



COM5940: NEW MEDIA BUSINESS MODEL & INNOVATION REGRESSION AND CLASSIFICATION MODELS IN SUPERVISED LEARNING

Bernard Suen

Center for Entrepreneurship

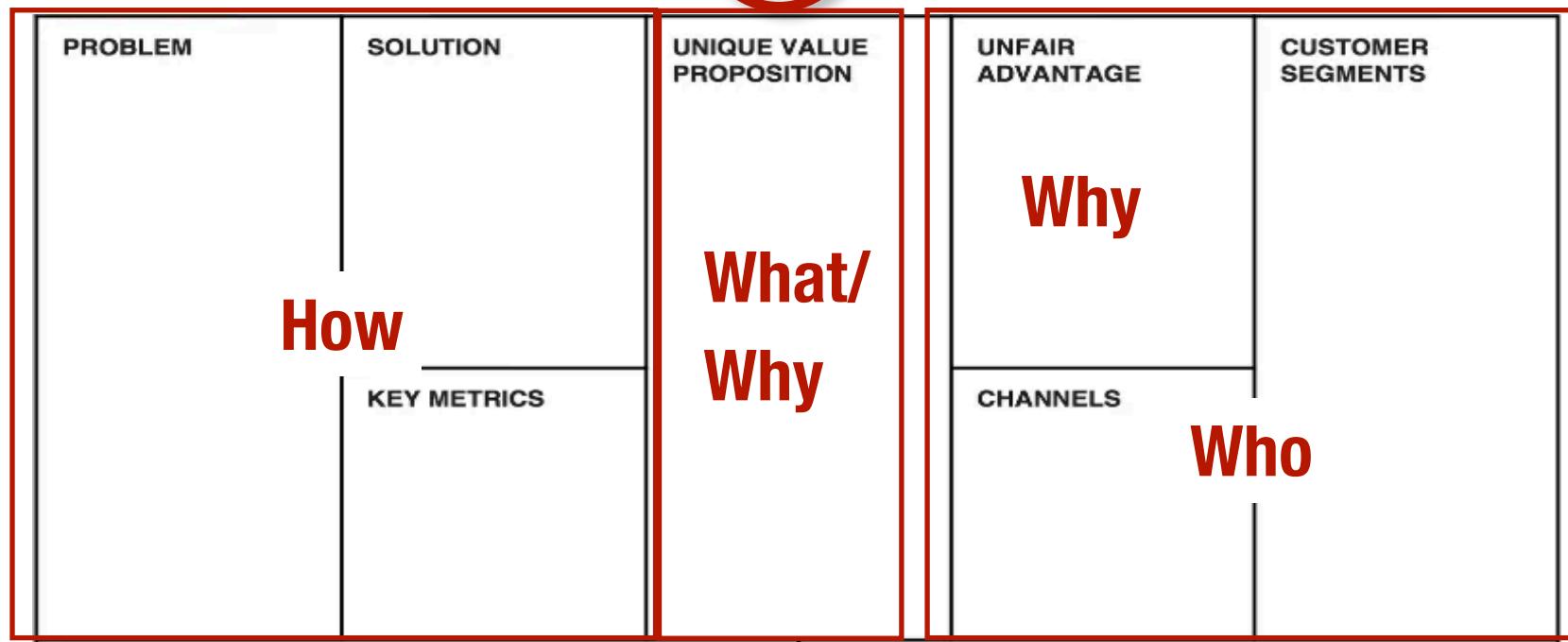
Chinese University of Hong Kong



Center for
Entrepreneurship

General Comments on 2nd Progress Report

Differences Between Business Model Canvas and Lean Canvas



Key Partners

Key Activities

Value Propositions

Customer Relationships

Customer Segments

Key Resources

Channels

Cost Structure

Revenue Streams

Problem

Solution

**Value
Propositions**

**Unfair
Advantage**

**Customer
Segments**

Key Metrics

Channels

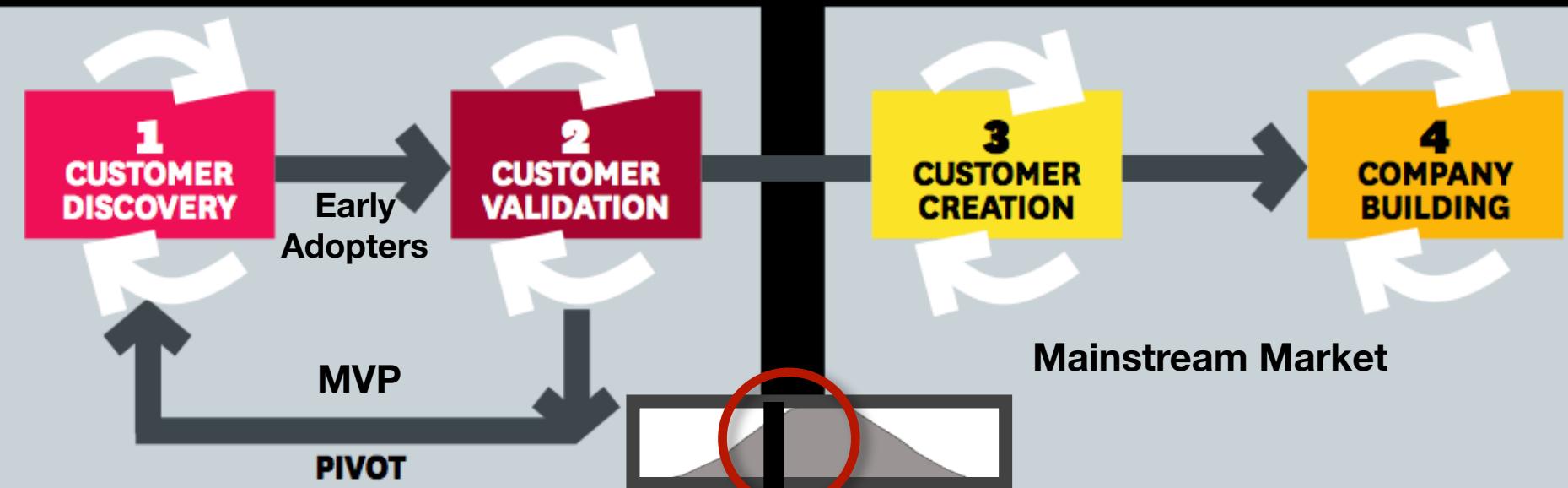
Cost Structure

Revenue Streams

“Why” the differences?

SEARCH

EXECUTION



Problem	Solution	Value Propositions	Unfair Advantage	Customer Segments	Key Partners	Key Activities	Value Propositions	Customer Relationships	Customer Segments
Key Metrics					Key Resources				
Cost Structure					Revenue Streams				

1. Some teams did not address the **2A3R** concerns.
2. Problem definition not **specific** enough.
3. **Data and content acquisition strategy** has not addressed the key issues (e.g. domain knowledge & expertise)
4. **Incentives** for KOL, brand and other potential business partners to work with you are not adequately addressed.
5. Only one team mentioned data scraping to address the **content shortage** problem.
6. Not addressing where **unfair advantage** can be exploited (except one team).

What has brought us here?

為什麼我们的學習會來到這一步!

API (e.g. REST)



BACK-END TECHNOLOGY

后台

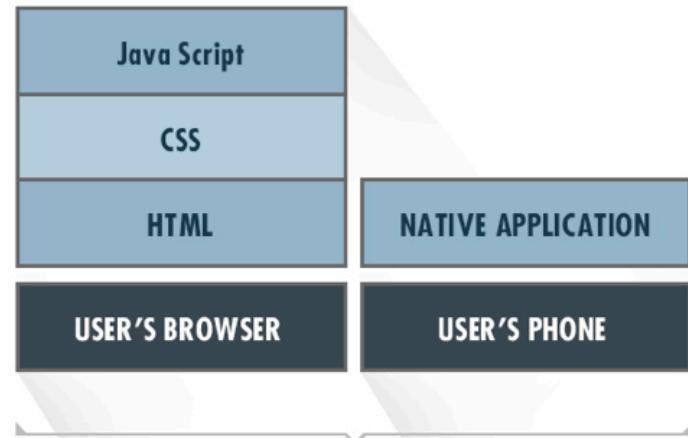
WHAT IS YOUR “CLOUD” AND “STACK” STRATEGY?

云架构及前端和后端的全栈策略

{ JSON }



THE INTERNET



CLIENT-SIDE

FRONT-END TECHNOLOGY

前台

云计算服务架构重组

Software as a Service

- SaaS
- Email, Map, Commerce
 - Group Productivity/SN
 - Matching/Lifestyle/Games

Function as a Service

- FaaS
- Serverless (e.g. AWS Lambda)
 - Process & scalability on demand
 - API/Event Driven/idle time no cost

Platform as a Service

- PaaS
- Web/App Server
 - Lang/Software Framework
 - Staging/API supports

Container as a Service

- CaaS
- Container Mgmt (e.g. Kubernetes)
 - Container (e.g. Docker)
 - Micro-service

Infrastructure as a Service

- IaaS
- Virtual Machines/OS
 - Compute/Storage/Network/Hardware

End-User Oriented

Clients

- Web Browser
- Mobile App
- IoT products

Developer Oriented

Public/Private/Hybrid Clouds

**Unsupervised
Learning**

**Supervised
Learning**

**Reinforcement
Learning**

What has brought us here?

為什麼我们的學習會來到這一步!



顶级程序员
319 篇文章 >

Python与人工智能的关系原来是这样的...

源 / AI时间 文 / 数据挖掘机

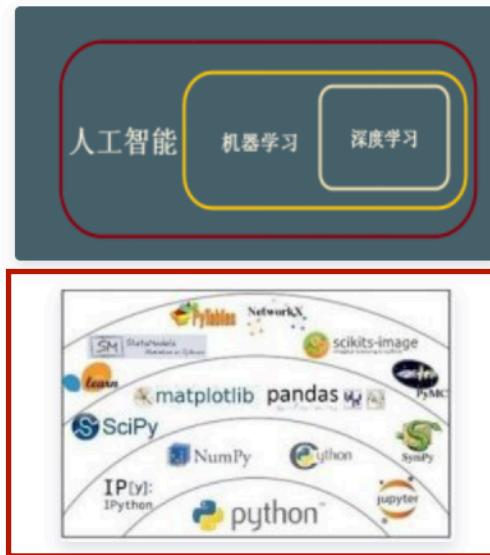
5

0

分享



人工智能掀起了世界的新一波科技浪潮，如今，你要要是不懂点AI、机器学习和python都不好意思说你是现代人，那么python究竟和人工智能什么关系 **为什么人工智能把python也给带火了** 今天就给大家简单介绍下python和人工智能的关系及应用，以及想要学人工智能的你，究竟需要学些什么python的知识，先来上两张图人工智能和python的图。



<https://cloud.tencent.com/developer/article/1161025>

从上图可以看出，人工智能包含常用机器学习和深度学习两个很重要的模块，而右图中python拥有matplotlib、Numpy、sklearn、keras等大量的库，像pandas、sklearn、matplotlib这些库都是做数据处理、数据分析、数据建模和绘图的库，基本上机器学习中对数据的爬取（scrapy）、对数据的处理和分析（pandas）、对数据的绘图（matplotlib）和对数据的建模（sklearn）在python中全都能找到对应的库来进行处理，所以，要想学习AI而不懂python，那就相当于想学英语而不认识单词，所以，python学起来吧。

PLATFORM (平台)





What data scientists do. We now know how data science works, at least in the tech industry. First, data scientists lay a solid data foundation in order to perform robust analytics. Then they use online experiments, among other methods, to achieve sustainable growth. Finally, they build machine learning pipelines and personalized data products to better understand their business and customers and to make better decisions. In other words, in tech, data science is about infrastructure, testing, machine learning for decision making, and data products. **数据变成产品和运营的一个重要构成部分**

Source: By Hugo Bowen-Anderson
August 15, 2018
Harvard Business Review

How does Machine Learning fit into the Problem Solving Framework (CRISP-DM)?

Source: Elements of User Experience
by Jesse James Garrett

Solution Space

how and
how much

Problem Space

who, what,
and why

The Surface Plane



The Skeleton Plane



The Structure Plane



The Scope Plane



The Strategy Plane



UX Elements

Concrete

Completion

Abstract

Conception

Visual Design

Interface Design Navigation Design
Information Design

Interaction Design Information Architecture

Functional Specifications Content Requirements

User Needs
Site Objectives

Deployment

Evaluation

Modeling

Data Preparation

Data Understanding

Business Understanding

time



生产环境中进行机器学习模型部署 (using Flask)



数据蛙datafrog
数据挖掘农民工

使用Flask REST API进行Python机器学习预测

2020-04-07 00:34:45

分类 : Python / Python Flask / 机器学习

阅读(11) 评论(0)

本文概述

- [通用机器学习架构](#)
- [机器学习训练与预测](#)
- [过度拟合的连接](#)
- [用于预测的REST API](#)
- [实施示例-自行车共享](#)
- [本文总结](#)

<http://www.srcmini.com/42615.html>



AI科技大本营

2.5K 篇文章 >

为什么说爱奇艺是一家技术公司？

AI无处不在

在近日举行的爱奇艺世界大会上，爱奇艺展示了AI技术和落地应用。爱奇艺首席技术官兼基础架构和智能内容分发事业群总裁刘文峰介绍称，在爱奇艺的整个运营流程中，AI发挥的作用无孔不入，贯穿视频内容的创作、生产、理解、分发、播放到变现以及客服整个流程中。

5
点赞0
收藏

分享

首先以智能生产环节为例。爱奇艺的内容主要来自三个部分，一部分是传统的采买，另外一部分是自制剧，以及爱奇艺号等原创用户产生的一些内容。对于自制部分，爱奇艺采用AI打造了一个爱创媒资系统，这个系统主要是针对导演、制作人的创作痛点。内容创作者在后期制作的过程中往往会被一些非原创性的工作所影响，以至于没有办法集中精力，而且拍摄素材繁多，导致1~2个小时的综艺后期制作的时间往往要花上8~10天才能最终看到成片。通过使用身份识别、表情识别、镜头识别等AI技术，类似工作可缩短至数分钟，剪辑师可以很方便地从繁杂的素材中挑选精彩片段，效率提高20%–30%。

比如，导演想要在成百上千个小时的素材里找到吴亦凡生气的特写，AI可能只需要三步就能解决：第一是身份识别，找出吴亦凡；第二是表情识别，识别出生气的表情；第三是镜头识别，识别特写镜头。这样就能很快把导演想要的片段找出来。

<https://cloud.tencent.com/developer/article/1430210>

三 菜单



云社区

专栏

问答 沙龙 快讯 团队主页 开发者手册 智能钛AI | 腾讯云大学 TVP



专栏首页 > 企鹅号快讯 > 号外！号外！Python纳入高考内容了！人工智能时代就要来临了！



10



0

分享



号外！号外！Python纳入高考内容了！人工智能时代就要来临了！<https://cloud.tencent.com/developer/article/1024699>

2018-01-19 阅读 445

就在前几天，和一位高校的信息技术老师聊天，我得到了一个震惊的消息：明年，浙江省信息技术教材将不会在使用晦涩难懂的VB语言，而是改学更简单易懂的Python语言。也就是说：

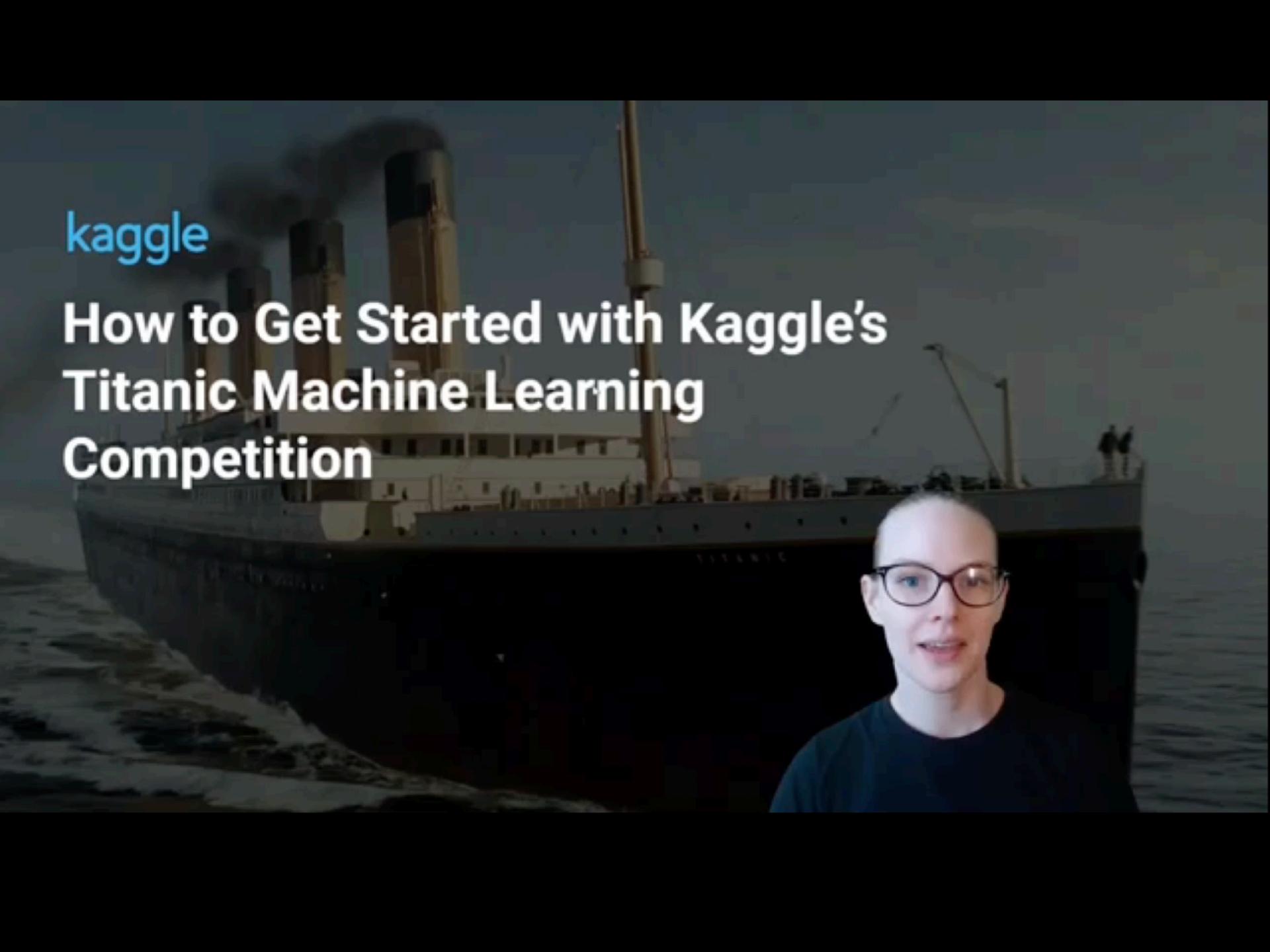
Python语言将纳入高考内容之一。

编程语言在升学中的比重逐渐加大，将要成为高考加分的一大利器。

Python将被纳入高考内容！

浙江省信息技术课程改革方案已经出台，Python确定进入浙江省信息技术高考，从2018年起浙江省信息技术教材编程语言将会从vb更换为Python。

其实不止浙江，教育大省北京和山东也确定要把Python编程基础纳入信息技术课程和高考的内容体系，Python语言课程化也将成为孩子学习的一种趋势。

A woman with short hair and glasses is speaking to the camera. She is positioned in the lower right foreground. In the background, there is a dark, grainy image of the RMS Titanic ship at night, with its name visible on the side.

kaggle

How to Get Started with Kaggle's Titanic Machine Learning Competition



- Home
- Compete
- Data
- Notebooks
- Discuss
- Courses
- More

Getting Started Prediction Competition

Titanic: Machine Learning from Disaster

Start here! Predict survival on the Titanic and get familiar with ML basics

Kaggle · 17,673 teams · Ongoing

[Overview](#) [Data](#) [Notebooks](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Join Competition](#)

Recently Viewed

- How Models Work
- Basic Data Exploration
- Exercise: Explore Your ...
- Air Pollution in Seoul
- Melbourne Housing Sn...

Overview

Description

👋 🛳 Ahoy, welcome to Kaggle! You're in the right place.

Evaluation

This is the legendary Titanic ML competition – the best, first challenge for you to dive into ML competitions and familiarize yourself with how the Kaggle platform works.

Frequently Asked Questions

The competition is simple: use machine learning to create a model that predicts which passengers survived the Titanic shipwreck.

Read on or watch the video below to explore more details. Once you're ready to start competing, click on the ["Join Competition button](#) to create an account and gain access to the [competition data](#).

<https://www.kaggle.com/c/titanic/>



为什么Kaggle对找工作有帮助？如何入门？



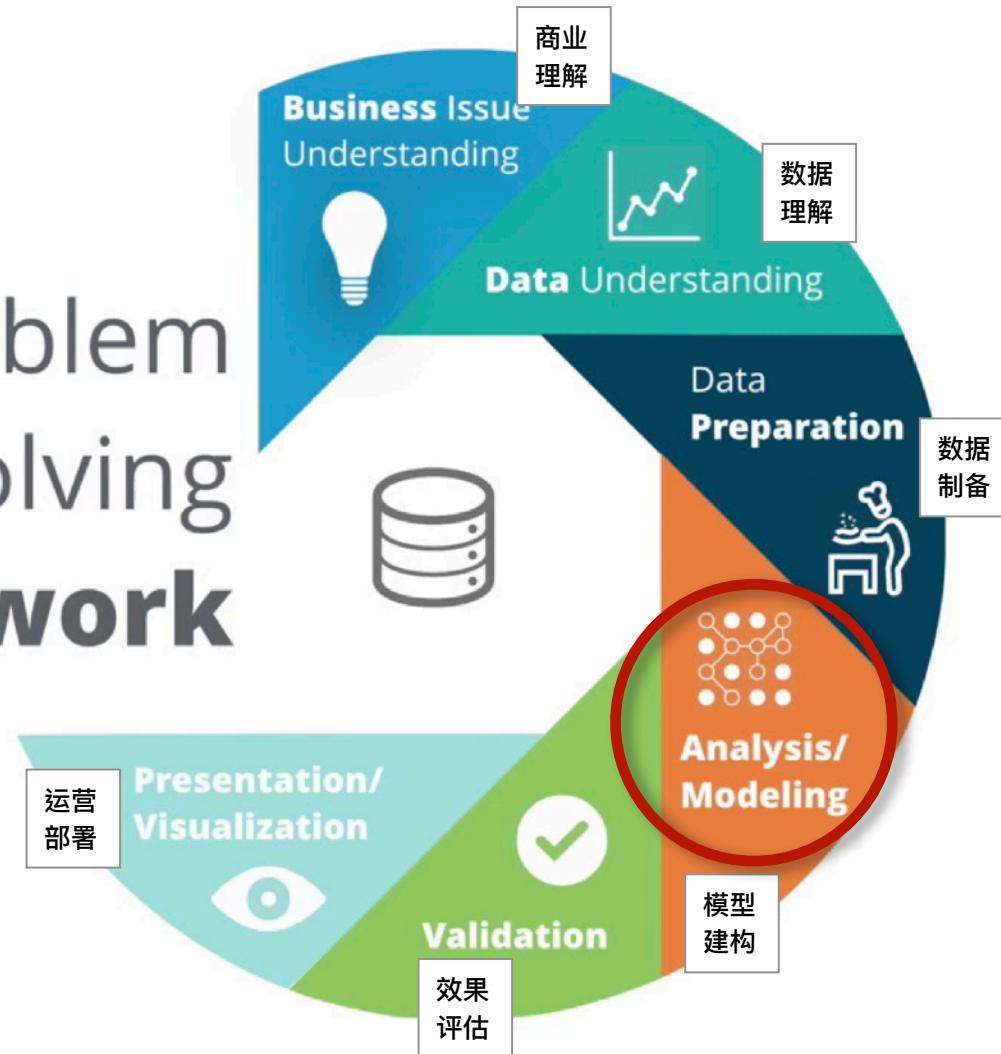
猴子

中国科学院大学 电子与通信工程硕士

148 人赞同了该文章



Problem Solving Framework



The Methodology Map

聚焦未来				聚焦过去和现在
Predict Outcome				Data Analysis
Data Rich		Data Poor		Geospatial
Numeric		Classification	A/B Test	Segmentation
Continuous	Count	Binary	Non Binary	Aggregation
<ul style="list-style-type: none"> • Linear Regression, • Multiple Regression 	<ul style="list-style-type: none"> • Count Regression 	<ul style="list-style-type: none"> • Logistic Regression • Decision Tree 	<ul style="list-style-type: none"> • Random Forest • K-Nearest Neighbour 	Descriptive
e.g. Scikit-learn, Tensorflow, Orange			e.g. Google Optimize	e.g. Google Analytics

Source: Udacity Model Selection Methodology Map

Predictive Analytics/Machine Learning/Deep Learning

预测分析/机器学习/深度学习

- Linear Regression
 - Decision Tree
 - Forest Model
 - Boosted Model
-
- Count Regression
 - Logistic Regression
 - Decision Tree
-
- Forest Model
 - Booted Model

Inferential Statistics & Hypothesis Testing

推论统计与假设检验

Explorative Data Analysis & Descriptive Statistics

探索性数据分析与
描述性统计

Source: Udacity Model Selection Methodology Map

Today's Agenda

1. Labelling of **numerical** and **categorical** data for supporting **Supervised Learning**.
2. Practical use of **Unsupervised Learning**
3. Separation of dataset into **training** and **testing** sets before building the model.
4. Use of **regression** or **classification** algorithms in training the model.
5. Model selection, training, testing, evaluation and deployment (as REST services)
6. Demonstration of regression and classification models in **Orange3**, **Scikit-learn** and **TensorFlow**.

**Patterns in
Unsupervised
Learning**

**Classification &
Regression in
Supervised
Learning**

**Demo in Orange,
Scikit-learn, &
Tensor Flow**

**Featuring Engineering & Dimension
Reduction**

<https://www.youtube.com/watch?v=N9fDIAfICMY>

With the possibility of having irrelevant and redundant features in confounding the outcome, how can we deal with them?



- ② Reduce the dimensions of the construct
- a) Principal Component Analysis
 - b) Exploratory Factor Analysis

<https://www.youtube.com/watch?v=TuSV2o3OcwI&list=PL7SDML5LXeEp6so26ZdKNZ6Ycvc55WnoG>

Unsupervised Learning is used in feature engineering and dimension reduction by leveraging its power in pattern recognition without labelling.

Features/Dimensions/Attributes/Independent Variables

年齢	性別	聘有家佣	家庭人数	居住状况	教育程度	职业	每週买菜次数
1							
2							
3							
4							
5							
6							
7							

Features/Dimensions/Attributes/Independent Variables

Observations/Instances

	年龄	性別	聘有家佣	家庭人数	居住状况	教育程度	职业	每週买菜次数	N+1	N+...
1										
2										
3										
4										
5										
6										
7										

Can anyone of the features be removed to focus on the key ones?

