

PROFESSIONAL & CONTINUING EDUCATION

UNIVERSITY *of* WASHINGTON

CLASSIFICATION SCHEMA



CLASSIFICATION SCHEMA

- > Modeling Dataset
 - Rectangular Dataset (aka table)
 - Schema
 - > Input columns
 - > Output column (target, outcome)
 - Classification: Category Column
 - Regression: Numeric Column
 - Horizontal partition of modeling data into training and test data
- > Incremental data has same schema as modeling data, except:
 - Incremental data does not have the output column (target, outcome)
 - Incremental data is not partitioned into training and test data



CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123	red	T	Yes
330-3150	Seaborg	No	0.987	green	T	No
330-3202	Seaborg	Yes	0.245	blue	F	Yes
415-2008	Seaborg	Yes	0.254	blue	T	Yes
415-2081	Seaborg	Bad	0.244	blue	F	No
415-2796	Seaborg		0.415	green	F	Maybe
415-2799	Seaborg	Yes	0.925	red	T	Yes
415-2913	Seaborg	Yes	0.376	green	F	Yes
415-3659	Seaborg	Bad	0.615	green	T	No
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No



CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123	red	T	Yes
330-3150	Seaborg	No	0.987	green	T	No
330-3202	Seaborg	Yes	0.245	blue	F	Yes
415-2008	Seaborg	Yes	0.254	blue	T	Yes
415-2081	Seaborg	Bad	0.244	blue	F	No
415-2796	Seaborg		0.415	green	F	Maybe
415-2799	Seaborg	Yes	0.925	red	T	Yes
415-2913	Seaborg	Yes	0.376	green	F	Yes
415-3659	Seaborg	Bad	0.615	green	T	No
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

Here is a rectangular dataset. The table has columns with headers and the data in each column have the same datatype. The data have been prepared and are ready for modeling.



CLASSIFICATION SCHEMA

Elsewhere, I have new data that do not contain the target outcome. I want to predict categorical values, like these, from this new data. For each row in the new data, I want to use the values from the other columns in the same row to predict the value in the missing column. This predicted value is called the “Target Outcome”.		Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
				Good	0.123	red	T	Yes
				No	0.987	green	T	No
				Yes	0.245	blue	F	Yes
				Yes	0.254	blue	T	Yes
				Bad	0.244	blue	F	No
					0.415	green	F	Maybe
				Yes	0.925	red	T	Yes
				Yes	0.376	green	F	Yes
				Bad	0.615	green	T	No
					0.321	blue	F	Maybe
				No	0.098	green	F	No
		598-2454	Seaborg	Bad	0.765	red	T	No

Target Outcome

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123	red	T	Yes
330-3150	Seaborg	No	0.987	green	T	No
330-3202	Seaborg	Yes	0.245	blue	F	Yes
415-2008	Seaborg	Yes	0.254	blue	T	Yes
415-2081	Seaborg	Bad	0.244	blue	F	No
415-2796	Seaborg		0.415	green	F	Maybe
415-2799	Seaborg	Yes	0.925	red	T	Yes
415-2913	Seaborg	Yes	0.376	green	F	Yes
415-3659	Seaborg	Bad	0.615	green	T	No
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

Target
Outcome

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	
330-3141	Seaborg	Good	0.123	
330-3150	Seaborg	No	0.987	
330-3202	Seaborg	Yes	0.245	
415-2008	Seaborg	Yes	0.254	
415-2081	Seaborg	Bad	0.244	
415-2796	Seaborg		0.415	
415-2799	Seaborg	Yes	0.925	
415-2913	Seaborg	Yes	0.376	
415-3659	Seaborg	Bad	0.615	
595-8413	Seaborg		0.321	blue
598-1243	Seaborg	No	0.098	green
598-2454	Seaborg	Bad	0.765	red

blue	F	Maybe
green	F	No
red	T	No

Keys and random data should not be used as inputs for predictive analytics. Random data may appear to have patterns, but those patterns are fortuitous and will not be available when needed for predictions. Keys may contain patterns, but these patterns are deceptive and may also not be available when needed.

