Open for Innovation ®

# KNIME

# Guided Analytics for Machine Learning Automation

**Christian Dietz, Simon Schmid**

KNIME

# Automating Everything?

**Human Input Needed!**
- Data Selection – Is this relevant?!
- Analysis Goal – What is interesting?
- Exploration – This looks weird?...

**Automating Data Integration:**
- Parsing
- Record Matching
- ...

**Legacy Data**

**In-house Data**

**Cloud**

**Data Blending**

**Cleaning & Transforming**

**Analysis**

**Explore**

**Deploy**

**Automating Data Proc:**
- Feature Selection
- Feature Construction
- Data Cleaning

**Automating Analytics:**
- Parameter Optimization
- Model Selection
- Ensemble Construction

# Automating Everything?

**Human Input Needed!**
- Data Selection – Is this relevant?!
- Analysis Goal – What is interesting?
- Exploration – This looks weird?...

**Automating Data Integration:**
- Parsing
- Record Matching
- ...



Legacy Data

In house Data

Cloud

Data Blending

Cleaning & Transforming

Analysis

Explore

Deploy

**Automating Data Proc:**
- Feature Selection
- Feature Construction
- Data Cleaning

**Automating Analytics:**
- Parameter Optimization
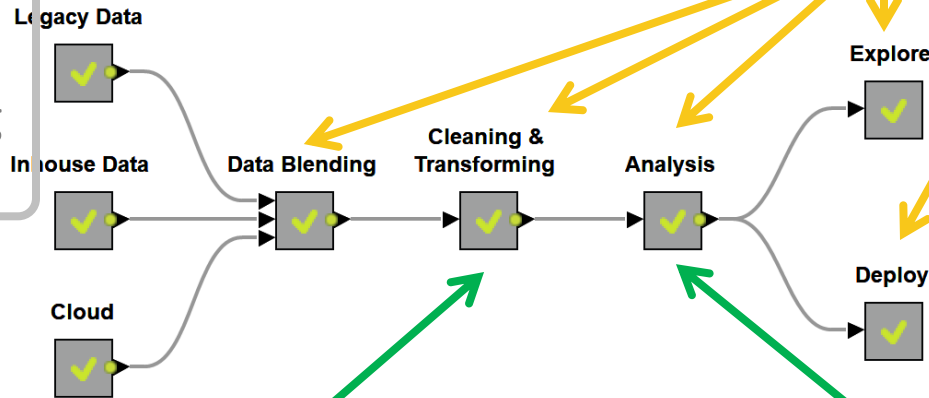- Model Selection
- Ensemble Construction

Open for Innovation ®
KNIME

# Building a Guided Automation Workflow

## Interaction Points



## Automated

# Guided Automation: Automation + Interaction

# Guided Automation: Automation + Interaction

# Guided Automation: Automation + Interaction

# Guided Automation: Automation + Interaction

# Guided Automation: Automation + Interaction

**WebPortal**



**Workflow** ✅



**KNIME Server**

Open for Innovation
KNIME ®

# Guided Automation on KNIME Server

**Live Demo**

# Scoring Workflow

# Scoring Workflow



## KNIME Server

# Customize the Blueprint for Text Processing



**Upload Data and Process Setup**

1. Upload your data / define file path
2. Select the target column for the prediction
3. Filter columns to exclude from the model
4. Select models and whether to fine-tune the model parameters

**Fine-tune Model Parameters or Use Automatic Settings**

If you selected to fine-tune the model parameters:
1. Define the parameter values used for optimization
2. Define feature engineering settings
Otherwise the settings will be automatically set.

**Execution Settings**

Decide which execution environment you want to use.

**ML Automation**

Model parameters are automatically optimized and features engineered.

**Download Models**

Compare and inspect the results of the models and download the desired ones.

**Added by in-house expert.**

# Thank You!



**Upload Data and Process Setup**

1. Upload your data / define file path
2. Select the target column for the prediction
3. Filter columns to exclude from the model
4. Select models and whether to fine-tune the model parameters

**Fine-tune Model Parameters or Use Automatic Settings**

If you selected to fine-tune the model parameters:
1. Define the parameter values used for optimization
2. Define feature engineering settings
Otherwise the settings will be automatically set.

**Execution Settings**

Decide which execution environment you want to use.

**ML Automation**

Model parameters are automatically optimized and features engineered.

**Download Models**

Compare and inspect the results of the models and download the desired ones.

