

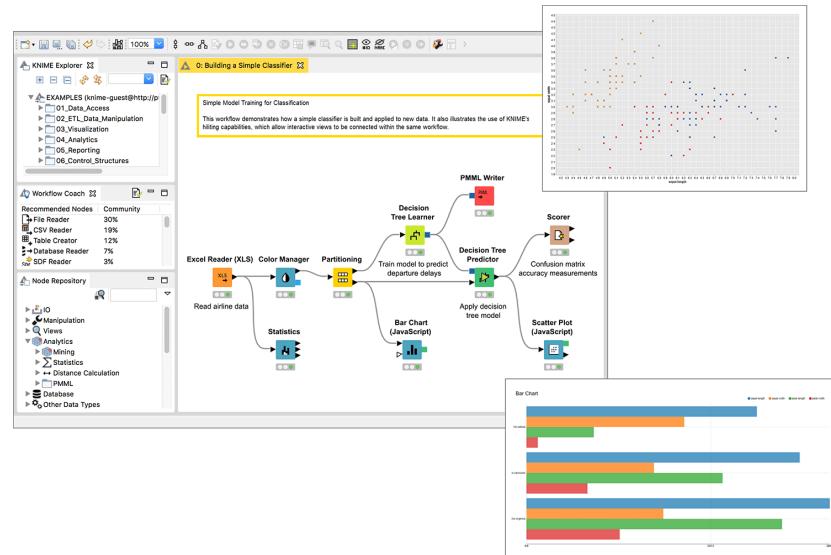
# Overview

# KNIME Analytics Platform



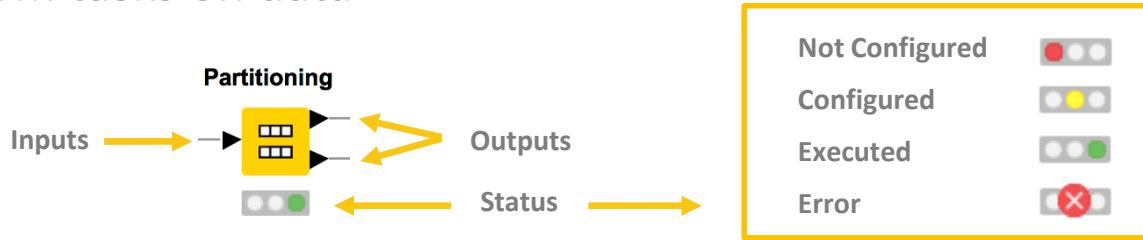
# What is KNIME Analytics Platform?

- A tool for data analysis, manipulation, visualization, and reporting
- Based on the graphical programming paradigm
- Provides a diverse array of extensions:
  - Text Mining
  - Network Mining
  - Cheminformatics
  - Many integrations, such as Java, R, Python, Weka, Keras, Plotly, H2O, etc.

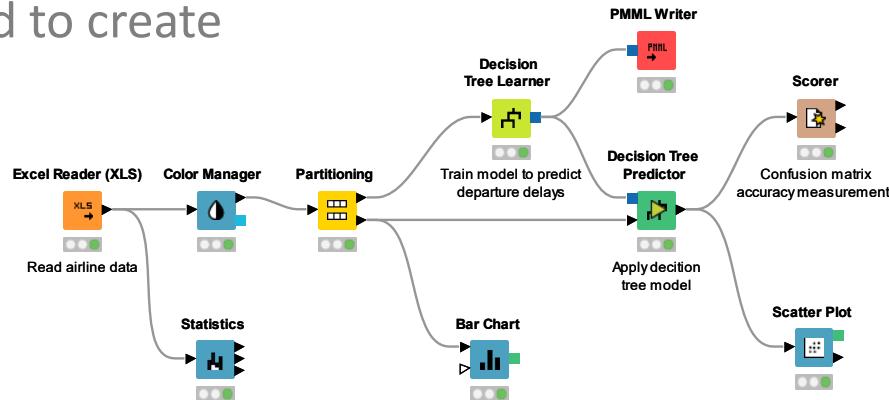


# Visual KNIME Workflows

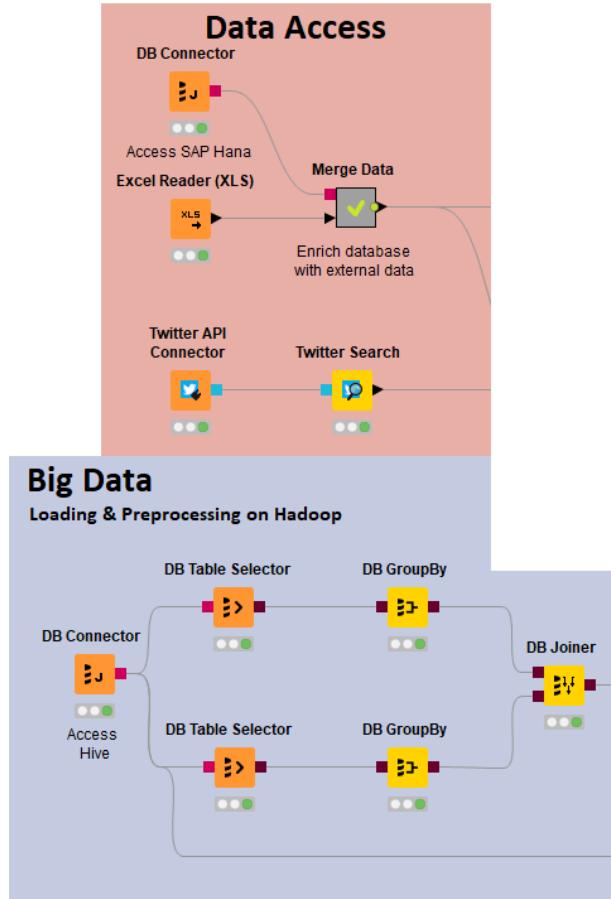
NODES perform tasks on data



Nodes are combined to create  
**WORKFLOWS**



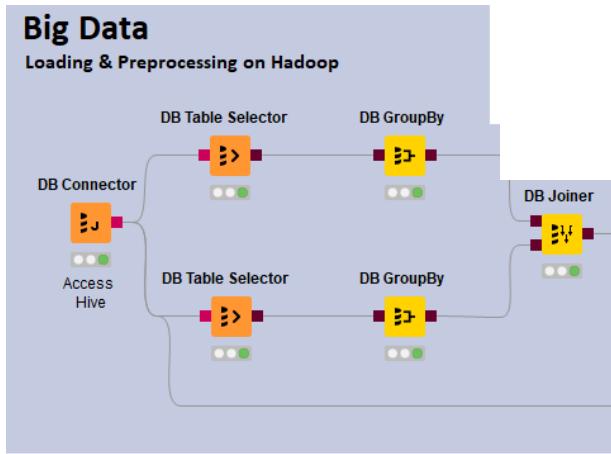
# Data Access



- Databases
  - MySQL, PostgreSQL
  - any JDBC (Oracle, DB2, MS SQL Server)
- Files
  - CSV, txt
  - Excel, Word, PDF
  - SAS, SPSS
  - XML
  - PMML
  - Images, texts, networks, chem
- Web, Cloud
  - REST, Web services
  - Twitter, Google

