

"The work contained and presented here is my work and my work alone."

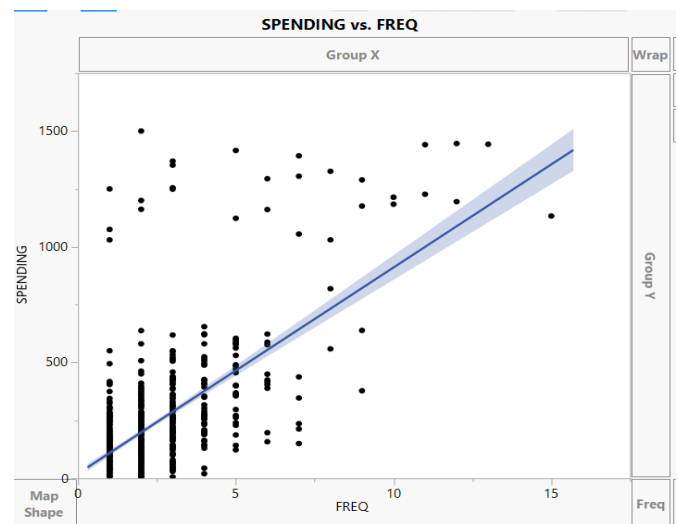
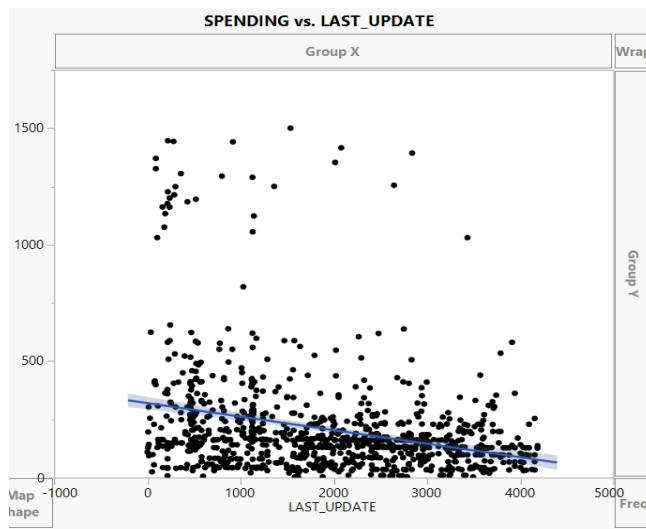
6.2 Predicting Software Reselling Profits.

a. Explore the spending amount by creating a tabular summary of the categorical variables and computing the average and standard deviation of spending in each category.

A-Following is the tabular summary

	WEB	GENDER	ADDRESS_RES	ADDRESS_US	N Rows	Mean(SPENDING)	Std Dev(SPENDING)
1	0	0	0	0	39	177.58974359	131.51740288
2	0	0	0	1	131	212.83206107	237.0980697
3	0	0	1	0	11	206.18181818	104.05269644
4	0	0	1	1	48	195.125	144.16267776
5	0	1	0	0	26	299.61538462	361.74997741
6	0	1	0	1	152	209.31578947	246.97913742
7	0	1	1	0	5	137.4	53.125323528
8	0	1	1	1	44	191.77272727	125.28736525
9	1	0	0	0	34	255.94117647	219.29818377
10	1	0	0	1	157	219.77070064	277.2668083
11	1	0	1	0	9	230	92.54458385
12	1	0	1	1	57	181.78947368	147.93673079
13	1	1	0	0	35	160.51428571	106.64437092
14	1	1	0	1	203	201.12807882	219.85454992
15	1	1	1	0	8	183.25	171.23897253
16	1	1	1	1	41	159.17073171	122.41096814

b. Explore the relationship between spending and each of the two continuous predictors by creating two scatterplots (SPENDING vs. FREQ, and SPENDING vs. LAST_UPDATE). Does there seem to be a linear relationship?



There seems to be a linear relationship. The linear equation has been partially fit. Not all the values are under the linear relation. First one is a negative equation. Second one is a equation with a positive slope.

c. To fit a predictive model for SPENDING:

i. Partition the 1000 records into training and validation sets. (Note that the dataset has a validation column—create your own validation column for this exercise.)

A-Validation set has been created.

	source_a	source_c	source_b	source_d	source_e	source_m	source_o	source_h	source_r	source_s	source_t	source_u	source_p	source_x	source_w	Validation
1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	Training
2	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	Training
3	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Training
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Training
5	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	Validation
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Training
7	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	Training
8	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	Training
9	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	Training
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	Training
1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	Training
2	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Training
3	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	Validation
4	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	Training
5	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	Training

ii. Run a multiple linear regression model for SPENDING vs. all six predictors. Give the estimated predictive equation.