Upload Dataset · Select Target · Feature Quality Calculation · Filter Columns · Select Models · IF Switch · Parameter Settings · Feature Engineering Settings · Automatic Parameter Settings · Automatic Feature Engineering Settings · End IF · Execution Settings · Training and Validation of Models · Download Models

# Download workflow from **knime.com** and get started!

Open for Innovation ®

KNIME

Christian

https://datascience1.knime.com/knime/#/Guided_Analytics_for_ML_Automation/01_Guided_Analytics_for_ML_Automation?exec=9b056c55-f35a-420e-b9b9-aaf7289fcf57&single

# KNIME WebPortal
Open for Innovation

Administration    Logout

## KNIME
Open for Innovation

Upload
Dataset

Select
Target

Filter
Columns

Select
Models

Parameter
Settings

Feature Eng.
Settings

Execution
Settings

Download
Models

## Upload Dataset

Upload the dataset to be used.

**Change File**    Selected file "airline.csv" (4 MB)

## Guide

### Upload Dataset

Upload the dataset to train the model. The dataset must be a representative sample of the past data history for the prediction problem at hand. Your file should be a KNIME table or a CSV file. The data will be uploaded onto the server for further processing.

REST API    © KNIME AG, Switzerland - Version 4.7.0

Open for Innovation
KNIME

https://datascience1.knime.com/knime/#/Guided_Analytics_for_ML_Automation/01_Guided_Analytics_for_ML_Automation?exec=9b056c55-f35a-420e-b9b9-aaf7289fcf57&single

**KNIME WebPortal** — Open for Innovation

Administration  Logout

**KNIME** — Open for Innovation

Upload Dataset → Select Target → Filter Columns → Select Models → Parameter Settings → Feature Eng. Settings → Execution Settings → Download Models

## Select Target

Select the target column whose values should be predicted.

**Select:**

IsDepDelayed ▼

| Row ID | Month | DayOfWeek | UniqueCarrier | Cancelled | CancellationCode | Diverted | WeatherDelay | SecurityDelay | IsArrDelayed | IsDepDelayed |
|--------|-------|-----------|---------------|-----------|------------------|----------|--------------|---------------|--------------|--------------|
| Row0 | 10 | 3 | PS | 0 | NA | 0 | ? | ? | YES | YES |
| Row1 | 10 | 4 | PS | 0 | NA | 0 | ? | ? | YES | NO |
| Row2 | 10 | 6 | PS | 0 | NA | 0 | ? | ? | YES | YES |
| Row3 | 10 | 7 | PS | 0 | NA | 0 | ? | ? | NO | NO |
| Row4 | 10 | 1 | PS | 0 | NA | 0 | ? | ? | YES | YES |
| Row5 | 10 | 3 | PS | 0 | NA | 0 | ? | ? | NO | NO |
| Row6 | 10 | 4 | PS | 0 | NA | 0 | ? | ? | YES | NO |
| Row7 | 10 | 5 | PS | 0 | NA | 0 | ? | ? | YES | YES |
| Row8 | 10 | 6 | PS | 0 | NA | 0 | ? | ? | YES | YES |
| Row9 | 10 | 7 | PS | 0 | NA | 0 | ? | ? | YES | NO |

## Guide

### Select Target

To train a model, a target column must be selected. In case of doubt, explore the preview (top 10 rows) on the left of all possible target columns for this prediction problem.

### Details for Experts

The possible columns can be of different types. Different Machine Learning methods can be used with different types of targets. Based on the type of the selected target column, only valid methods for the prediction task will be surfaced.

List of prediction tasks:

- If the target column is a categorical feature with 2 possible values, then it is a *Binary Classification* task.
- If the target column is a categorical feature with more than 2 possible values, then it is a *Multiclass Classification* task.
- If the target column is a numerical feature with only 2 different numbers, then it is a *Binary Classification* task.
- If the target column is a numerical feature with more than 2 and not more than 20 different numbers, then it is a *Multiclass Classification* task.

REST API   © KNIME AG, Switzerland - Version 4.7.0

**KNIME** WebPortal

Open for Innovation

Administration    Logout

Open for Innovation
**KNIME**

Upload Dataset — Select Target — Filter Columns — Select Models — Parameter Settings — Feature Eng. Settings — Execution Settings — Download Models

## Filter Columns

### Set Column Relevance Filter

Use the slider to select a subset of columns based on their relevance. If in doubt, do not change.

