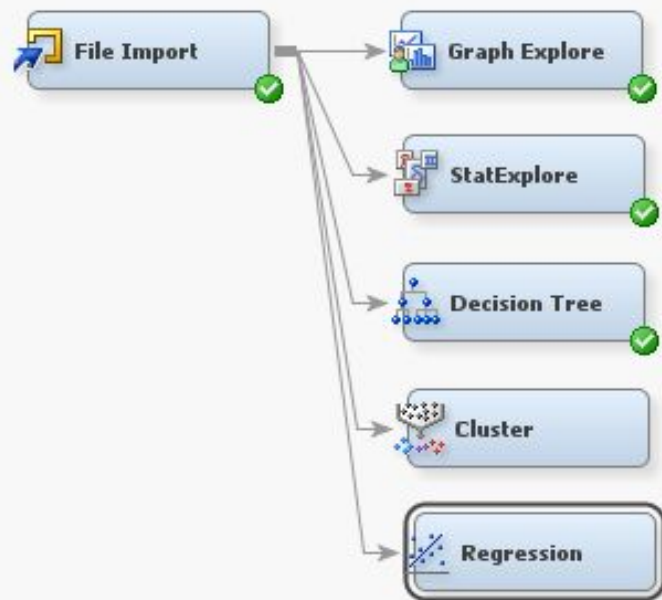


Using SAS Enterprise Miner to Analyze Netflix Data

By: Ashley Krause

- Secondary Source:
Kaggle
 - Column 1
 - Title
 - Year (1905 - 2022)
 - Kind
 - Genre
 - Rating
 - Vote
 - Country
 - Language
 - Cast
 - Director
 - Composer
 - Writer
 - Runtime





Variables - FIMPORT

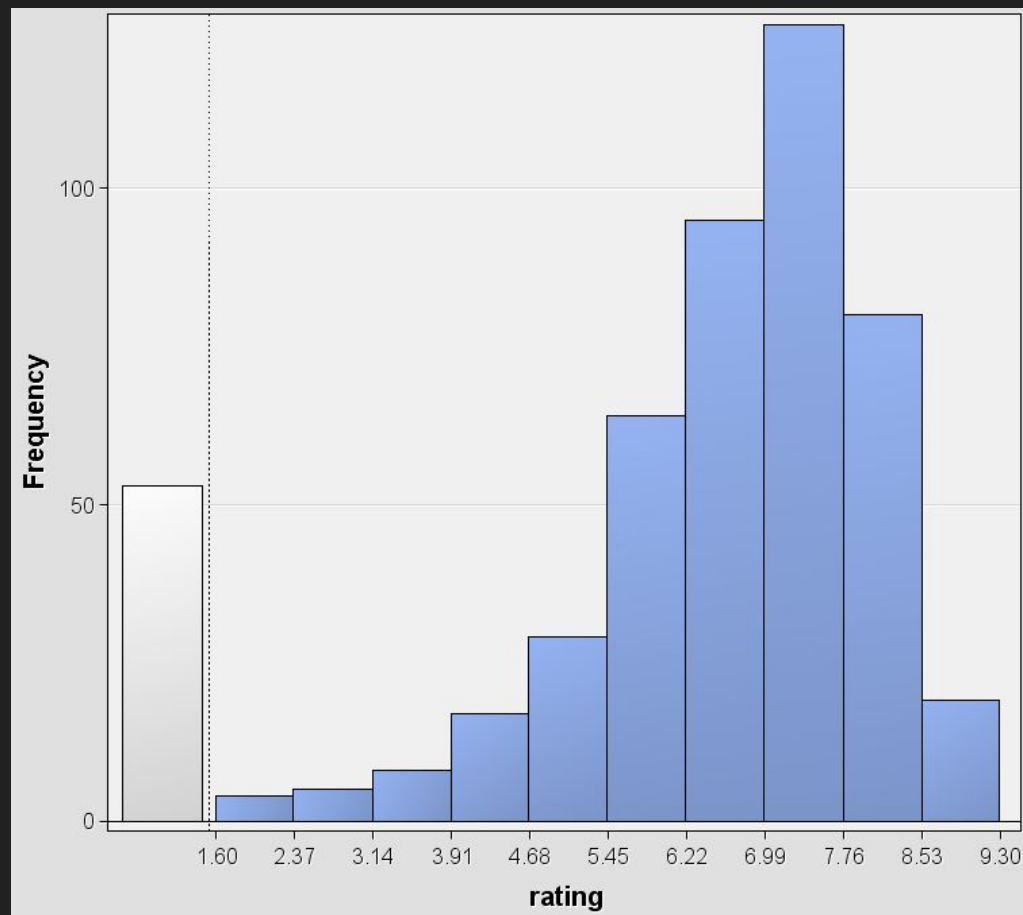
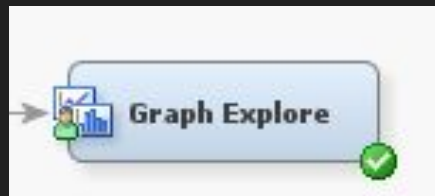
(none) ☐ not Equal to ☐ ...

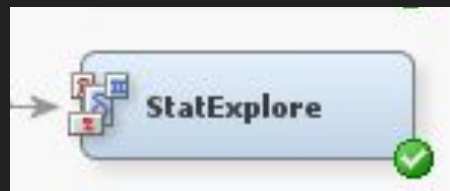
Columns: ☐ Label

☐ Mining

☐ Basic

Name	Role	Level	Report	Order	Drop	Lower Limit	Upper Limit
cast	Text	Nominal	No		No	.	.
Column1	ID	Interval	No		No	.	.
composer	Text	Nominal	No		No	.	.
country	Text	Nominal	No		No	.	.
director	Input	Nominal	No		No	.	.
genre	Text	Nominal	No		No	.	.
kind	Input	Nominal	No		No	.	.
language	Text	Nominal	No		No	.	.
rating	Input	Interval	No		No	.	.
runtime	Input	Interval	No		No	.	.
title	Text	Nominal	No		No	.	.
vote	Input	Interval	No		No	.	.
writer	Text	Nominal	No		No	.	.
year	Input	Interval	No		No	.	.