# Big Data

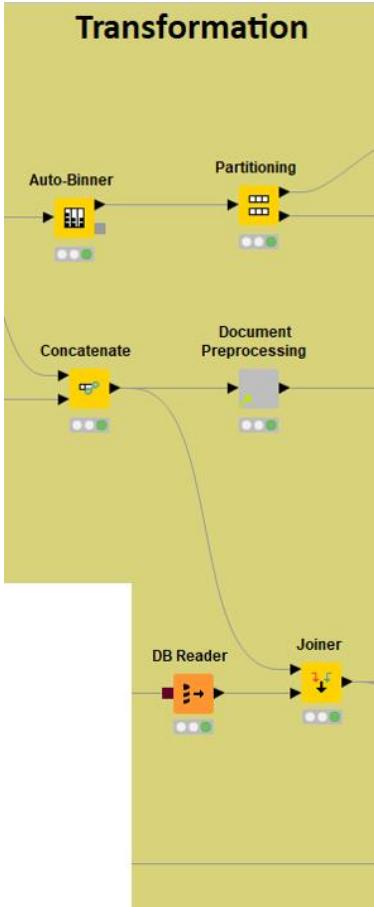
---



- Spark
- HDFS support
- Hive
- Impala
- In-database processing

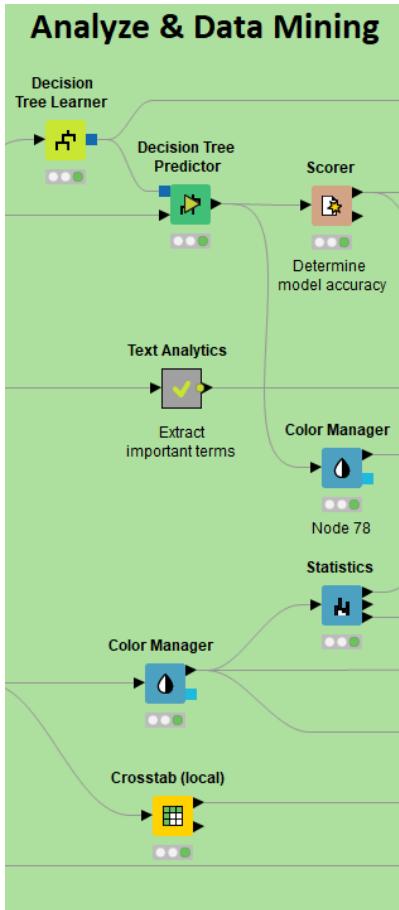


# Transformation



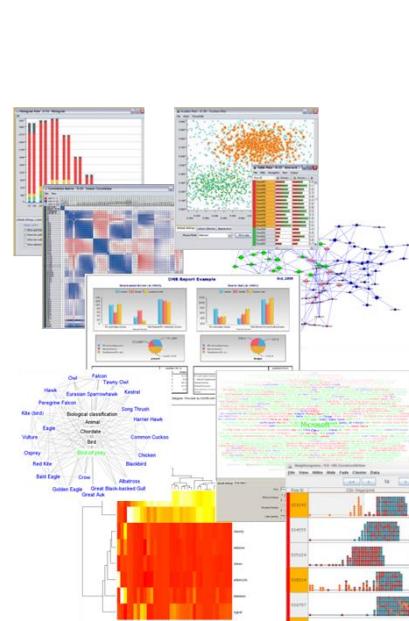
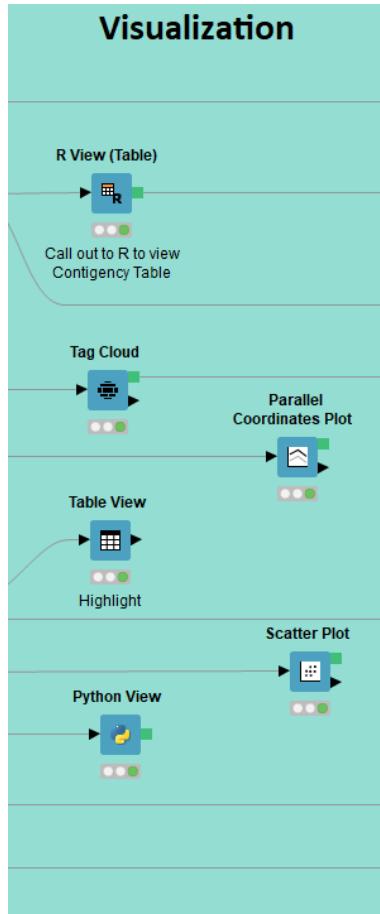
- Preprocessing
  - Row, column, matrix based
- Data blending
  - Join, concatenate, append
- Aggregation
  - Grouping, pivoting, binning
- Feature Creation and Selection

# Analysis & Data Mining



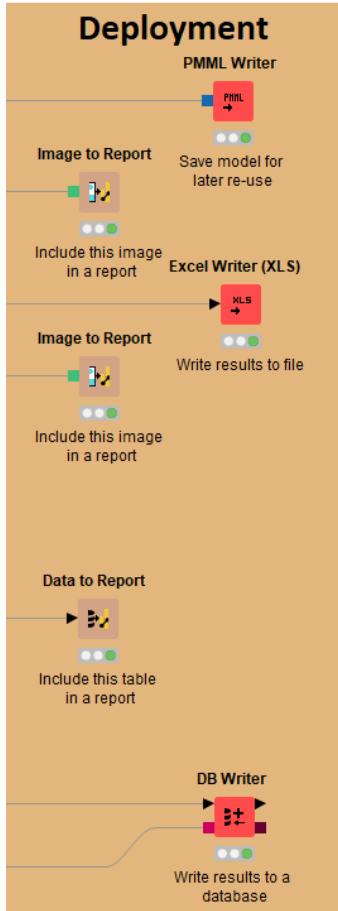
- Regression
  - Linear, logistic
- Classification
  - Decision tree, ensembles, SVM, MLP, Naïve Bayes
- Clustering
  - k-means, DBSCAN, hierarchical
- Validation
  - Cross-validation, scoring, ROC
- Deep Learning
  - Keras, DL4J
- External
  - R, Python, Weka, H2O, Keras

# Visualization



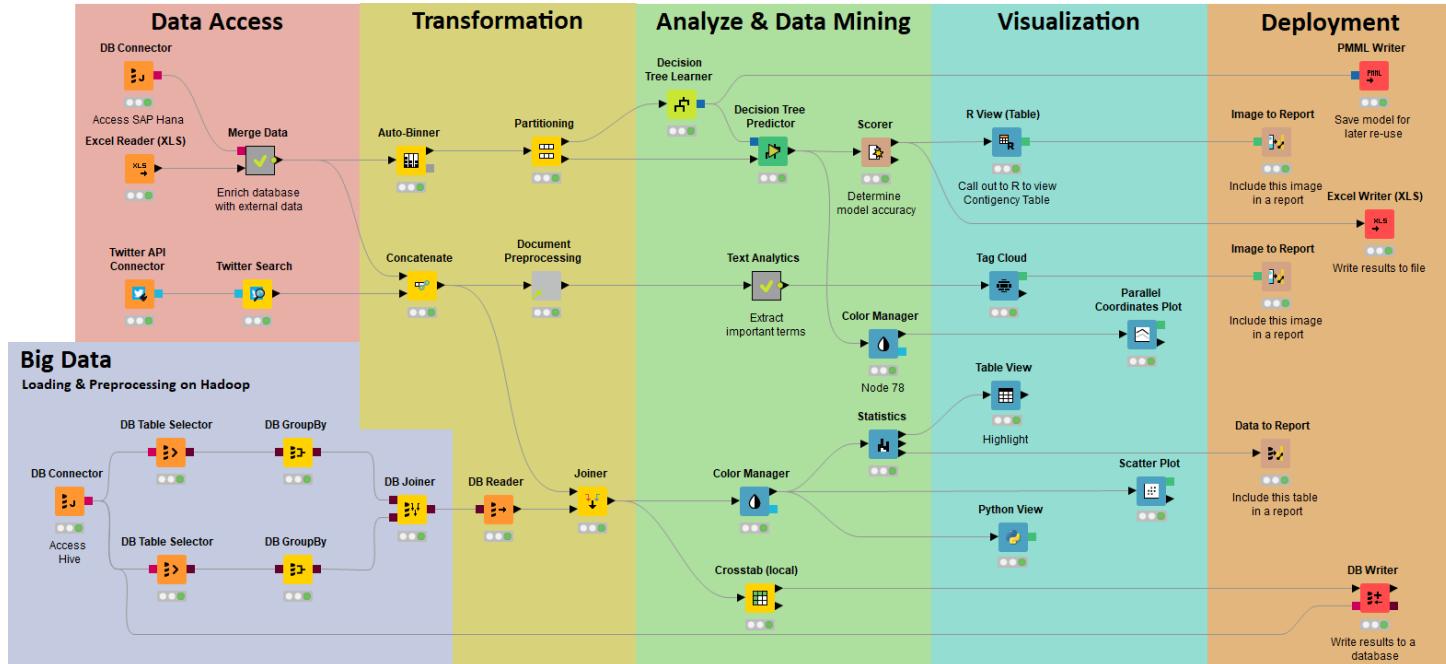
- Interactive Visualizations
- JavaScript-based nodes
  - Scatter Plot, Box Plot, Line Plot
  - Networks, ROC Curve, Decision Tree
  - Plotly Integration
  - Adding more with each release!
- Misc
  - Tag cloud, open street map, molecules
- Script-based visualizations
  - R, Python

# Deployment



- Database
- Files
  - Excel, CSV, txt
  - XML
  - PMML
  - to: local, KNIME Server, SSH-, FTP-Server
- BIRT Reporting

# Over 2000 Native and Embedded Nodes Included:



## Data Access

MySQL, Oracle, ...  
SAS, SPSS, ...  
Excel, Flat, ...  
Hive, Impala, ...  
XML, JSON, PMML  
Text, Doc, Image, ...  
Web Crawlers  
Industry Specific  
Community / 3rd

## Transformation

Row  
Column  
Matrix  
Text, Image  
Time Series  
Java  
Python  
Community / 3rd

## Analysis & Mining

Statistics  
Data Mining  
Machine Learning  
Web Analytics  
Text Mining  
Network Analysis  
Social Media Analysis  
R, Weka, Python  
Community / 3rd