**Patterns in
Unsupervised
Learning**

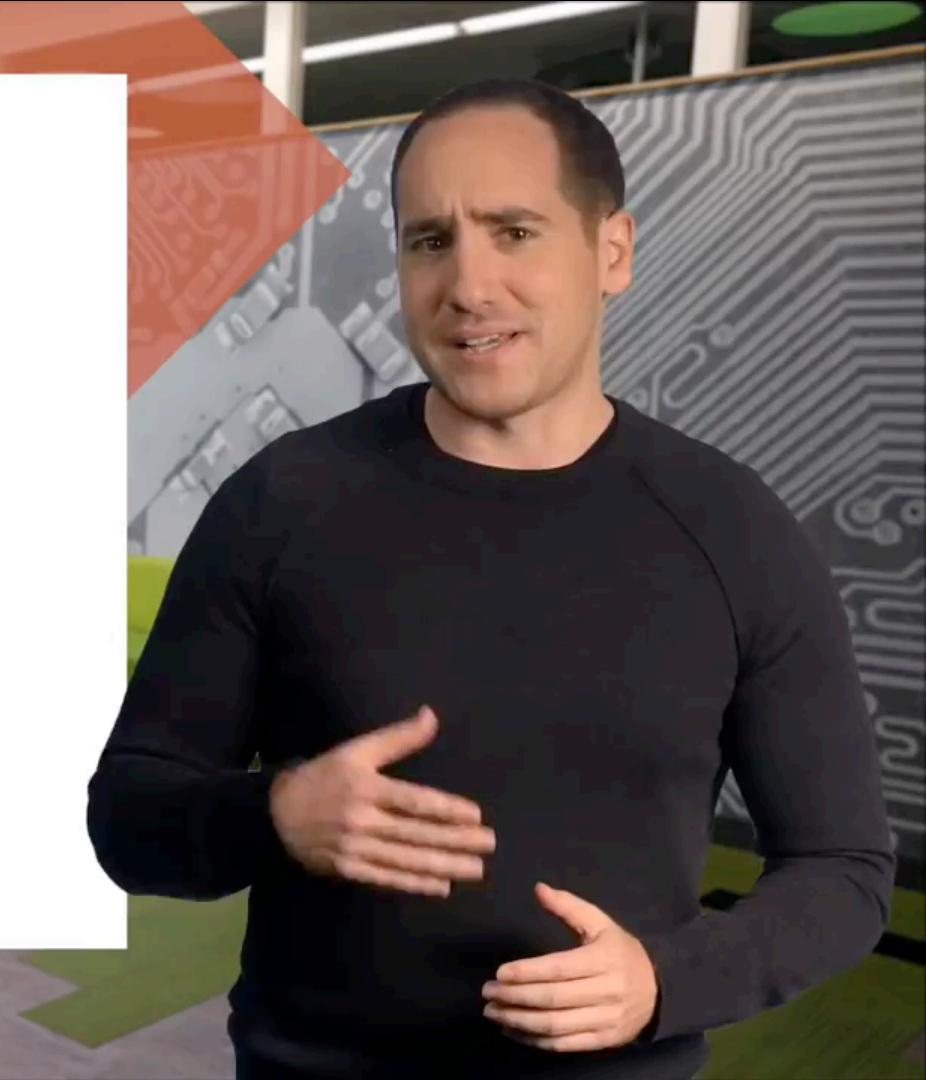
**Classification &
Regression in
Supervised
Learning**

**Demo in Orange,
Scikit-learn, &
Tensor Flow**

{ML} Recipes

Josh Gordon presents:

Let's Write a Pipeline

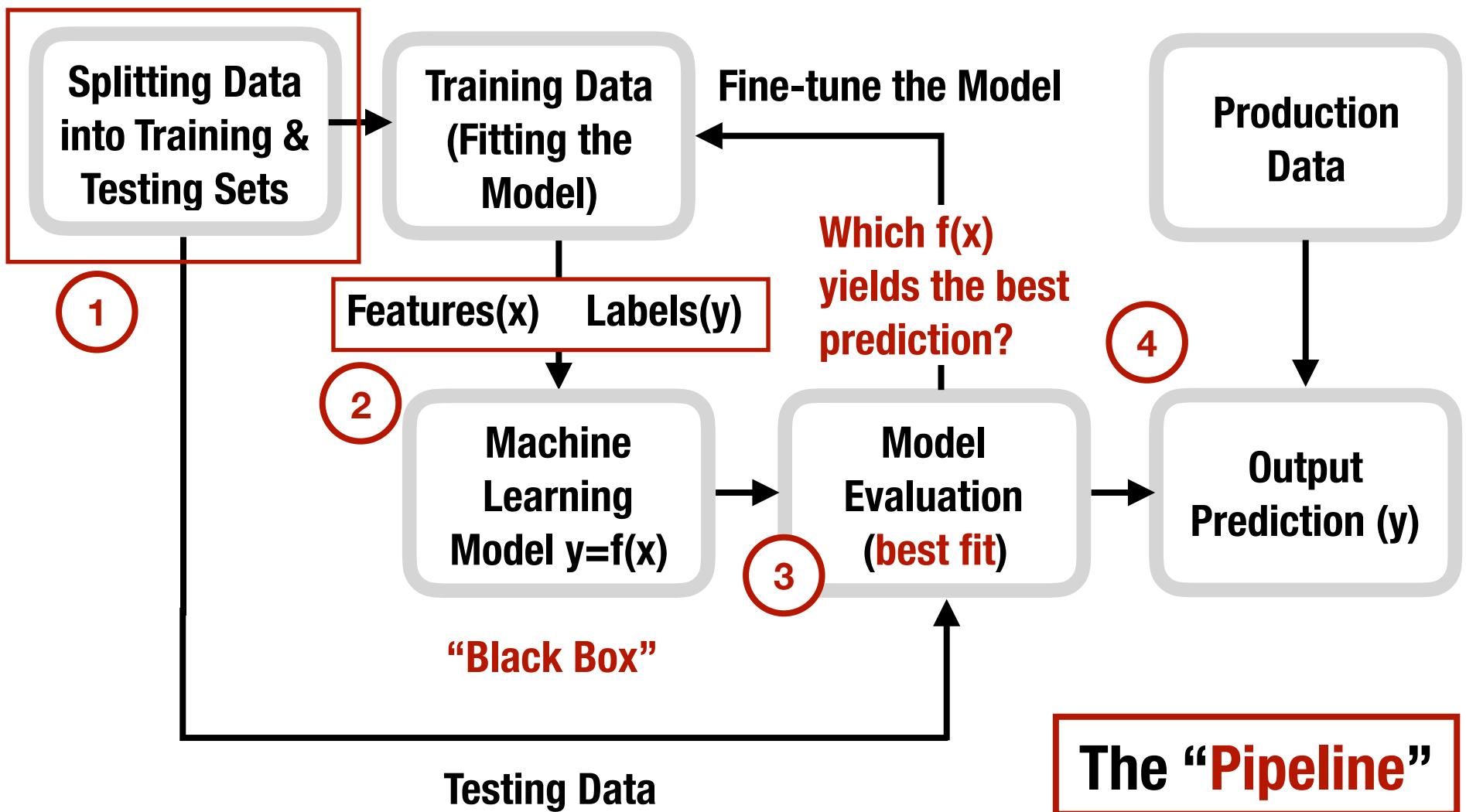


<https://www.youtube.com/watch?v=84gqSbLcBFE>

“A machine learning pipeline is used to help automate machine learning workflows.

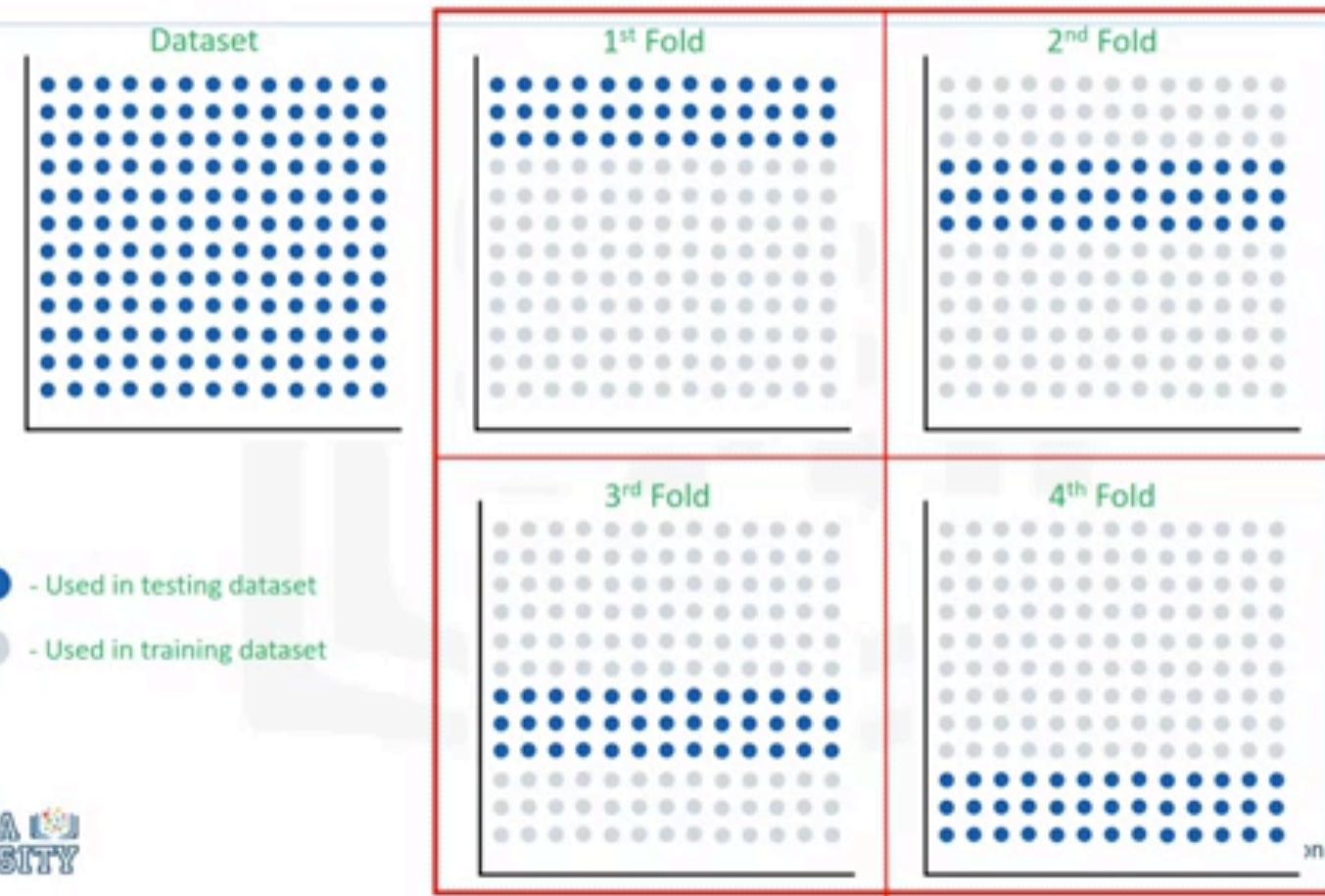
They operate by enabling a sequence of data to be transformed and correlated together in a model that can be tested and evaluated to achieve an outcome, whether positive or negative.”

Source: What is a Pipeline in Machine Learning? How to create one? by Shashanka M



Why Splitting the Data into Training and Testing Sets in Prediction?

K-fold Cross Validation



一个比喻

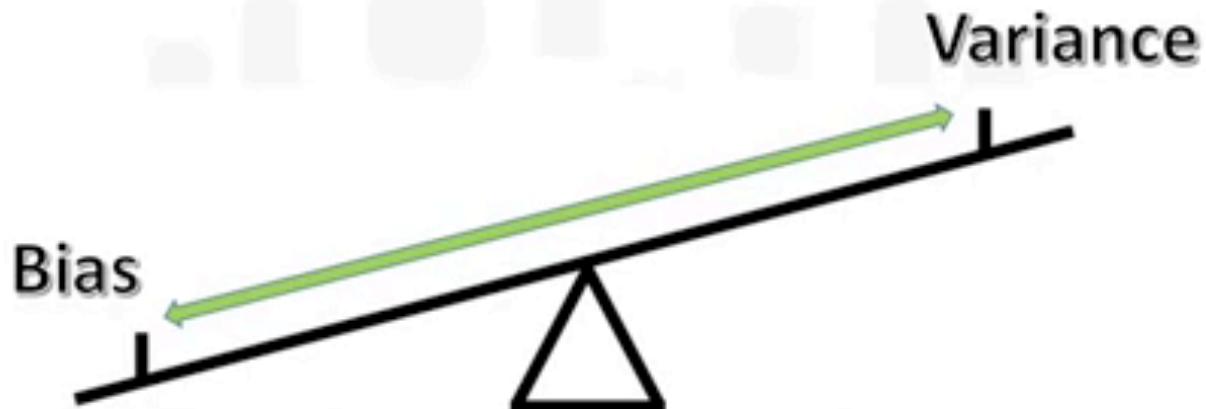
以考试成绩来预测学生表现

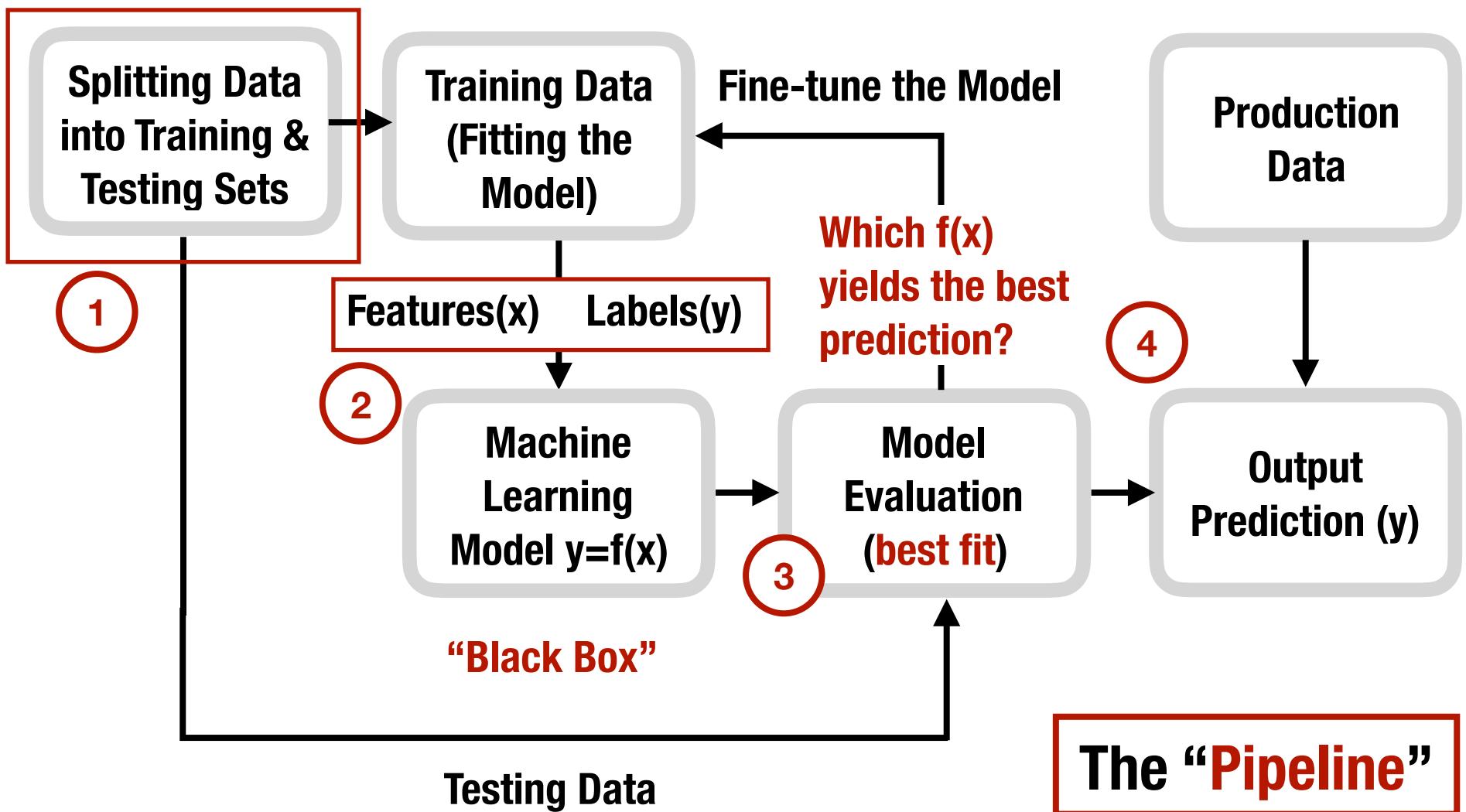
以旧试题来训练学生 → 重用旧试题来做评估

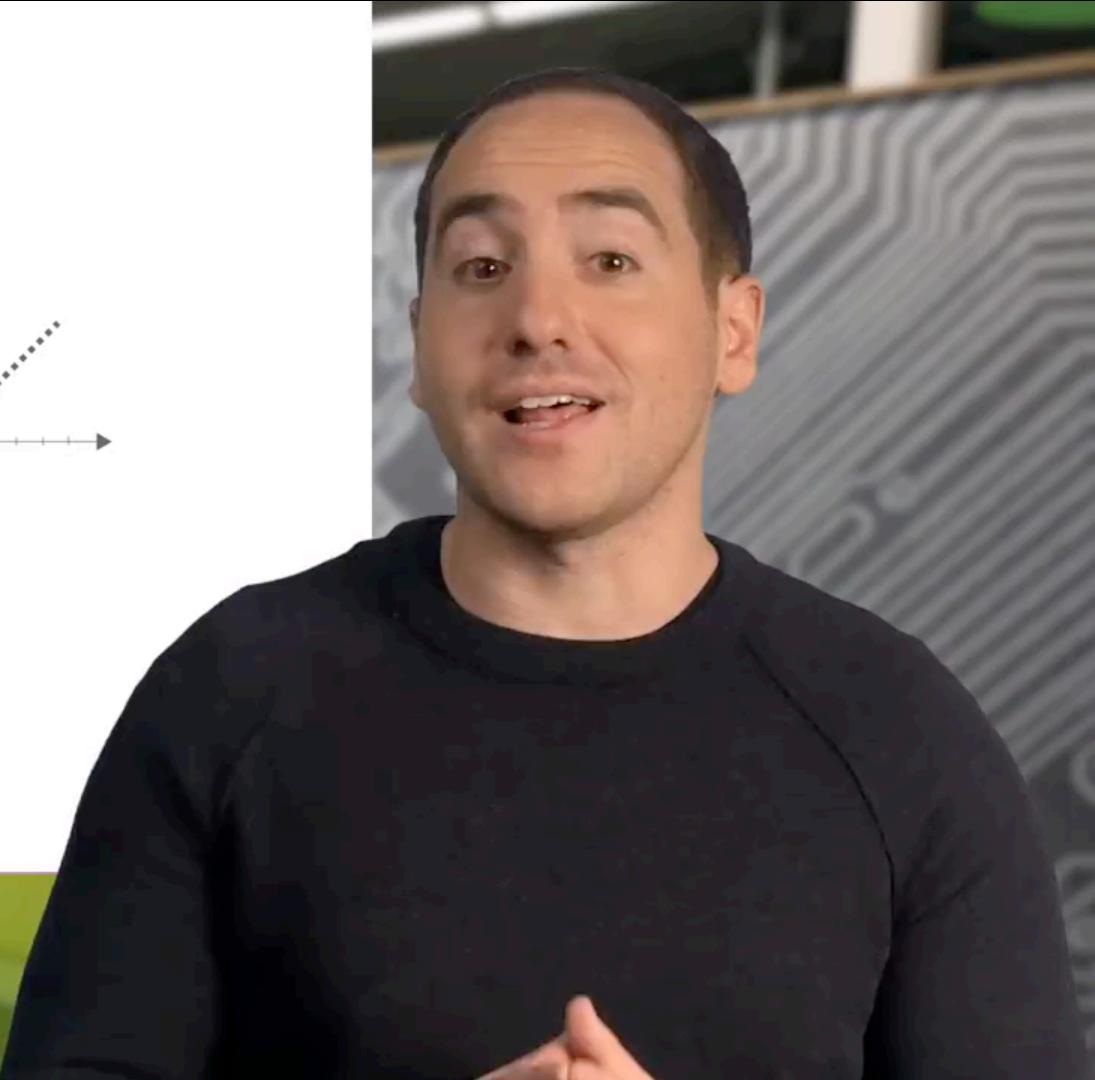
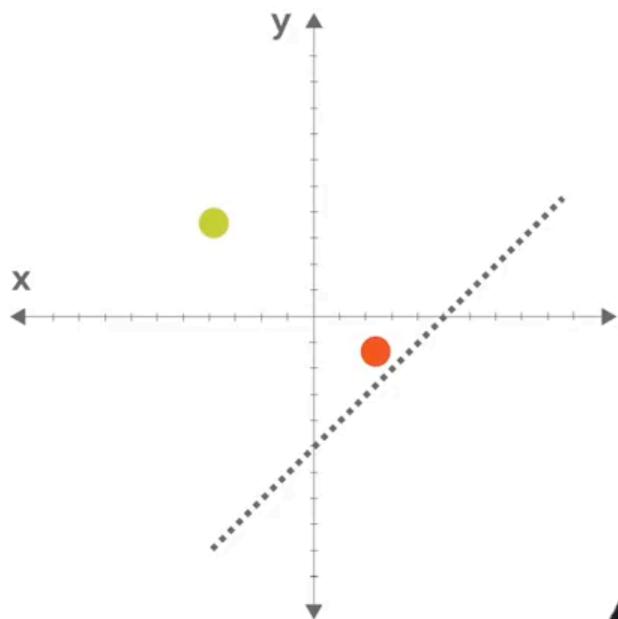
结果出来成绩优异能代表甚么？

够竟我们要培训会考试的学生或是会思考的学生？同样道理也适用于以Supervised Learning来训练ML模型。

Bias Variance Trade-Off







<https://www.youtube.com/watch?v=84gqSbLcBFE>

Which Model to Choose?

聚焦未来				聚焦过去和现在	
Predict Outcome				Data Analysis	
Data Rich				Data Poor	Geospatial
Numeric		Classification		A/B Test	Segmentation
Continuous	Count	Binary	Non Binary		Aggregation
<ul style="list-style-type: none"> • Linear Regression, • Multiple Regression 	<ul style="list-style-type: none"> • Count Regression 	<ul style="list-style-type: none"> • Logistic Regression • Decision Tree 	<ul style="list-style-type: none"> • Random Forest • K-Nearest Neighbour 		Descriptive
e.g. Scikit-learn, Tensorflow, Orange				e.g. Google Optimize	e.g. Google Analytics

Source: Udacity Model Selection Methodology Map

Business Problem

Predict Outcome

Data Rich

Data Poor

<https://www.youtube.com/watch?v=Lkpkl7KmGLA>

<https://www.youtube.com/watch?v=cAlPn6YRPc0>

Let's learn a few basic ML algorithms (models)

K-Nearest Neighbour (Classification Model)

If K=11 and the new cell is between two (or more) categories, we simply pick the category that “gets the most votes”.

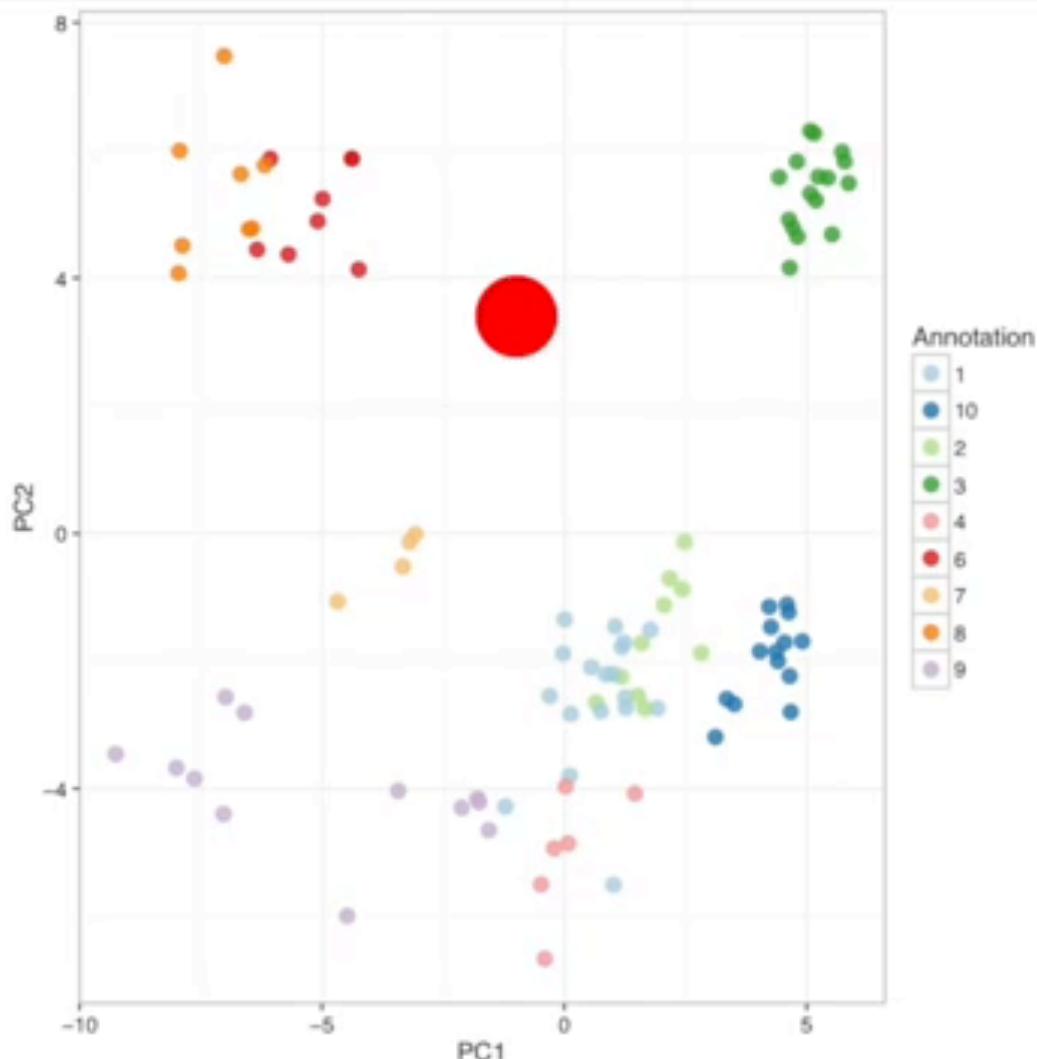
In this case....

7 nearest neighbors are **RED**.

3 nearest neighbors are **ORANGE**.

1 nearest neighbor is **GREEN**.

Since **RED** got the most votes, the final assignment is **RED**.



Linear Regression (Regression Model)

Linear and Polynomial Regression



Decision Tree (Classification Model)

Decision Tree

Gender	Age	App
F	15	
F	25	
M	32	
F	40	
M	12	
M	14	



Examples in Orange, Scikit-learn, & Tensor Flow

**Patterns in
Unsupervised
Learning**

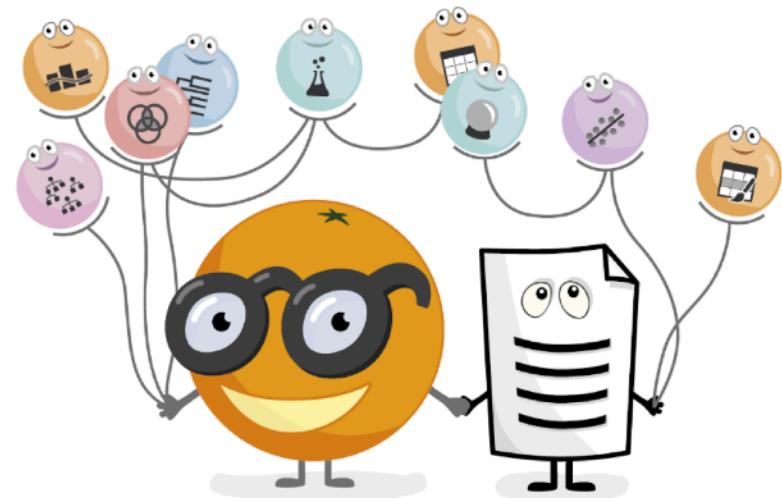
**Classification &
Regression in
Supervised
Learning**

**Demo in Orange,
Scikit-learn, &
Tensor Flow**

[Features](#)[Screenshots](#)[Workflows](#)[Download](#)[Blog](#)[Docs](#)[Training](#)[Donate](#)

Data Mining Fruitful and Fun

Open source machine learning and data visualization for novice and expert. Interactive data analysis workflows with a large toolbox.

[Download Orange](#)



- Home
- Compete
- Data
- Notebooks
- Discuss
- Courses
- More

Getting Started Prediction Competition

Titanic: Machine Learning from Disaster

Start here! Predict survival on the Titanic and get familiar with ML basics

Kaggle · 17,673 teams · Ongoing

[Overview](#) [Data](#) [Notebooks](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Join Competition](#)

Recently Viewed

- How Models Work
- Basic Data Exploration
- Exercise: Explore Your ...
- Air Pollution in Seoul
- Melbourne Housing Sn...

Overview

Description

👋 🛳 Ahoy, welcome to Kaggle! You're in the right place.

Evaluation

This is the legendary Titanic ML competition – the best, first challenge for you to dive into ML competitions and familiarize yourself with how the Kaggle platform works.

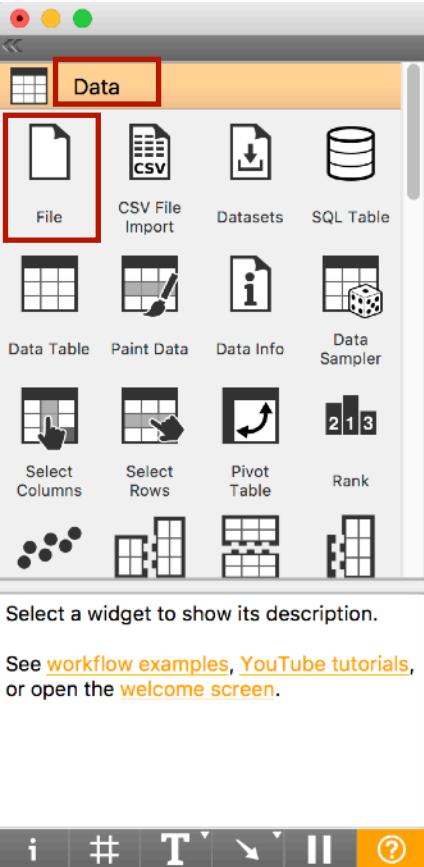
Frequently Asked Questions

The competition is simple: use machine learning to create a model that predicts which passengers survived the Titanic shipwreck.