41.07

0.00                                                                                          100.00

Overall Column Relevance

| | Feature Name | Overall Column Relevance | Correlation with Target (%) | ID/Noise Test (%) | Constant Value Test (%) | Missing Value Test (%) |
|---|---|---|---|---|---|---|
| | CRSElapsedTime | 96.44 | 28.756 | 3.04 | 3.56 | 0.04 |
| | ArrDelay | 96.39 | 54.009 | 2.42 | 3.61 | 2.69 |
| | ActualElapsedTime | 96.3 | 26.456 | 3.7 | 2.69 | 2.69 |
| | CRSDepTime | 95.19 | 18.621 | 4.81 | 2.49 | 0 |
| | Year | 94.96 | 29.379 | 0.21 | 5.04 | 0 |
| | Distance | 94.239 | 18.965 | 5.761 | 4.69 | 0.07 |
| | DayofMonth | 92.74 | 21.197 | 0.3 | 7.26 | 0 |
| | Origin | 91.66 | 26.04 | 1.27 | 8.34 | 0 |
| | CRSArrTime | 90.749 | 19.645 | 9.251 | 1.38 | 0 |

### Guide

#### Set Column Relevance Filter

By default, all columns will be used to train the model that creates the prediction. However, not all columns contribute with the same importance or relevance to the final prediction. In some cases, columns are not informative or contain spurious information. To help you decide, the overall column relevance towards the final prediction is measured.

- **Column Relevance** is an overall metric summarizing the metrics belows. Use the slider to select the input features based on their *Overall Column Relevance*.

The additional metrics calculated automatically and used to determine *Overall Column Relevance* include:

- **ID/Noise Test** measures how likely the column is a representation used to identify each row in your table. Row identifiers are uninformative for your model and should be removed.

- **Constant Value Test** measures how often the column contains the exact same value. Columns with just a constant value also carry no information. You should avoid using them.

- **Missing Value Test** measures the percentage of missing values in a column over the entire dataset. You should remove features with a percentage of missing values too high.

By using the slider, columns can be excluded from model training based on their column relevance.

Furthermore you can use the linear correlation between each column and the column to predict to refine your input set.

- **Correlation with Target** measures the linear correlation with the column the model will predict: *IsDepDelayed*. It is important to keep in mind if a feature is highly or poorly correlated. If you have high correlation (close to + or - 100%) this will help the model to achieve a good performance, unless the column has too many unique values (e.g. on high *ID/Noise Test*). If instead you have low

REST API          © KNIME AG, Switzerland - Version 4.7.0

Open for Innovation
**KNIME**

# Manually Select Columns

In addition, columns can be visually examined and then manually selected for exclusion below. If in doubt, do nothing.

Numeric    Nominal    Data Preview

Search: 

| Column | Exclude Column | Minimum | Maximum | Mean | Median | Standard Deviation | Variance | Skewness | Kurtosis | Overall Sum | No. zeros | No. missings |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ⊕ Year | ☐ | 1987 | 2008 | 1997.465 | 1997 | 6.316 | 39.894 | 0.007 | -1.193 | 19974654 | 0 | 0 |
| ⊕ Month | ☐ | 1 | 10 | 1.399 | 1 | 1.852 | 3.430 | 4.430 | 17.629 | 13987 | 0 | 0 |
| ⊕ DayofMonth | ☐ | 1 | 31 | 14.695 | 14 | 9.202 | 84.673 | 0.167 | -1.222 | 146948 | 0 | 0 |
| ⊕ DayOfWeek | ☐ | 1 | 7 | 3.843 | 4 | 1.910 | 3.648 | 0.154 | -1.092 | 38433 | 0 | 0 |
| ⊖ DepTime | ☐ | 1 | 2400 | 1345.779 | 1330 | 466.991 | 218080.269 | 0.085 | -1.117 | 13129419 | 0 | 244 |

No. NaN      0

No. +∞      0

No. -∞      0

**Histogram**

**KNIME** WebPortal

Open for Innovation
**KNIME**

Upload Dataset — Select Target — Filter Columns — Select Models — Parameter Settings — Feature Eng. Settings — Execution Settings — Download Models

## Select Models

Choose one or more machine learning models to train for your prediction task.

**Simple models**

☑ Naive Bayes

☑ Decision Tree

☑ Logistic Regression

**Complex models**

☐ Support Vector Machine

☑ Random Forest

☑ Generalized Linear Models

☑ Gradient Boosted Trees

☐ Deep Learning

## Fine-tune Model Parameters

By default (unchecked), all parameters for the selected models and for feature engineering are automatically fully optimized. However, by checking this option you can guide the automatic optimization process.

☑ **Finetune Model Parameters**

## Outlier Treatment

### Guide

#### Select Models

Choose which models you want to train. The available models have different levels of complexity. Less complex models are simpler to interpret and understand and generally faster to train and more efficient to use in production. In contrast, more complex models are capable of solving more complicated problems at possibly a finer level of detail but possibly at the cost of longer training times and less efficient usage in production. If the time is not an issue and you simply want to see the best performing model, use all of the proposed models, especially the ones with higher complexity. You can later compare model performances as well as runtime to choose the model that best solves your task. If a convenient solution is what you are aiming for, enabling only the simpler models will save you training time and will allow for a more efficient execution.