Random
or Keys

Target
Outcome

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123	red	T	Yes
330-3150	Seaborg	No	0.987	green	T	No
330-3202	Seaborg	Yes	0.245	blue	F	Yes
415-2008	Seaborg	Yes	0.254	blue	T	Yes
415-2081	Seaborg	Bad	0.244	blue	F	No
415-2796	Seaborg		0.415	green	F	Maybe
415-2799	Seaborg	Yes	0.925	red	T	Yes
415-2913	Seaborg	Yes	0.376	green	F	Yes
415-3659	Seaborg	Bad	0.615	green	T	No
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

**Random
or Keys**

**Target
Outcome**

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123			
330-3150	Seaborg	No	0.987			
330-3202	Seaborg	Yes	0.245			
415-2008	Seaborg	Yes	0.254			
415-2081	Seaborg	Bad	0.244			
415-2796	Seaborg		0.415			
415-2799	Seaborg	Yes	0.925			
415-2913	Seaborg	Yes	0.376			
415-3659	Seaborg	Bad	0.615			
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

Columns with constant data are unnecessary. In general, they will not affect the algorithm and therefore the model will be the same. But, they distract from the task. Also, they may increase memory and processing requirements.

Random
or Keys

Constant

Target
Outcome

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123	red	T	Yes
330-3150	Seaborg	No	0.987	green	T	No
330-3202	Seaborg	Yes	0.245	blue	F	Yes
415-2008	Seaborg	Yes	0.254	blue	T	Yes
415-2081	Seaborg	Bad	0.244	blue	F	No
415-2796	Seaborg		0.415	green	F	Maybe
415-2799	Seaborg	Yes	0.925	red	T	Yes
415-2913	Seaborg	Yes	0.376	green	F	Yes
415-3659	Seaborg	Bad	0.615	green	T	No
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

Random
or Keys

Constant

Target
Outcome

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123			
330-3150	Seaborg	No	0.987			
330-3202	Seaborg	Yes	0.245			
415-2008	Seaborg	Yes	0.254			
415-2081	Seaborg	Bad	0.244			
415-2796	Seaborg		0.415			
415-2799	Seaborg	Yes	0.925			
415-2913	Seaborg	Yes	0.376			
415-3659	Seaborg	Bad	0.615			
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

A proxy column is a column that was created after the “target” was observed. The proxy contains information that would not be available for predictions. The proxy column correlates well with the target .

Random
or Keys

Constant

Proxy

Target
Outcome

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123	red	T	Yes
330-3150	Seaborg	No	0.987	green	T	No
330-3202	Seaborg	Yes	0.245	blue	F	Yes
415-2008	Seaborg	Yes	0.254	blue	T	Yes
415-2081	Seaborg	Bad	0.244	blue	F	No
415-2796	Seaborg		0.415	green	F	Maybe
415-2799	Seaborg	Yes	0.925	red	T	Yes
415-2913	Seaborg	Yes	0.376	green	F	Yes
415-3659	Seaborg	Bad	0.615	green	T	No
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

Random
or Keys

Constant

Proxy

Target
Outcome

W

CLASSIFICATION SCHEMA

Some inputs to supervised learning are continuous attributes, like integers, floats and time.

Some inputs to supervised learning are categories, like strings, binned numbers, and factors.

Some inputs to supervised learning are binary attributes, like categories with only two states and binarized multi-state categories.

<p>inputs to supervised learning are continuous attributes, like integers, floats and time.</p> <p>inputs to supervised learning are categories, like strings, binned numbers, and factors.</p> <p>inputs to supervised learning are binary attributes, like series with only two states and serialized multi-state categories.</p>			Column 3	Column 4	Column 5	Column 6	Column 7
			Good	0.123	red	T	Yes
			Good	0.987	green	T	No
			Good	0.245	blue	F	Yes
			Good	0.254	blue	T	Yes
			Good	0.244	blue	F	No
			Good	0.415	green	F	Maybe
			Good	0.925	red	T	Yes
			Good	0.376	green	F	Yes
			Good	0.615	green	T	No
			0.321	blue	F	Maybe	
598-1243	Seaborg	No	0.098	green	F	No	
598-2454	Seaborg	Bad	0.765	red	T	No	