Read on or watch the video below to explore more details. Once you're ready to start competing, click on the ["Join Competition button](#) to create an account and gain access to the [competition data](#).

<https://www.kaggle.com/c/titanic/>

**Let's train/build a Titanic survivor model
pipeline in Orange to get you started in
joining the Kaggle community.**



File: titanic_train.csv

URL:

Info

891 instance(s)
10 feature(s) (2.0% missing values)
Data has no target variable.
3 meta attribute(s)

Columns (Double click to edit)

	Name	Type	Role	Values
1	PassengerId	N numeric	feature	
2	Survived	C categorical	✓ feature target	0, 1
3	Pclass	N numeric	meta	
4	Sex	C categorical	skip	female, male
5	Sex-Binary	C categorical	feature	0, 1
6	Age	N numeric	feature	
7	SibSp	N numeric	feature	

Browse documentation datasets

Reset

Apply

?

File

File: titanic_train.csv

URL:

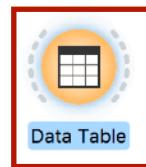
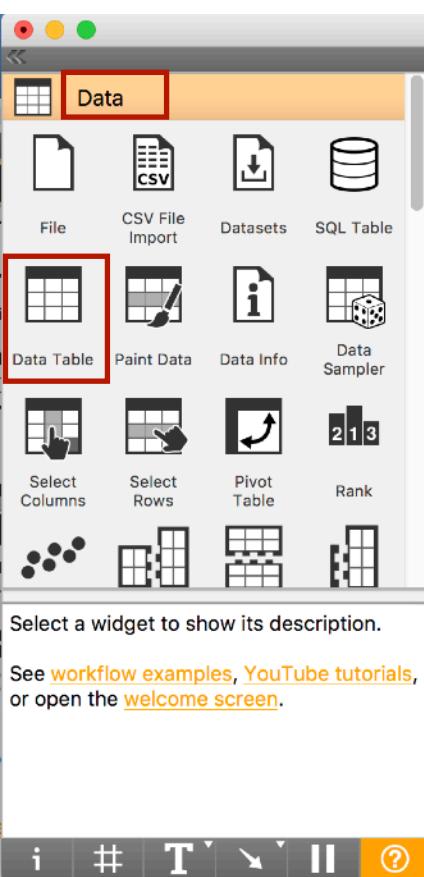
Info

891 instance(s)
10 feature(s) (2.0% missing values)
Data has no target variable.
3 meta attribute(s)

Columns (Double click to edit)

	Name	Type	Role	Values
1	PassengerId	N numeric	feature	
2	Survived	C categorical	target	0, 1
3	Pclass	N numeric	feature	
4	Sex	C categorical	feature	female, male
5	Sex-Binary	C categorical	feature	0, 1
6	Age	N numeric	feature	
7	SibSp	N numeric	feature	

?



Data Table

File CSV File Import Datasets SQL Table

Data Table Paint Data Data Info Data Sampler

Select Columns Select Rows Pivot Table Rank

Correlati... Merge Data Concate... Select by Data Index

Select a widget to show its description.

See [workflow examples](#), [YouTube tutorials](#), or open the [welcome screen](#).

File Data Table

Data

File CSV File Import Datasets SQL Table

Data Table Paint Data Data Info Data Sampler

Select Columns Select Rows Pivot Table Rank

Data Table

View the dataset in a spreadsheet.

more...

Info

891 instances
10 features (2.0% missing values)
No target variable.
3 meta attributes (25.7% missing values)

Variables

Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

Selection

Select full rows

Restore Original Order

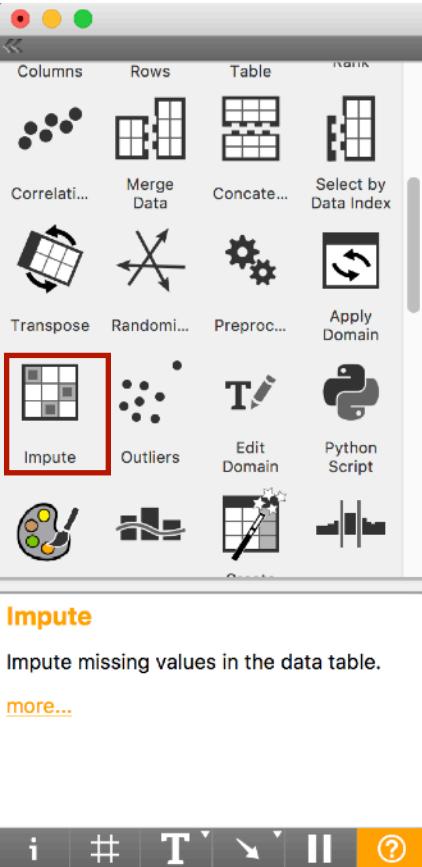
Send Automatically

sibsp - Number of Siblings/Spouses Aboard. **parch** - Number of Parents/Children Aboard. **ticket** - Ticket Number. **fare** - Passenger Fare. Jan 18, 2018

embarked - Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton)
boat - Lifeboat (if survived) **body** - Body number (if did not survive and body was recovered) May this help you: <https://data.world/nrippner/titanic-disaster-dataset>.

parch - Number of Parents/Children Aboard. **ticket** - Ticket Number. **fare** - Passenger Fare. **cabin** - Cabin. **embarked** - Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton) Jan 18, 2018

Survived	Age	SibSp	Parch	Fare	Embarked
1	22.00	1	0	7.2500	S
2	38.00	1	0	71.2833	C
3	26.00	0	0	7.9250	S
4	35.00	1	0	53.1000	S
5	35.00	0	0	8.0500	S
6	?	0	0	8.4583	Q
7	54.00	0	0	51.8625	S
8	2.00	3	1	21.0750	S
9	27.00	0	2	11.1333	S
10	14.00	1	0	30.0708	C
11	4.00	1	1	16.7000	S
12	58.00	0	0	26.5500	S
13	20.00	0	0	8.0500	S
14	39.00	1	5	31.2750	S
15	14.00	0	0	7.8542	S
16	55.00	0	0	16.0000	S
17	2.00	4	1	29.1250	Q
18	?	0	0	13.0000	S
19	31.00	1	0	18.0000	S
20	?	0	0	7.2250	C

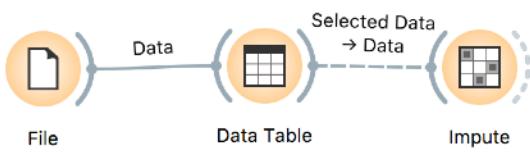
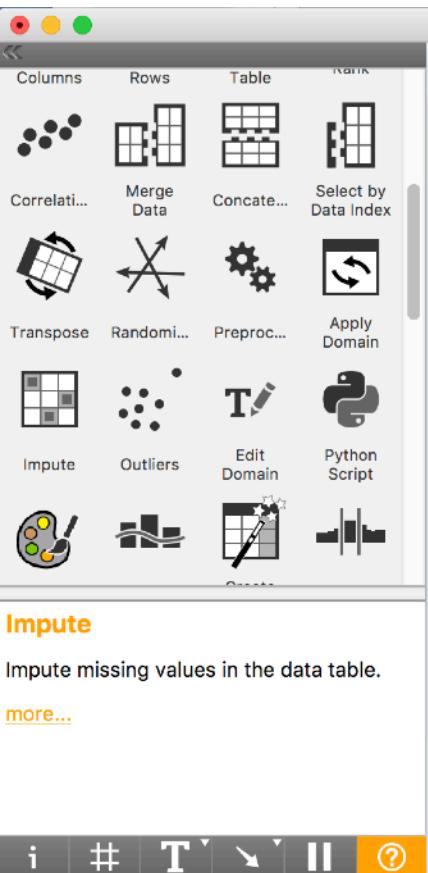


Impute

Impute missing values in the data table.

[more...](#)





Impute

Impute missing values in the data table.

more...

i # T ↻ || ?

Impute

Default Method

- Don't impute
- Average/Most frequent
- As a distinct value
- Model-based imputer (simple tree)
- Random values
- Remove instances with unknown values

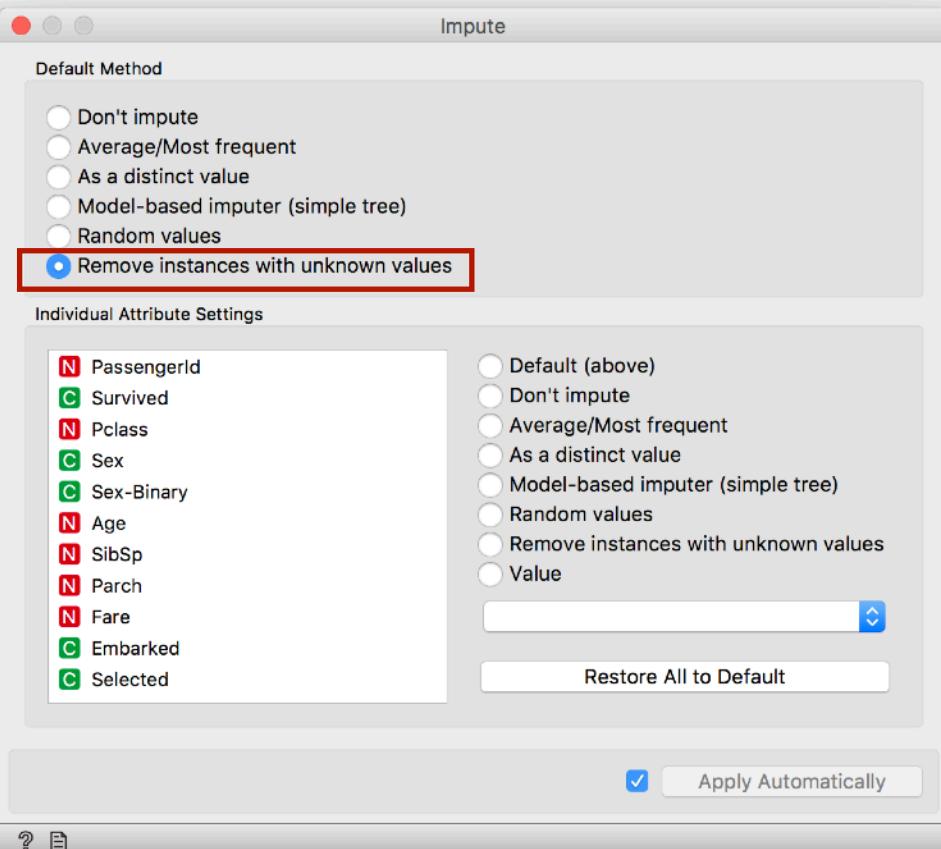
Individual Attribute Settings

<input type="radio"/> N	PassengerId	<input type="radio"/> Default (above)
<input type="radio"/> C	Survived	<input type="radio"/> Don't impute
<input type="radio"/> N	Pclass	<input type="radio"/> Average/Most frequent
<input type="radio"/> C	Sex	<input type="radio"/> As a distinct value
<input type="radio"/> C	Sex-Binary	<input type="radio"/> Model-based imputer (simple tree)
<input type="radio"/> N	Age	<input type="radio"/> Random values
<input type="radio"/> N	SibSp	<input type="radio"/> Remove instances with unknown values
<input type="radio"/> N	Parch	<input type="radio"/> Value
<input type="radio"/> N	Fare	
<input type="radio"/> C	Embarked	
<input type="radio"/> C	Selected	

Restore All to Default

Apply Automatically

? ⌘



com5940-2020-11 — Edited

View Zoom Add Slide Ungroup Group Masters Play Keynote Live Mask Table Chart Text Shape Media Comment Collaborate Format Animate Document

Columns Rows Table Rank

Correlat... Merge Data Concat... Select by Data Index

Transpose Randomi... Preproc... Apply Domain

Impute Outliers Edit Domain Python Script

Clear All Cancel OK

Edit Links

Data Table **Impute**

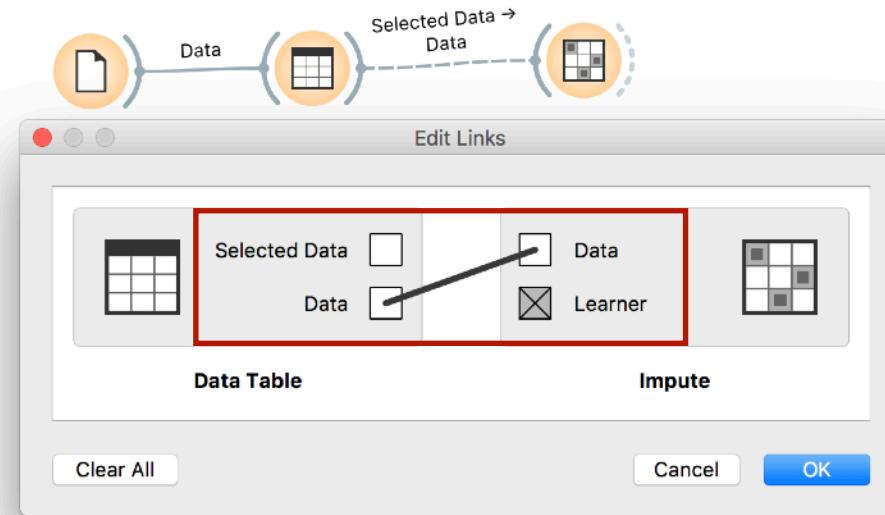
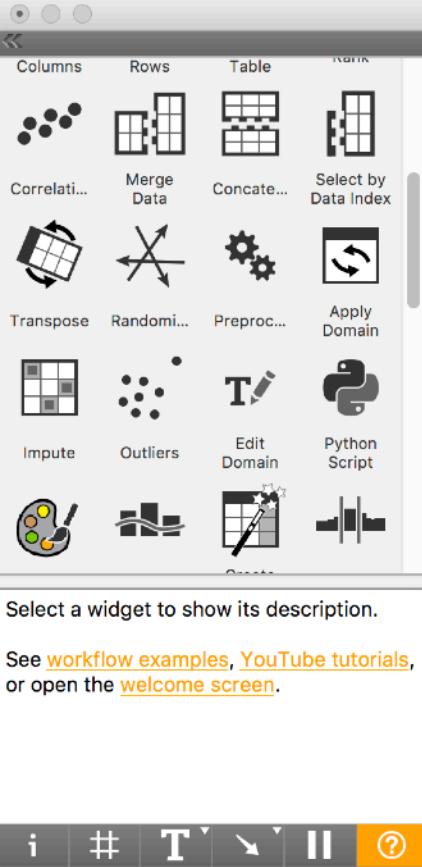
Selected Data → Data Learner

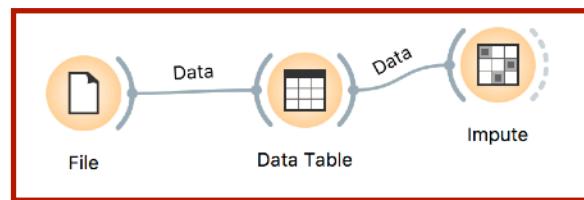
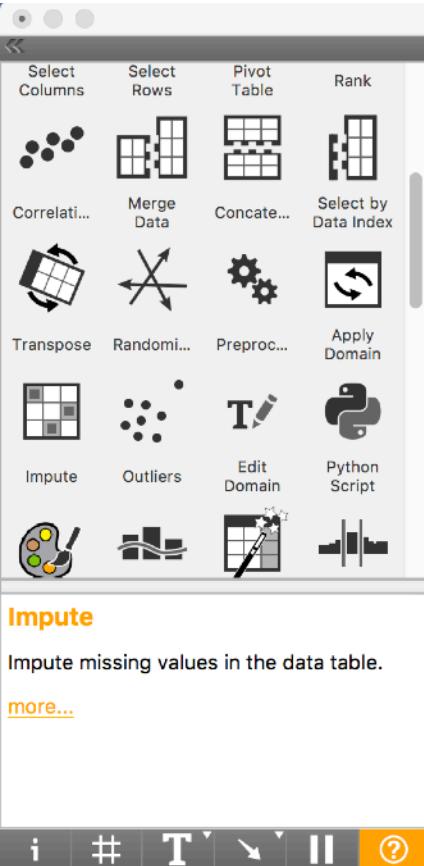
Impute

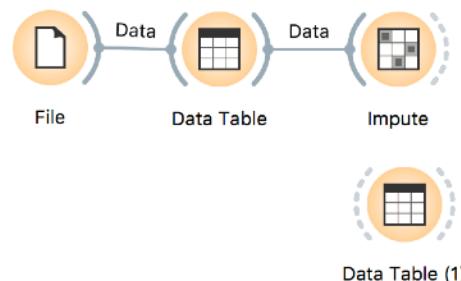
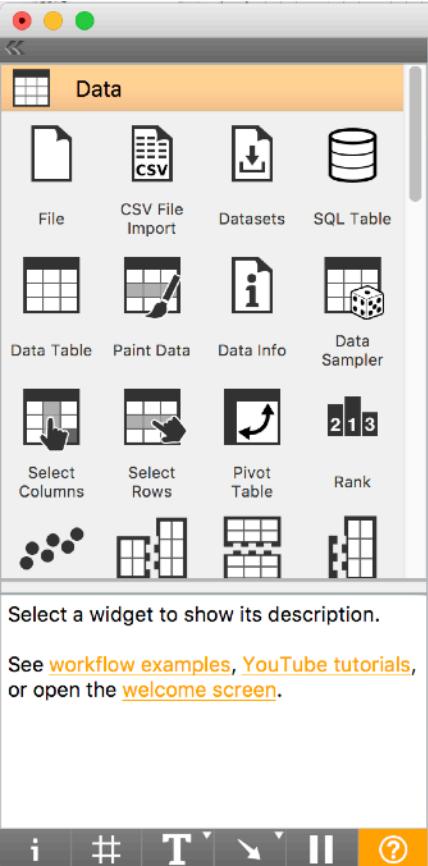
Impute missing values in the data table.

[more...](#)

i # T ↻ || ?







The screenshot shows the KNIME interface with the following components:

- Data Palette (Left):** A sidebar titled "Data" containing various icons for data manipulation nodes. The nodes include:
 - File, CSV File Import, Datasets, SQL Table
 - Data Table, Paint Data, Data Info, Data Sampler
 - Select Columns, Select Rows, Pivot Table, Rank
 - Scatter Plot, Data Table, Data Table
- Data Table Node (Bottom Left):** A node titled "Data Table" with the sub-label "(1)" below it. It has three ports: "Data" (input), "Data" (output), and "Data" (output) from an "Impute" node.
- Flow Diagram (Center):** A flow diagram showing the connection between a "File" node, a "Data Table" node, an "Impute" node, and a final "Data Table" node labeled "(1)".
- Toolbar (Bottom):** A horizontal toolbar with icons for information, help, and other system functions.

Data Table

View the dataset in a spreadsheet.

[more...](#)

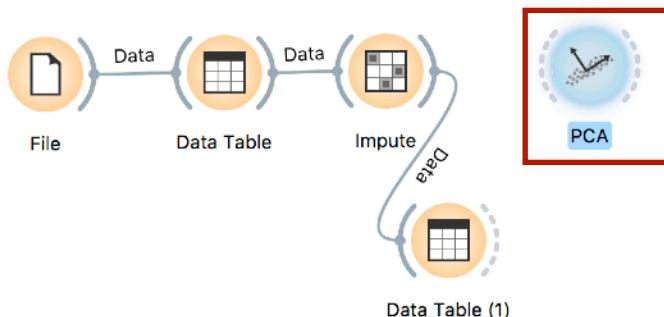
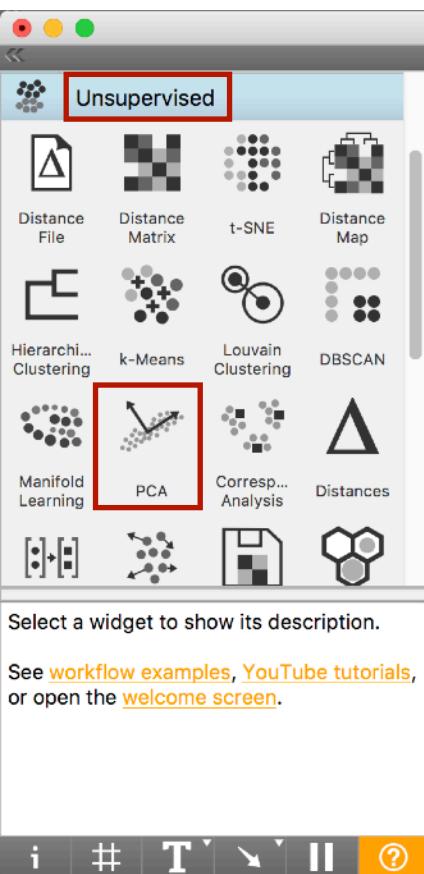
The screenshot shows the Orange data mining software interface. On the left, the 'Data' tab is selected, indicated by an orange bar. Below it, there are several icons for data import and manipulation: File, CSV File Import, Datasets, SQL Table, Data Table, Paint Data, Data Info, Data Sampler, Select Columns, Select Rows, Pivot Table, Rank, and three visualization icons. A 'Data Table' section provides a brief overview of the dataset: 712 instances, 10 features (no missing values), Discrete class with 2 values (no missing values), and 3 meta attributes (24.8% missing values). Below this are sections for Variables (checkboxes for Show variable labels (if present), Visualize numeric values, and Color by instance classes), Selection (checkbox for Select full rows), and buttons for Restore Original Order and Send Automatically.

No more zeros

Data Table (1)

	Age	SibSp	Parch	Fare	Embarked
1	22.00	1	0	7.2500	S
2	38.00	1	0	71.2833	C
3	26.00	0	0	7.9250	S
4	35.00	1	0	53.1000	S
5	35.00	0	0	8.0500	S
6	54.00	0	0	51.8625	S
7	2.00	3	1	21.0750	S
8	27.00	0	2	11.1333	S
9	14.00	1	0	30.0708	C
10	4.00	1	1	16.7000	S
11	58.00	0	0	26.5500	S
12	20.00	0	0	8.0500	S
13	39.00	1	5	31.2750	S
14	14.00	0	0	7.8542	S
15	55.00	0	0	16.0000	S
16	2.00	4	1	29.1250	Q
17	31.00	1	0	18.0000	S
18	35.00	0	0	26.0000	S
19	34.00	0	0	13.0000	S
20	15.00	0	0	8.0500	C

**Dimension reduction and feature engineering
in Orange using Unsupervised Learning for
finding patterns.**



Unsupervised

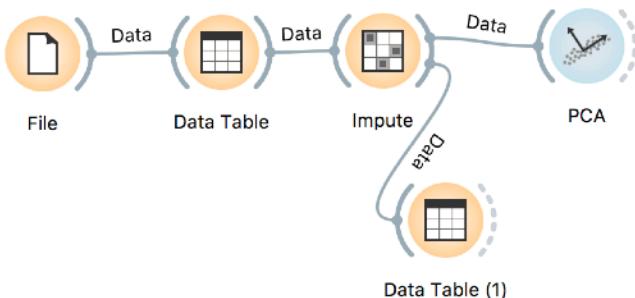
- Distance File
- Distance Matrix
- t-SNE
- Distance Map
- Hierarch... Clustering
- k-Means
- Louvain Clustering
- DBSCAN
- Manifold Learning
- PCA
- Corresp... Analysis
- Distances

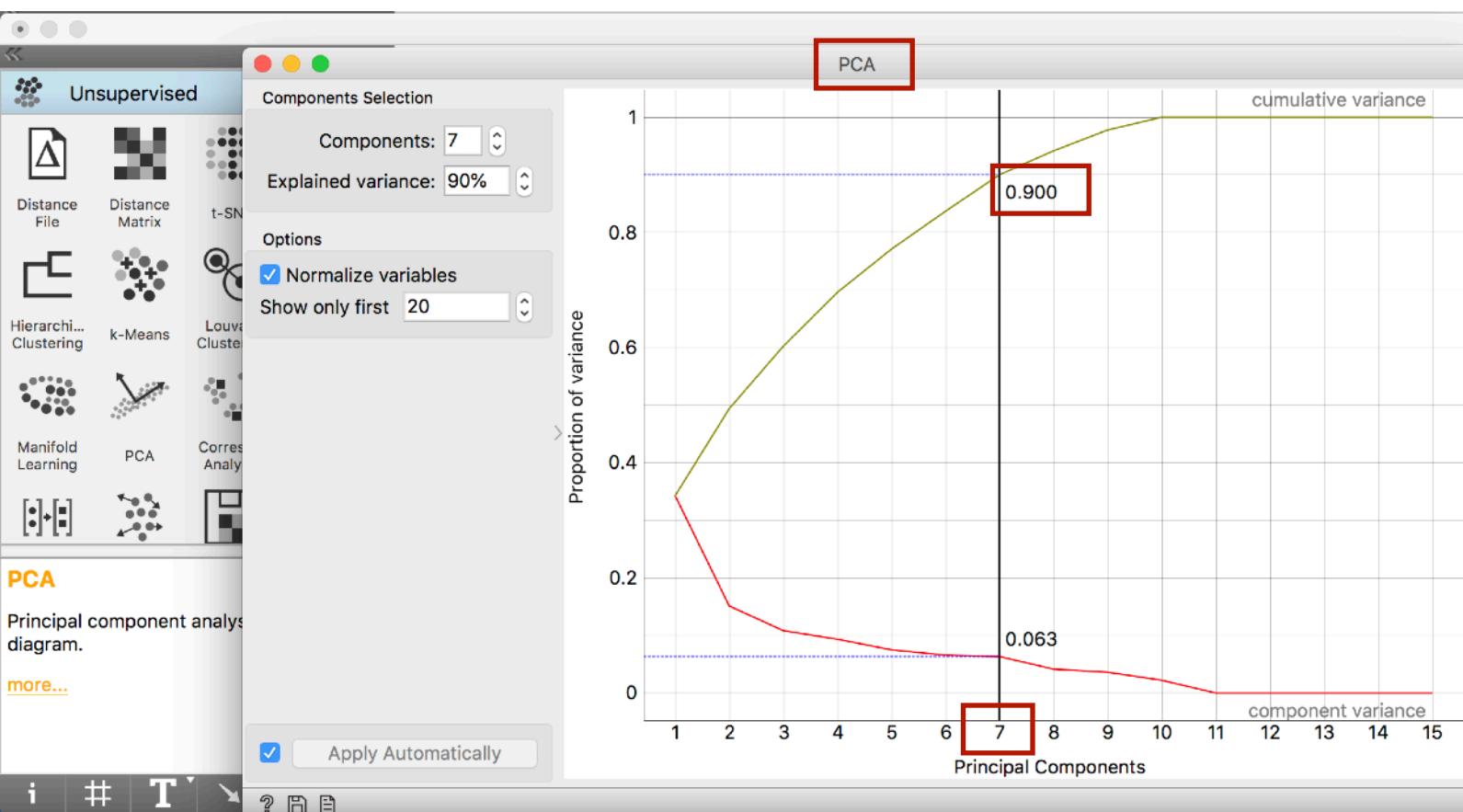
PCA

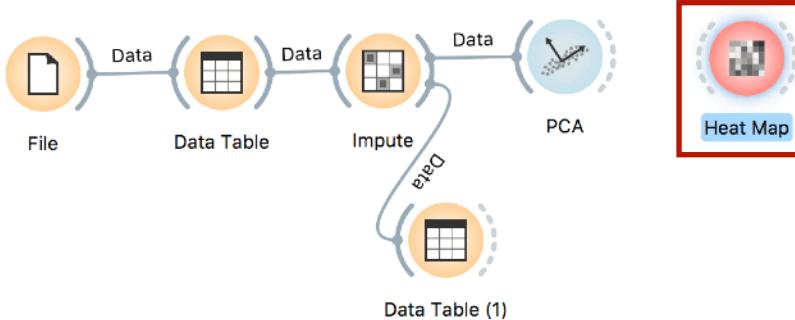
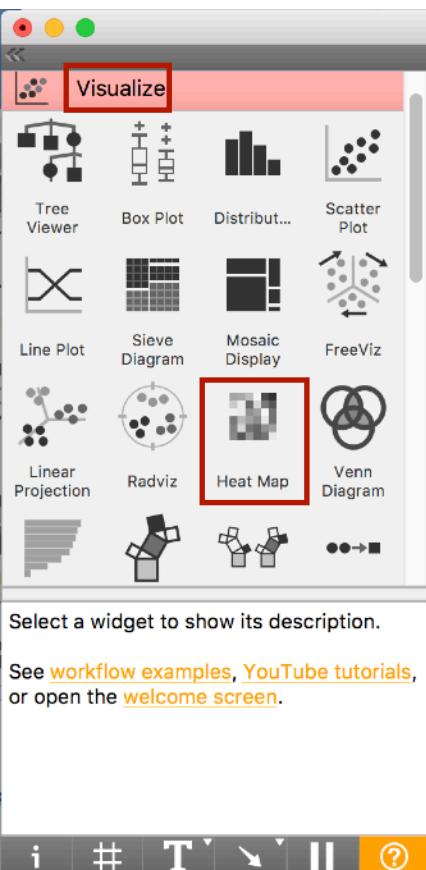
Principal component analysis with a scree-diagram.