A-The formula for prediction is stated beside.

iii. Based on this model, what type of purchaser is most likely to spend a large amount of money?

A-Based on the formula, we can see that frequency has the highest impact and hence, purchaser is most likely to spend a large amount if the frequency is high.

Prediction Expression

```

43.5552808034149
+ 91.747540897913 * FREQ
+ -0.0212438274742 * LAST_UPDATE
+ Match ( WEB ) { 0 ⇒ 5.18759802407864
                  1 ⇒ -5.1875980240786
                  else ⇒ .
}
+ Match ( GENDER ) { 0 ⇒ -0.6729549583002
                    1 ⇒ 0.6729549583002
                    else ⇒ .
}
+ Match ( ADDRESS_RES ) { 0 ⇒ 43.0881461279384
                          1 ⇒ -43.088146127938
                          else ⇒ .
}
+ Match ( ADDRESS_US ) { 0 ⇒ 9.98756191229101
                        1 ⇒ -9.987561912291
                        else ⇒ .
}

```

iv. If we used backward elimination to reduce the number of predictors, which predictor would be dropped first from the model?

A-For my model, Gender was the first variable that was removed. As shown in picture beside.

Stepwise Fit for SPENDING

Stepwise Regression Control

Stopping Rule: Max Validation RSquare

Direction: Backward

Rules: Combine

Go Stop Step

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

SSE	DFE	RMSE	RSquare	RSquare Adj	Cp	p	AICc	BIC	RSquare Validation	RMSE Validation
21156247	734	169.77404	0.4744	0.4723	3.1726687	4	9678.9	9701.837	0.3115	148.1296

Current Estimates

Lock	Entered	Parameter	Estimate	nDF	SS	"F Ratio"	"Prob > F"
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Intercept	38.8994192	1	0	0.000	1
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	FREQ	91.4027132	1	14816514	514.048	1.1e-86
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	LAST_UPDATE	-0.0219893	1	376950.3	13.078	0.00032
<input type="checkbox"/>	<input type="checkbox"/>	WEB(1-0)	0	1	21475.4	0.745	0.38841
<input type="checkbox"/>	<input type="checkbox"/>	GENDER(1-0)	0	1	0.426257	0.000	0.99693
<input type="checkbox"/>	<input checked="" type="checkbox"/>	ADDRESS_RES(1-0)	-42.942637	1	893303.6	30.992	3.63e-8
<input type="checkbox"/>	<input type="checkbox"/>	ADDRESS_US(1-0)	0	1	42765	1.485	0.22344

Step History

Step	Parameter	Action	"Sig Prob"	Seq SS	RSquare	Cp	p	AICc	BIC	RSquare Validation
1	All	Entered	.	.	0.4760	7	7	9682.83	9719.46	0.3038
2	GENDER(1-0)	Removed	0.9149	330.0008	0.4760	5.0114	6	9680.79	9712.87	0.3041
3	WEB(1-0)	Removed	0.4098	19598.98	0.4755	3.6906	5	9679.44	9706.95	0.3092

v. Show how the prediction and the prediction error are computed for the first purchase in the validation set.

Validation	Pred Formula SPENDING	Residual SPENDING
ng	90.319205366	37.680794634

SPENDING	FREQ	LAST_UPDATE
128	2	3662

After the regression model, I have saved the Pred formula Spending, Residual Spending. It can be seen in save columns, save Predicting Formula and Residuals.

vi. Evaluate the predictive accuracy of the model by examining its performance on the validation set.

Effect Tests

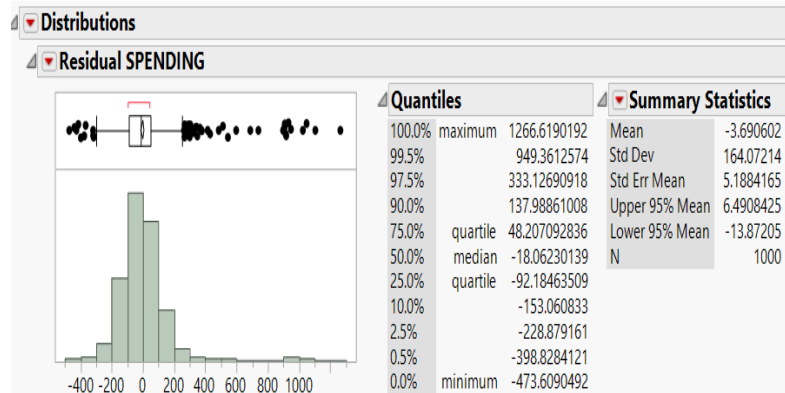
Crossvalidation

Source	RSquare	RASE	Freq
Training Set	0.4760	169.06	738
Validation Set	0.3038	148.95	262

We can see that the Rsquare value is more for Training set rather than Validation set. So, that model is slightly better,

vii. Create a histogram of the model residuals. Do they appear to follow a normal distribution? How does this affect the predictive performance of the model?