If a convenient solution is what you are aiming for, enabling only the simpler models will save you both time to create and be more efficient in executing.

#### Levels of Complexity

- **Simple models**
    - *Naive Bayes* is a simple probabilistic classifier based on the Bayes' Theorem.
    - *Decision Tree* is a simple to understand tree-like model which makes predicitons based on rules.
    - *Logistic Regression* is a statistical model which maximizes a likelihood function.
- **Complex models**
    - *Support Vector Machine* is a non-probabilistic linear classifier.
    - *Random Forest* is an ensemble learning method which constructs multiple decision trees.
    - *Generalized Linear Model* is a flexible generalization

**Complex models**

☐ Support Vector Machine

☑ Random Forest

☑ Generalized Linear Models

☑ Gradient Boosted Trees

☐ Deep Learning

## Fine-tune Model Parameters

By default (unchecked), all parameters for the selected models and for feature engineering are automatically fully optimized. However, by checking this option you can guide the automatic optimization process.

☑ **Finetune Model Parameters**

## Outlier Treatment

By default (checked), outliers are removed automatically. By unchecking this option, outliers are not removed.

☑ **Automatically Remove Outliers**

---

well as runtime to choose the model that best solves your task. If a convenient solution is what you are aiming for, enabling only the simpler models will save you training time and will allow for a more efficient execution.

If a convenient solution is what you are aiming for, enabling only the simpler models will save you both time to create and be more efficient in executing.

**Levels of Complexity**

- **Simple models**

  ○ *Naive Bayes* is a simple probabilistic classifier based on the Bayes' Theorem.

  ○ *Decision Tree* is a simple to understand tree-like model which makes predicitons based on rules.

  ○ *Logistic Regression* is a statistical model which maximizes a likelihood function.

- **Complex models**

  ○ *Support Vector Machine* is a non-probabilistic linear classifier.

  ○ *Random Forest* is an ensemble learning method which constructs multiple decision trees.

  ○ *Generalized Linear Model* is a flexible generalization of linear regression models.

  ○ *Gradient Boosted Trees* is a complex ensemble learning method which constructs multiple decision trees.

  ○ *Deep Learning* is a complex non-linear multi-level neural network.

---

**Fine-tune Model Parameters**

By default (unchecked), all parameters for the selected models and for feature engineering are automatically fully optimized. However, by checking this option you can guide the automatic optimization process.

When checked, you will be presented with optimization options for the parameters of the selected models. In addition, you will also be presented with options for the creation of additional feature columns.

**Outlier Treatment**

By default (checked), outliers are removed automatically. By unchecking this option, outliers are not removed.

Open for Innovation ®
**KNIME** WebPortal

Administration    Logout

Open for Innovation ®
**KNIME**

Upload Dataset — Select Target — Filter Columns — Select Models — Parameter Settings — Feature Eng. Settings — Execution Settings — Download Models

## Parameter Settings

### Gradient Boosted Machine

Set the parameter ranges for the Gradient Boosted Machine.

30 ——————————— 70

5                                                    100

Number of Trees

12 —————————————————————— 20

8                                                    20

Maximal Depth

0.08 —————————————————————— 0.10

0                                                    0

Learning Rate

## Guide

### Parameter Settings

Each machine learning model has its own unique set of parameters. In keeping with the selected models, set the appropriate ranges for the optimization of the exposed parameters. The larger the set of values, the better the model. In contrast, smaller sets of values lead to a faster runtime.

Changing these settings is optional.

Listed here you can find info regarding each model available parameters.

- Random Forest
  - **Number of Trees** controls the number of trees to learn.
  - **Maximal Depth** limits the number of tree levels to be learned.

- Gradient Boosted Machine
  - **Number of Trees** controls the number of trees to learn.
  - **Maximal Depth** limits the number of tree levels to be learned.
  - **Learning Rate** specifies the rate at which the model is learned.

- Decision Tree
  - **Maximal Depth** limits the number of tree levels to be learned.

- Logistic Regression
  - Regularization controls which regularization prior to use:
  - **Uniform**: this prior corresponds to no regularization

REST API        © KNIME AG, Switzerland - Version 4.7.0

Open for Innovation ®
**KNIME**

Christian

https://datascience1.knime.com/knime/#/Guided_Analytics_for_ML_Automation/01_Guided_Analytics_for_ML_Automation?exec=9b056c55-f35a-420e-b9b9-aaf7289fcf57&single

# KNIME WebPortal

Administration    Logout

## Open for Innovation KNIME

Upload Dataset — Select Target — Filter Columns — Select Models — Parameter Settings — Feature Eng. Settings — Execution Settings — Download Models

# Feature Engineering Settings

## Select Techniques

Please select the feature engineering techniques you want to use.