Random
or Keys

Constant

Proxy

Continuous
Input

Categorical
Input

Binary
Input

Target
Outcome

W

CLASSIFICATION SCHEMA

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7
330-3141	Seaborg	Good	0.123	red	T	Yes
330-3150	Seaborg	No	0.987	green	T	No
330-3202	Seaborg	Yes	0.245	blue	F	Yes
415-2008	Seaborg	Yes	0.254	blue	T	Yes
415-2081	Seaborg	Bad	0.244	blue	F	No
415-2796	Seaborg		0.415	green	F	Maybe
415-2799	Seaborg	Yes	0.925	red	T	Yes
415-2913	Seaborg	Yes	0.376	green	F	Yes
415-3659	Seaborg	Bad	0.615	green	T	No
595-8413	Seaborg		0.321	blue	F	Maybe
598-1243	Seaborg	No	0.098	green	F	No
598-2454	Seaborg	Bad	0.765	red	T	No

Random
or Keys

Constant

Proxy

Continuous
Input

Categorical
Input

Binary
Input

Target
Outcome

W

CLASSIFICATION SCHEMA

	Input 1	Input 2	Input 3	Outcome
	0.123	red	T	Yes
	0.987	green	T	No
	0.245	blue	F	Yes
	0.254	blue	T	Yes
	0.244	blue	F	No
	0.415	green	F	Maybe
	0.925	red	T	Yes
	0.376	green	F	Yes
	0.615	green	T	No
	0.321	blue	F	Maybe
	0.098	green	F	No
	0.765	red	T	No

Continuous
Input

Binary
Input

Categorical
Input

Target
Outcome

W

CLASSIFICATION SCHEMA

Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	No
0.321	blue	F	Maybe
0.098	green	F	No
0.765	red	T	No

**Continuous
Input**

**Binary
Input**

**Categorical
Input**

**Target
Outcome**

W

CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3

Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	No
0.321	blue	F	Maybe
0.098	green	F	No
0.765	red	T	No



CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3

Modeling Data
(300-100000 rows)

Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	No
0.321	blue	F	Maybe
0.098	green	F	No
0.765	red	T	No

W

CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3

Training Data
(200-50000 rows)

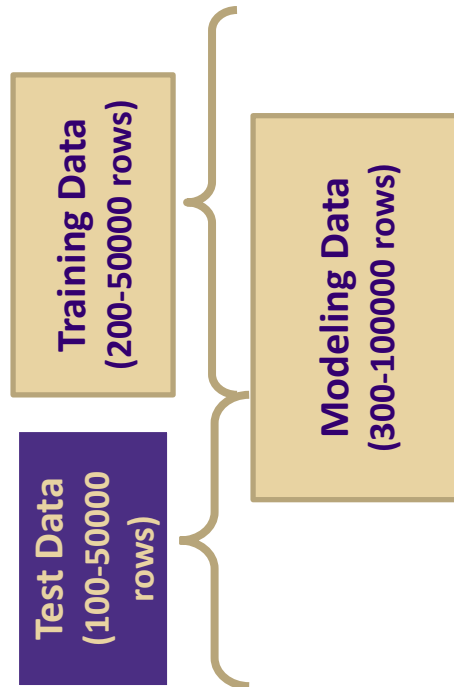
Modeling Data
(300-100000 rows)

Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	No
0.321	blue	F	Maybe
0.098	green	F	No
0.765	red	T	No