[more...](#)

i | # | T | < | || | ?







Visualize

The interface features a sidebar on the left containing a grid of icons, each representing a different type of visualization tool. The tools include:

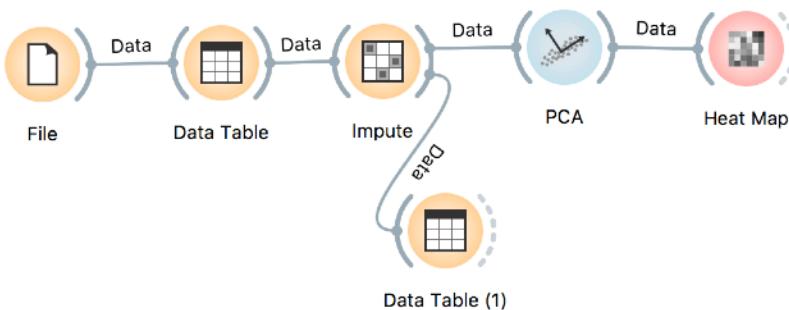
- Tree Viewer
- Box Plot
- Distribut...
- Scatter Plot
- Line Plot
- Sieve Diagram
- Mosaic Display
- FreeViz
- Linear Projection
- Radviz
- Heat Map
- Venn Diagram
- Heat Map
-

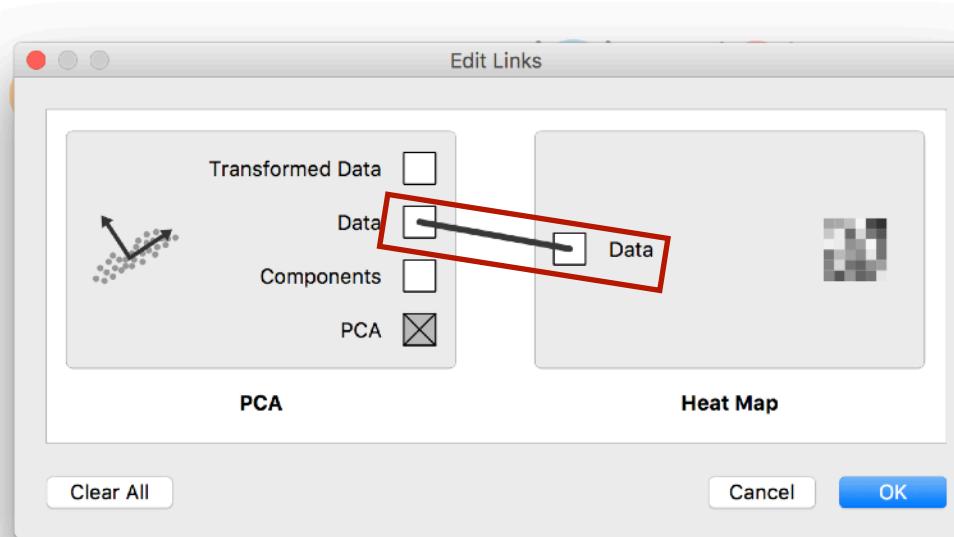
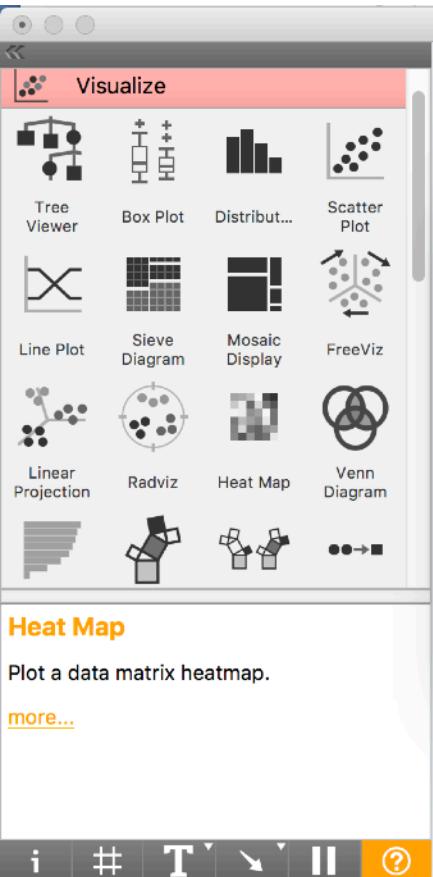
Heat Map

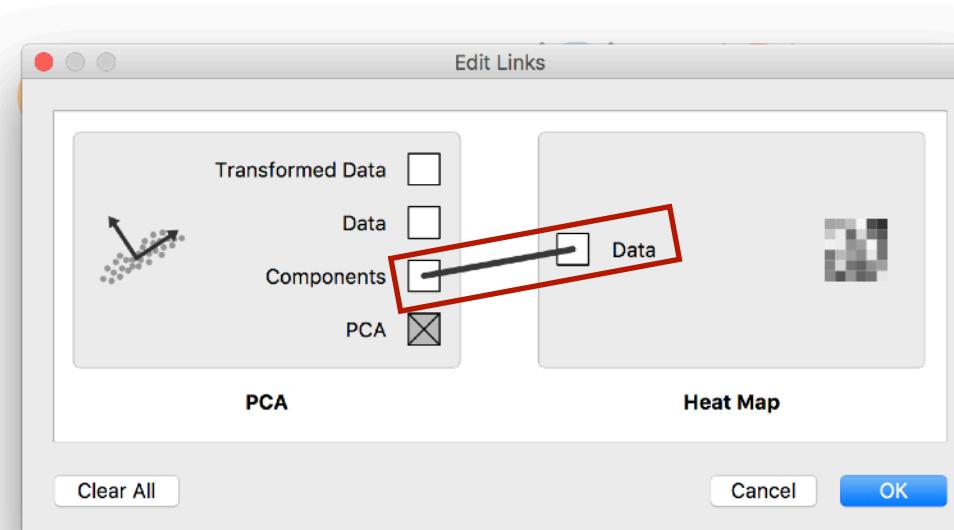
Plot a data matrix heatmap.

[more...](#)

Navigation icons at the bottom: i, #, T, ×, ||, ?







Visualize

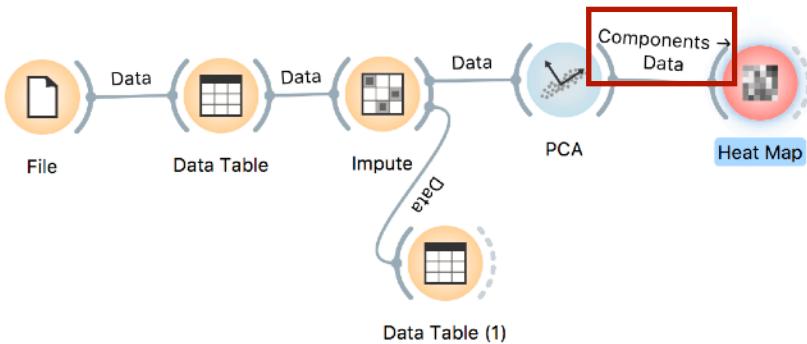
Tree Viewer
Box Plot
Distribution
Scatter Plot
Line Plot
Sieve Diagram
Mosaic Display
FreeViz
Linear Projection
Radviz
Heat Map
Venn Diagram

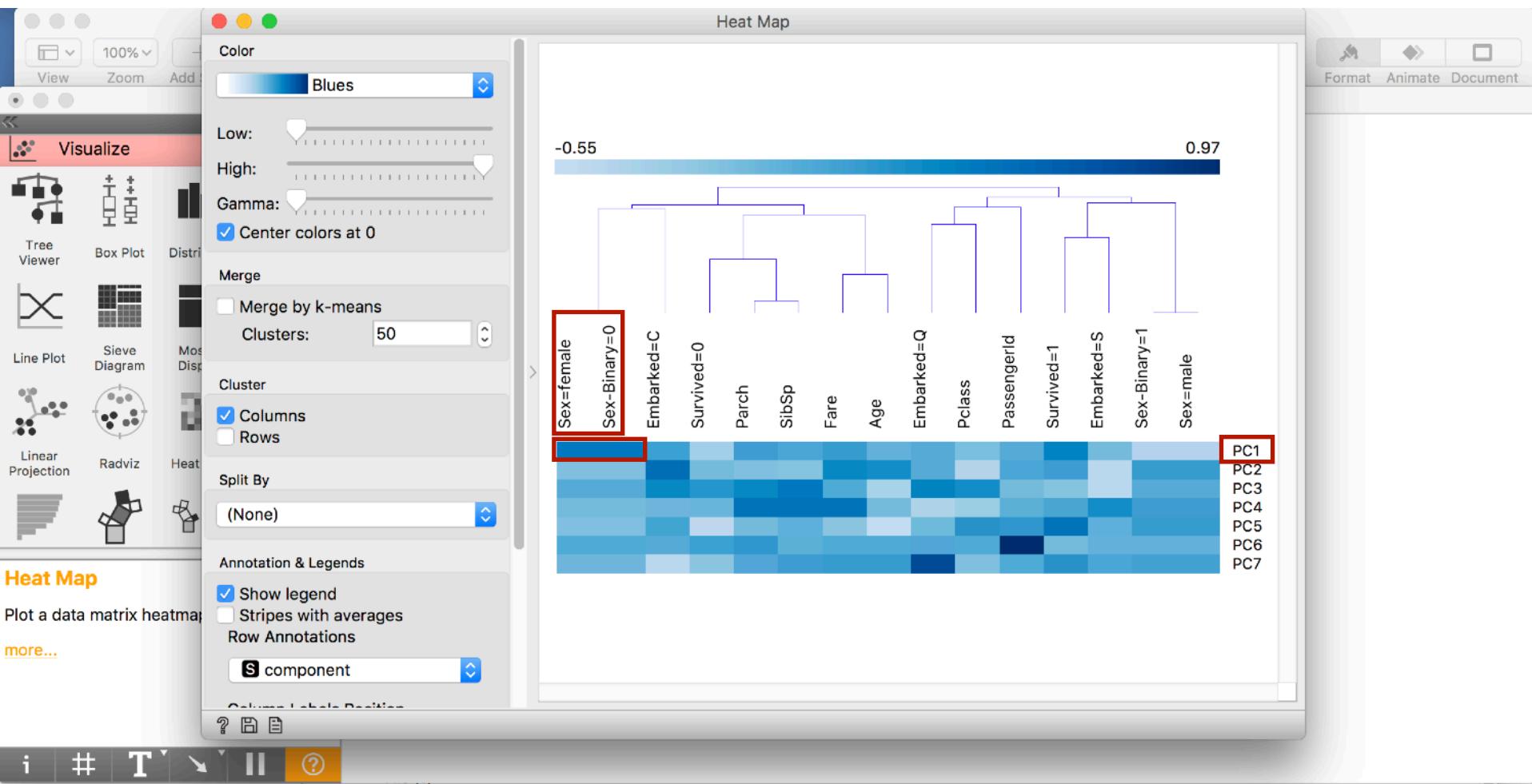
Heat Map

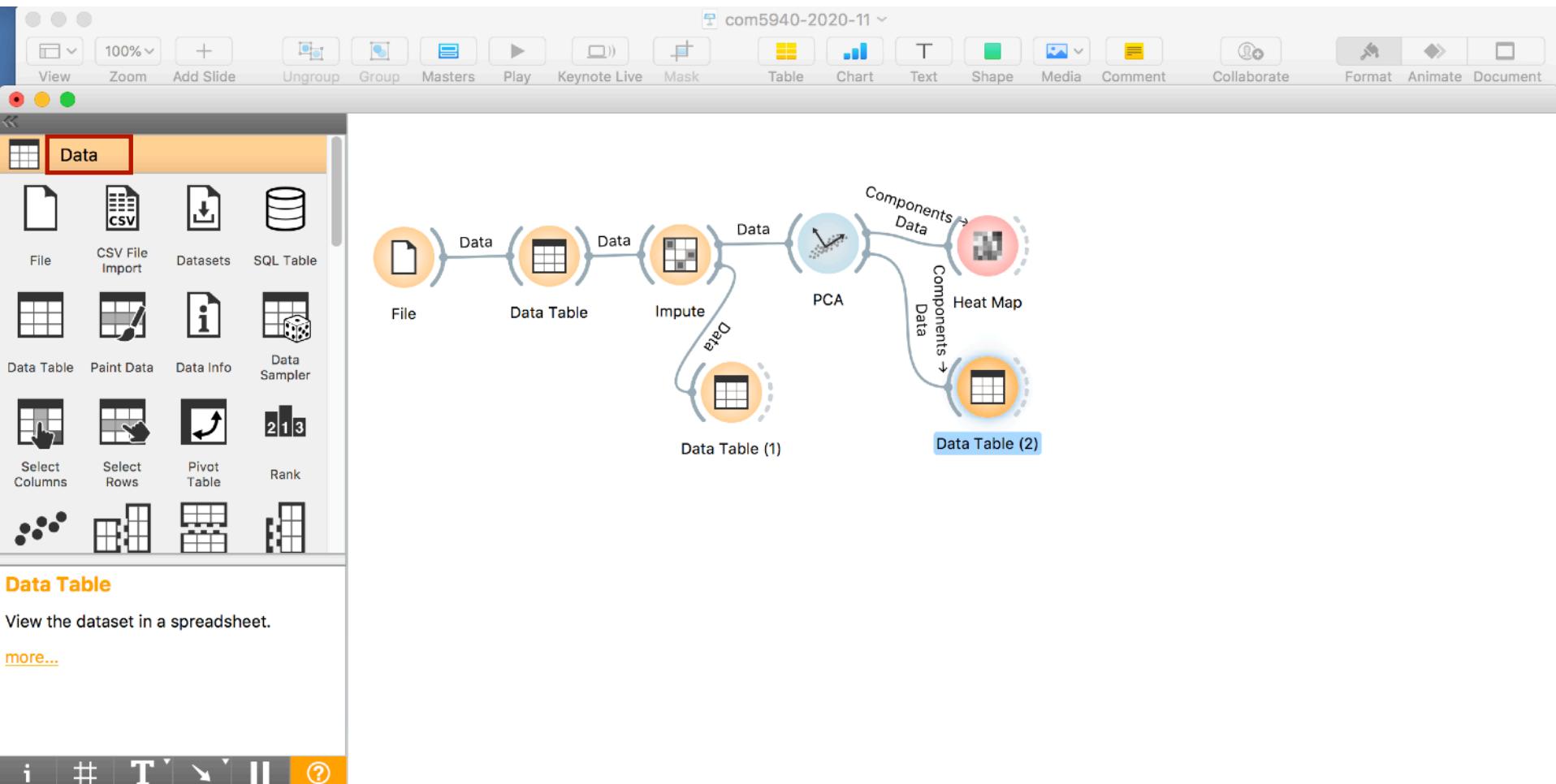
Plot a data matrix heatmap.

[more...](#)

i | # | T | ↺ | II | ⓘ







titanic_clasification_model_demo1.ows

Data Table (2)

component	PassengerId	Survived=0	Survived=1	Pclass	Sex=female
1 PC1	-0.00531213	-0.33085	0.33085	-0.145724	0.414614
2 PC2	0.0542579	-0.126653	0.126653	-0.410663	-0.156404
3 PC3	-0.126597	0.192752	-0.192752	0.259059	-0.00550433
4 PC4	-0.0131455	-0.0675409	0.0675409	-0.285139	-0.123198
5 PC5	0.176347	-0.470361	0.470361	0.213315	-0.159174
6 PC6	0.966783	0.115323	-0.115323	0.0335903	0.0345148
7 PC7	0.0293683	-0.219769	0.219769	-0.168819	-0.0869588

Info

7 instances (no missing values)
15 features (no missing values)
No target variable.
1 meta attribute (no missing values)

Variables

Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

Selection

Select full rows

Restore Original Order
 Send Automatically

Data Table
View the dataset in a spreadsheet.
[more...](#)

i # T ↵ II ?

Data Table Paint Data Data Info Data Sampler

Select Columns Select Rows Pivot Table Rank

Correlat... Merge Data Concat... Select by Data Index

Transpose Randomi... Preproc... Apply Domain

Select a widget to show its description.

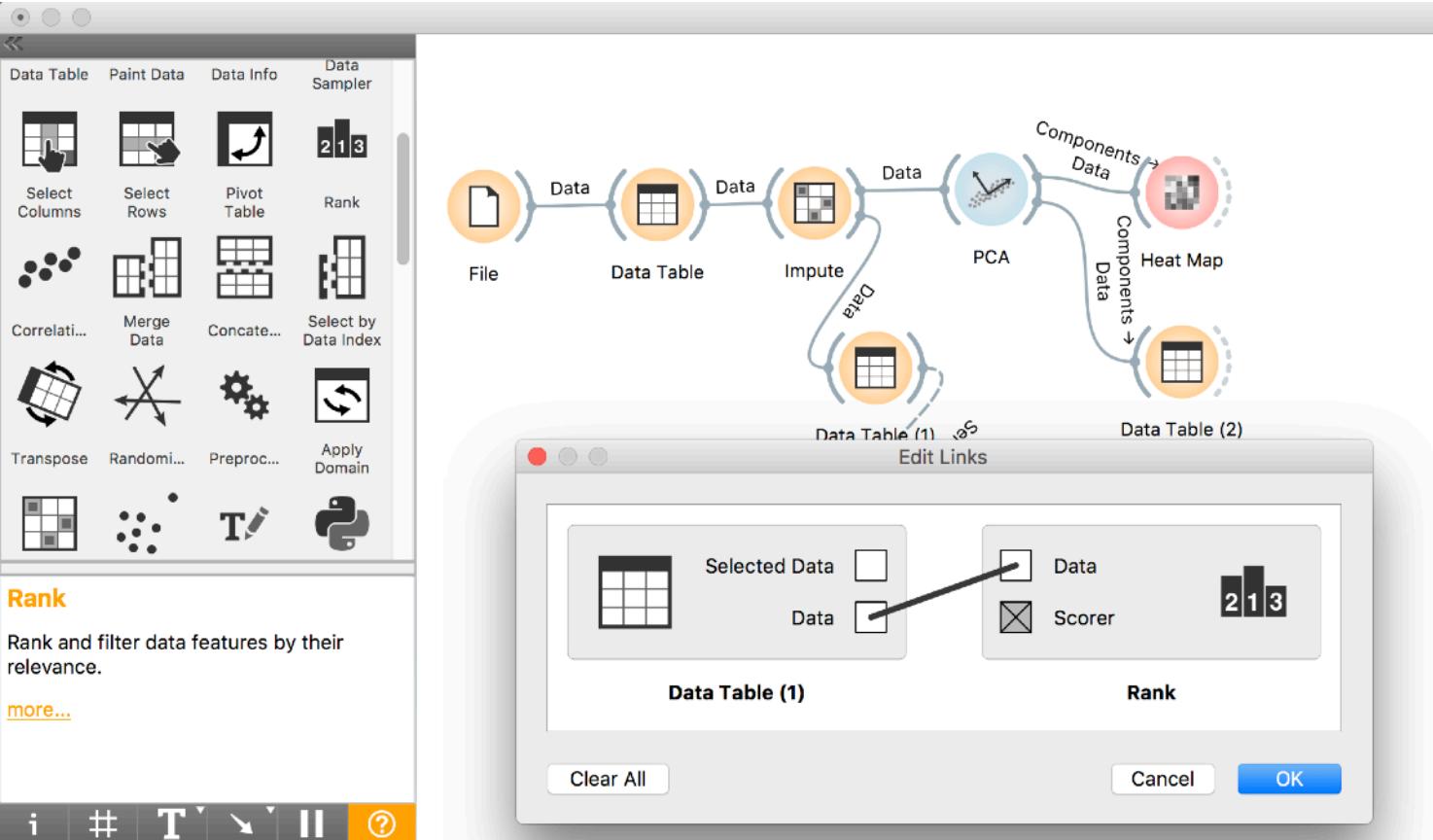
See [workflow examples](#), [YouTube tutorials](#), or open the [welcome screen](#).

The screenshot shows a data analysis application window. On the left, a sidebar lists various data manipulation tools: Select Columns, Select Rows, Pivot Table, Rank (which is highlighted with a red box), Correlat..., Merge Data, Concat..., Select by Data Index, Transpose, Randomi..., Preproc..., Apply Domain, and a few icons at the bottom. Below this is a message: "Select a widget to show its description." and links to "workflow examples", "YouTube tutorials", and the "welcome screen". At the bottom are standard window controls (minimize, maximize, close) and a toolbar with icons for information, hash, text, back, forward, double back, double forward, and help.

File Data Data Table Data Impute Data PCA Components Data Heat Map Components Data Data Table (1) Data Table (2)

Rank

The main area displays a workflow diagram. It starts with a "File" node connected to a "Data Table" node. This is followed by an "Impute" node, then a "PCA" node which generates a "Heat Map" node. The "PCA" node also outputs "Components Data" to two "Data Table" nodes labeled "(1)" and "(2)". A separate "Rank" node, also highlighted with a red box, is shown at the bottom, likely representing a step in the process.



Rank

Rank and filter data features by their relevance.

[more...](#)

Scoring Methods

- Information Gain
- Information Gain Ratio
- Gini Decrease
- ANOVA
- χ^2
- ReliefF
- FCBF

	#	Gain ratio	Gini
N PassengerId		0.000	0.000
C Survived	2	0.000	0.000
N Pclass		0.000	0.000
C Sex	2	0.000	0.000
C Sex-Binary	2	0.000	0.000
N Age		0.000	0.000
N SibSp		0.000	0.000
N Parch		0.000	0.000
N Fare		0.000	0.000
C Embarked	3	0.000	0.000

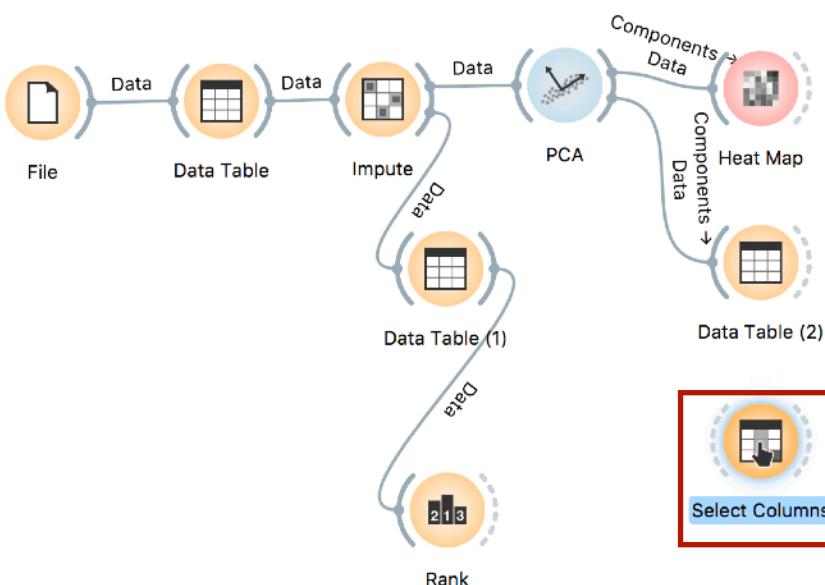
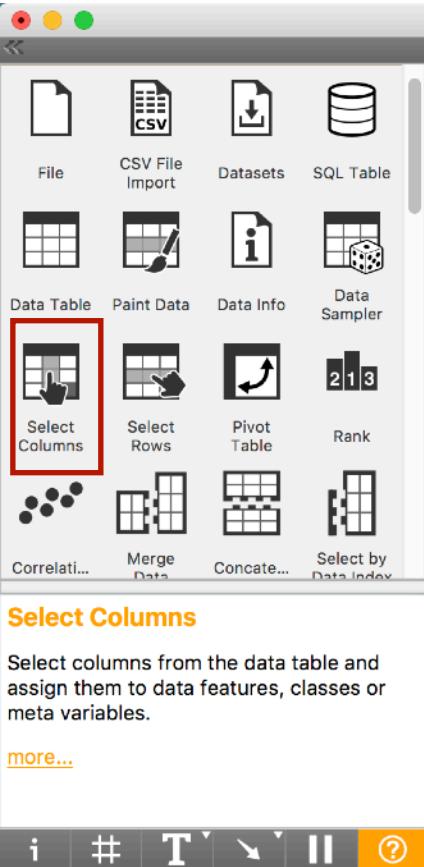
Select Attributes

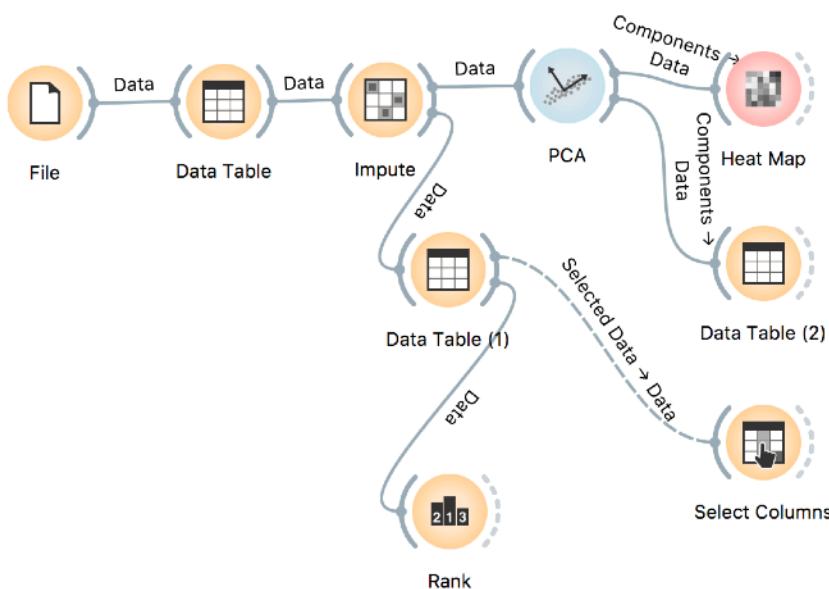
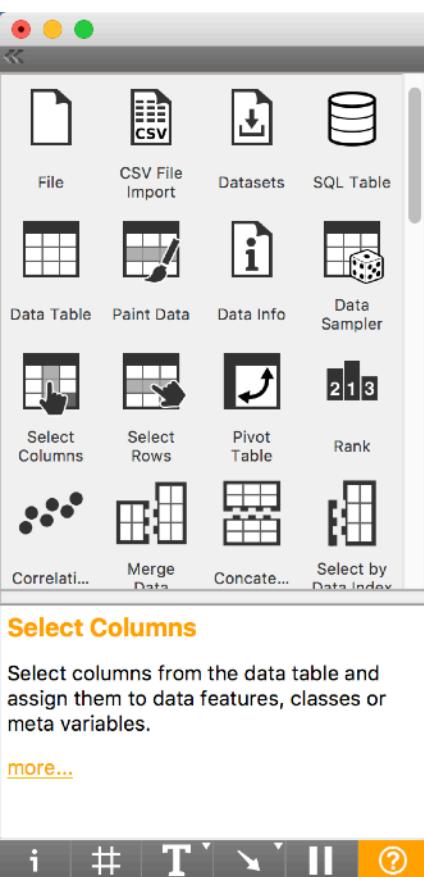
- None
- All
- Manual

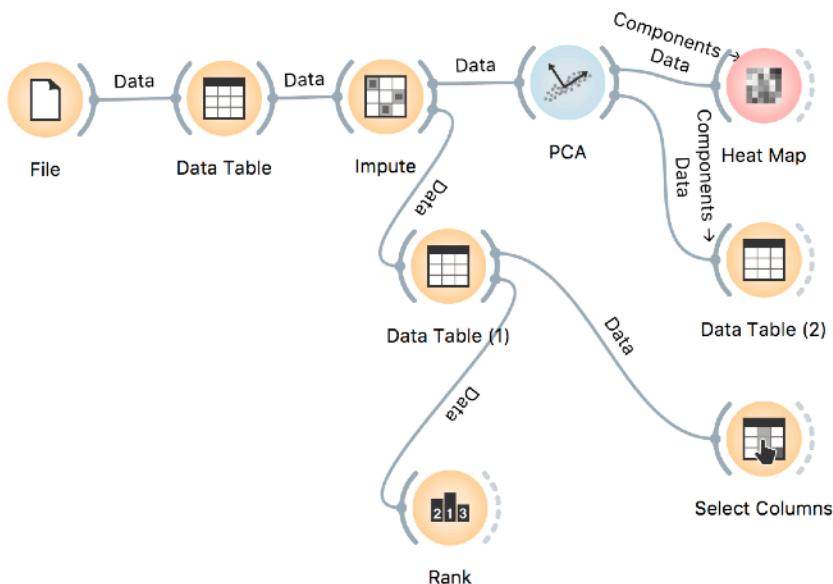
Best ranked:

Send Automatically

?