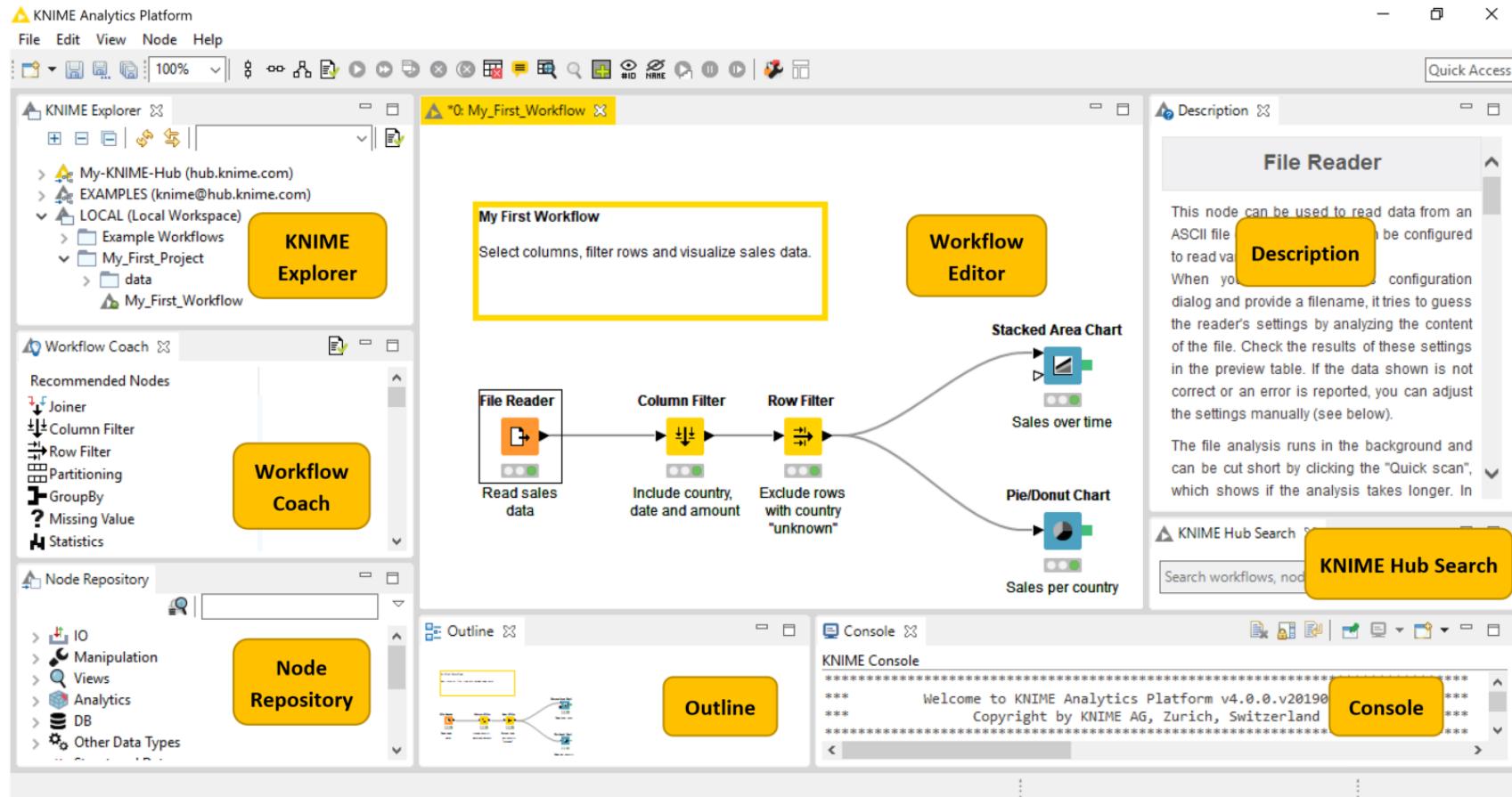
## Visualization

R  
JFreeChart  
JavaScript  
Plotly  
Community / 3rd

## Deployment

via BIRT  
PMML  
XML, JSON  
Databases  
Excel, Flat, etc.  
Text, Doc, Image  
Industry Specific  
Community / 3rd

# The KNIME Workbench

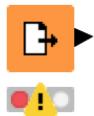


# More on Nodes...

---

A node can have 3 states:

**File Reader**



## Not Configured:

The node is waiting for configuration or incoming data.

**File Reader**



## Configured:

The node has been configured correctly, and can be executed.

**File Reader**

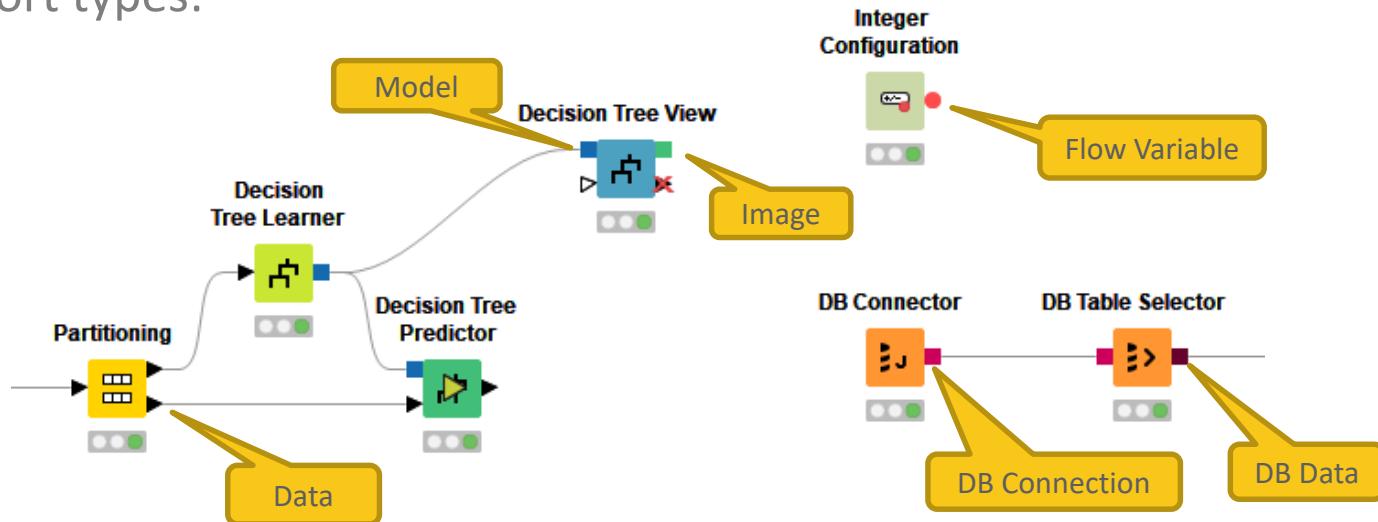


## Executed:

The node has been successfully executed. Results may be viewed and used in downstream nodes.

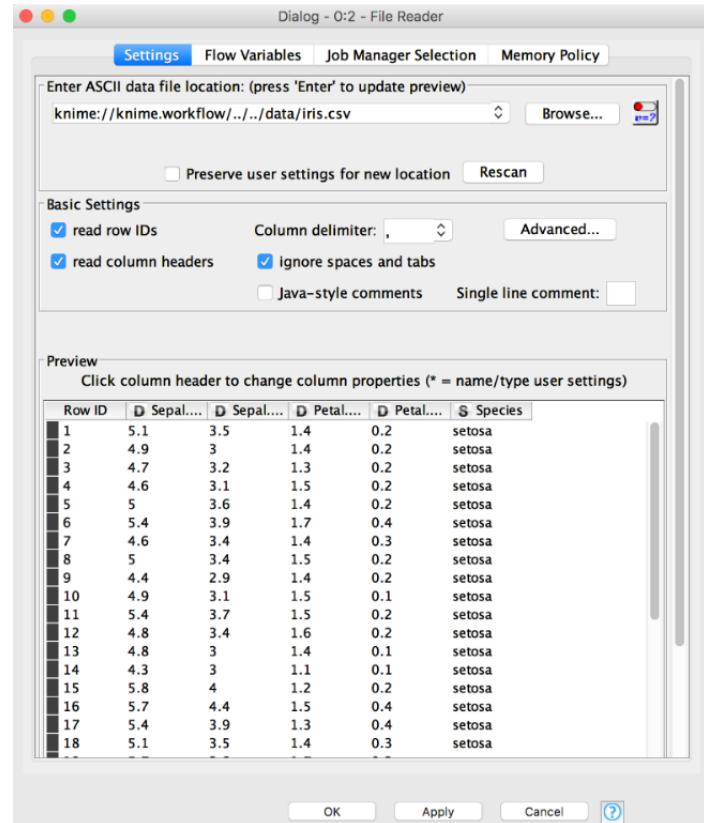
# Inserting and Connecting Nodes

- Insert nodes into workspace by dragging them from Node Repository or by double-clicking in Node Repository
- Connect nodes by left-clicking output port of Node A and dragging the cursor to (matching) input port of Node B
- Common port types:



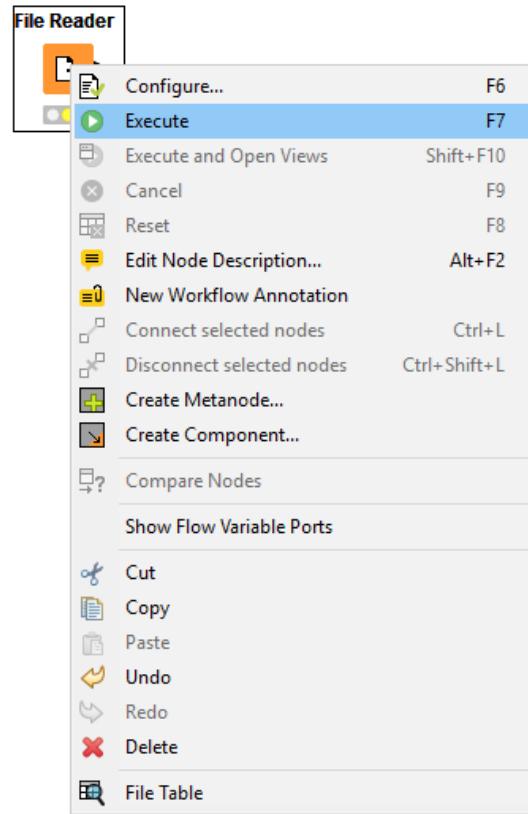
# Node Configuration

- Most nodes require configuration
- To access a node configuration window:
  - Double-click the node
  - Right-click -> Configure



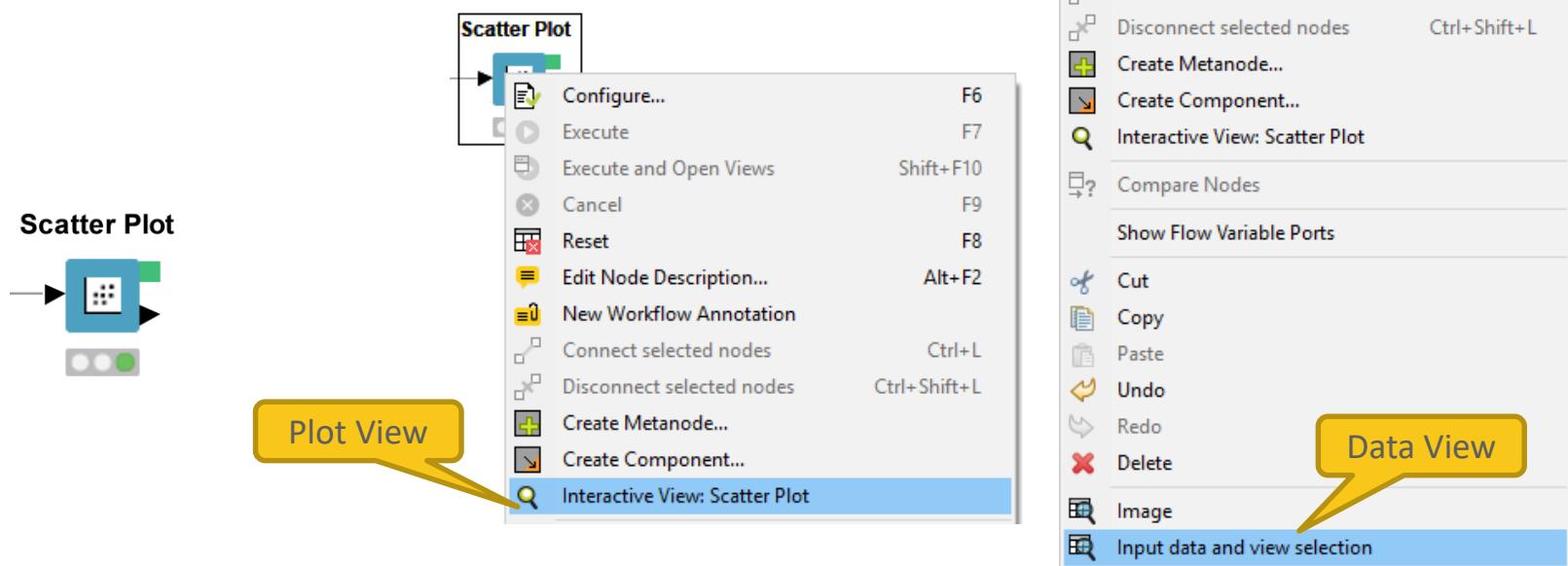
# Node Execution

- Right-click node
- Select Execute in the context menu
- If execution is successful, status shows green light
- If execution encounters errors, status shows red light



# Node Views

- Right-click node
- Select Views in context menu
- Select output port to inspect execution results