**Select:**

☑ Simple Transformations

☑ Feature Combinations

☑ Dimensionality Reduction

☑ Cluster Distance Transformation

## Aggressiveness

Please select the level of *aggressiveness* of the feature engineering.

| 0.35 |

0.0          0.5          1.0

## Guide

### Feature Engineering Settings

#### Select Techniques

In many cases, a model can be improved by creating new data columns from existing ones. This is called feature engineering. There is a number of techniques available and you can select here the ones to be used:

- **Simple Transformations:** mathematical transformations are applied on numerical features (e.g., *logarithm*, *exponential*, *square*, *tanh*, etc.).

- **Feature Combinations:** features are combined by either adding, subtracting, multiplying, or dividing two numerical features.

- **Dimensionality Reduction:** a *Principal Component Analysis (PCA)* is applied on the selected features.

- **Cluster Distance Transformation:** the data is clustered by the selected features and for each data point the distance to a chosen cluster center is calculated.

**If you do not choose any techniques, no feature engineering will be performed.**

#### Aggressiveness

You can also determine how aggressively feature engineering should be practiced. Setting the slider higher means that more different subsets of features are used to train the models. Each subset is evaluated based on model performance. An increased level of aggressiveness leads to an increased runtime, but also to a wider exploration of the possible feature combinations. If set to zero, no feature engineering is performed at all. For simple prediction tasks optimizing just the model parameters without any - or less - feature engineering might already be sufficient to produce a good model.

**If you do not set the slider, the default value of 0.2 is applied.**

REST API    © KNIME AG, Switzerland - Version 4.7.0

Secure | https://datascience1.knime.com/knime/#/Users/simon.schmid/Guided_Analytics_for_ML_Automation/01_Guided_Analytics_for_ML_Automation_executed_95_prettyROC?exec=2ea9b0af-f9b7-4cd8-8c5e-679bb4413aa8&single

**KNIME** WebPortal

Administration    Logout

**KNIME** Open for Innovation

Upload Dataset — Select Target — Filter Columns — Select Models — Parameter Settings — Feature Eng. Settings — Execution Settings — Download Models

# Download Models

Here is a summary of information (performances) about the models trained based on your specifications. The first chart compares the accuracy and Area under the Curve of each model. The second chart compares the training times. The third chart compares the prediction time on a new record. The fourth chart shows the ROC (AUC). After the table to download the model parameters, a performance summary for each model is shown.

### Compare Model Metric Performance

This bar chart visualize different performance metrics to assess the quality of each model.

**Main Performance Metrics**

● Area Under Curve (%)   ● Accuracy (%)       zoom

Metric value (%)

100.0

80.0

60.0

40.0

20.0

0.0

Naive Bayes    Generalized Linear Model    Logistic Regression    Random Forest    Gradient Boosted Trees    Decision Tree

Model

## Guide

### Download Models

The models shown were trained based on your specifications.

Each model has its own hyperparameter optimization, feature engineering and feature selection based on either the automatic settings or the manual settings. By means of the visualizations below, compare the selected models to decide which model to deploy.

The first chart shows model accuracy and AUC. A higher accuracy is better than a lower accuracy. The amount of time needed to train a model is in the next chart. If a model is only trained occasionally then the amount of time may be irrelevant. If a new model needs to be re-trained more often (daily, hourly, etc.) then the training time may be important. The next chart shows the relative time it would take to apply the model to (or score on) a new record. When you have many records that need to be scored in a short amount of time and possibly at a frequent rate then this time is important.

The fourth chart shows the ROC curve. It provides an additional way of looking at model performance.

The information provided by the four charts should help you decide which model is most suitable. The model with the highest accuracy may take much longer to be applied to (or scored on) new data than less accurate models. These last ones however might execute faster. The "right" model for your situation will depend on a combination of all these factors.

### Details for Experts

You can compare the models now by different metrics. Above you can see a bar chart with two main measures of performance. *Accuracy* is the percentage of correct predictions among all predictions. *Area Under Curve (AUC)* measures the area under the Receiver Operating Characteristic (ROC) curve.

The ROC curve plot describes the *Receiver Operating Characteristic* curves, one for each model. On the y-axis you have the true positive rate, on the x-axis you have the false posive rate.