Select Columns

Available Variables

Filter

Up > Down

Features

Filter

- N PassengerId
- C Survived
- N Pclass
- C Sex
- C Sex-Binary
- A

Target Variable

Up > Down

- C Selected

Meta Attributes

Up > Down

- S Name
- S Ticket
- S Cabin
- C Selected (1)

Reset Send Automatically

The diagram illustrates a data processing workflow. It starts with a 'Data' input node, which feeds into a 'PCA' component. The 'PCA' component has two outputs: one labeled 'Components Data' leading to a 'Heat Map' visualization, and another labeled 'Data' leading to a 'Data Table (1)' component. From 'Data Table (1)', the data flows through a 'Select Columns' component (indicated by a hand cursor icon) to 'Data Table (2)'. Finally, 'Data Table (2)' outputs 'Components Data' to a 'Heat Map' visualization.

Select Columns

Available Variables

Features

- Pclass
- Sex
- Age

Target Variable

- Survived

Meta Attributes

Data Table (1)

Data Table (2)

Select Columns

PCA

Heat Map

Data

Components Data

Components Data

Components Data

Data

Data

Data

Rank

Send Automatically

The screenshot shows a data analysis application with a central workspace and a left sidebar. The sidebar contains icons for File, Data Table, Select Columns, Correlation, and a 'Select Columns' section which is currently active. The main workspace displays a 'Select Columns' dialog and a data flow diagram.

The 'Select Columns' dialog lists available variables: PassengerId, Sex-Binary, SibSp, Parch, Fare, Embarked, Selected, Name, Ticket, Cabin, and Selected (1). It also lists selected features (Pclass, Sex, Age) and the target variable (Survived). A 'Send Automatically' checkbox is checked at the bottom.

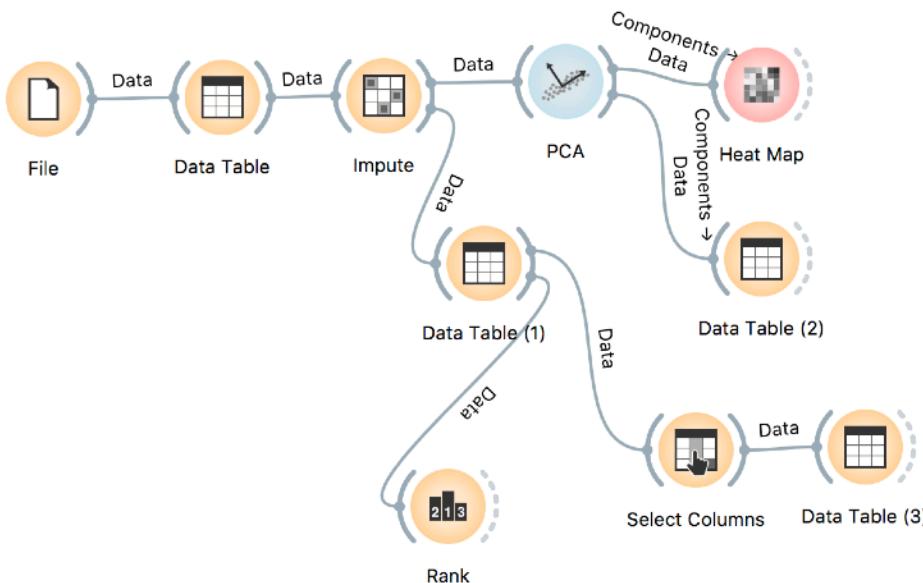
The data flow diagram on the right shows a 'PCA' node connected to a 'Heat Map' node. Both nodes receive 'Data' from a 'Data Table (1)' node and 'Components Data' from a 'Select Columns' node. The 'Heat Map' node also receives 'Components Data' from the 'PCA' node. A 'Data Table (2)' node receives 'Data' from the 'Heat Map' node. A 'Rank' node receives 'Data' from the 'Data Table (2)' node. A 'Select Columns' node is also present in the diagram.

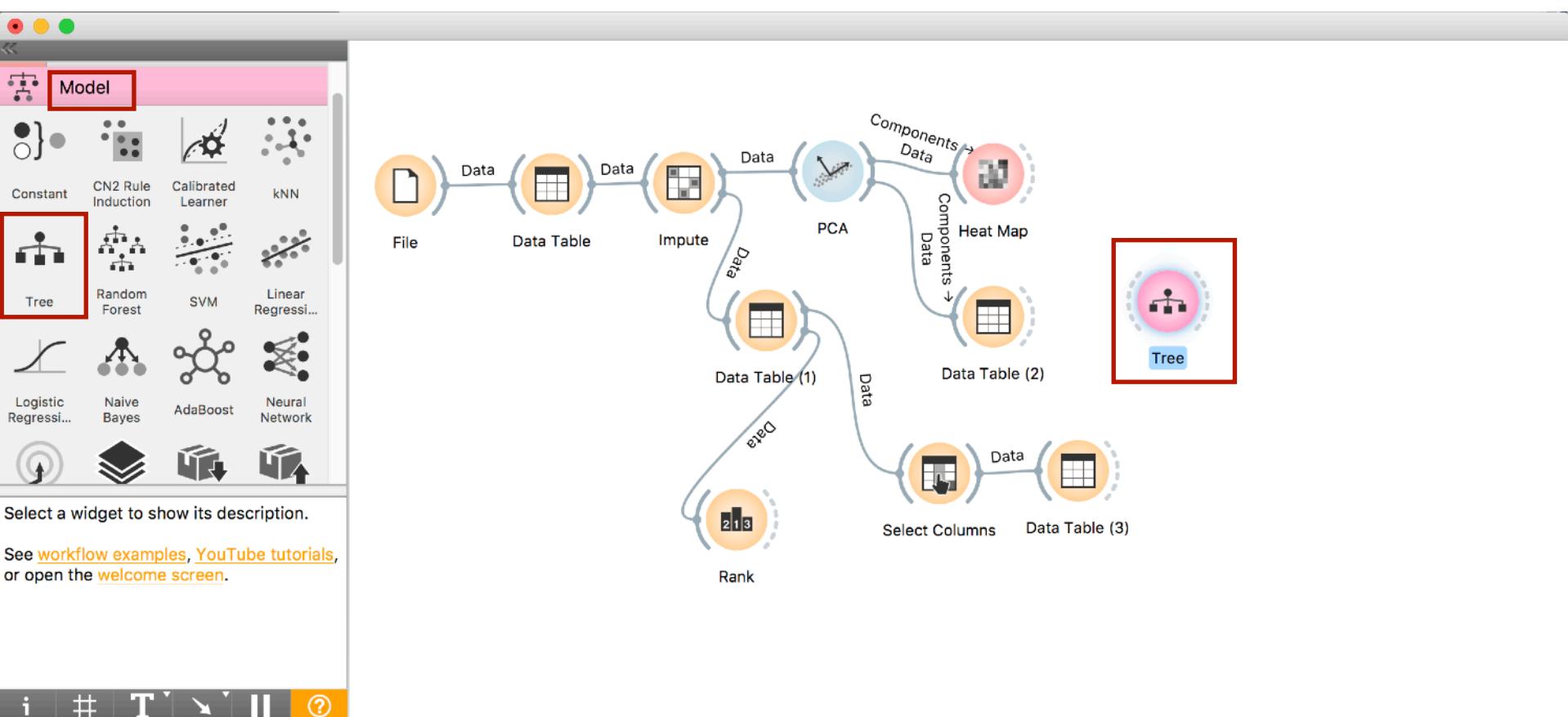
Model

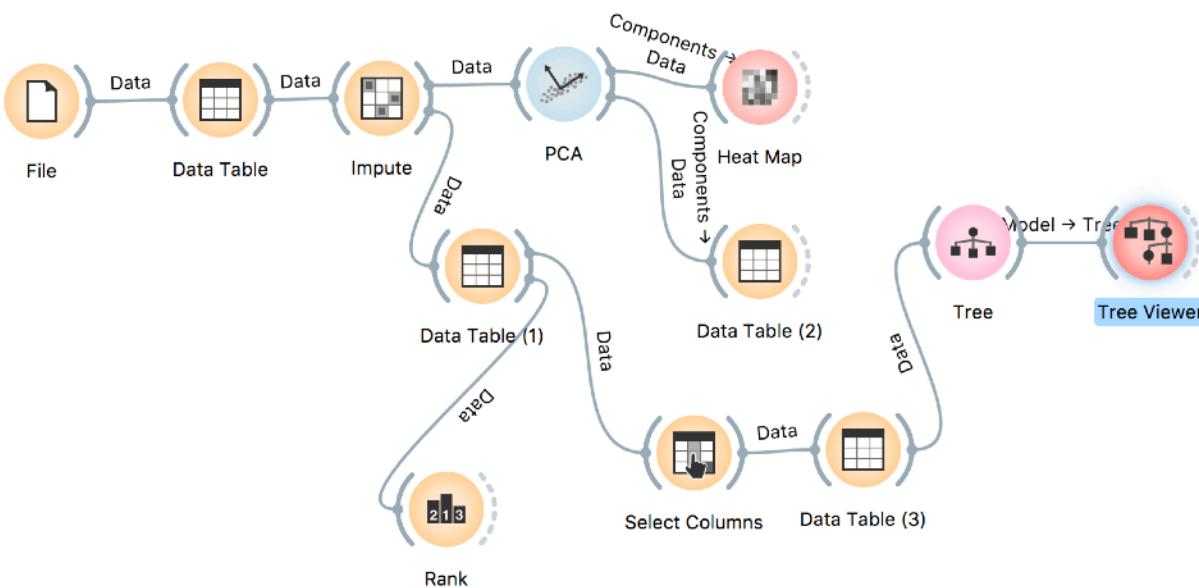
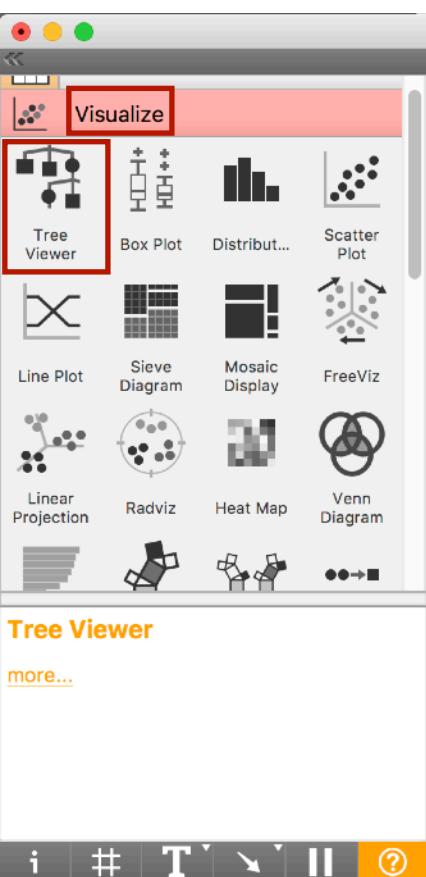
- Constant
- CN2 Rule Induction
- Calibrated Learner
- kNN
- Tree
- Random Forest
- SVM
- Linear Regressi...
- Logistic Regressi...
- Naive Bayes
- AdaBoost
- Neural Network

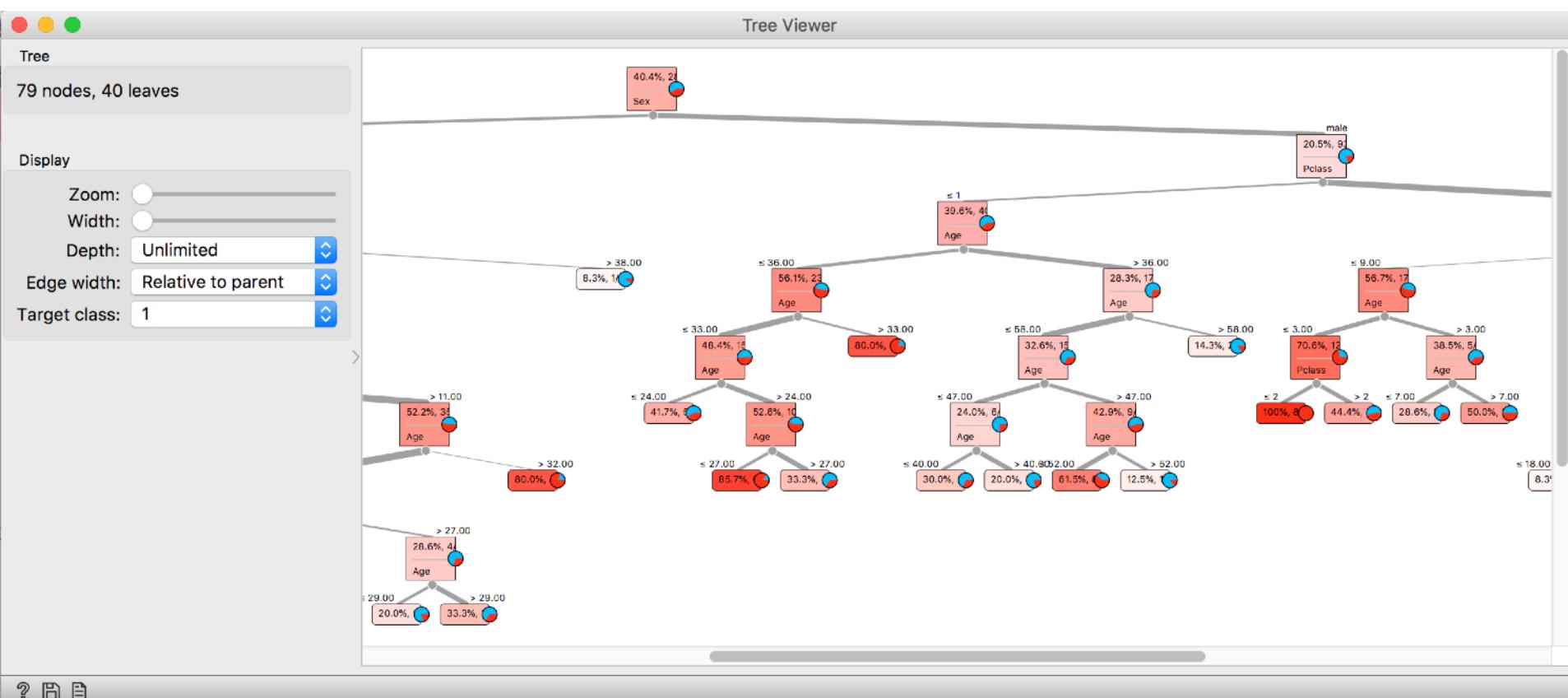
Data Table
View the dataset in a spreadsheet.
[more...](#)

i # T ⌛ || ?





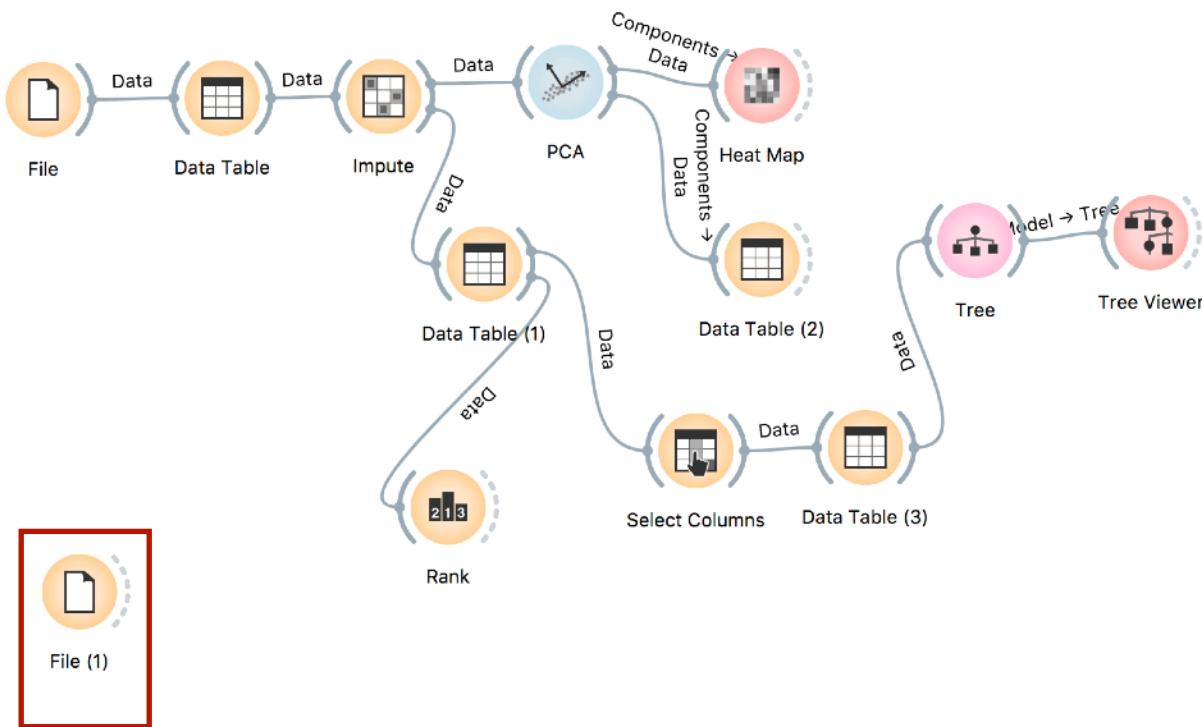




After training the model, it's time for testing it.

The screenshot shows the KNIME interface with the following components:

- Toolbar:** Includes icons for File, CSV File Import, Datasets, SQL Table, Data Table, Paint Data, Data Info, Data Sampler, Select Columns, Select Rows, Pivot Table, Rank, and a scatter plot icon.
- Data Tab:** Selected tab, showing options for File, CSV File Import, Datasets, and SQL Table.
- File Section:** Describes reading data from an input file or network and sending a data table to the output.
- more...:** A link to additional information.
- Bottom Bar:** Includes icons for Help, View, Tools, Run, Stop, and a question mark.



File (1)

File: **titanic_test.csv**

URL:

Info

891 instance(s)
 10 feature(s) (2.0% missing values)
 Data has no target variable.
 3 meta attribute(s)

Columns (Double click to edit)

	Name	Type	Role	Values
1	PassengerId	N numeric	feature	
2	Survived	C categorical	feature	0, 1
3	Pclass	N numeric	feature	
4	Sex	C categorical	feature	female, male
5	Sex-Binary	C categorical	feature	0, 1
6	Age	N numeric	feature	
7	SibSp	N numeric	feature	

Components → Data

```

graph LR
    CA((CA)) -- Components → Data --> HeatMap((Heat Map))
    CA -- Components → Data --> DT2((Data Table (2)))
    CA -- Components → Data --> Tree((Tree))
    CA -- Components → Data --> DT3((Data Table (3)))
    DT2 -- Data --> Tree
    Tree -- Model → Tree --> TV((Tree Viewer))
    Tree -- Data --> DT3
    DT3 -- Data --> TV
    
```

CA
 Heat Map
 Data Table (2)
 Tree
 Tree Viewer
 Select Columns
 Data Table (3)

Data

File CSV File Import Datasets SQL Table

Data Table Paint Data Data Info Data Sampler

Select Columns Select Rows Pivot Table Rank

Data Table
View the dataset in a spreadsheet.
[more...](#)

891 instances
10 features (2.0% missing values)
No target variable.
3 meta attributes (25.7% missing values)

Variables

Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

Selection

Select full rows

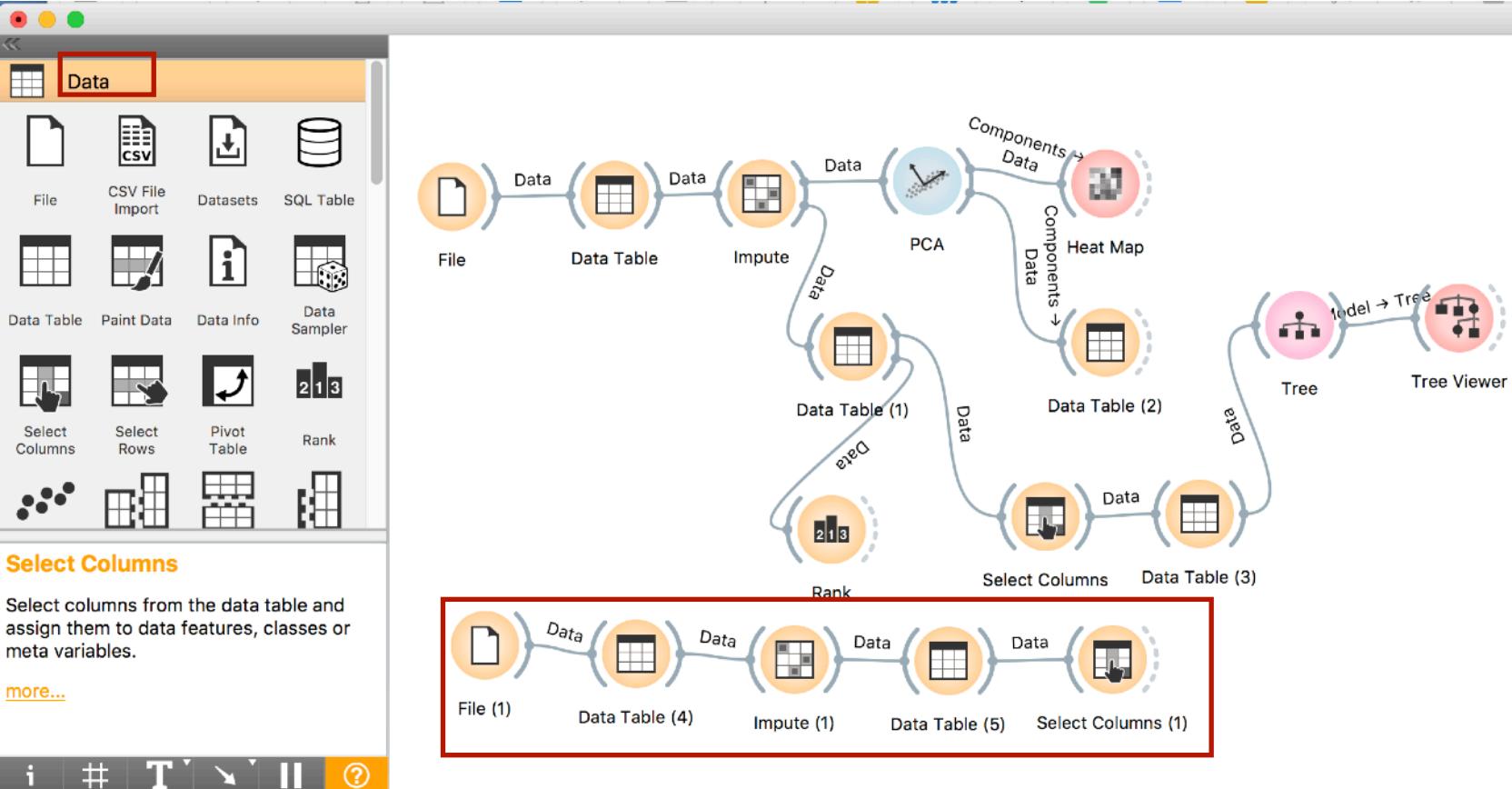
Restore Original Order

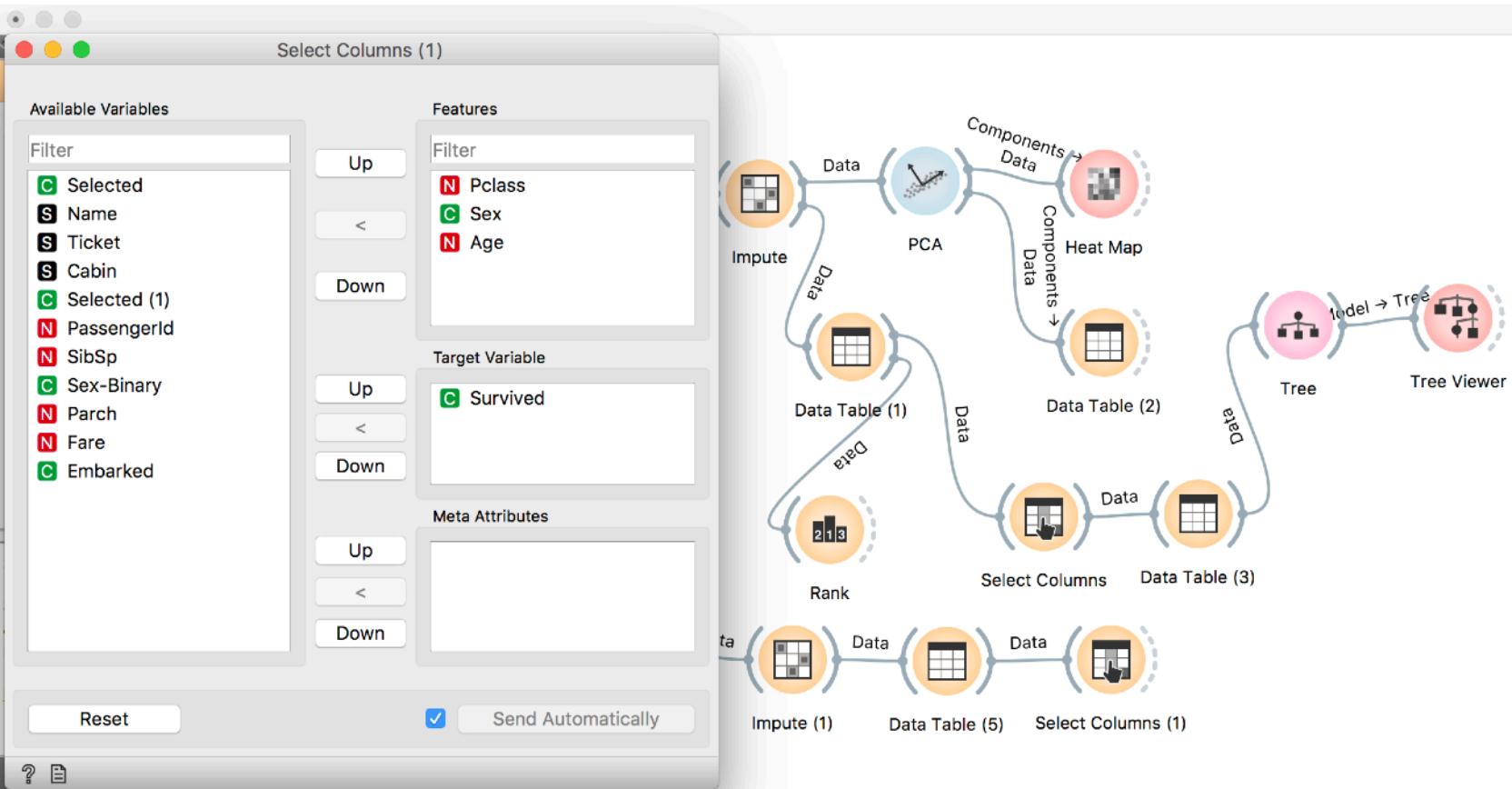
Send Automatically

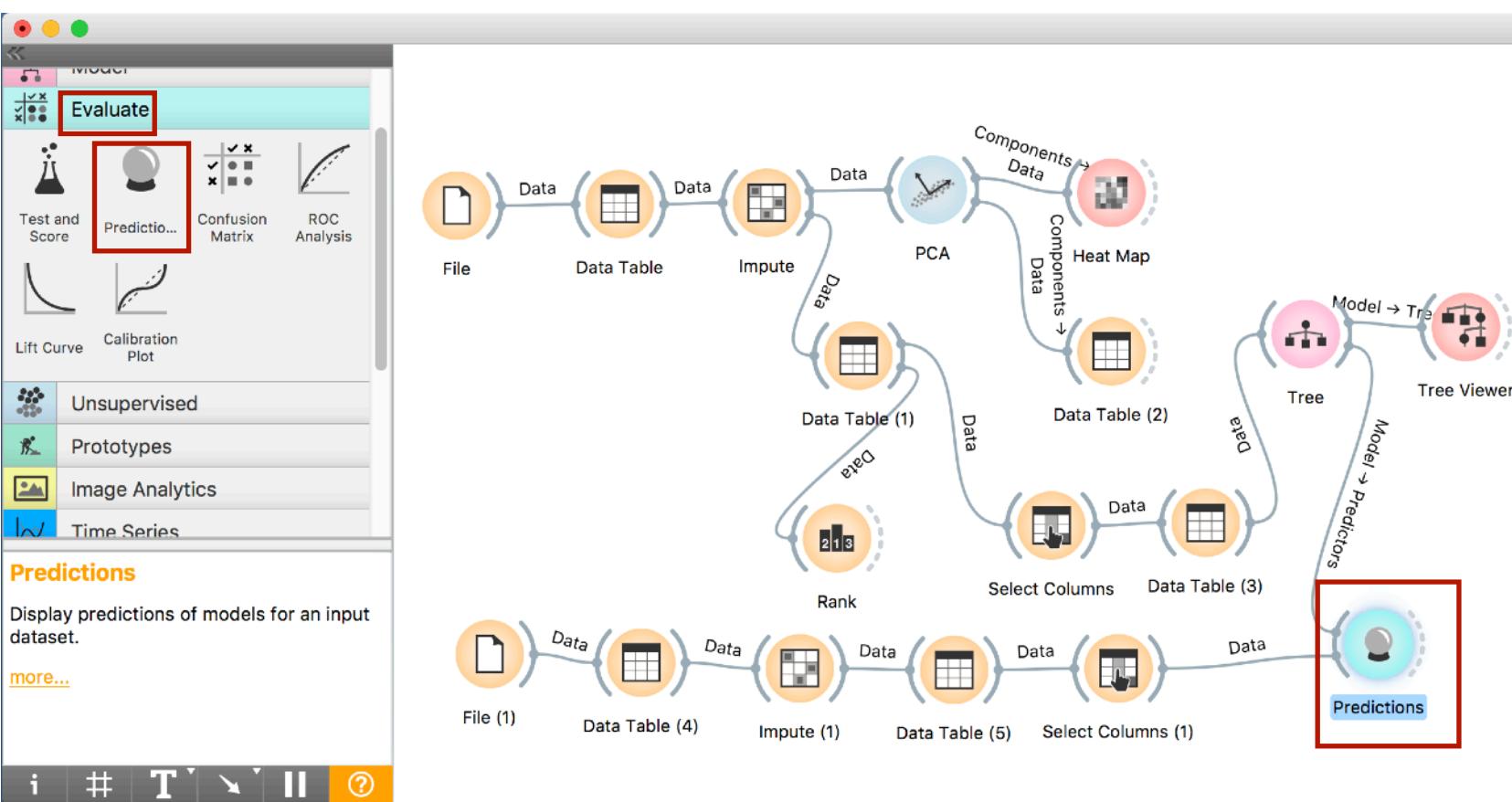
?