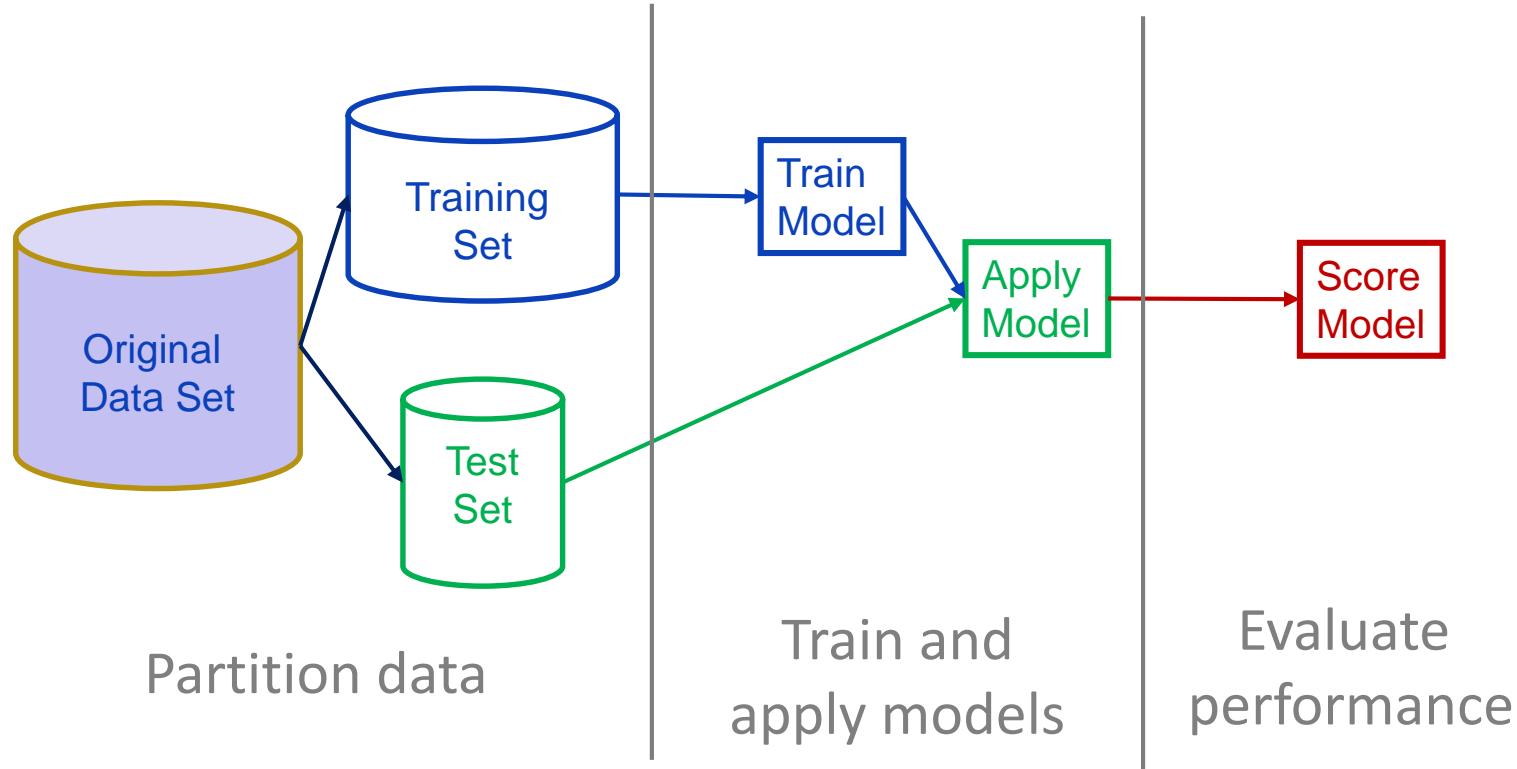


# Data Mining

Partition, Learn, Predict, Score

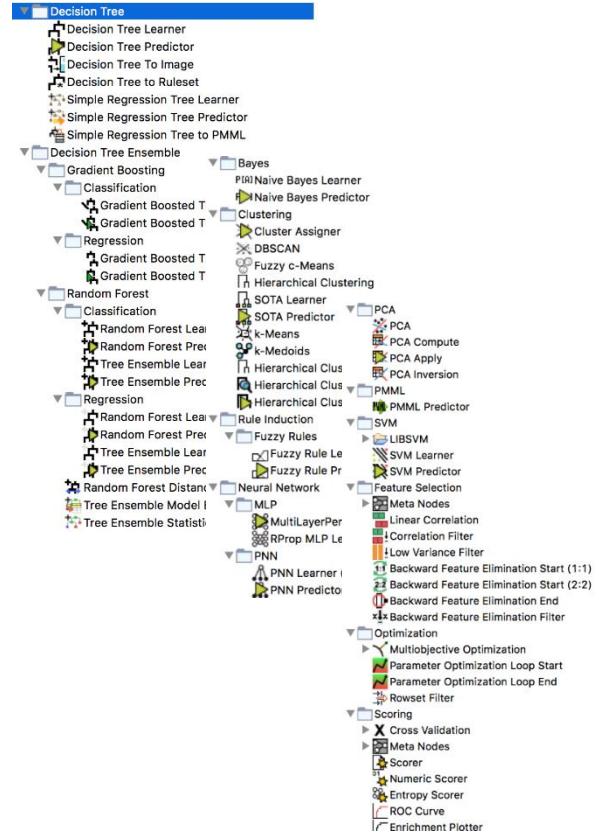


# Data Mining: Process Overview



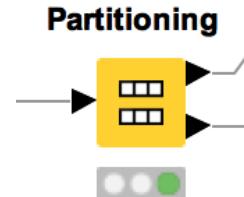
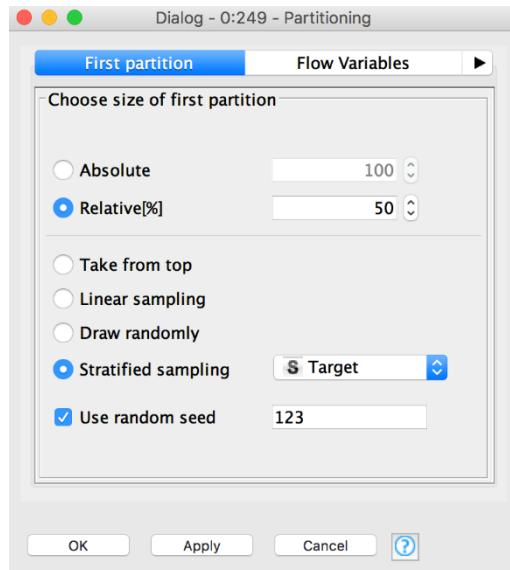
# Data Mining in KNIME

- KNIME has many modeling tools!
  - Decision tree, random forest, SVM, regression, neural networks, clustering, ...
  - and integrations with other libraries: R, Python, H2O, WEKA, libSVM, etc.
- And many model evaluation nodes
  - ROC, standard, numeric and entropy scorers
  - Feature elimination
  - Cross validation



# New Node: Partitioning

- Use to split data into training and evaluation sets
  - Partition by count (e.g. 10 rows) or fraction (e.g. 10%)
  - Sample by a variety of methods; random, linear, stratified



Two tables representing partitions of the 'default' dataset (Rows: 5775, Spec - Columns: 13). The columns are Row ID, Marita..., Gender, Estim..., Numbr..., and Age.

**First partition (as defined in dialog) - 0:249 - Partitioning**

Row ID	Marita...	Gender	Estim...	Numbr...	Age
Row0	M	M	90000	0	44
Row7	M	M	60000	2	46
Row9	S	M	70000	1	46
Row10	S	F	70000	1	46
Row13	M	M	100000	3	42
Row14	S	F	100000	3	42
Row15	S	F	30000	1	31
Row17	S	F	20000	2	66
Row18	S	M	30000	2	66
Row20	S	M	40000	2	32
Row21	S	F	40000	1	32

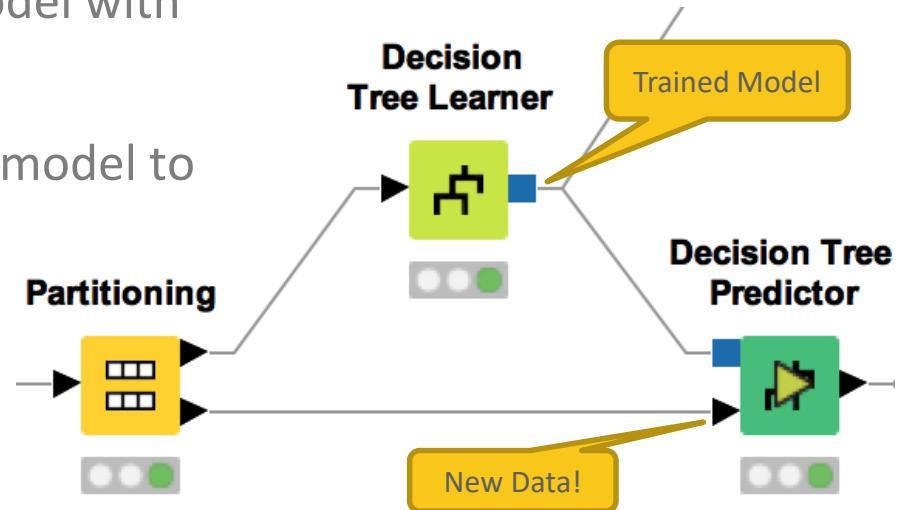
**Second partition (remaining rows) - 0:249 - Partitioning**

Row ID	Marita...	Gender	Estim...	Numbr...	Age
Row1	S	M	60000	1	45
Row2	M	M	60000	1	45
Row3	S	F	70000	1	42
Row4	S	F	80000	4	42
Row5	S	M	70000	1	45
Row6	S	F	70000	1	44
Row8	S	F	60000	3	46
Row11	M	M	60000	4	46
Row12	M	F	100000	2	42
Row16	M	M	30000	1	31
Row17	S	M	40000	2	32
Row19	S	M	30000	2	64

# Learner-Predictor Motif

---

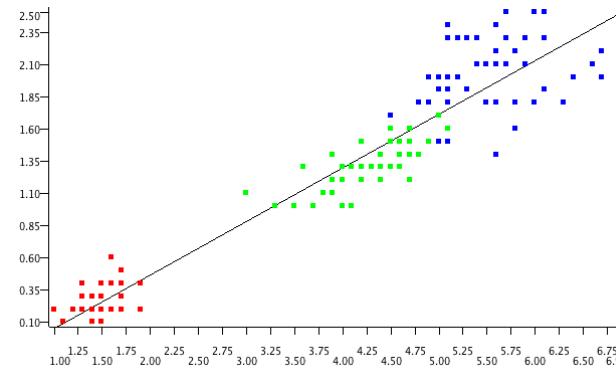
- Most data mining approaches in KNIME use a Learner-predictor motif.
- The Learner node trains the model with its input data.
- The Predictor node applies the model to a different subset of data.



# Regression

Predict *numeric* outcomes on existing data (supervised)

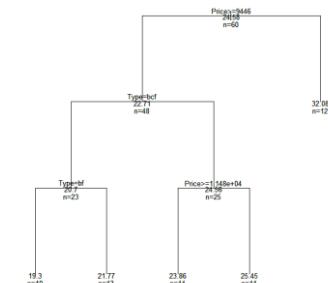
- Applications
  - Forecasting
  - Quantitative Analysis
- Methods
  - Linear
  - Polynomial
  - Regression Trees
  - Partial Least Squares



Statistics on Linear Regression

Variable	Coeff.	Std. Err.	t-value	P> t
Petal.Length	0.4158	0.0096	43.3872	0.0
Intercept	-0.3631	0.0398	-9.1312	4.44E-16

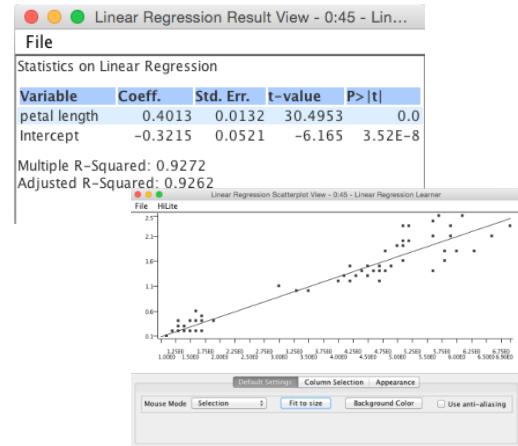
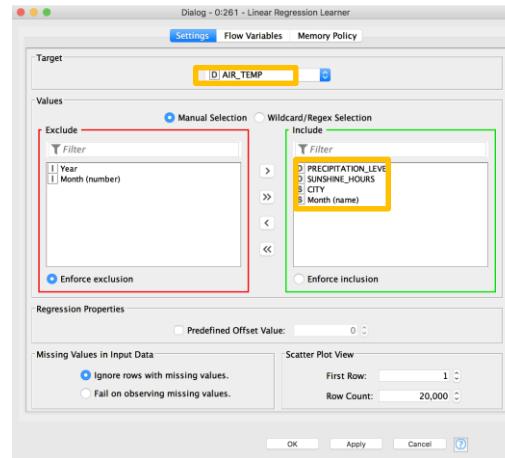
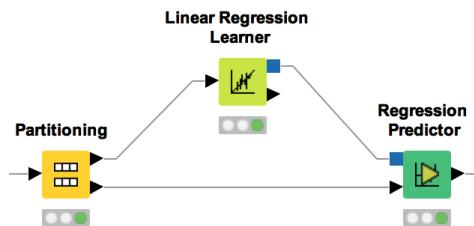
Multiple R-Squared: 0.9271  
Adjusted R-Squared: 0.9266