REST API       © KNIME AG, Switzerland - Version 4.7.0

Open for Innovation **KNIME**

Secure | https://datascience1.knime.com/knime/#/Users/simon.schmid/Guided_Analytics_for_ML_Automation/01_Guided_Analytics_for_ML_Automation_executed_95_prettyROC?exec=2ea9b0af-f9b7-4cd8-8c5e-679bb4413aa8&single

# KNIME WebPortal
Open for Innovation

Administration    Logout

## Compare Training and Prediction Times

The first bar chart compares the training times of all models. The second bar chart compares the prediction time for one single sample.

### Training Time

● Naive Bayes   ● Generalized Linear Model   ● Random Forest   ● Logistic Regression   ● Gradient Boosted Trees   ● Decision Tree

Time (s)

22.0
20.0

0.0

Models               Training Time (s)

### Prediction Time per Sample

● Random Forest   ● Naive Bayes   ● Logistic Regression   ● Generalized Linear Model   ● Gradient Boosted Trees   ● Decision Tree

Time (ms)

0.1
0.1

Prediction Time (ms)
● Random Forest   0.1

## Advanced Assessment of Models

The advanced assessment of models sections shows four additional charts per model.

- **1. Performance Metrics Bar Charts**

  For this visualization we measured the following metrics:

  - *Recall* (or True Positive Rate) (% of "NO" rows correctly classified)
  - *Precision* (or Positive Predicted Value) (% of predicted "NO" rows correctly classified)
  - *Specificity* (or True Negative Rate) (% of not "NO" rows correctly classified)
  - *F-measure* (harmonic average between Recall and Precision)

- **2. Cumulative Gain Chart and Lift Chart**

  This chart can display two different charts: the Cumulative Gain Chart and the Lift Chart. By default the cumulative gain chart is displayed. This chart is drawing a curve that reflects how well the model is doing compared to a random classifier. You are selecting rows from the test ranked by the probability of class "NO". On the x-axis you have the percentage of top ranked rows by the model that define the partition of rows you are considering. On the y-axis you measure the response as the percentage of "NO" rows over their total number in your partition of top ranked rows. If the model is bad, the curve will be close to the black line (random classifier), where the percentage of original "NO" rows is exactly equal the percentage of selected rows. The cumulative gain curve should be above the bisector line and the greater the area between the cumulative gain curve and the bisector line is, the better the model is.

  If you click on the top right corner of this chart, you will be able to visualize the relative lift chart as well. The lift on the y-axis measure the different between the cumulative gain chart curve and the bisector line.

- **3. Global Feature Importance Bar Chart**

  This chart shows the global feature importance. A surrogate random forest model is trained overfitting the test set predicted classes. From such a model it is possible to measure how often each feature is useful to outcome a prediction. In the chart the six most important features are shown whereby only features of the original data set are considered. More information at this link.

- **4. Confusion Matrix Heatmap**

  This chart shows a confusion matrix. A confusion matrix is summarizing all the predictions on the test set by considering how many instances fall in each cell according to prediction and ground truth. The heatmap is encoding with shades of blue the number of instances in each cell. A

REST API    © KNIME AG, Switzerland - Version 4.7.0

Open for Innovation
KNIME

Secure | https://datascience1.knime.com/knime/#/Users/simon.schmid/Guided_Analytics_for_ML_Automation/01_Guided_Analytics_for_ML_Automation_executed_95_prettyROC?exec=2ea9b0af-f9b7-4cd8-8c5e-679bb4413aa8&single
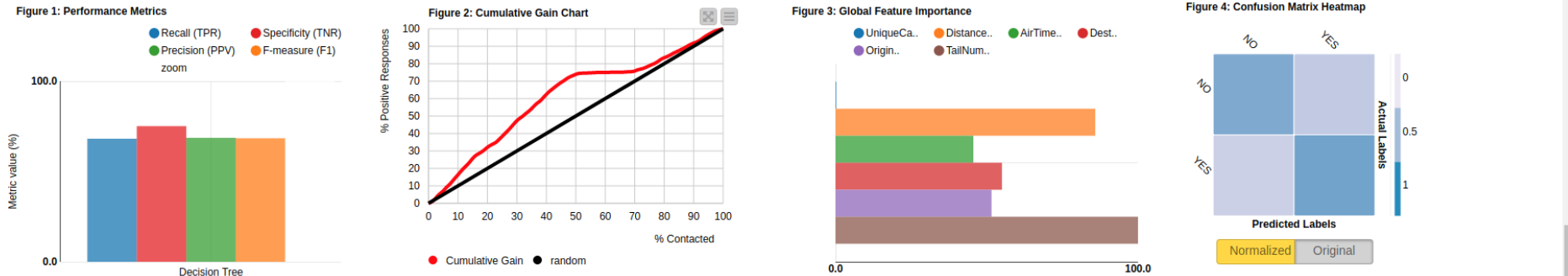
# KNIME WebPortal
Open for Innovation