Data Table (4)

	Age	SibSp	Parch	Fare	Embarked
1	22.00	1	0	7.2500	S
2	38.00	1	0	71.2833	C
3	26.00	0	0	7.9250	S
4	35.00	1	0	53.1000	S
5	35.00	0	0	8.0500	S
6	?	0	0	8.4583	Q
7	54.00	0	0	51.8625	S
8	2.00	3	1	21.0750	S
9	27.00	0	2	11.1333	S
10	14.00	1	0	30.0708	C
11	4.00	1	1	16.7000	S
12	58.00	0	0	26.5500	S
13	20.00	0	0	8.0500	S
14	39.00	1	5	31.2750	S
15	14.00	0	0	7.8542	S
16	55.00	0	0	16.0000	S
17	2.00	4	1	29.1250	Q
18	?	0	0	13.0000	S
19	31.00	1	0	18.0000	S
~	~	~	~	~	~







Predictions

Show probabilities for

	Tree	Survived	Pclass	Sex	Age
0	0.91 : 0.09	0	3	male	22.00
1	0.04 : 0.96	1	1	female	38.00
2	0.17 : 0.83	1	3	female	26.00
3	0.04 : 0.96	1	1	female	35.00
4	0.96 : 0.04	0	3	male	35.00
5	0.88 : 0.12	0	1	male	54.00
6	0.56 : 0.44	0	3	male	2.00
7	0.17 : 0.83	1	3	female	27.00
8	0.07 : 0.93	1	2	female	14.00
9	0.00 : 1.00	1	3	female	4.00
10	0.17 : 0.83	1	1	female	58.00
11	0.83 : 0.17	0	3	male	20.00
12	0.96 : 0.04	0	3	male	39.00
13	0.30 : 0.70	0	3	female	14.00
14	0.04 : 0.96	1	2	female	55.00
15	0.56 : 0.44	0	3	male	2.00
16	0.67 : 0.33	0	3	female	31.00
17	0.96 : 0.04	0	2	male	35.00

Model AUC CA F1 Precision Recall

Tree 0.913 0.851 0.848 0.854 0.851

Restore Original Order

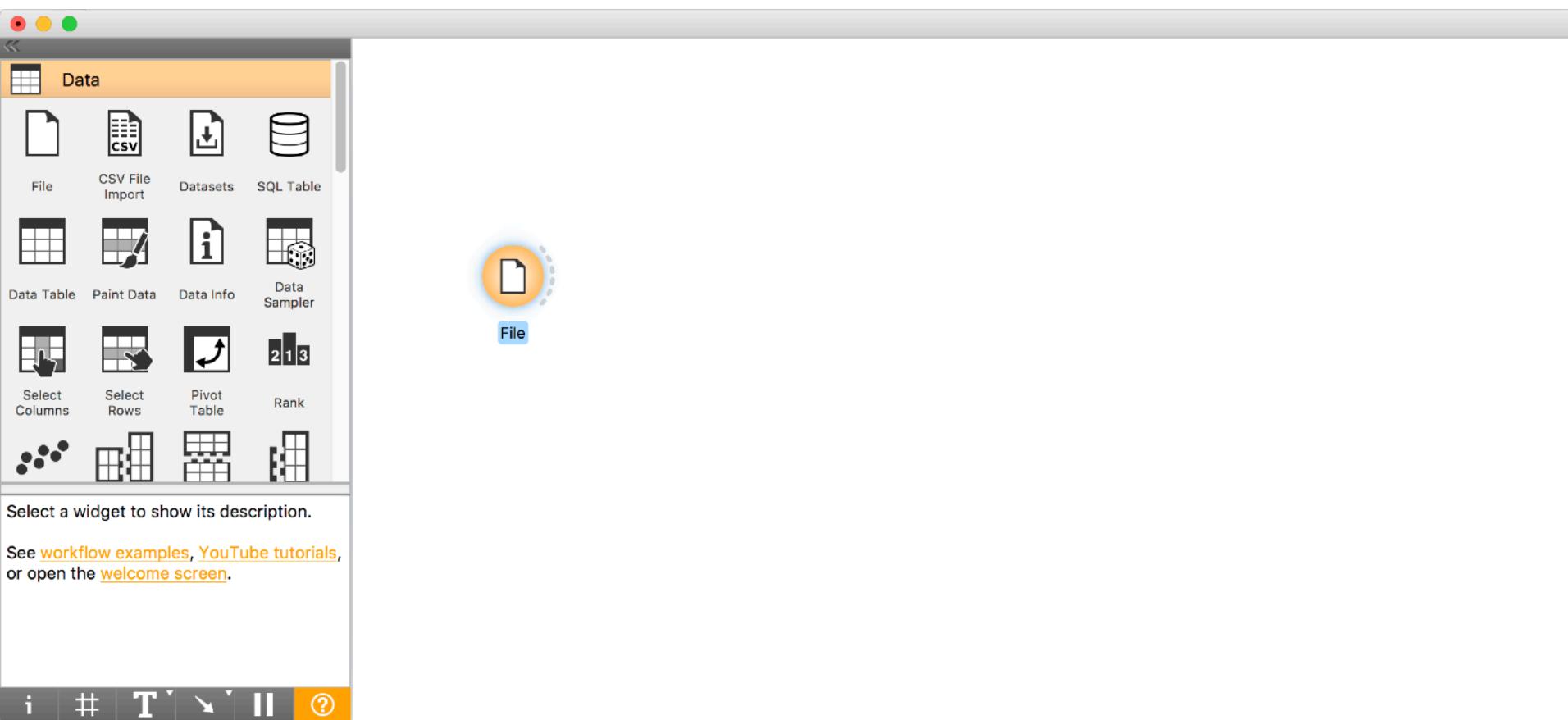
712

```

graph LR
    DT1[Data Table (1)] -- Data --> SC1[Select Columns (1)]
    DT2[Data Table (2)] -- Components Data --> HeatMap[Heat Map]
    DT2 -- Components Data --> Tree[Tree]
    DT3[Data Table (3)] -- Components Data --> Tree
    DT3 -- Data --> SC2[Select Columns (2)]
    DT3 -- Data --> Predictions[Predictions]
    SC1 -- Model > Tree --> Tree
    Tree -- Model > Predictors --> Predictions
    SC2 -- Model > Tree --> Tree
    Tree -- Model > Tree --> TV[Tree Viewer]
    TV -- Model > Predictors --> Predictions
  
```

The diagram illustrates a data science workflow. It starts with three data tables (DT1, DT2, DT3) which feed into two selection components (Select Columns (1) and Select Columns (2)). The output from Select Columns (1) goes to a 'Tree' component, which then feeds into a 'Tree Viewer'. The output from Select Columns (2) also goes to the 'Tree' component. The 'Tree' component then outputs to a 'Predictions' component. Additionally, DT2 feeds into a 'Heat Map' component and a 'Tree' component. DT3 feeds into both a 'Heat Map' component and a 'Tree' component. The 'Tree' component from DT3 also feeds into the 'Tree Viewer'.

**Another way of training and testing the model
(when you only have one dataset not separated
into training and testing sets).**



File: titanic.csv

URL: http://localhost/datasets/iris.csv

Info

891 instance(s)
3 feature(s) (6.6% missing values)
Classification; categorical class with 2 values (no missing values)
4 meta attribute(s)

Columns (Double click to edit)

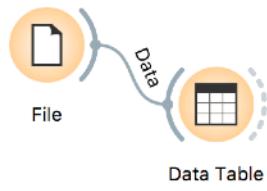
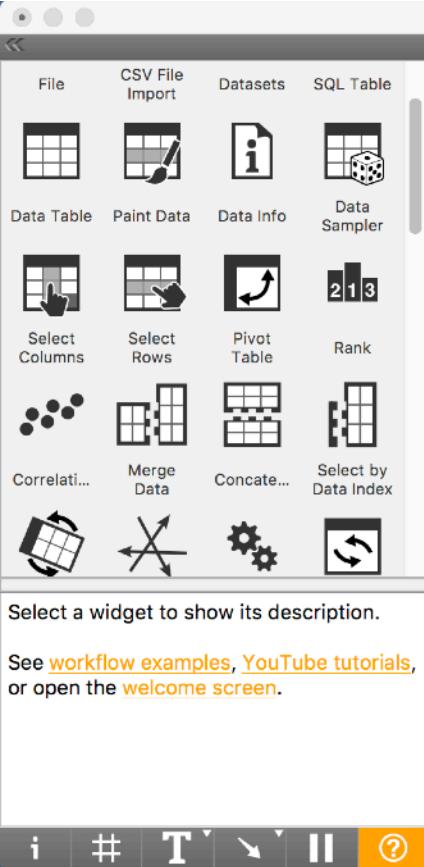
	Name	Type	Role	Values
1	Pclass	N numeric	feature	
2	Sex	C categorical	feature	female, male
3	Age	N numeric	feature	
4	Survived	C categorical	target	0, 1
5	Name	S text	meta	
6	Ticket	S text	meta	
7	Cabin	S text	meta	

Browse documentation datasets

Reset Apply

Tree Viewer

Predictors Predictions



Data Table

View the dataset in a spreadsheet-like view.

more...

CSV File Import

Info

891 instances

3 features (6.6% missing values)

Discrete class with 2 values (no missing values)

4 meta attributes (19.3% missing values)

Variables

Show variable labels (if present)

Visualize numeric values

Color by instance classes

Selection

Select full rows

Restore Original Order

Send Automatically

Cabin Selected Pclass Sex Age

	Cabin	Selected	Pclass	Sex	Age
3	?	No	3.0	female	26.00
4	C123	No	1.0	female	35.00
5	?	No	3.0	male	35.00
6	?	No	3.0	male	?
7	E46	No	1.0	male	54.00
8	?	No	3.0	male	2.00
9	?	No	3.0	female	27.00
10	?	No	2.0	female	14.00
11	G6	No	3.0	female	4.00
12	C103	No	1.0	female	58.00
13	?	No	3.0	male	20.00
14	?	No	3.0	male	39.00
15	?	No	3.0	female	14.00
16	?	No	2.0	female	55.00
17	?	No	3.0	male	2.00
18	?	No	2.0	male	?
19	?	No	3.0	female	31.00
20	?	No	3.0	female	?
21	?	No	2.0	male	35.00

The screenshot shows a software interface for data processing. On the left is a vertical toolbar with various icons and labels:

- File
- CSV File Import
- Datasets
- SQL Table

Below these are several data manipulation icons:

- Data Table
- Paint Data
- Data Info
- Data Sampler

Then there are more complex operations:

- Select Columns
- Select Rows
- Pivot Table
- Rank

Statistical and correlation tools:

- Correlat...
- Merge Data
- Concaten...
- Select by Data Index

Finally, there are utility icons:

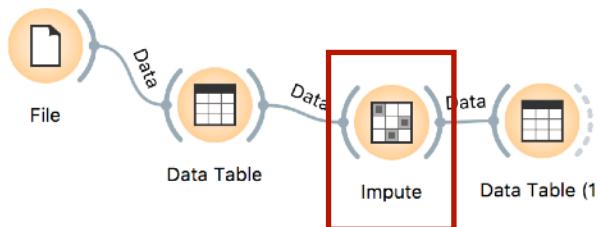
- Transpose
- Swap
- Settings
- Refresh

Data Table

View the dataset in a spreadsheet.

[more...](#)

At the bottom are standard window controls (Minimize, Maximize, Close) and a toolbar with icons for search, filter, column width, and help.



Data Table (1)

Info

714 instances
3 features (no missing values)
Discrete class with 2 values (no missing values)
5 meta attributes (14.8% missing values)

Variables

Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

Selection

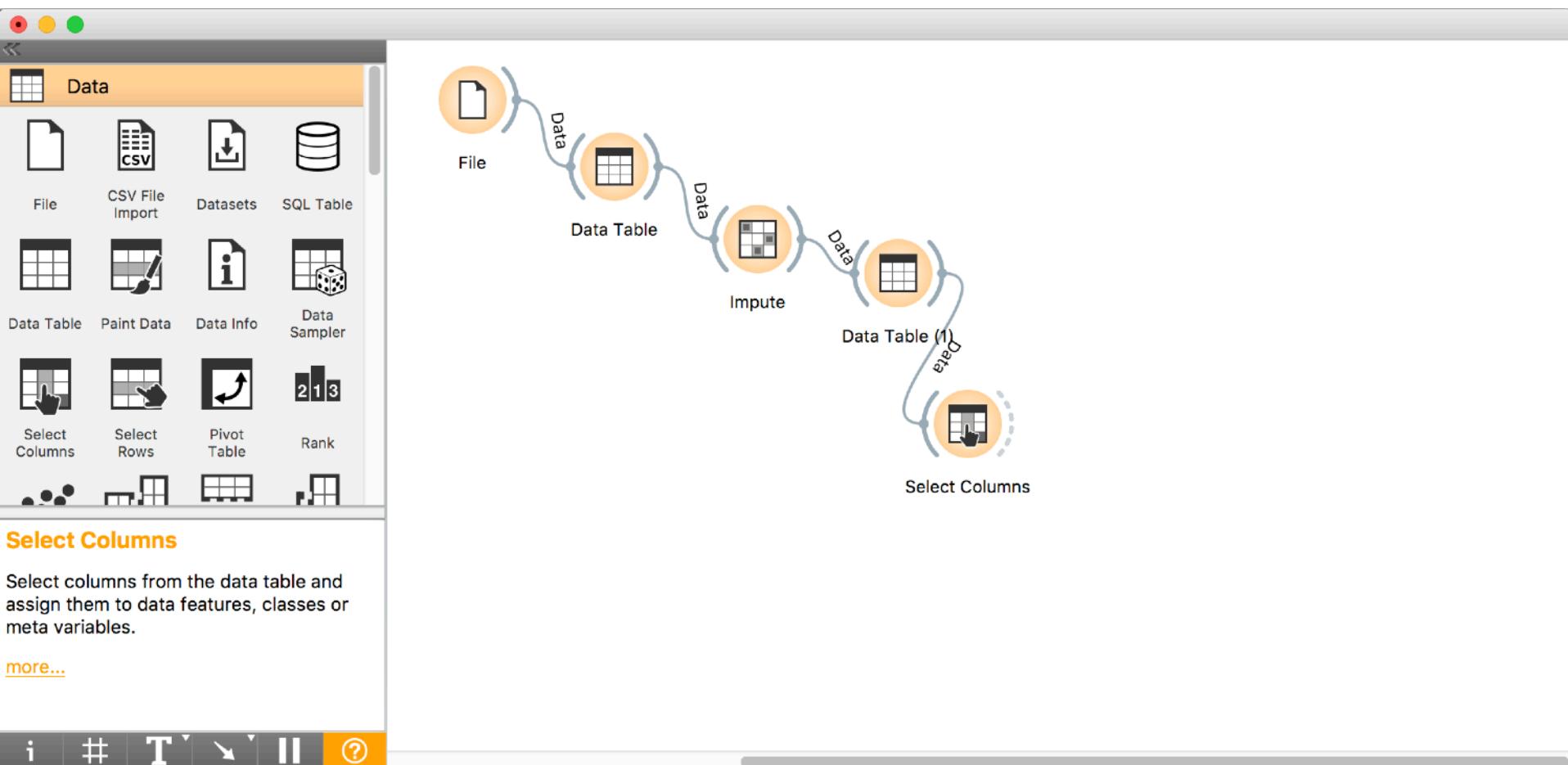
Select full rows

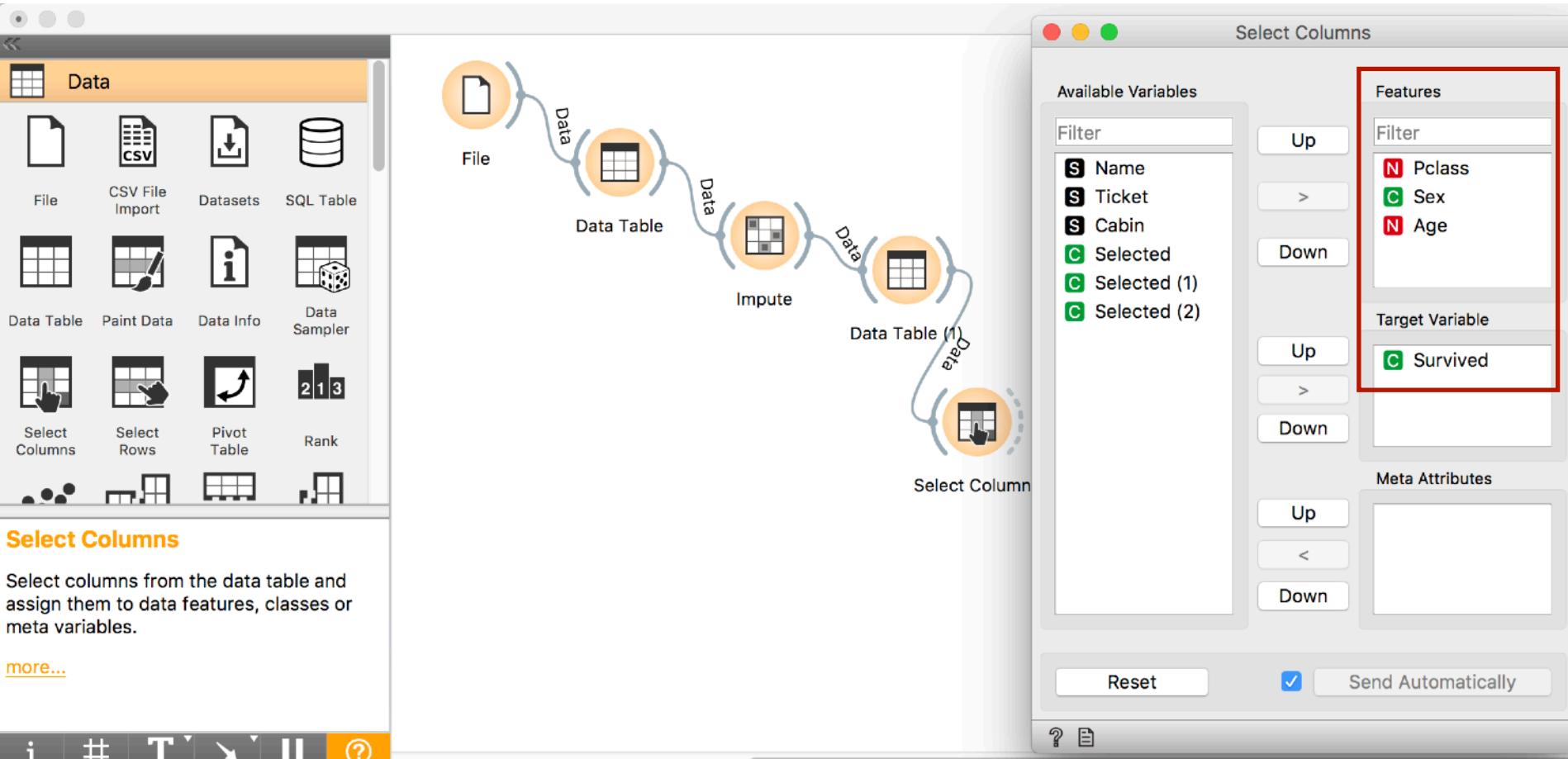
Restore Original Order

Send Automatically

Age

	Selected	Selected (1)	Pclass	Sex	Age
48	No	No	1.0	male	45.00
49	No	No	3.0	male	4.00
50	No	No	2.0	female	29.00
51	No	No	3.0	male	19.00
52	No	No	3.0	female	17.00
53	No	No	3.0	male	26.00
54	No	No	2.0	male	32.00
55	No	No	3.0	female	16.00
56	No	No	2.0	male	21.00
57	No	No	3.0	male	26.00
58	No	No	3.0	male	32.00
59	No	No	3.0	male	25.00
60	No	No	2.0	male	0.83
61	No	No	3.0	female	30.00
62	No	No	3.0	male	22.00
63	No	No	3.0	male	29.00
64	No	No	1.0	male	28.00
65	No	No	2.0	female	17.00
66	No	No	3.0	female	33.00





Model

Evaluate

- Test and Score
- Prediction...
- Confusion Matrix
- ROC Analysis

- Lift Curve
- Calibration Plot

Unsupervised

Prototypes

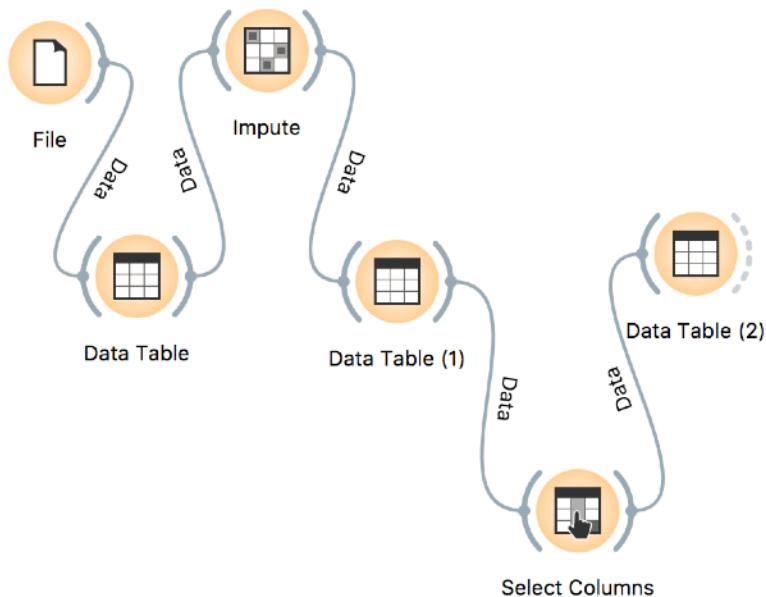
Image Analytics

Data Table

View the dataset in a spreadsheet.

[more...](#)

i # T ↻ II ?