# New Nodes: Linear Regression Learner & Regression Predictor

A linear model relating a dependent variable to 1 or more independent variables

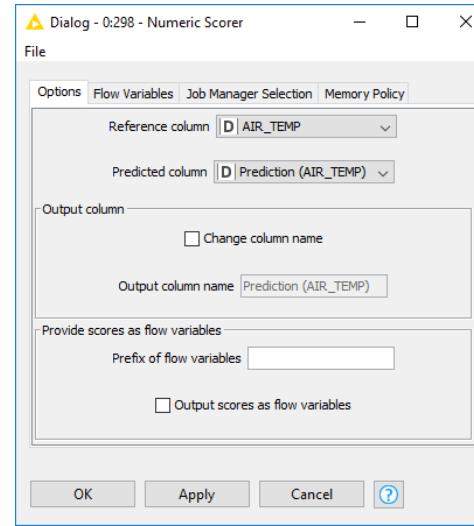
- Model coefficients provided in 2nd output port
- Also available: Polynomial and Tree Ensemble Regression nodes



# New Node: Numeric Scorer

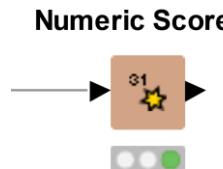
Similar to scorer node, but for nodes with *numeric* predictions (e.g. linear/polynomial regression)

- Compare dependent variable values to predicted values to evaluate goodness of fit.
- Report R<sup>2</sup>, RMSD, SEM etc.



The window title is 'Statistics - 0:298 - Numeric Scorer'. It shows a table titled 'Table "Scores" - Rows: 6 Spec - Column: 1 Properties Flow Variables'. The table has two columns: 'Row ID' and 'Prediction (AIR\_TEMP)'. The data rows are:

Row ID	Prediction (AIR_TEMP)
R^2	0.333
mean absolute error	3.574
mean squared error	21.329
root mean squared error	4.618
mean signed difference	1.048
mean absolute percentage error	NaN

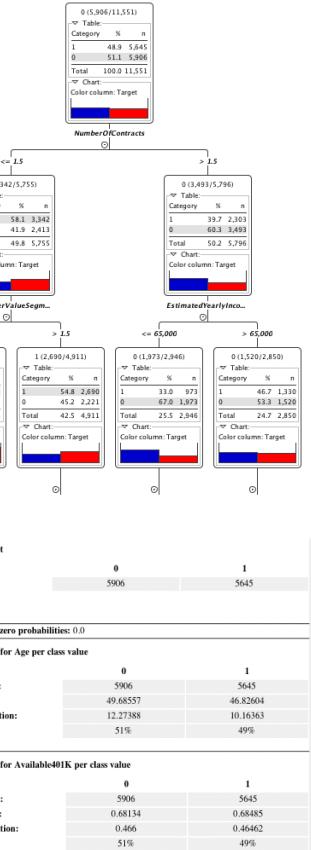
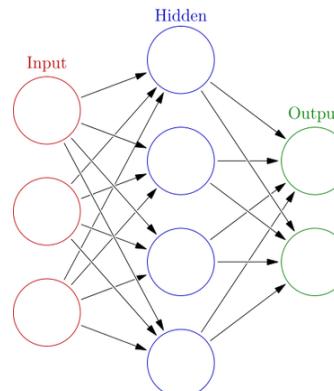


# Classification

Predict *nominal* outcomes on existing data (supervised)

- Applications
  - Churn analysis (yes/no)
  - Chemical activity (active/inactive)
  - Spam detection (spam/not spam)
  - Optical character recognition (A-Z)

- Methods
  - Decision Trees
  - Neural Networks
  - Naïve Bayes
  - Logistic Regression



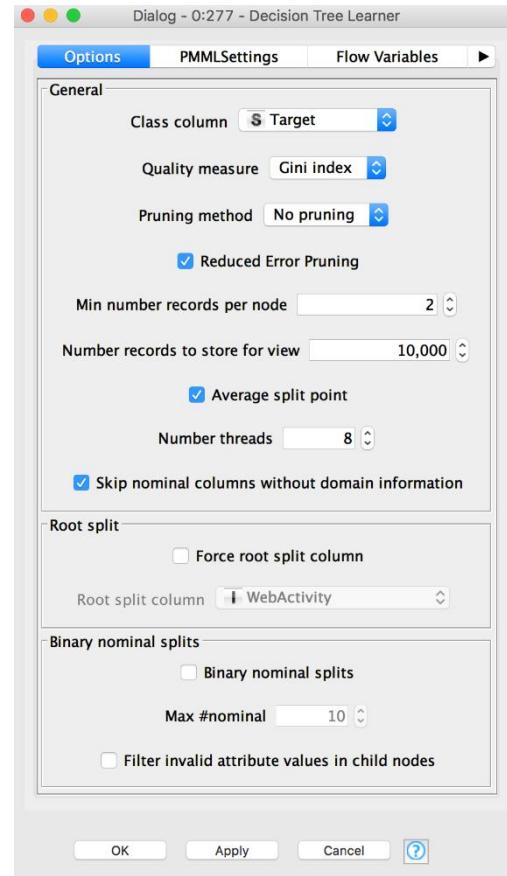
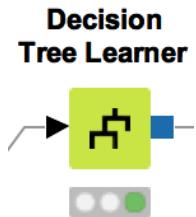
# Target Column

- Target column contains values that are predicted by the classification model
- Binomial target values are often encoded to 1 and 0

Application	Target Column	Target Values
Churn analysis	Churn	Yes/No or 1/0
Chemical activity	Active	Yes/No or 1/0
Spam Detection	Spam	Yes/No or 1/0
Optical Character Recognition	Character	A-Z

Output data - 0:311 - Column Resorter						
File Hilite Navigation View						
Table "default" - Rows: 5776 Spec - Columns: 17 Properties Flow Variables						
R...	I CustomerKey	S Marital...	S Gender	S Target	S Prediction (Target)	
...	11001	S	M	1	0	
...	11002	M	M	1	0	
...	11003	S	F	1	1	
...	11004	S	F	1	0	
...	11005	S	M	1	1	
...	11006	S	F	1	1	
...	11008	S	F	1	0	
...	11011	M	M	1	0	
...	11012	M	F	0	1	
...	11016	M	M	1	1	

# New Node: Decision Tree Learner



# Decision Tree View

Decision Tree View - 0:277 - Decision Tree Learner

File HiLite Tree

The decision tree starts at the root node (0 (2,953/5,775)) which splits based on the feature *NumberOfContracts*. The left branch ( $\leq 1.5$ ) leads to a node with 1 (1,694/2,924) and a table showing Category 1 (48.9%) and Category 0 (51.1%). The right branch ( $> 1.5$ ) leads to a node with 0 (1,723/2,851) and a table showing Category 1 (39.6%) and Category 0 (60.4%). A yellow callout points to the right branch with the text: "Most of the people who don't churn have more than one contract". The tree structure is shown on the right side of the interface.

0 (2,953/5,775)

Table:

Category	%	n
1	48.9	2,822
0	51.1	2,953
Total	100.0	5,775

*NumberOfContracts*

$\leq 1.5$

$> 1.5$

1 (1,694/2,924)

Table:

Category	%	n
1	57.9	1,694
0	42.1	1,230
Total	50.6	2,924

0 (1,723/2,851)

Table:

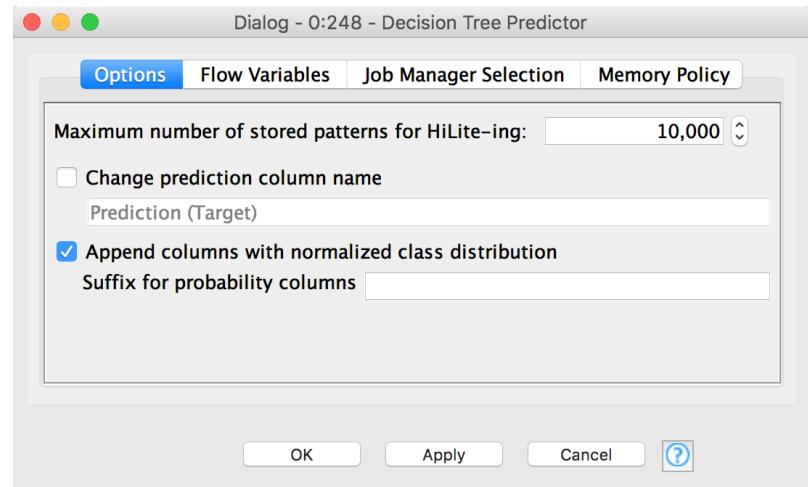
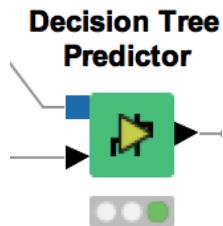
Category	%	n
1	39.6	1,128
0	60.4	1,723
Total	49.4	2,851