Class Variable Summary Statistics (maximum 500 observations printed)

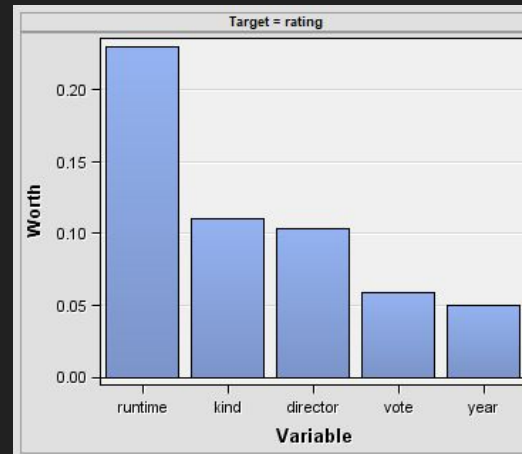
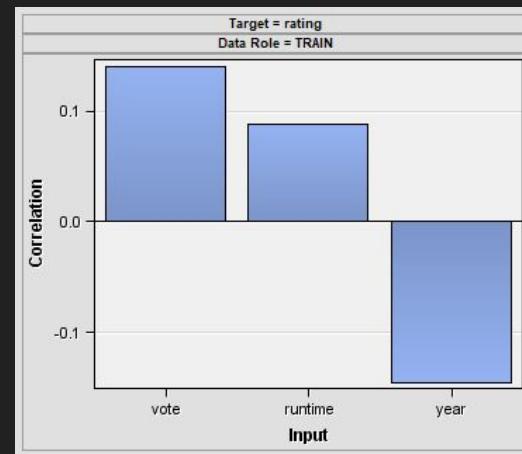
Data Role=TRAIN

Data Role	Variable Name	Role	Number of Levels	Missing	Mode	Mode Percentage	Mode2	Mode2 Percentage
TRAIN	director	INPUT	513	138		20.26	['Akira Kurosawa']	0.59
TRAIN	kind	INPUT	8	0	movie	56.80	video movie	14.54

Interval Variable Summary Statistics (maximum 500 observations printed)

Data Role=TRAIN

Variable	Role	Mean	Standard Deviation	Non Missing	Missing	Minimum	Median	Maximum	Skewness	Kurtosis
rating	INPUT	6.680635	1.285113	8949	807	1	6.9	9.6	-0.81107	0.757476
runtime	INPUT	98.11594	63.68129	8763	993	1	94	1620	8.410935	141.9943
vote	INPUT	21218.21	98048.73	8949	807	5	1535	2462087	12.20194	201.5009
year	INPUT	1994.74	16.24509	9756	0	1905	1999	2023	-1.46315	3.023835





```

34
35
36             Cluster Summary for 1 Cluster
37
38
39             Cluster    Variation    Proportion    Second
40             Cluster    Members    Variation    Explained    Explained    Eigenvalue
41             -----
42             1          4          4      1.096628      0.2742      1.0181
43
44             Total variation explained = 1.096628 Proportion = 0.2742
45
46             Cluster 1 will be split because it has the largest second eigenvalue, 1.018102, which is greater than the MAXEIGEN=1 value.
47
48             Clustering algorithm converged.
49
50
51             Cluster Summary for 2 Clusters
52
53
54             Cluster    Variation    Proportion    Second
55             Cluster    Members    Variation    Explained    Explained    Eigenvalue
56             -----
57             1          2          2      1.061685      0.5308      0.9383
58             2          2          2      1.031453      0.5157      0.9685
59
60             Total variation explained = 2.093138 Proportion = 0.5233
61
62
63             R-squared with
64             2 Clusters
65             -----
66             Cluster    Variable    Own    Next    1-R**2
67             Cluster    Variable    Cluster    Closest    Ratio
68             -----
69             Cluster 1    runtime    0.5308    0.0010    0.4696
70             Cluster 1    vote      0.5308    0.0007    0.4695
71             Cluster 2    VARI      0.5157    0.0025    0.4855
72             Cluster 2    year      0.5157    0.0000    0.4843
73
74             No cluster meets the criterion for splitting.
75
76
77             Total    Proportion    Minimum    Maximum
78             Variation    of    Proportion    Second
79             Explained    by    Variation    Explained
80             by    Explained    by a
81             Clusters    Clusters    by Clusters    Cluster
82             -----
83             1      1.096628      0.2742      0.2742      1.018102      0.0445
84             2      2.093138      0.5233      0.5157      0.968547      0.5157      0.4855
85
86             *-----*

```




Output

The DMREG Procedure

Model Information

Training Data Set WORK.EM_DMREG.VIEW
DMDB Catalog WORK.REG_DMDB
Target Variable rating
Target Measurement Level Interval
Error Normal
Link Function Identity
Number of Model Parameters 2
Number of Observations 8453

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	108.089400	108.089400	65.89	<.0001
Error	8451	13863	1.640360		
Corrected Total	8452	13971			

Model Fit Statistics

R-Square	0.0077	Adj R-Sq	0.0076
AIC	4185.5230	BIC	4187.5239
SEC	4199.6075	C(p)	2.0000

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	t Value	Pr > t
Intercept	1	6.4701	0.0260	249.29	<.0001
runtime	1	0.00178	0.000219	8.12	<.0001

* Score Output

* Report Output