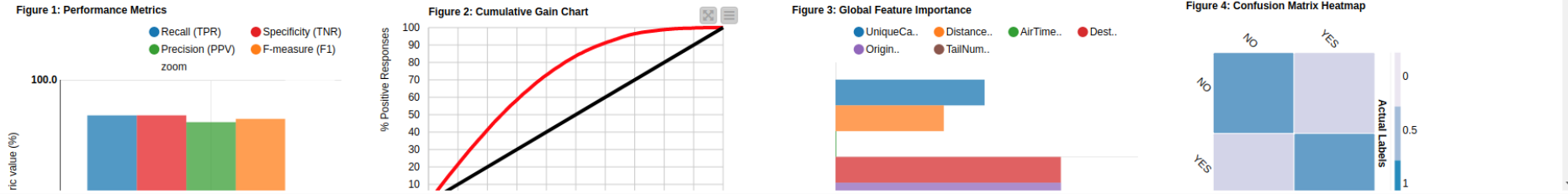Administration | Logout

## Advanced Assessment of Models

Each row represents a series of additional information about each trained model.

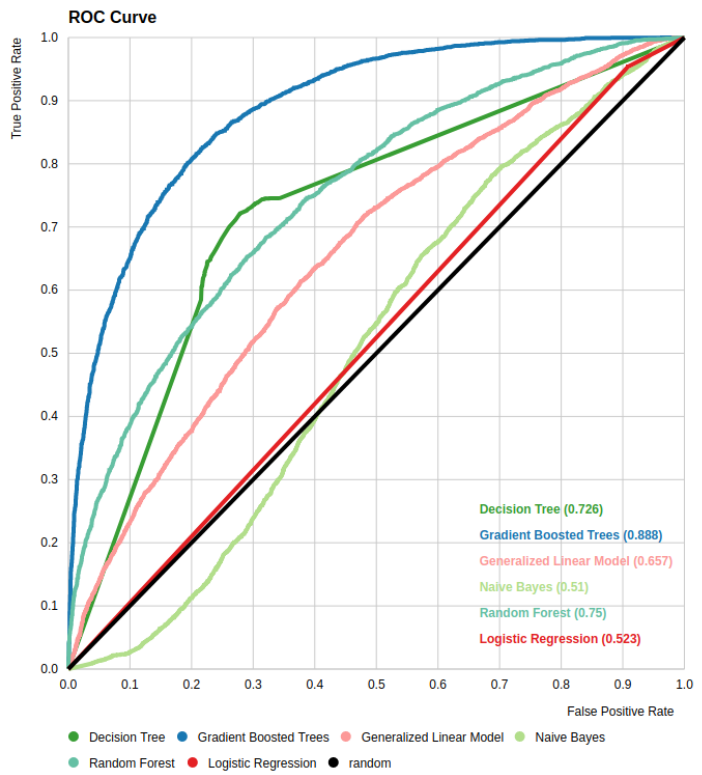- target feature: *IsArrDelayed*
- positive class: *NO*

### Decision Tree



**Figure 1: Performance Metrics**
- Recall (TPR)
- Specificity (TNR)
- Precision (PPV)
- F-measure (F1)

zoom

**Figure 2: Cumulative Gain Chart**
- Cumulative Gain
- random

**Figure 3: Global Feature Importance**
- UniqueCa..
- Distance..
- AirTime..
- Dest..
- Origin..
- TailNum..

**Figure 4: Confusion Matrix Heatmap**

Normalized | Original

### Gradient Boosted Trees



**Figure 1: Performance Metrics**
- Recall (TPR)
- Specificity (TNR)
- Precision (PPV)
- F-measure (F1)

zoom

**Figure 2: Cumulative Gain Chart**

**Figure 3: Global Feature Importance**
- UniqueCa..
- Distance..
- AirTime..
- Dest..
- Origin..
- TailNum..

**Figure 4: Confusion Matrix Heatmap**

REST API | © KNIME AG, Switzerland - Version 4.7.0

Open for Innovation
KNIME

Plots the ROC curves, one for each model. The greater the area under a curve the better the model is. To plot this chart the following settings for the target *IsArrDelayed* were automatically defined:

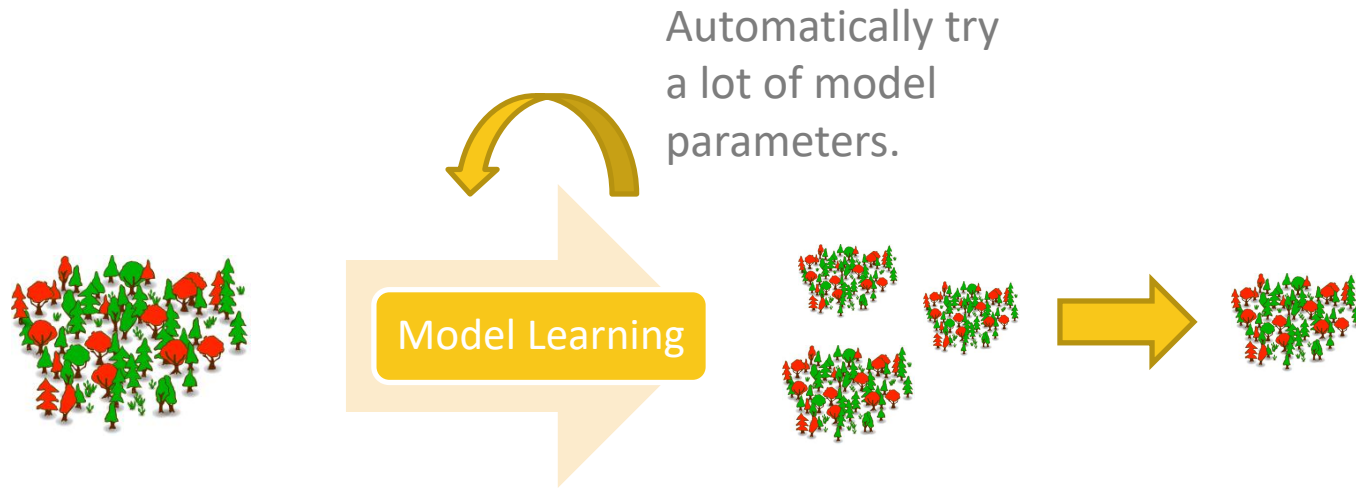- positive class: *NO*
- negative class: *YES*

**Download Model**

The following table summarizes the information in the charts. Please select the model you would like to download and use for predictions.

| Model | Accuracy (%) | Area Under Curve (%) | Prediction Time (ms) | Training Time (s) |
|---|---|---|---|---|
| Decision Tree | 86.542 | 72.644 | 0.1 | 22 |
| Gradient Boosted Trees | 83.996 | 88.813 | 0.1 | 6.8 |
| Random Forest | 72.433 | 75.017 | 0.1 | 6.4 |
| Logistic Regression | 63.345 | 52.267 | 0.1 | 6.6 |
| Generalized Linear Model | 61.754 | 65.704 | 0.1 | 3.4 |
| Naive Bayes | 52.152 | 50.969 | 0.1 | 1.9 |

Showing 1 to 6 of 6 entries

**Decision Tree**

Download Workflow

**Gradient Boosted Trees**

Download Workflow

**Generalized Linear Model**

Download Workflow

**Naive Bayes**

Download Workflow

**Random Forest**

Download Workflow

**Logistic Regression**

Download Workflow

## Advanced Assessment of Models

Each row represents a series of additional information about each trained model.

- target feature: *IsArrDelayed*
- positive class: *NO*

# What do we Automate?

- Data Cleansing
  - Missing value handling, calculate statistics, outlier detection

- Feature Engineering
  - Mathematical transformations, feature combinations, dimensionality reduction and more

- Feature Selection
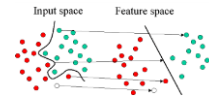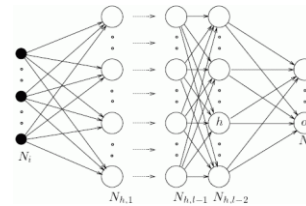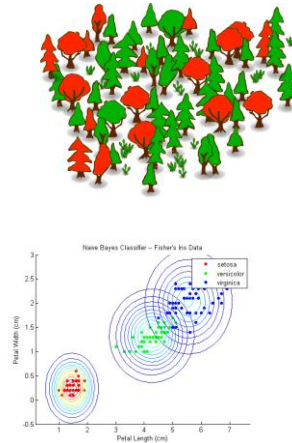  - Forward feature selection, backward feature elimination, genetic algorithm etc.

Open for Innovation ®
KNIME

# What do we Automate?

- Parameter Optimization



Automatically try a lot of model parameters.

Model Learning

Open for Innovation ®
KNIME

# What do we Automate?

- Model Selection
  - Try many models, but in an automated way.

# What do we Automate?

- Model Selection and Parametrization



Model Learning

**Check out Daria's Blogpost "Stuck in the Nine Circles of Hell? Try Parameter Optimization & A Cup Of Tea"** (05/28/18)

Model Learning