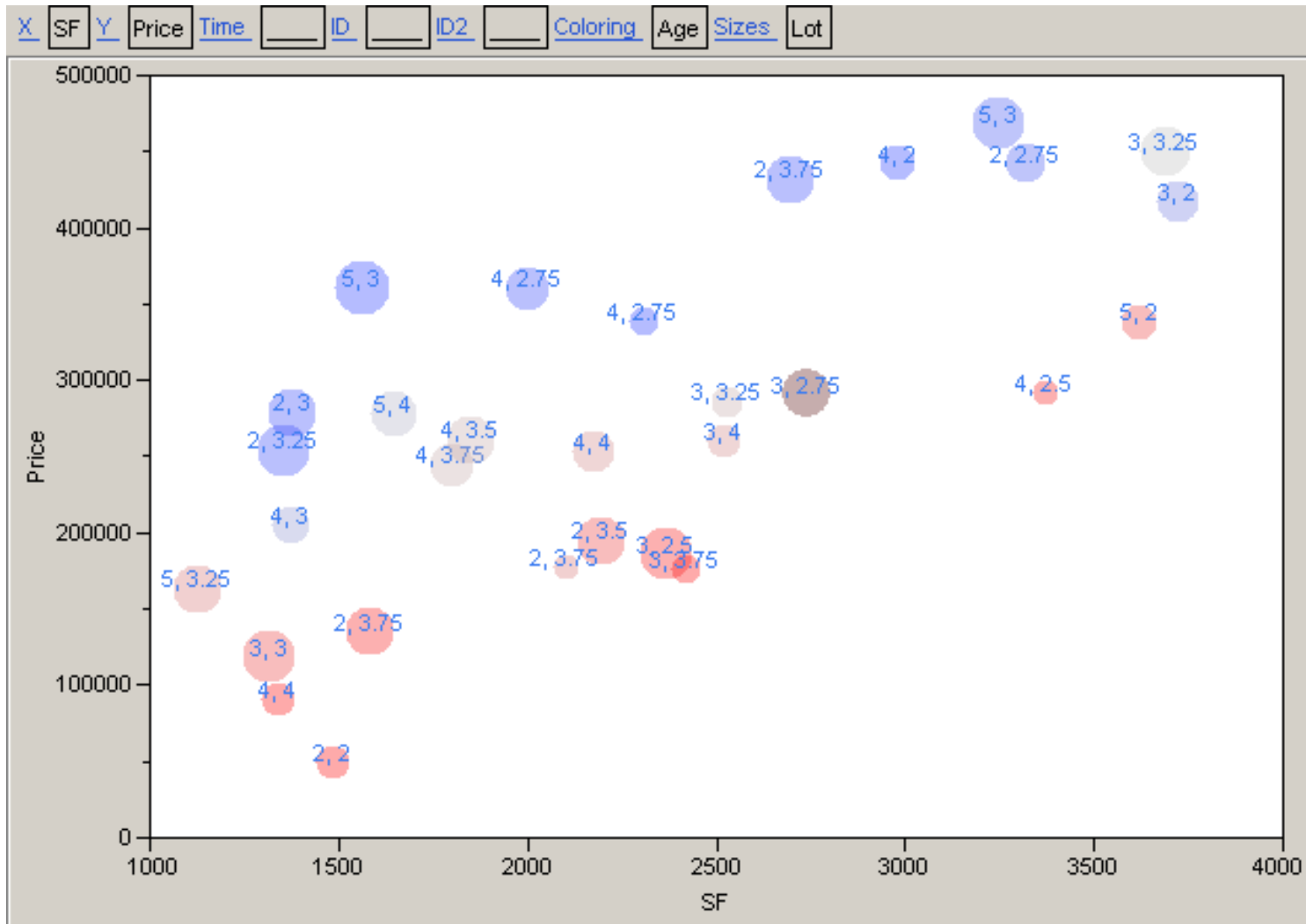
Use Prediction Profiler



What is the predicted price now?

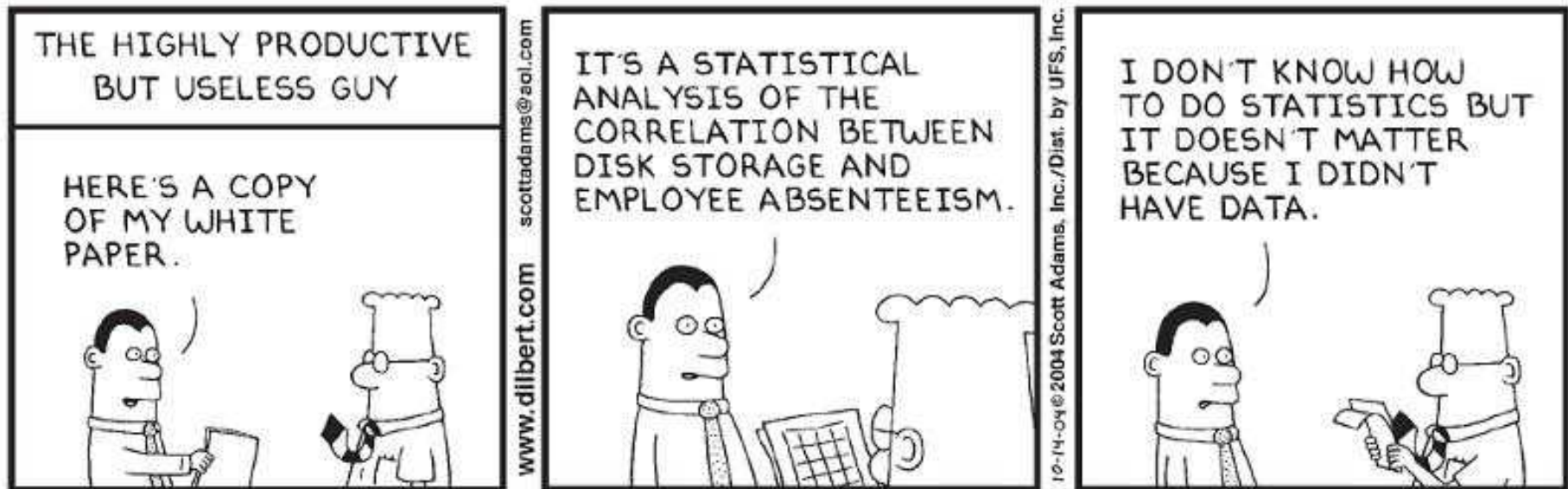
One Last Tip: Visually Display your Data!

How many dimensions are shown in this single graph?



Questions?

Dilbert





"... all models are wrong; the practical question is how wrong do they have to be to not be useful ..."

George Box and Norman Draper, Empirical Model Building and Response Surfaces, John Wiley, 1987, pg. 74