W

CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3

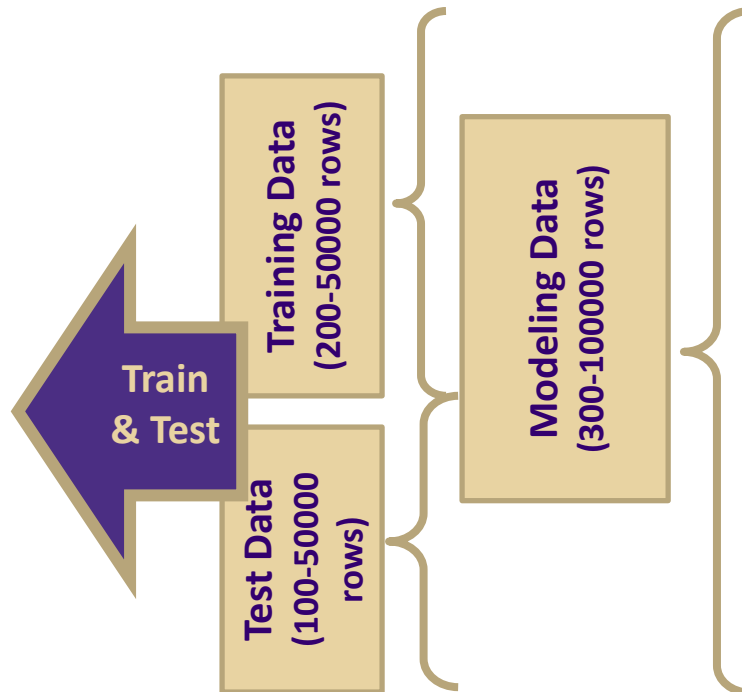


Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	No
0.321	blue	F	Maybe
0.098	green	F	No
0.765	red	T	No

W

CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3

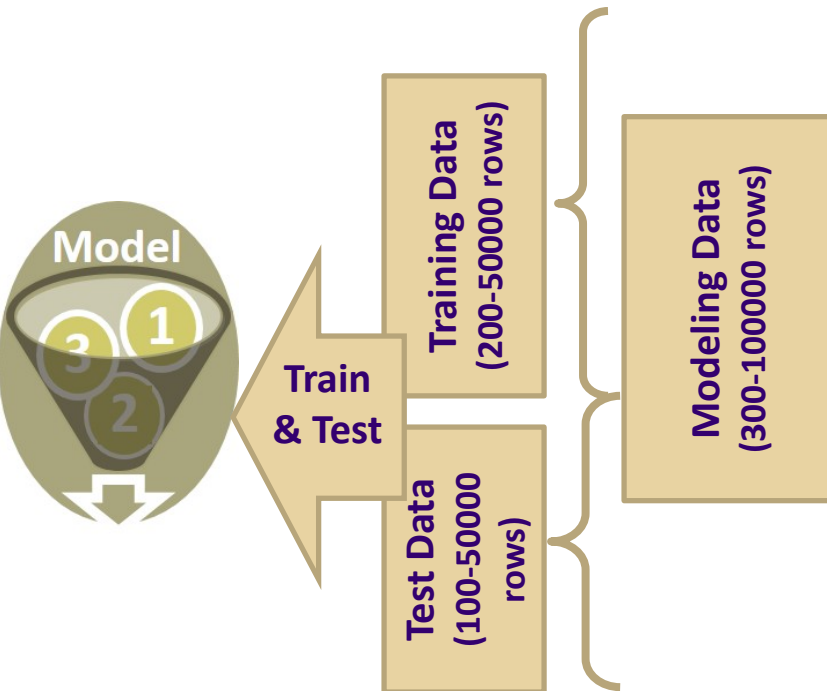


Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	No
0.321	blue	F	Maybe
0.098	green	F	No
0.765	red	T	No

W

CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3

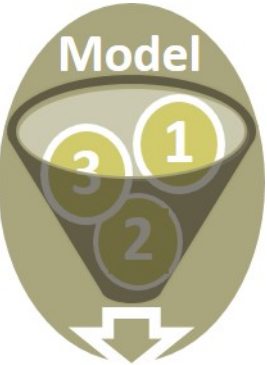


Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	No
0.321	blue	F	Maybe
0.098	green	F	No
0.765	red	T	No

W

CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3



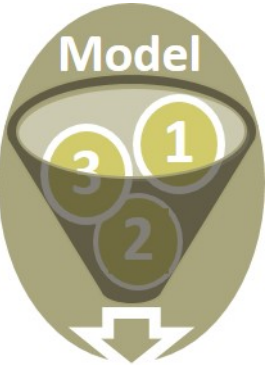
Elsewhere, I have new data that do not contain the target outcome. I want to predict categorical values, like these, from this new data. For each row in the new data, I want to use the values from the other columns in the same row to predict the value in the missing column. This predicted value is called the “Target Outcome”.