Most of the people who  
don't churn have more  
than one contract

Zoom: 100.0%

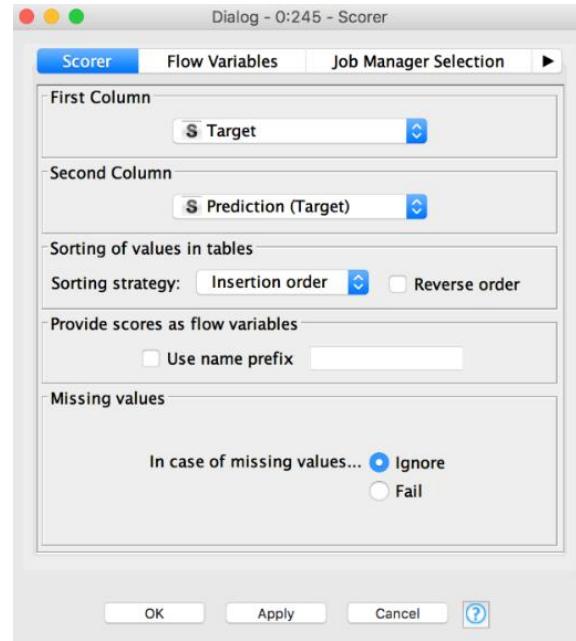
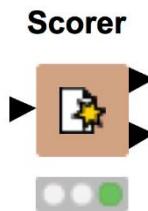
# New Node: Decision Tree Predictor

- Takes a decision tree model & applies it to new data
- Check the box to append class probabilities



# New Node: Scorer

Compare predicted results to known truth in order to evaluate model quality



# New Node: Scorer

Confusion matrix shows the distribution of model errors

Confusion Matrix - 0:297 - Scorer		
File	Hilite	
Target \ Prediction (Target)	1	0
1	2073	750
0	759	2193

Correct classified: 4,266

Accuracy: 73.87 %

Cohen's kappa ( $\kappa$ ) 0.477

Wrong classified: 1,509

Error: 26.13 %

An accuracy statistics table provides a detailed analysis of model quality

Accuracy statistics - 0:297 - Scorer													
File	Hilite	Navigation	View	Table "default" – Rows: 3 Spec – Columns: 11 Properties Flow Variables									
Row ID	TruePositives	FalsePositives	TrueNegatives	FalseNegatives	Recall	Precision	Sensitivity	Specificity	F-measure	Accuracy	Cohen's kappa		
1	2073	759	2193	750	0.734	0.732	0.734	0.743	0.733	?	?		
0	2193	750	2073	759	0.743	0.745	0.743	0.734	0.744	?	?		
Overall	?	?	?	?	?	?	?	?	?	0.739	0.477		

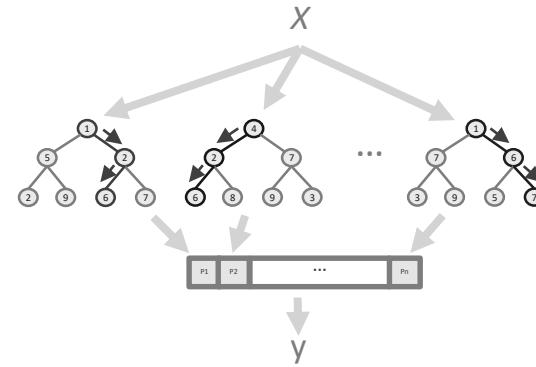
# Confusion Matrix

---

	Predicted class <b>POSITIVE</b> (churn)	Predicted class <b>NEGATIVE</b> (no churn)
Actual class <b>POSITIVE</b> (churn)	TRUE POSITIVE (TP)  2073	FALSE NEGATIVE (FN)  750
Actual class <b>NEGATIVE</b> (no churn)	FALSE POSITIVE (FP)  759	TRUE NEGATIVE (TN)  2193

# KNIME's Tree Ensemble Models

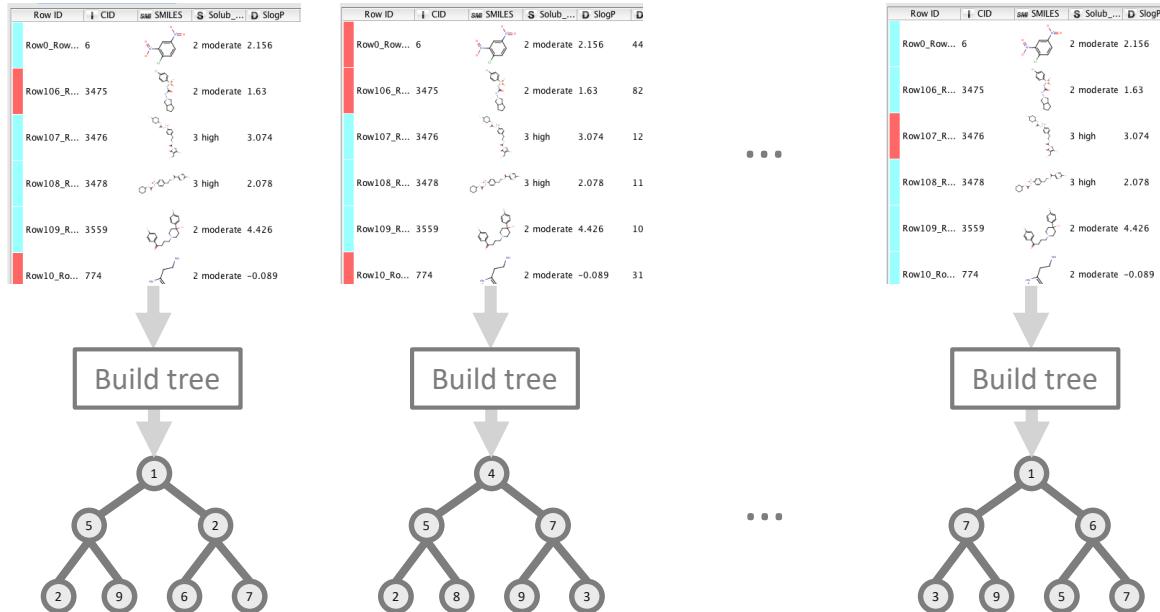
- The general idea is to take advantage of the “wisdom of the crowd”
- Ensemble models: Combining predictions from a large number of weak predictors, e.g. decision trees
- Leads to a more accurate and robust model
- This is called “bagging”



Typically: for classification the individual models vote and the majority wins; for regression, the individual predictions are averaged

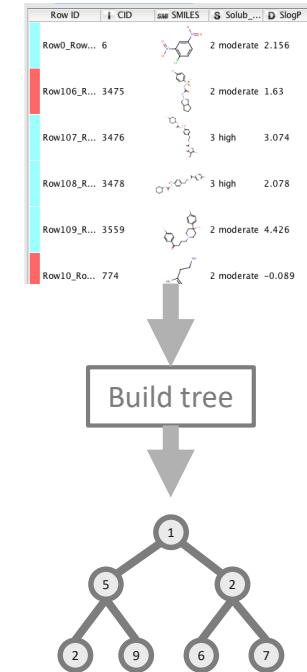
# How Does Bagging Work?

- Pick a different random subset of the training data for each model in the ensemble (bag)



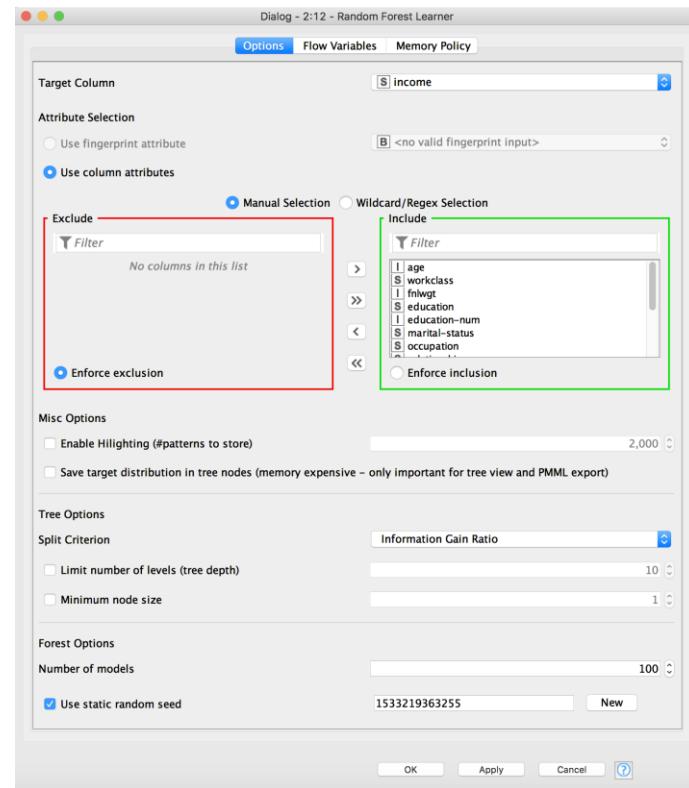
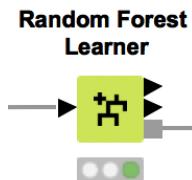
# Random Forest