The histogram is skewed on right. With a lot of outliers. There is no difference in the predictive performance since the residuals do not affect the predictive capabilities. They only mention the variance from the predicted value from the actual value.



6.3 Predicting Airfare on New Routes

6.3 a. Explore the numerical predictors and response (FARE) by creating histograms, a correlation table and a scatterplot matrix. Examining potential relationships between FARE and those predictors. What seems to be the best single predictor of FARE?





After examining the relationships between all the predictors and fare, I can conclude that Distance has the highest correlation with FARE. Hence, it is the best single predictor.

b. Explore the categorical predictors (excluding the first four) by computing the percentage of flights in each category. Create a tabular summary with the average fare in each category. Which categorical predictor seems best for predicting FARE?

	VACATION	SW	SLOT	GATE	N Rows	Mean(FARE)
1	No	No	Controlled	Constrained	18	\$206.798
2	No	No	Controlled	Free	114	\$208.330
3	No	No	Free	Constrained	78	\$210.582
4	No	No	Free	Free	119	\$196.184
5	Yes	No	Controlled	Free	28	\$139.902
6	Yes	No	Free	Constrained	20	\$137.314
7	Yes	No	Free	Free	67	\$143.977
8	No	Yes	Controlled	Constrained	1	\$74.280
9	No	Yes	Controlled	Free	17	\$110.171
10	No	Yes	Free	Constrained	7	\$139.960
11	No	Yes	Free	Free	114	\$96.951
12	Yes	Yes	Controlled	Free	4	\$131.605
13	Yes	Yes	Free	Free	51	\$89.811

South west is the best since for all the fares, when south west is "NO", the fares are higher and they follow a trend.

c. Find a model for predicting the average fare on a new route:

i. Partition the data into training and validation sets. The model will be fit to the training data and evaluated on the validation set.

As we can see, the one in blue are the optimal regression models with best Rsquare value. But I would choose the third one, with 11 variables. Since the Rsquare values are similar, it is a little less expensive as it has less number of variables. The one with 13 values are is also a good model with high RSquare value.

So, the result in model ii has only 9 variables using mixed stepwise. Model iii has 11 variables which is more optimum and with high RSquare value.

iv. Compare the predictive accuracy of both models (ii) and (iii) using measures such as RMSE, *Cp*, *AICc*, and Validation RSquare.

MODEL ii (With 9 variables)

SSE	DFE	RMSE	RSquare	RSquare Adj	Cp	p	AICc	BIC	RSquare Validation	RMSE Validation
408638.01	368	33.323106	0.7978	0.7928	12.307385	10	3736.03	3778.592	0.7518	39.74746

MODEL iii (with 11 variables)

SSE	DFE	RMSE	RSquare	RSquare Adj	Cp	p	AICc	BIC	RSquare Validation	RMSE Validation
402054.5	366	33.143771	0.8010	0.7951	10.341411	12	3734.169	3784.323	0.7593	39.14315

There is a very slight difference in Rsquare and Rsquare ADJ. There is a larger difference in BIC, SIC p, Cp value.

v. Using model (iii), predict the average fare on a route with the following characteristics: COUPON = 1.202, NEW = 3, VACATION = No, SW = No, HI = 4442.141, S_INCOME = \$28,760, E_INCOME = \$27,664, S_POP = 4,557,004, E_POP = 3,195,503, SLOT = Free, GATE = Free, PAX = 12,782, DISTANCE = 1976 miles.

Parameter	
Intercept	0
COUPON	1.202
NEW	3
VACATION[No]	1
SW[No]	1
HI	4442.141
S_INCOME	28760
E_INCOME	27664
S_POP	4557004
E_POP	3195503
SLOT[Controlled]	0
GATE[Constrained]	0
DISTANCE	1976
PAX	12782
=	0
Value	290.5966461
Std Error	29.44861105
t Ratio	9.8679236724
Prob> t	1.326407e-20
SS	113468.88682
Sum of Squares	113468.88682
Numerator DF	1
F Ratio	97.375917603
Prob > F	1.326407e-20

The average Fare is - \$290.59

vi. Using model (iii), predict the reduction in average fare on the route in (v) if Southwest decides to cover this route.