Operational Data
(1-∞ rows)

Input 1	Input 2	Input 3	Outcome
0.123	red	T	Yes
0.987	green	T	No
0.245	blue	F	Yes
0.254	blue	T	Yes
0.244	blue	F	No
0.415	green	F	Maybe
0.925	red	T	Yes
0.376	green	F	Yes
0.615	green	T	Target Outcome
0.321	blue	F	
0.098	green	F	
0.765	red	T	
0.234	green	T	
0.567	blue	F	
0.890	green	T	
0.314	red	T	

CLASSIFICATION SCHEMA

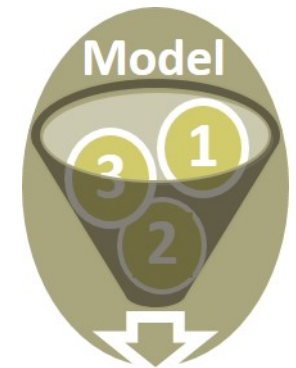
Outcome from Input 1, Input 2, Input 3



Operational Data (1-∞ rows)	Input 1	Input 2	Input 3	Target Outcome
	0.234	green	T	
	0.567	blue	F	
	0.890	green	T	
	0.314	red	T	

CLASSIFICATION SCHEMA

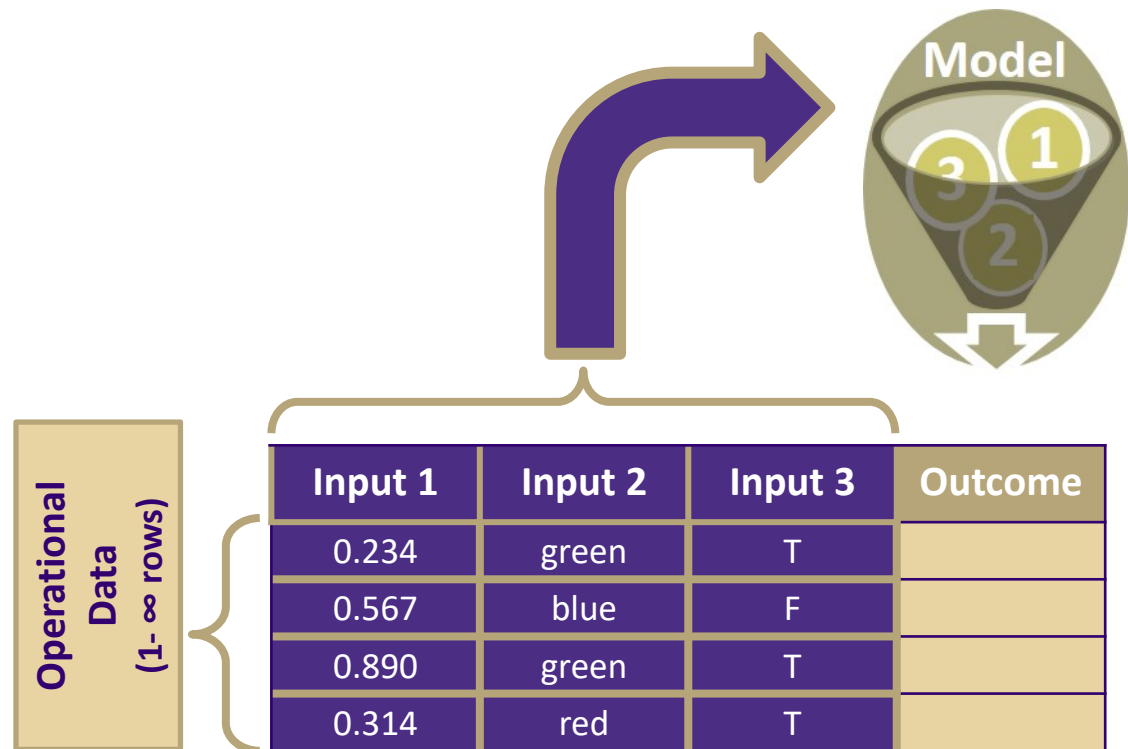
Outcome from Input 1, Input 2, Input 3



Operational Data (1-∞ rows)	Input 1	Input 2	Input 3	Outcome
	0.234	green	T	
	0.567	blue	F	
	0.890	green	T	
	0.314	red	T	