- Bag of decision trees, with an extra element of randomization when building the trees: **each node in the decision tree only “sees” a subset of the input columns**, typically  $\sqrt{N}$
- Random forests tend to be very robust w.r.t. overfitting (though the individual trees are almost certainly overfit)
- Extra benefit: training tends to be much faster



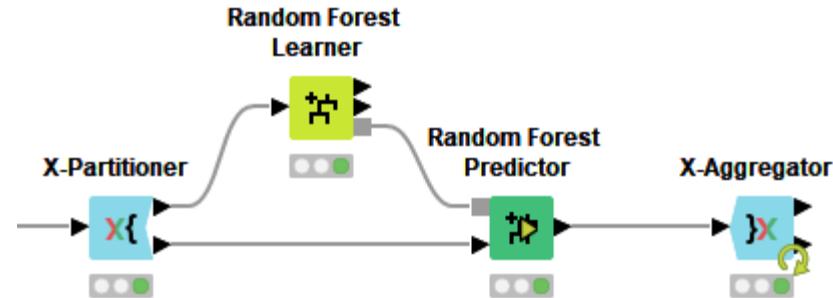
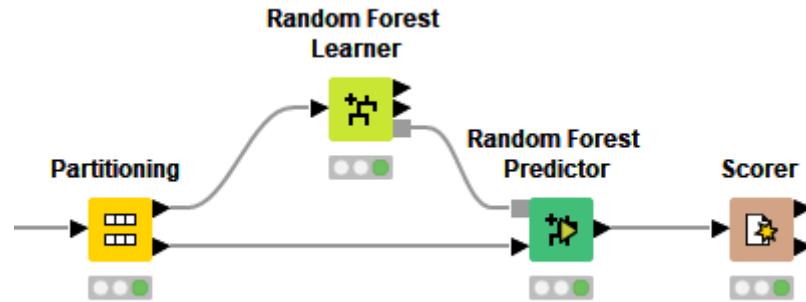
# New Nodes: Random Forest Learner

- The output model describes a random forest and is applied in the corresponding predictor node using a simple majority vote
- The statistics table on the attributes tells how often each attribute...
  - ... is used in the first three splits
  - ... was a possible candidate in the first three splits



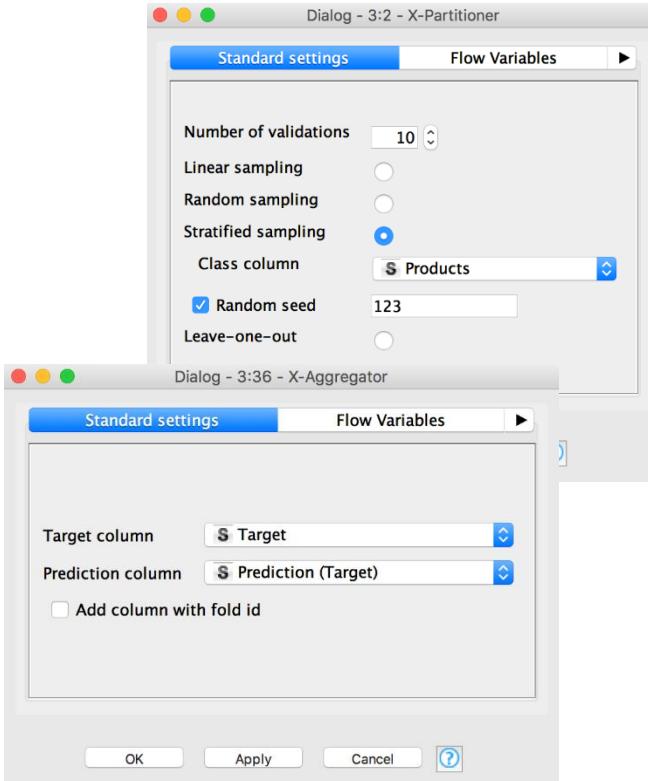
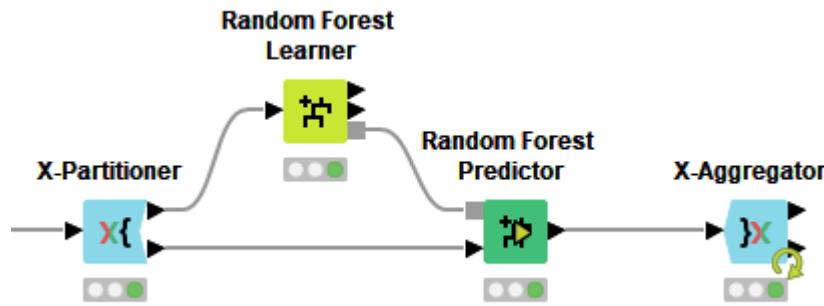
# Cross Validation

- Used to evaluate model stability
- Re-execute the modeling process many times using different data partitions
- Collect aggregated statistics on model accuracy



# Example: Cross Validation

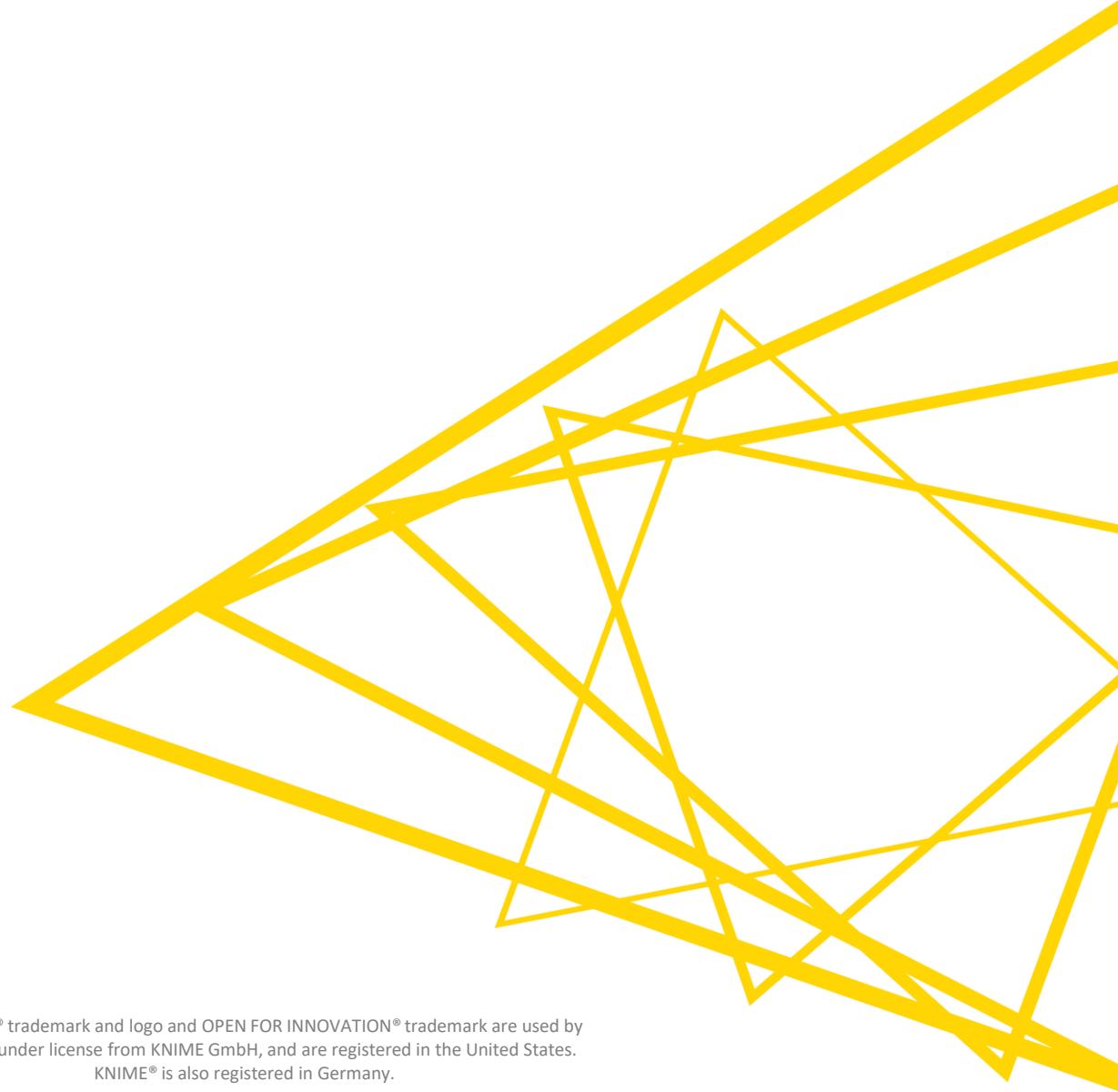
- X-Partitioner → X-Aggregator
- X-Partitioner replaces Partition
- X-Aggregator replaces Scorer
- Can be used with any learner/predictor





# Thank You!

[education@knime.com](mailto:education@knime.com)



The KNIME® trademark and logo and OPEN FOR INNOVATION® trademark are used by KNIME AG under license from KNIME GmbH, and are registered in the United States. KNIME® is also registered in Germany.