

## Linear Regression - Predicting Airfares on New Routes

### 1. Perform Typcasting. Which variables are categorical?

Categorical Variables:

- S\_CODE: starting airport's code
- S\_CITY: starting city
- E\_CODE: ending airport's code
- E\_CITY: ending city
- VACATION: whether a vacation route (Yes) or not (No)
- SW: whether Southwest Airlines serves that route (Yes) or not (No)
- SLOT: whether either endpoint airport is slot controlled or not; this is a measure of airport congestion
- GATE: whether either endpoint airport has gate constraints or not; this is another measure of airport congestion

Numerical Variables:

- COUPON: average number of coupons
- NEW: number of new carriers
- HI: Herfindel Index – measure of market concentration
- S\_INCOME: starting city's average personal income
- E\_INCOME: ending city's average personal income
- S\_POP: starting city's population
- E\_POP: ending city's population
- DISTANCE: distance between two endpoint airports in miles
- PAX: number of passengers on that route during the period of data collection
- FARE: average fare on that route

- Partition the data into training (60%) and validation sets (40%). The model will be fit to the training data and evaluated on the validation set.

When we partition the data using standard partition and the 60/40 split, we get the following:

**Output Navigator**

Inputs	Partition Summary	Partitioned Data

**Elapsed Times in Milliseconds**

Data Reading Time	Algorithm
24	

**Inputs**

Data	
Workbook	Airfares (4).xlsx
Worksheet	data
Range	SAS1:\$R\$639
# Records in the input data	638

Variables	
# Selected Variables	18
Selected Variables	S_CODE S_CITY E_CODE E_CITY COUPON NEW

Partitioning Parameters	
Partitioning type	RANDOM
Random seed	12345
Ratio - Training	0.6
Ratio - Validation	0.4

**Partition Summary**

Partition	# Records
Training	383
Validation	255

**Partitioned Data**

Record ID	S_CODE	S_CITY	E_CODE	E_CITY	COUPON	NEW	VACATION	SW
Record 1								

data codes data **STDPartition**

### 3. What is the Fare value for the last three records included in validation set?

Unfortunately, when the partition is run it does not order the records in order, so as per the example in class the first 383 records (60%) belong to the training and the rest of the 255 (40%) records belong to the validation set. Therefore, the fare value for the last three records included in the validation set are:

- 52.92
- 154.73
- 102.95

1421287	2549844	Free	Free	869	6082	142.98
3036732	1021733	Free	Free	180	7646	66.14
8621121	113091	Controlled	Free	254	8469	127.38
9056076	1197234	Free	Free	226	16845	52.92
8621121	1021830	Free	Unstrained	426	12883	154.73
4459144	1421287	Free	Unstrained	963	16414	102.95

All  
Smooth  
Linear  
Recourse  
? Model T  
If Unknown  
model.

Count: 255

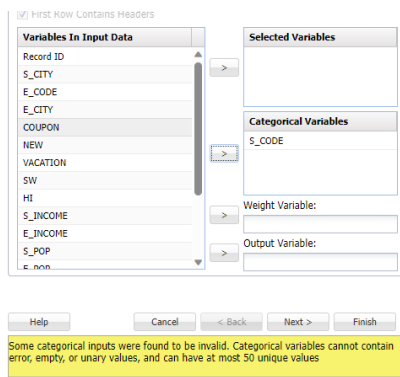
4. Try to add all categorical variables to create a linear regression model. Which variable cannot be added (not allowed by the software) and why is that?

Analytical solver we have two types:

- Classify: target variable is a **categorical variable**
- Predict: price, numerical variable

Steps: We click on data mining and then click on classify.

We then select linear regression and proceed to select the categorical values. However, when S\_CITY is selected it doesn't work - we get the below error, because there are blanks in the column for these.



The following categorical variables are not added:

- S\_CITY
- E\_CITY

5. Create a linear regression model only based numerical variables to predict the Fare variable. For scoring Training and Validation data, have detailed and summary reports created. What is RMSE for the validation?

**Data Mining: Linear Regression** Date: 01-Oct-2023 20:28:07

**Output Navigator**

Inputs	Regression Summary	Predictor Screen	Coefficients
PHML Model	Training Prediction Sum	Validation Predi	

**Elapsed Times in Milliseconds**

Data Reading Time	Algorithm Time	Report Time	Total
14	6	4	24

**Inputs**

**Data**

Workbook	Airfares (4).xlsx
Worksheet	STDPartition
Training data used for building the model	SC537:SU5419
# Records in the training data	383
Validation data	SC5420:SU5674
# Records in the validation data	255

**Variables**

# Variables	9
Scale Variables	COUPON NEW HI S_INCOME E_INCOME S_POP E_POP DISTANCE FARE
Categorical Variables	
Output Variable	FARE

**Rescaling: Fitting Parameters**

Rescale Data?	FALSE
---------------	-------

**Regression Model: Fitting Parameters**

Fit Intercept	TRUE
---------------	------

**Regression Model: Reporting Parameters**

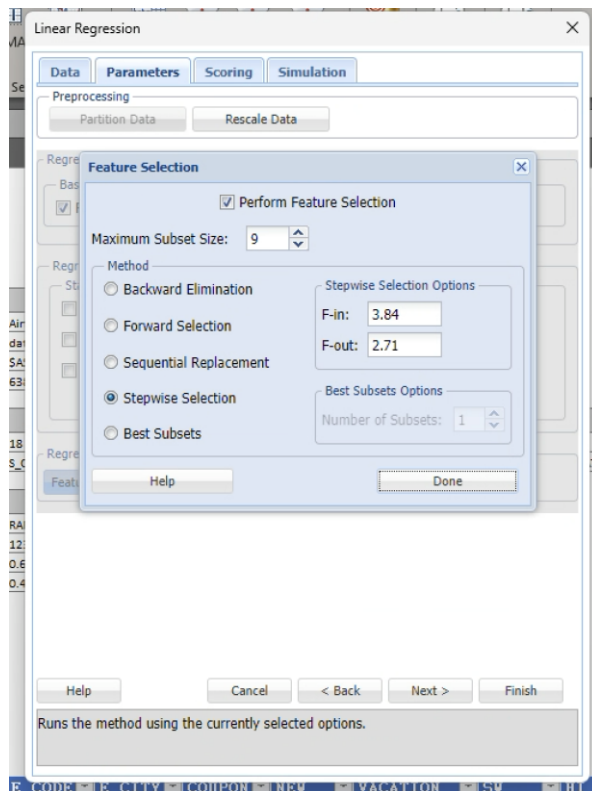
Variance-Covariance Matrix	FALSE
Multicollinearity Diagnostic	FALSE

data codes | data | STDPartition | **LinReg\_Output** | LinReg\_TrainingScore | LinReg\_ValidationScore | LinReg\_Stc ...

## Validation: Prediction Summary

Metric	Value
SSE	487875.6
MSE	1913.238
RMSE	43.74057
MAD	35.29559
R2	0.665976

6. Redo step 5 and this time use stepwise variable selection feature to reduce the number of predictors. How many subsets are created, and which variable(s) are omitted at the end?



See excel spreadsheet attached to assignment.