Custom Test	
Parameter	
Intercept	0
COUPON	1.202
NEW	3
VACATION[No]	1
SW[No]	0
HI	4442.141
S_INCOME	28760
E_INCOME	27664
S_POP	4557004
E_POP	3195503
SLOT[Controlled]	0
GATE[Constrained]	0
DISTANCE	1976
PAX	12782
=	0
Value	269.55818949
Std Err...	30.155143932
t Ratio	8.9390450299
Prob> t	1.715606e-17
SS	93112.391472
Sum of Squares	93112.391472
Numerator DF	1
F Ratio	79.906526046
Prob > F	1.715606e-17

If southwest decides to cover this, the reduction is approximately \$21.
The average fare now is \$269.55

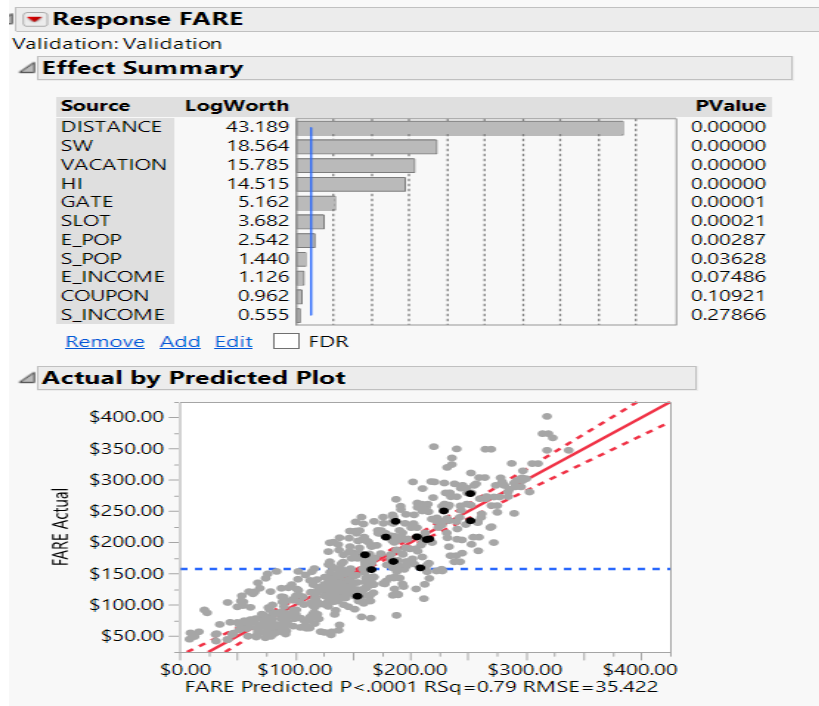
vii. In reality, which of the factors will be available for predicting the average from a new airport (i.e., before flights start operating on those routes)? Which ones can be estimated? How?

COUPON, NEW- not available

E_INCOME, S_INCOME, GATE, SLOT, PAX- estimated.

not
fare

viii. Select a model that includes only factors that are available before flights begin to operate on the new route. Use an exhaustive search (All Possible Models) to find such a model.



Following is the model with only available and predictors which can be estimated values.

ix. Use model (viii) to predict the average fare on a route with characteristics COUPON = 1.202, NEW = 3, VACATION = No, SW = No, HI = 4442.141, S_INCOME = \$28,760, E_INCOME = \$27,664, S_POP = 4,557,004, E_POP = 3,195,503, SLOT = Free, GATE = Free, PAX = 12,782, DISTANCE = 1976 miles.