Data Table (2)

	Survived	Pclass	Sex	Age
1	0	3.0	male	22.00
2	1	1.0	female	38.00
3	1	3.0	female	26.00
4	1	1.0	female	35.00
5	0	3.0	male	35.00
6	0	1.0	male	54.00
7	0	3.0	male	2.00
8	1	3.0	female	27.00
9	1	2.0	female	14.00
10	1	3.0	female	4.00
11	1	1.0	female	58.00
12	0	3.0	male	20.00
13	0	3.0	male	39.00
14	0	3.0	female	14.00
15	1	2.0	female	55.00
16	0	3.0	male	2.00
17	0	3.0	female	31.00
18	0	2.0	male	35.00
19	1	2.0	male	34.00
20	1	3.0	female	15.00

Info

714 instances (no missing values)
3 features (no missing values)
Discrete class with 2 values (no missing values)
No meta attributes

Variables

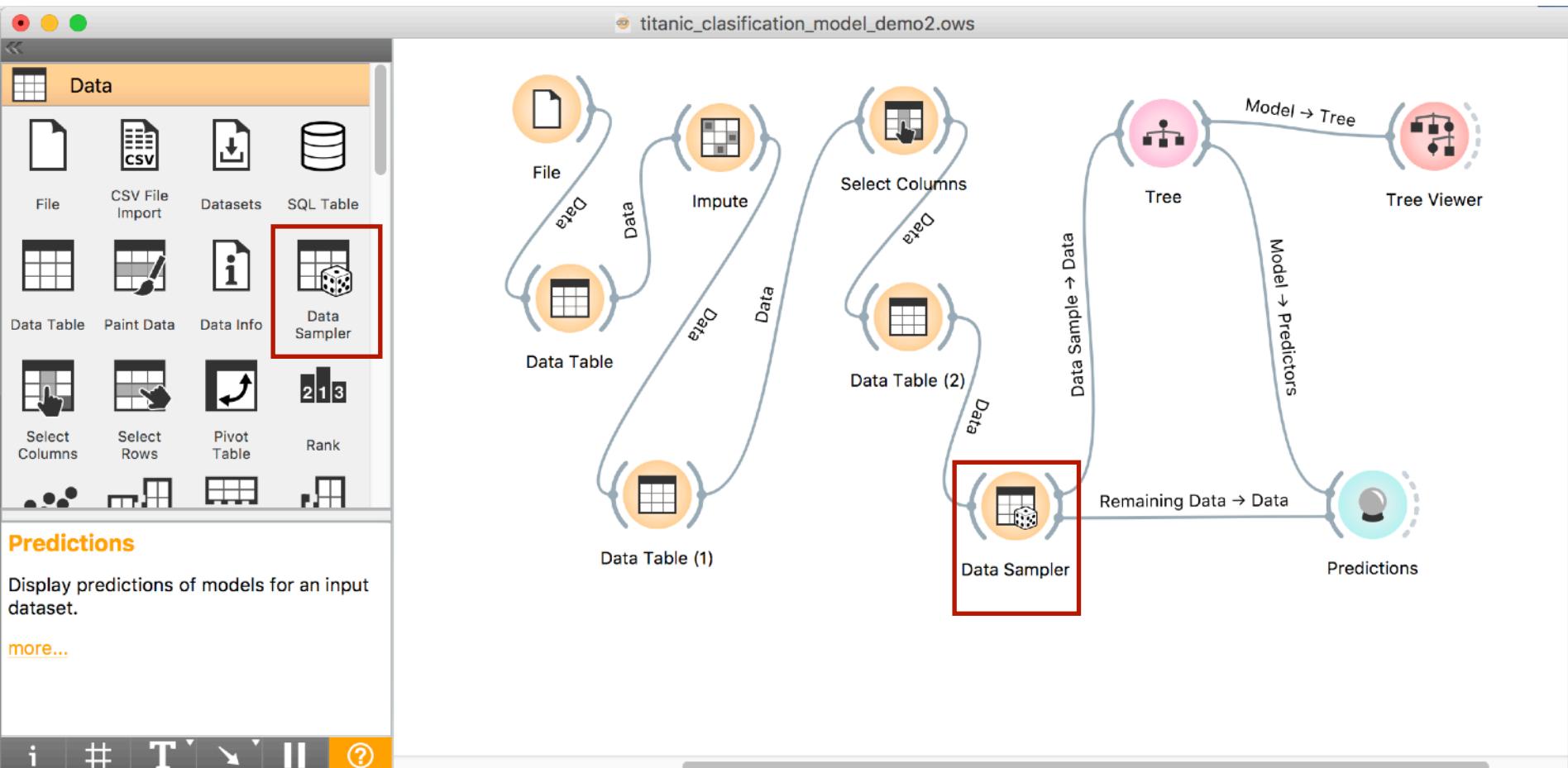
Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

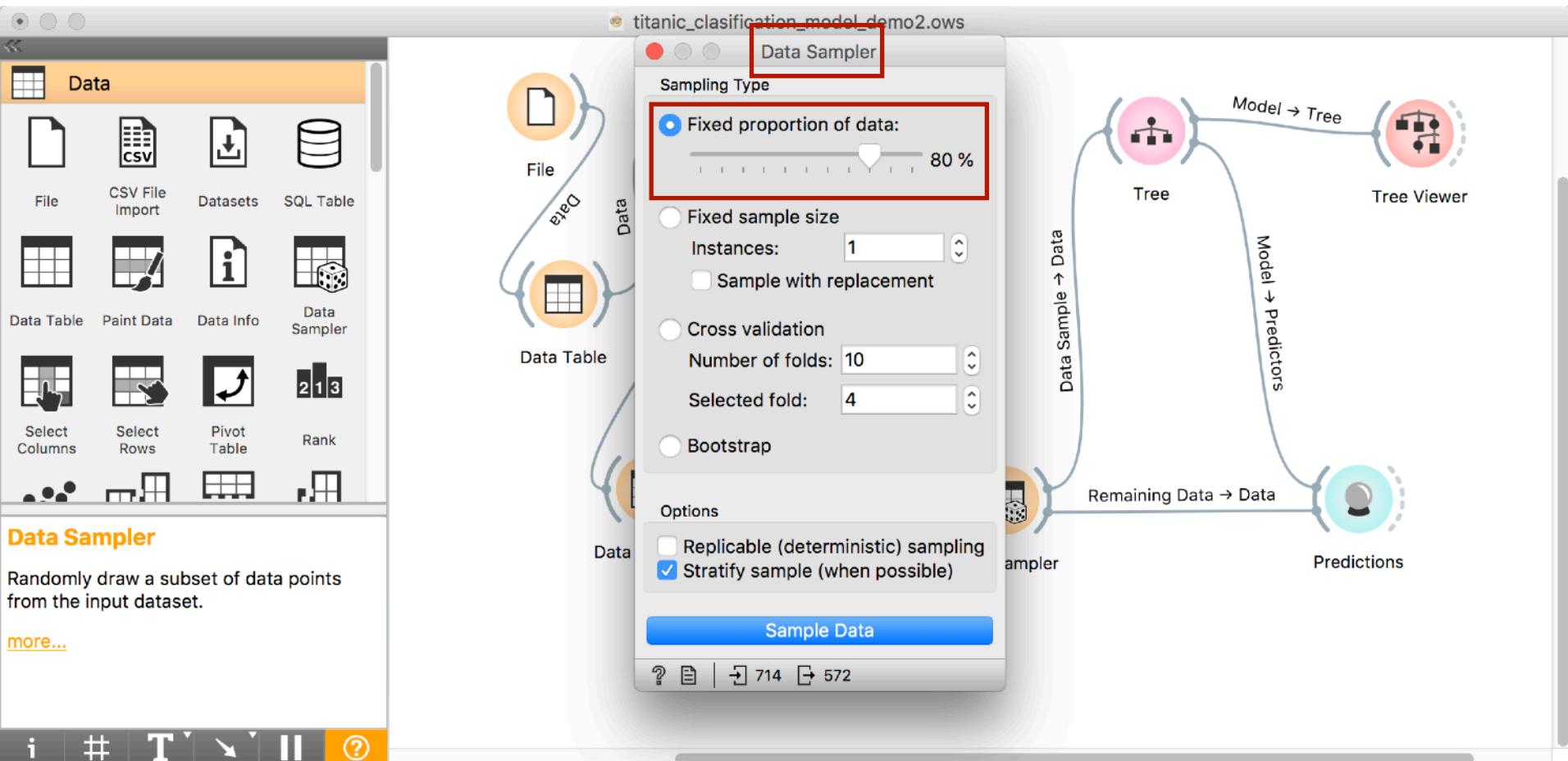
Selection

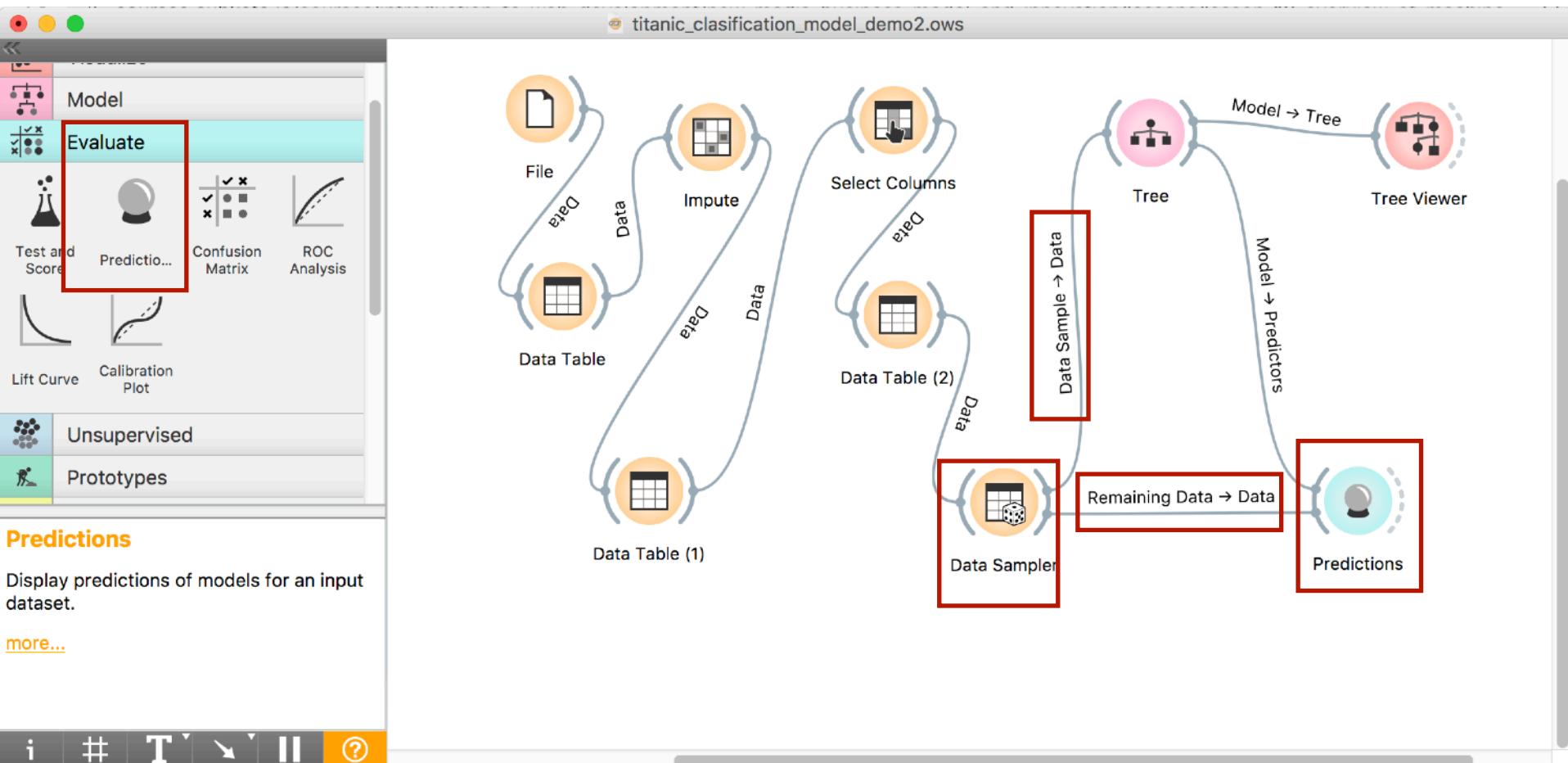
Select full rows

Restore Original Order

Send Automatically







titanic_clasification_model_demo2.ows

Show probabilities for

Tree	Survived	Selected	Pclass	Sex	Age
0	1	No	3.0	male	20.00
1	0.75 : 0.25 → 0	0	3.0	male	35.00
2	0.97 : 0.03 → 0	1	1.0	female	35.00
3	0.04 : 0.96 → 1	0	3.0	male	19.00
4	0.85 : 0.15 → 0	1	1.0	female	21.00
5	0.09 : 0.91 → 1	1	1.0	female	38.00
6	0.04 : 0.96 → 1	1	2.0	male	0.67
7	0.00 : 1.00 → 1	1	2.0	female	25.00
8	0.09 : 0.91 → 1	0	2.0	male	44.00
9	0.71 : 0.29 → 0	1	1.0	female	30.00
10	0.04 : 0.96 → 1	1	3.0	male	32.00
11	0.55 : 0.45 → 0	0	2.0	male	25.00
12	0.75 : 0.25 → 0	1	3.0	female	22.00
13	0.40 : 0.60 → 1	0	3.0	male	14.00
14	0.89 : 0.11 → 0	0	1.0	male	19.00
15	0.85 : 0.15 → 0	0	2.0	male	24.00
16	0.00 : 1.00 → 1	No	3.0	male	24.00

Predictions

Display predictions of model on dataset.

more...

Restore Original Order

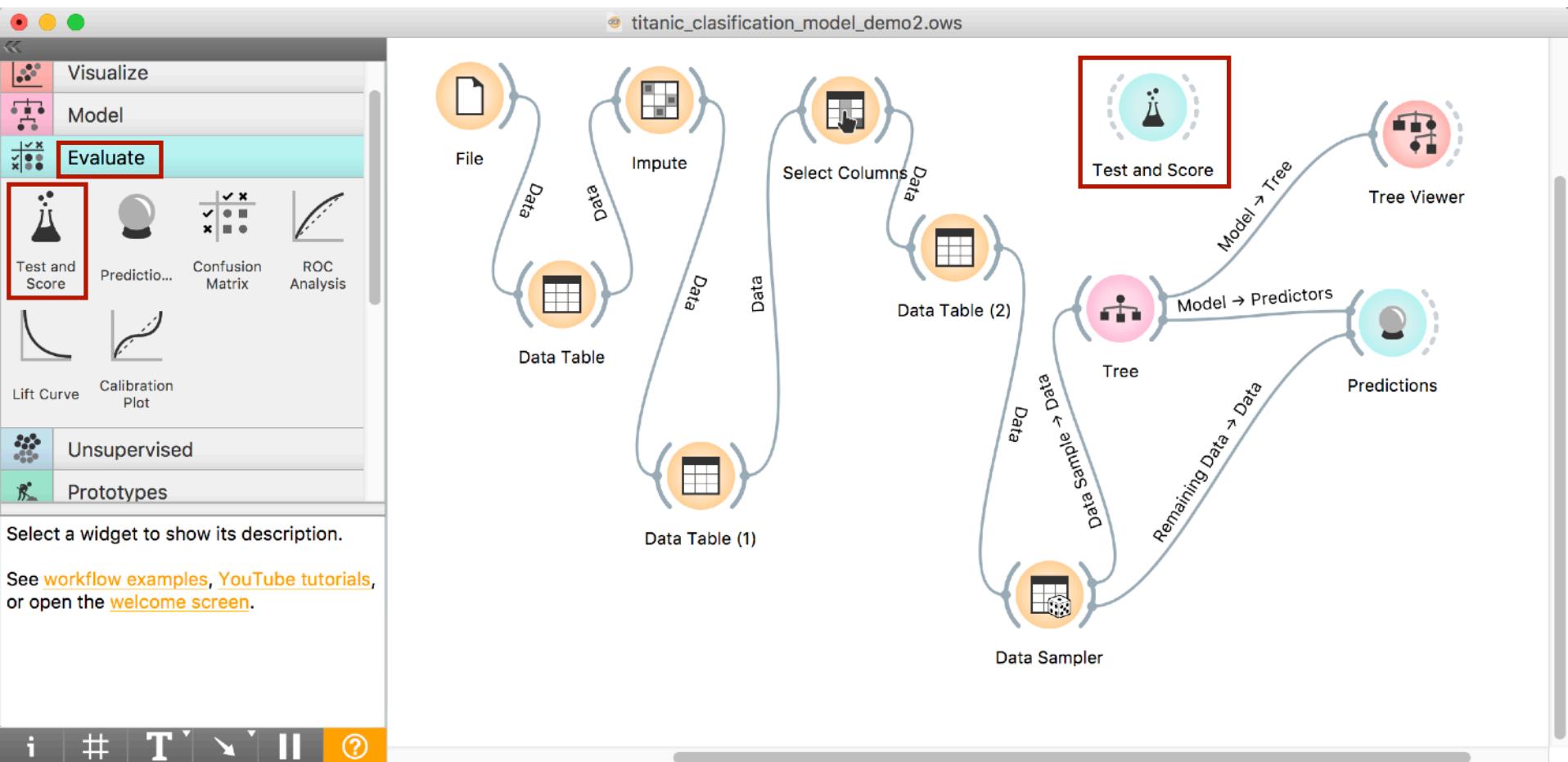
Model AUC CA F1 Precision Recall

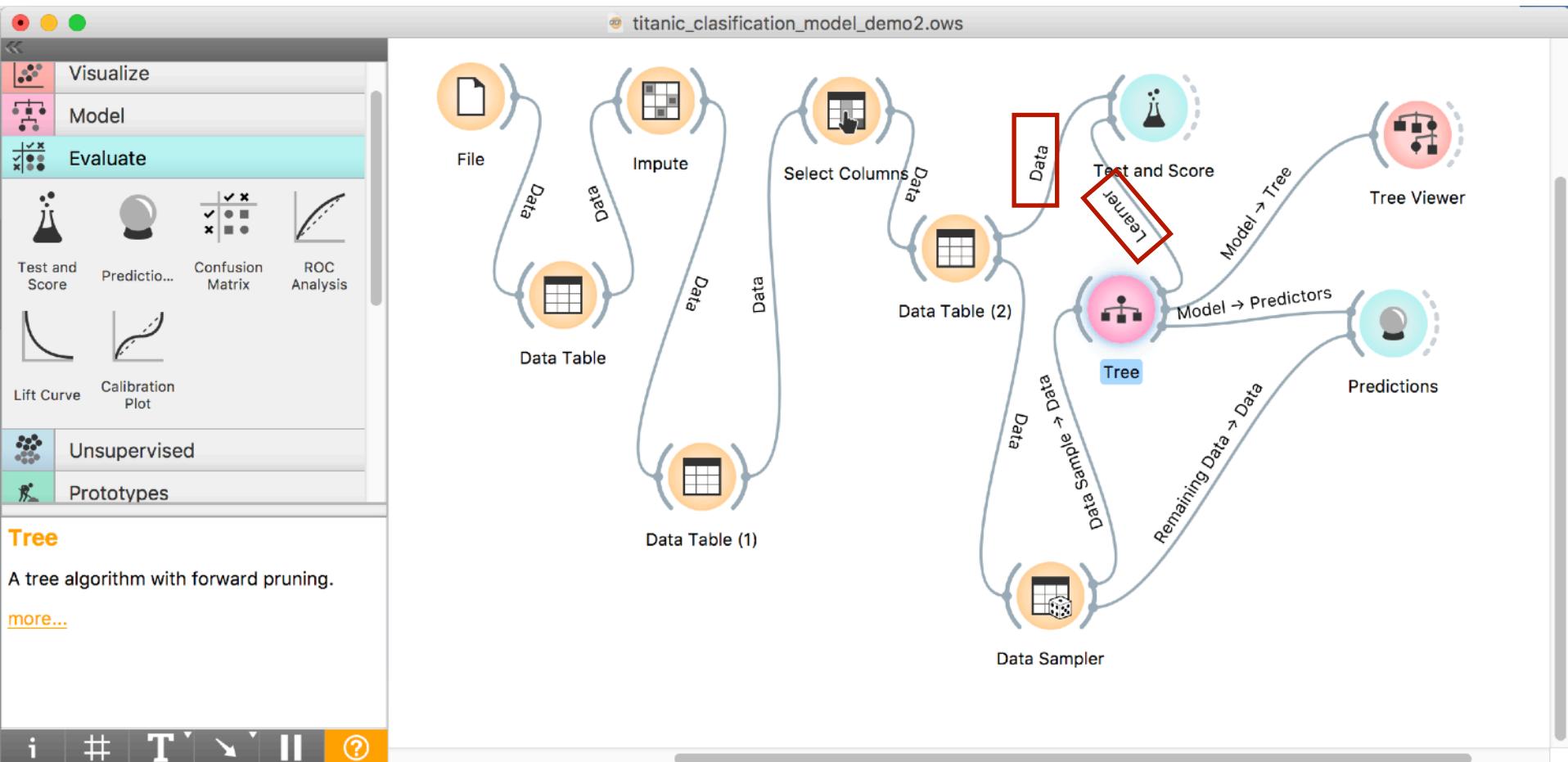
Tree 0.940 0.901 0.901 0.902 0.901

142

i # T ↻ II ?

One more method





titanic_classification_model_demo2.ows

Test and Score

Sampling

Cross validation
Number of folds: 10
 Stratified

Cross validation by feature
Selected

Random sampling
Repeat train/test: 100
Training set size: 80 %
 Stratified

Leave one out

Test on train data

Test on test data

Evaluation Results

Model	AUC	CA	F1	Precision	Recall
Tree	0.851	0.804	0.802	0.804	0.804

Visualize

Model

Evaluate

Test and Score

Prediction...

Confusion Matrix

ROC Analysis

Lift Curve

Calibration Plot

Unsupervised

Prototypes

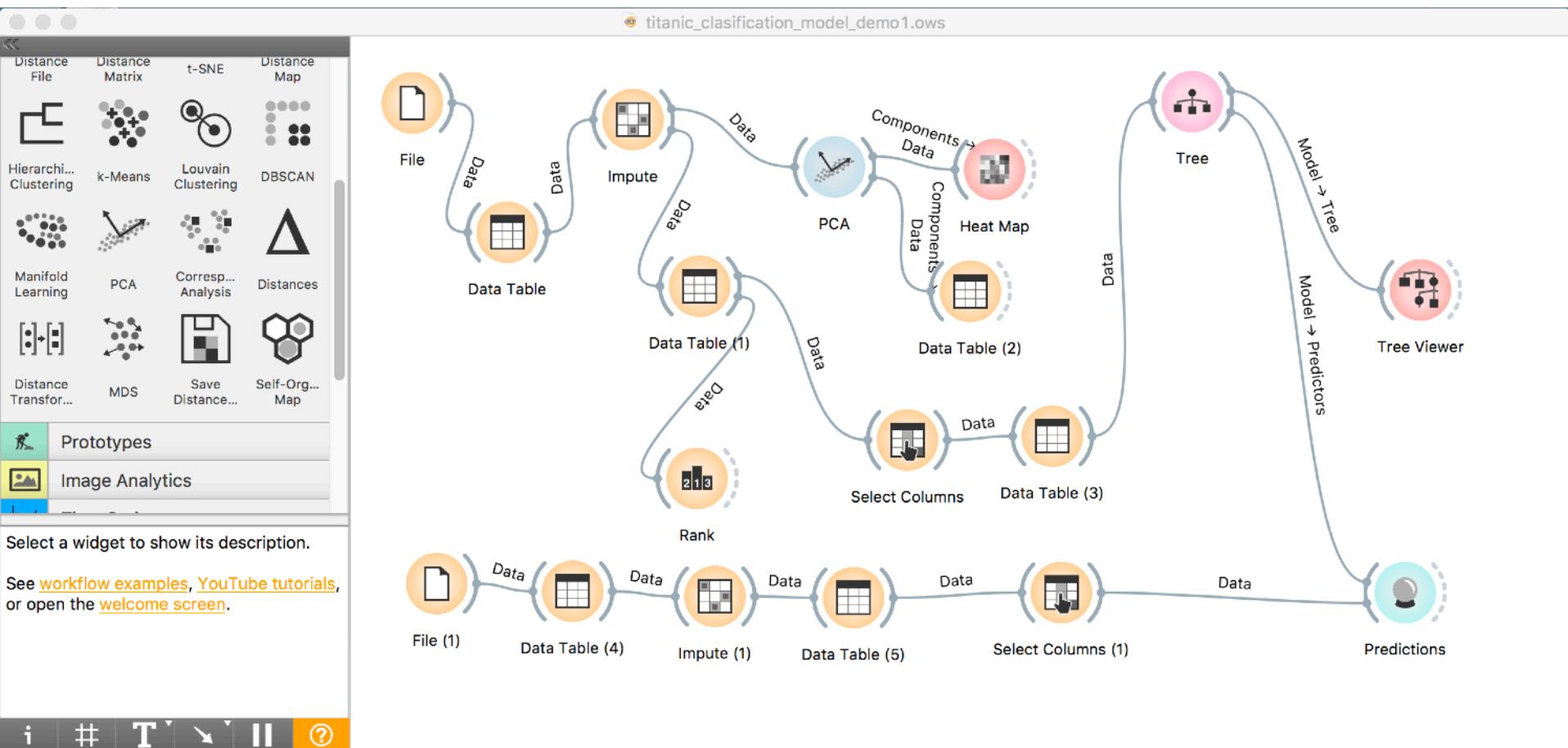
Test and Score

Cross-validation accuracy estimation.

more...

i # T ↻ II ?

Try out other models for choosing the best fit.



Predictions

	Tree	Survived	Pclass	Sex	Age
0	0.91 : 0.09	0	3	male	22.00
1	0.04 : 0.96	1	1	female	38.00
2	0.17 : 0.83	1	3	female	26.00
3	0.04 : 0.96	1	1	female	35.00
4	0.96 : 0.04	0	3	male	35.00
5	0.88 : 0.12	0	1	male	54.00
6	0.56 : 0.44	0	3	male	2.00
7	0.17 : 0.83	1	3	female	27.00
8	0.07 : 0.93	1	2	female	14.00
9	0.00 : 1.00	1	3	female	4.00
10	0.17 : 0.83	1	1	female	58.00
11	0.83 : 0.17	0	3	male	20.00
12	0.96 : 0.04	0	3	male	39.00
13	0.30 : 0.70	0	3	female	14.00
14	0.04 : 0.96	1	2	female	55.00
15	0.56 : 0.44	0	3	male	2.00
16	0.67 : 0.33	0	3	female	31.00
17	0.96 : 0.04	0	2	male	35.00

Model AUC CA F1 Precision Recall

Tree 0.913 0.851 0.848 0.854 0.851

Show probabilities for 0 1

Restore Original Order

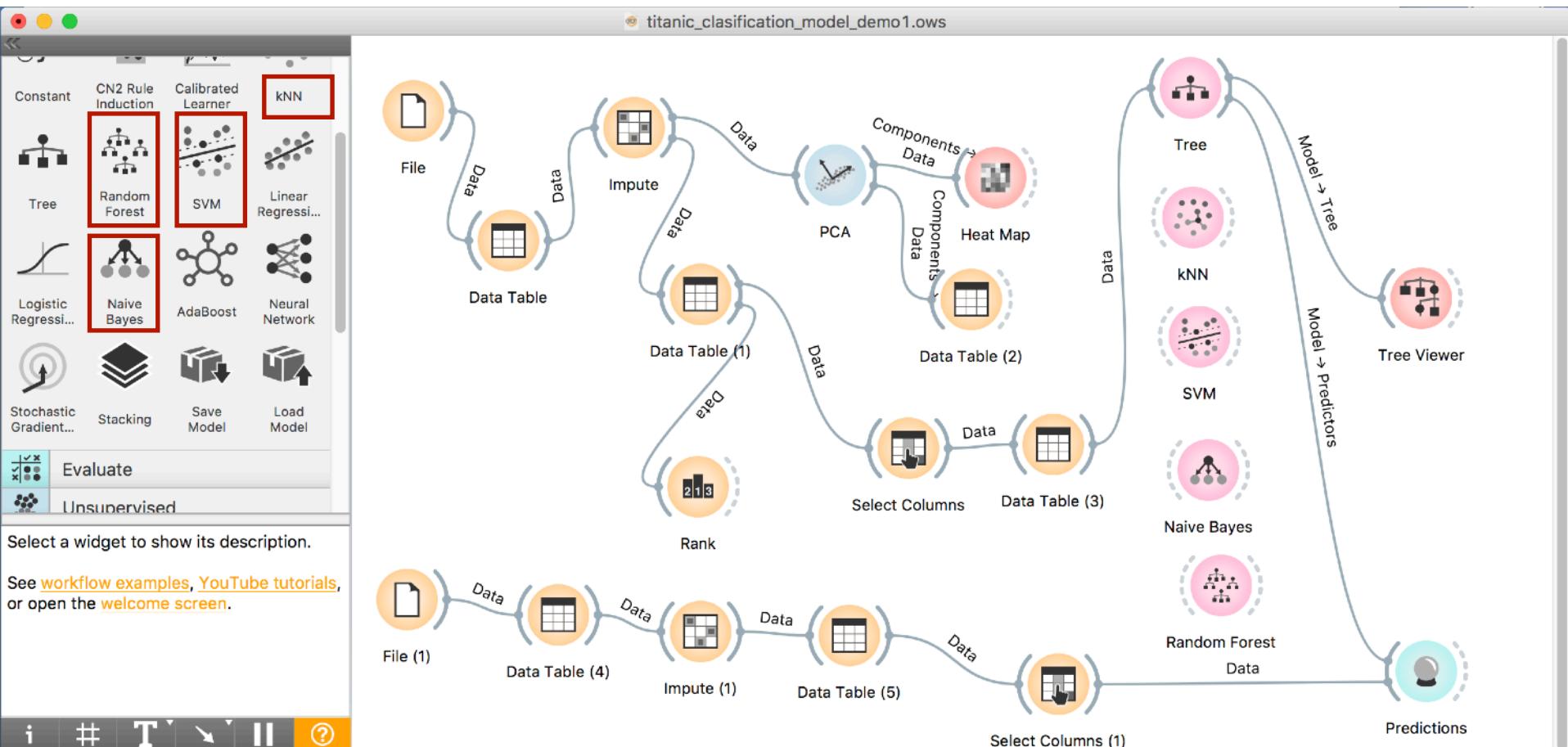
712

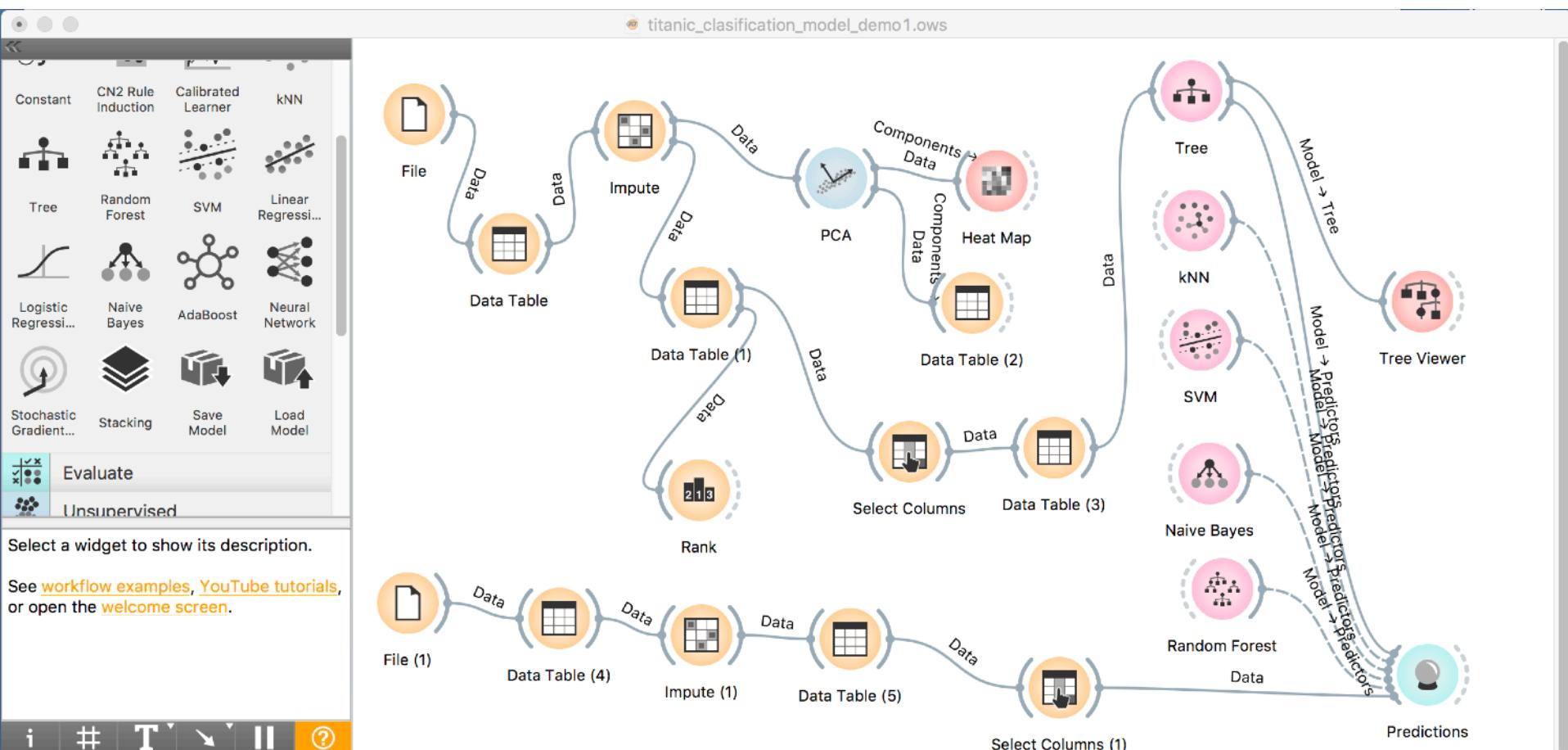
i # T ↴ || ?