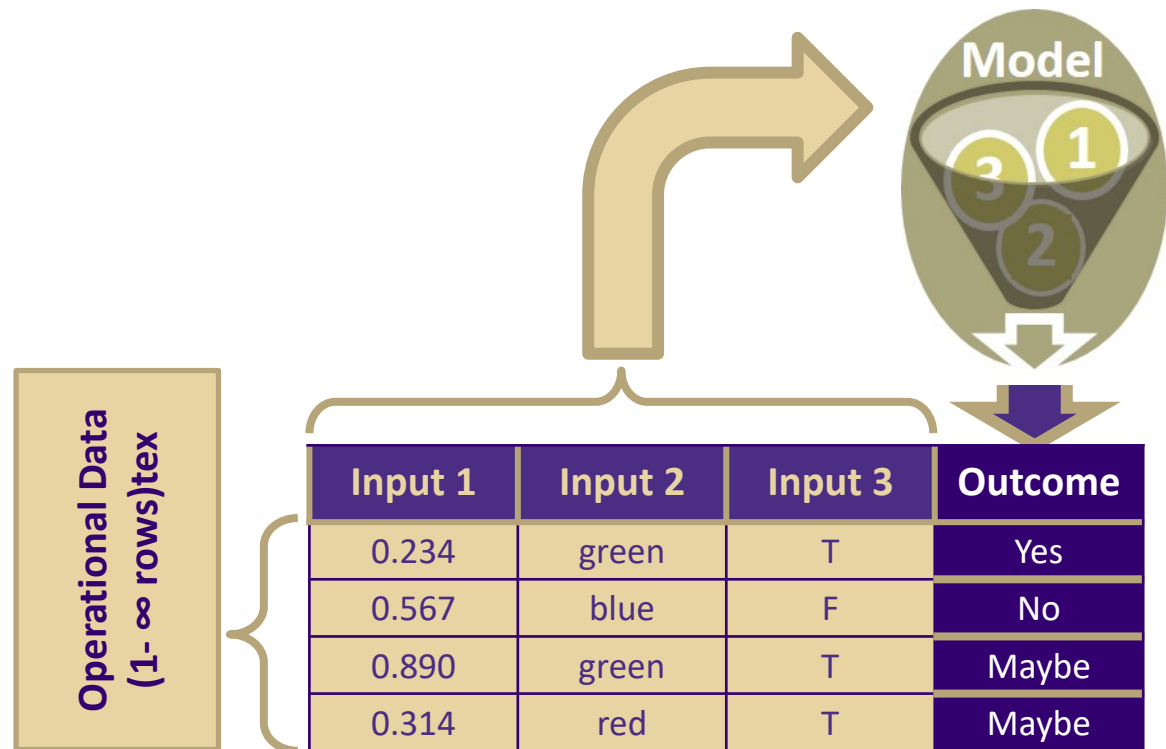
CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3



CLASSIFICATION SCHEMA

Outcome from Input 1, Input 2, Input 3



CLASSIFICATION SCHEMA

- > Attributes
 - All the columns are attributes
- > Input Column
 - Input columns are columns that can help predict the outcome. Input columns can be of type binary, ordinal, or category.
- > Target Outcome
 - The term "Target Outcome" is redundant. The outcome is the target and vice versa. The target or outcome is the output of a predict function. Providing target or outcome values during modeling makes the process supervised. Creating a model using a outcome is called supervised learning.
- > Proxy Column
 - A proxy column is a column that predicts too well. It is too good to be true. Something from the target leaked. This is also called target leakage. The leaked information is "not fair" to use in modeling. Values for that attribute will not be available when you want to predict the target outcome from operational data.



CLASSIFICATION SCHEMA

> Key Column

- In principle, a key column should not affect the model's prediction. The relationship between a key and any other attribute should be random. In practice, the algorithm will find a pattern in the key column and train on this pattern. This pattern is likely to be fortuitous, that means: random. The pattern will not hold for test data or when the model is applied. As a consequence, the key column will affect the model in a bad way.

> Constant Column

- A constant column should have no affect on the model's predictions. The constant column may increase computation time and cause other problems. It is standard practice to remove all constant columns prior to modeling.