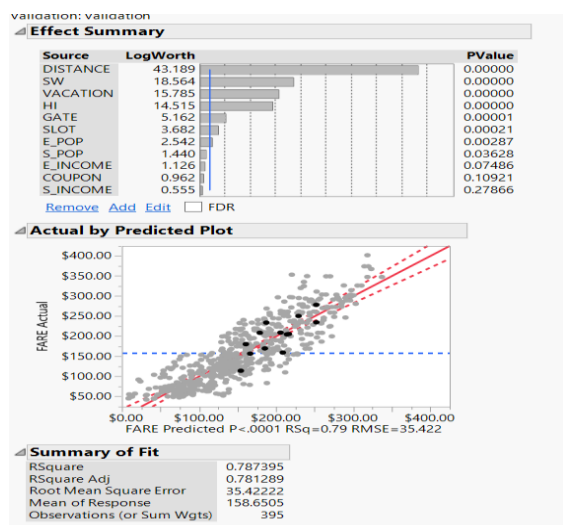
Custom Test

Parameter	
Intercept	0
COUPON	1.202
VACATION[No]	1
SW[No]	1
HI	4442.141
S_INCOME	28760
E_INCOME	27664
S_POP	4557004
E_POP	3195503
SLOT[Controlled]	0
GATE[Constrained]	0
DISTANCE	1976
=	0
Value	318.95685438
Std Error	29.367683088
t Ratio	10.860810961
Prob > t	3.901664e-24
SS	148004.93011

Sum of Squares	148004.93011
Numerator DF	1
F Ratio	117.95721472
Prob > F	3.901664e-24

The picture beside is the one with the values and the value of fare is roughly \$318.95.

x. Compare the predictive accuracy of this model with model (iii). Is this model good enough, or is it worthwhile reevaluating the model once flights begin on the new route?



Stepwise Regression Control

Go Stop Step

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

SSE	DFE	RMSE	RSquare	RSquare Adj	Cp	p	AICc	BIC	RSquare Validation	RMSE Validation
402054.5	366	33.143771	0.8010	0.7951	10.341411	12	3734.169	3784.323	0.7593	39.14315

Current Estimates

Lock	Entered	Parameter	Estimate	nDF	SS	F Ratio	Prob > F
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Intercept	-24.631075	1	0	0.000	1
<input type="checkbox"/>	<input type="checkbox"/>	COUPON	0	1	304.2553	0.276	0.59937
<input type="checkbox"/>	<input type="checkbox"/>	NEW	0	1	63.49901	0.058	0.81037
<input type="checkbox"/>	<input type="checkbox"/>	VACATION(Yes-No)	-15.728501	1	57805.28	52.622	2.4e-12
<input type="checkbox"/>	<input type="checkbox"/>	SW(Yes-No)	-17.275174	1	66652.52	60.675	7e-14
<input type="checkbox"/>	<input type="checkbox"/>	HI	0.00793237	1	51494.49	46.877	3.2e-11
<input type="checkbox"/>	<input type="checkbox"/>	S_INCOME	0.00118813	1	4234.746	3.855	0.05035
<input type="checkbox"/>	<input type="checkbox"/>	E_INCOME	0.00085118	1	3723.311	3.389	0.06643
<input type="checkbox"/>	<input type="checkbox"/>	S_POP	3.1009e-6	1	17264	15.716	8.85e-5
<input type="checkbox"/>	<input type="checkbox"/>	E_POP	4.72744e-6	1	30601.94	27.858	2.24e-7
<input type="checkbox"/>	<input type="checkbox"/>	SLOT(Free-Controlled)	-11.144471	1	25238.45	22.975	2.39e-6
<input type="checkbox"/>	<input type="checkbox"/>	GATE(Free-Constrained)	-12.831672	1	28814.13	26.230	4.91e-7
<input type="checkbox"/>	<input type="checkbox"/>	DISTANCE	0.07591607	1	623154.3	567.273	2.2e-76
<input type="checkbox"/>	<input type="checkbox"/>	PAX	-0.0006196	1	24076.73	21.918	4.02e-6

Step History

Step	Parameter	Action	Sig Prob	Seq SS	RSquare	Cp	p	AICc	BIC	RSquare Validation
1	Intercept	Entered								

All Possible Models

Ordered up to best 56 models up to 13 terms per model.

Model

COUPON,NEW,VACATION(Yes-No),SW(Yes-No),HLS_INCOME,E_INCOME,S_POP,E_POP,SLOT(Free-Controlled),GATE(Free-Constrained),DISTANCE,PAX

Model 3 is better due to more r square value. Yes, it is worthwhile.

d. In competitive industries, a new entrant with a novel business plan can have a disruptive effect on existing firms. If a new entrant's business model is sustainable, other players are forced to respond by changing their business practices. If the goal of the analysis was to evaluate the effect of Southwest Airlines' presence on the airline industry rather than predicting fares on new routes, how would the analysis be different? Describe technical and conceptual aspects.

Here, we can notice, the only variable that we can consider where can predict is the coupons. We can use the business acumen and predict and modulate the coupon prices. Hence, it's the only other thing which can predicted using the other predictors.