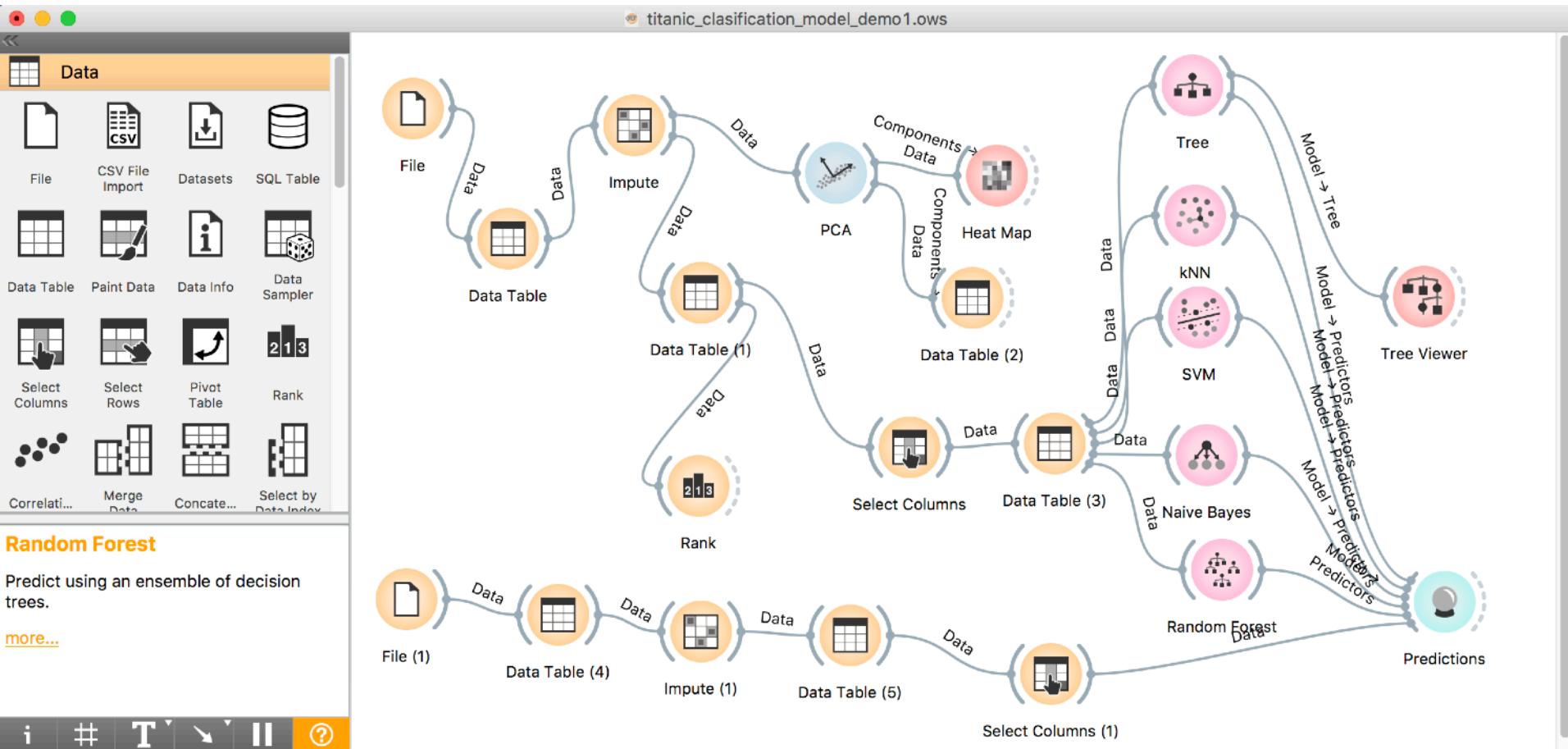
```

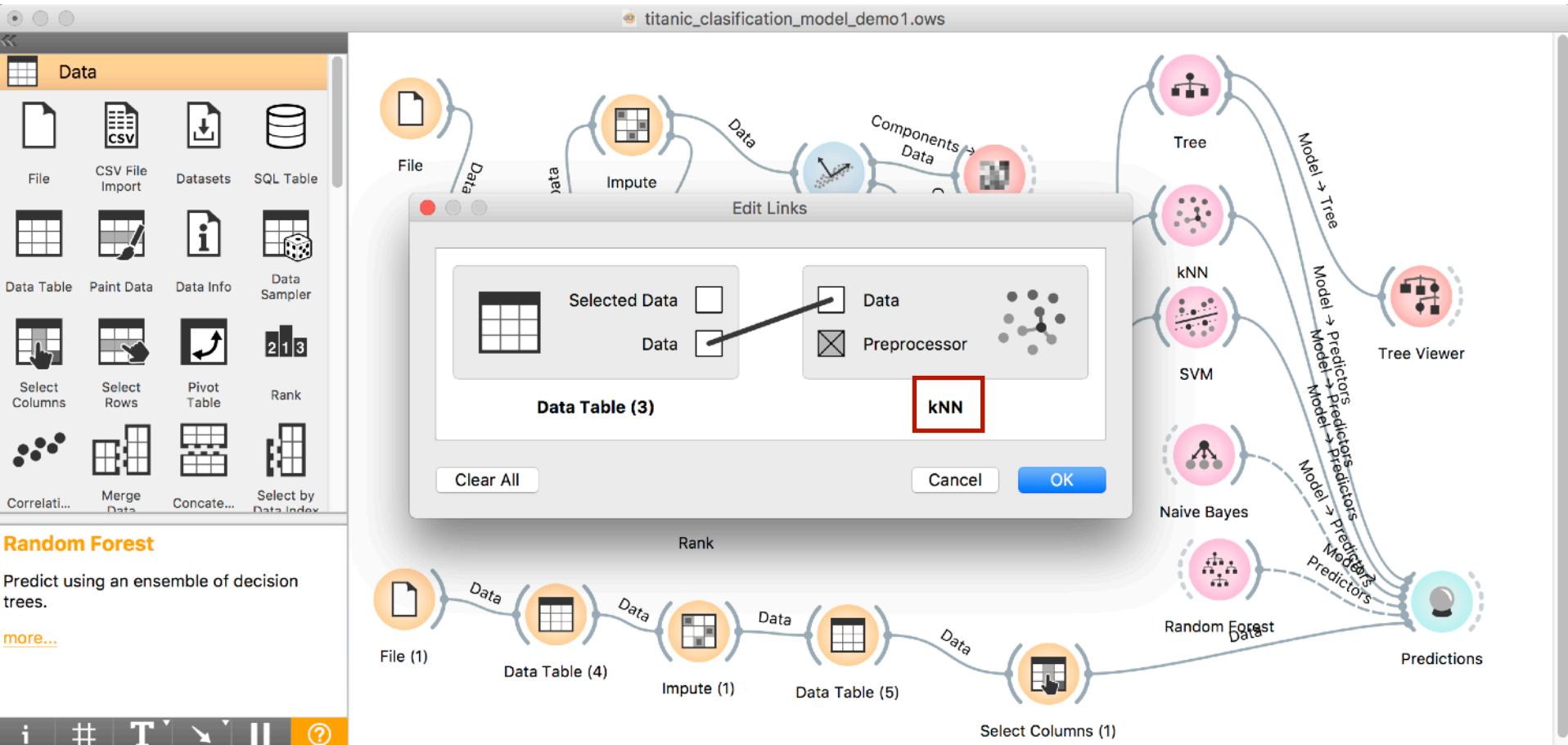
graph TD
    DT2[Data Table (2)] -- Data --> Tree((Tree))
    Tree -- Model --> TV[Tree Viewer]
    TV -- Model --> P[Predictions]
    SC[Select Columns (1)] -- Data --> Tree
    subgraph Components [Components]
        direction TB
        H[Heat Map] --- DT2
        DT2 --- DT3[Data Table (3)]
        DT3 --- SC
    end
    subgraph DataFlow [Data Flow]
        direction TB
        DT2 -- Data --> Tree
        Tree -- Model --> TV
        TV -- Model --> P
        SC -- Data --> Tree
    end

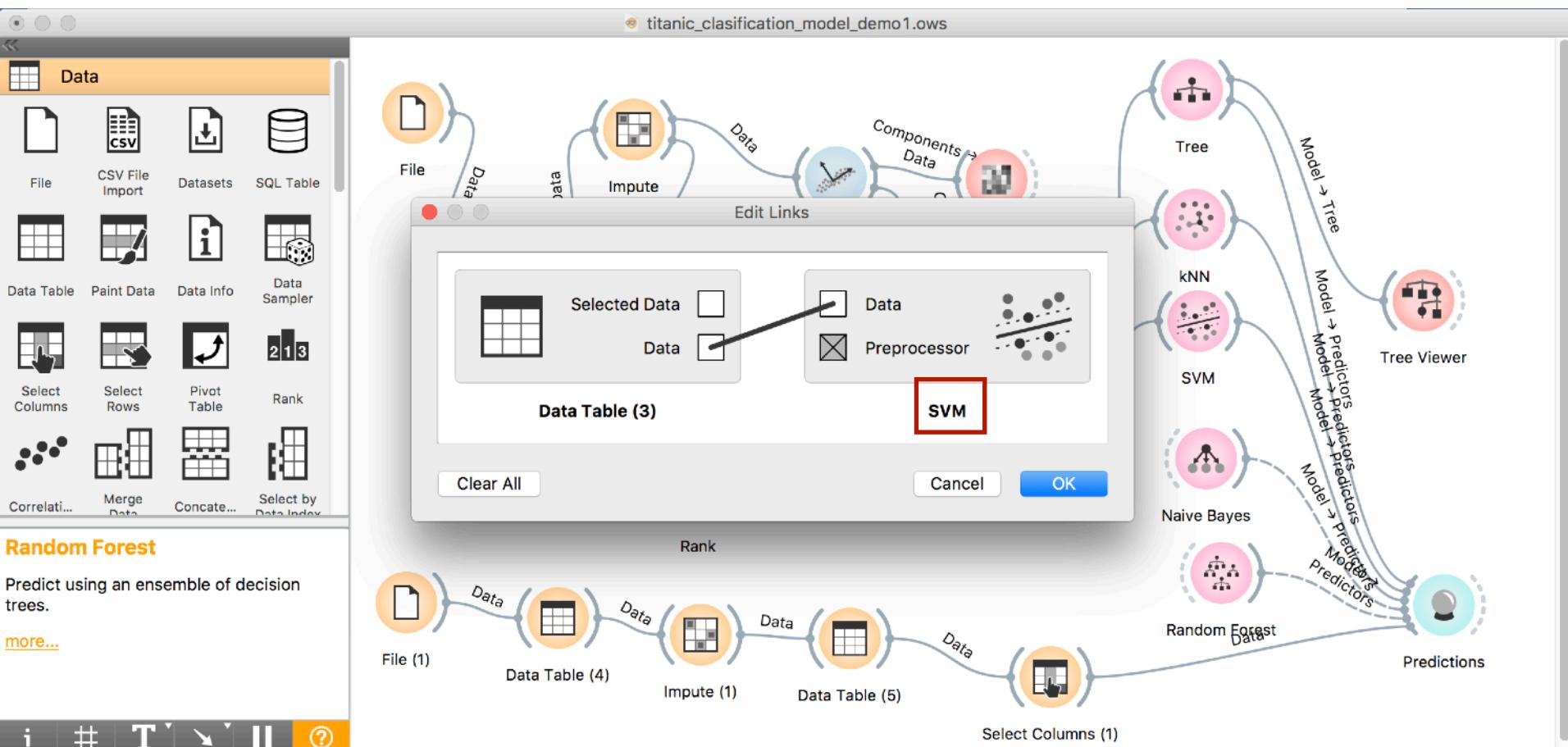
```

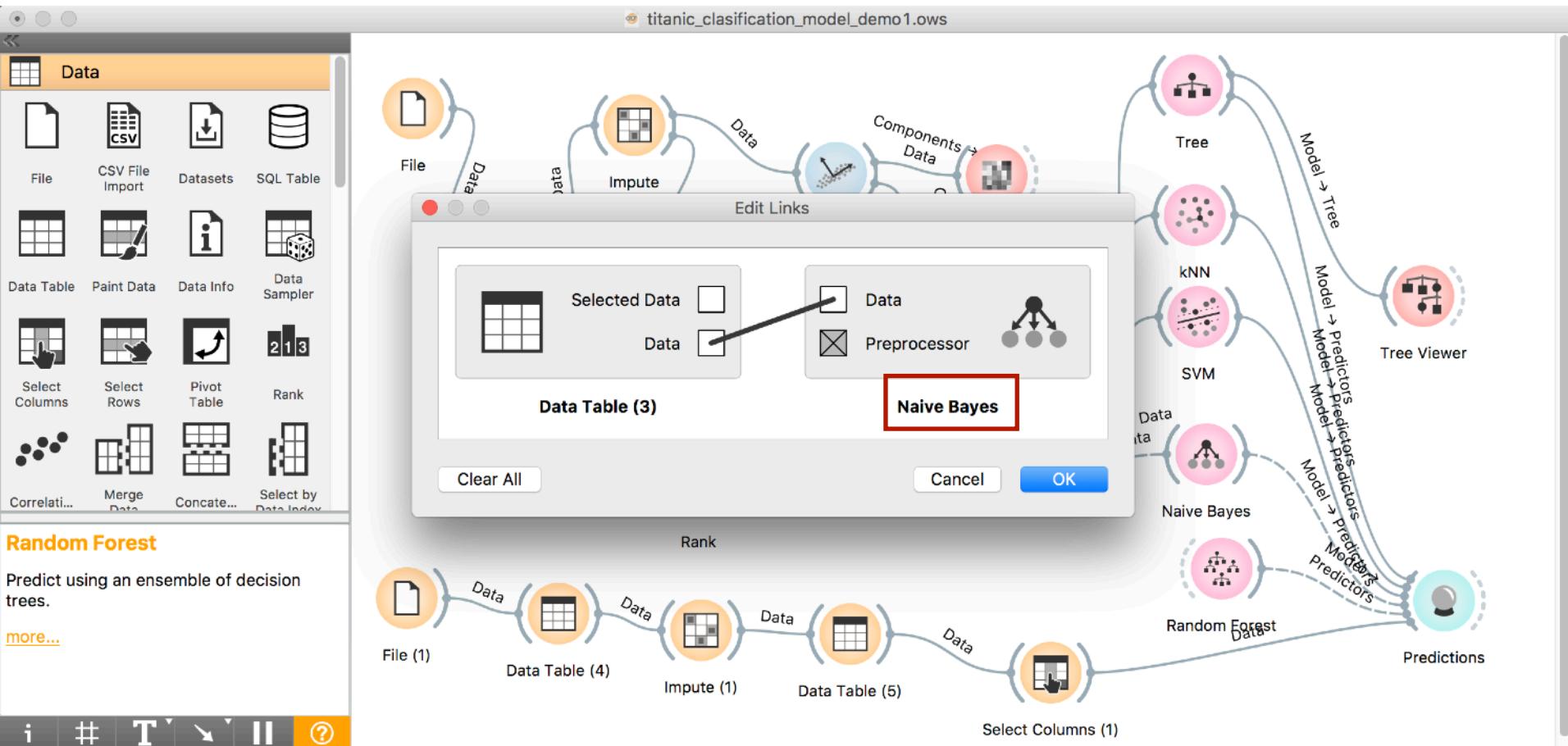


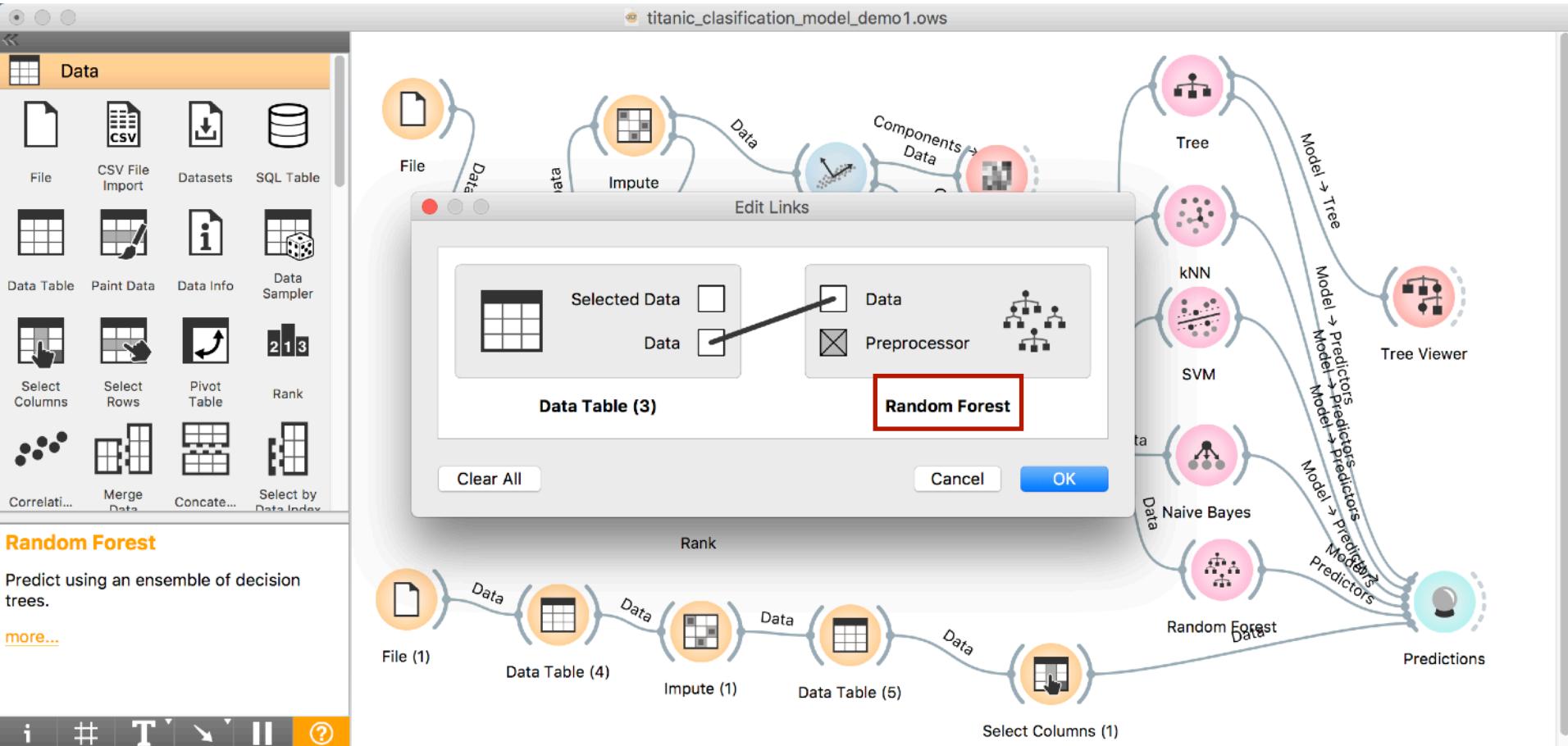






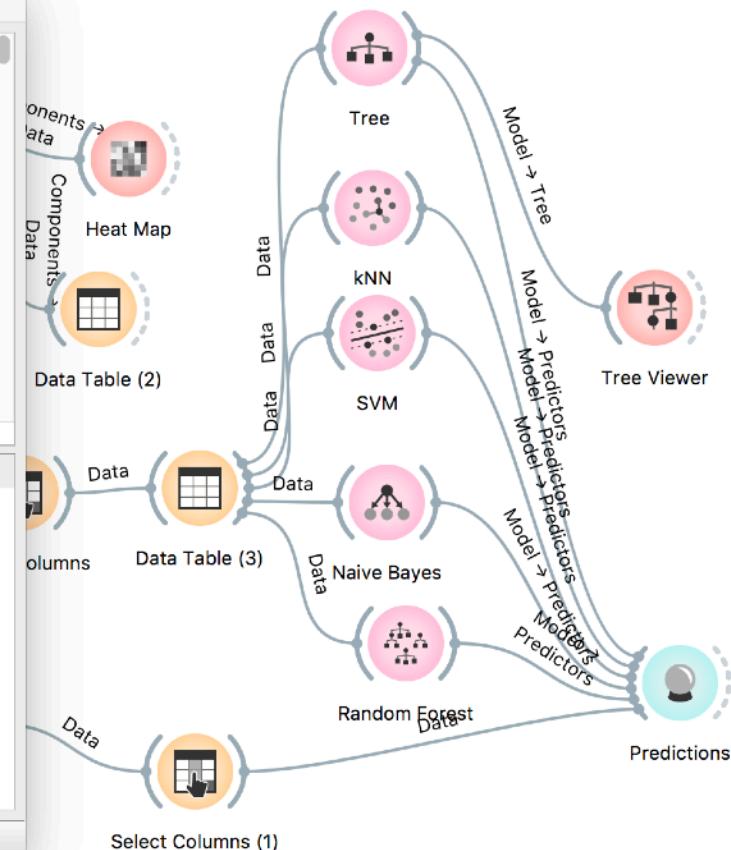






Predictions						
Show probabilities for		Tree	kNN	SVM	Naive Bay	Survived
0		0.91 : 0.09 → 0	1.00 : 0.00 → 0	0.80 : 0.20 → 0	0.91 : 0.09	0
1		0.04 : 0.96 → 1	0.20 : 0.80 → 1	0.26 : 0.74 → 1	0.10 : 0.90	1
		0.17 : 0.83 → 1	0.60 : 0.40 → 0	0.52 : 0.48 → 1	0.46 : 0.54	1
		0.04 : 0.96 → 1	0.00 : 1.00 → 1	0.27 : 0.73 → 1	0.10 : 0.90	1
		0.96 : 0.04 → 0	1.00 : 0.00 → 0	0.80 : 0.20 → 0	0.88 : 0.12	0
		0.88 : 0.12 → 0	1.00 : 0.00 → 0	0.88 : 0.12 → 0	0.62 : 0.38	0
		0.56 : 0.44 → 0	0.80 : 0.20 → 0	0.81 : 0.19 → 0	0.87 : 0.13	0
		0.17 : 0.83 → 1	0.20 : 0.80 → 1	0.52 : 0.48 → 1	0.46 : 0.54	1
		0.07 : 0.93 → 1	0.20 : 0.80 → 1	0.50 : 0.50 → 1	0.16 : 0.84	1
		0.00 : 1.00 → 1	0.20 : 0.80 → 1	0.65 : 0.35 → 0	0.36 : 0.64	1

Model	AUC	CA	F1	Precision	Recall
Naive Bayes	0.842	0.779	0.777	0.778	0.779
SVM	0.827	0.789	0.785	0.789	0.789
kNN	0.915	0.847	0.845	0.847	0.847
Tree	0.913	0.851	0.848	0.854	0.851
Random Forest	0.954	0.885	0.884	0.886	0.885





scikit-learn

Machine Learning in Python

[Getting Started](#)[What's New in 0.22.2](#)[GitHub](#)

- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying which category an object belongs to.

Applications: Spam detection, image recognition.

Algorithms: SVM, nearest neighbors, random forest, and more...

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, nearest neighbors, random forest, and more...

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, and more...



Please [cite us](#) if you use
the software.

7. Dataset loading utilities

- 7.1. General dataset API
- 7.2. Toy datasets
- 7.3. Real world datasets
- 7.4. Generated datasets
- 7.5. Loading other datasets

7.2. Toy datasets

scikit-learn comes with a few small standard datasets that do not require to download any file from some external website.

They can be loaded using the following functions:

<code>load_boston</code> ([return_X_y])	Load and return the boston house-prices dataset (regression).
<code>load_iris</code> ([return_X_y])	Load and return the iris dataset (classification).
<code>load_diabetes</code> ([return_X_y])	Load and return the diabetes dataset (regression).
<code>load_digits</code> ([n_class, return_X_y])	Load and return the digits dataset (classification).
<code>load_linnerud</code> ([return_X_y])	Load and return the linnerud dataset (multivariate regression).
<code>load_wine</code> ([return_X_y])	Load and return the wine dataset (classification).
<code>load_breast_cancer</code> ([return_X_y])	Load and return the breast cancer wisconsin dataset (classification).

These datasets are useful to quickly illustrate the behavior of the various algorithms implemented in scikit-learn. They are however often too small to be representative of real world machine learning tasks.

7.2.1. Boston house prices dataset

Data Set Characteristics:

**Let's train the Titanic classifier in Python Codes
with Scikit-learn.**

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

Building and Visualizing a Decision Tree Classifier for the Titanic Dataset

```
In [68]: 1 # Import Pandas, Scikit-learn, Graphviz modules
2 from sklearn import tree
3 from sklearn.model_selection import train_test_split
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.tree import export_graphviz
6 from sklearn.metrics import classification_report
7 from sklearn.metrics import accuracy_score
8 import pandas as pd
9 import pydotplus
10 from IPython.display import Image
```

```
In [69]: 1 # Import Titanic dataset
2 data = pd.read_csv('http://localhost/datasets/titanic.csv')
3 # Check the first 5 entries
4 data.head(5)
```

Out[69]:

	Pclass	Sex	Age	Survived	Name	Ticket	Cabin	Selected
0	continuous	discrete	continuous	discrete	string	string	string	discrete
1	NaN	NaN	NaN	class	meta	meta	meta	meta
2	3.0	male	22.0	0	Braund, Mr. Owen Harris	A/5 21171	Nan	No
3	1.0	female	38.0	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	PC 17599	C85 STON/O2. 3101282	No NaN
4	3.0	female	26.0	1				

```
In [70]: 1 # Convert Male and Female into binary values of 0 and 1.
```



```
In [70]: 1 # Convert Male and Female into binary values of 0 and 1.
2 data.Sex[data.Sex == 'male'] = 1
3 data.Sex[data.Sex == 'female'] = 0
4 print(data)
```

```
Pclass      Sex     Age  Survived \
0 continuous discrete continuous discrete
1          NaN      NaN      NaN    class
2          3.0      1       22.0      0
3          1.0      0       38.0      1
4          3.0      0       26.0      1
..        ...
888         2.0      1       27.0      0
889         1.0      0       19.0      1
890         3.0      0       NaN       0
891         1.0      1       26.0      1
892         3.0      1       32.0      0

                           Name      Ticket \
0                         string      string
1                         meta      meta
2           Braund, Mr. Owen Harris      A/5 21171
3  Cumings, Mrs. John Bradley (Florence Briggs Th...      PC 17599
4           Heikkinen, Miss. Laina  STON/O2. 3101282
..                     ...
888           Montvila, Rev. Juozas      211536
889           Graham, Miss. Margaret Edith      112053
890  Johnston, Miss. Catherine Helen "Carrie"      W./C. 6607
891           Behr, Mr. Karl Howell      111369
892           Dooley, Mr. Patrick      370376

Cabin Selected
0   string  discrete
```



```
In [71]: 1 # Check the first 5 items.  
2 data.head(5)
```

Out[71]:

	Pclass	Sex	Age	Survived	Name	Ticket	Cabin	Selected
0	continuous	discrete	continuous	discrete	string	string	string	discrete
1	NaN	NaN	NaN	class	meta	meta	meta	meta
2	3.0	1	22.0	0	Braund, Mr. Owen Harris	A/5 21171	Nan	No
3	1.0	0	38.0	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	PC 17599 STON/O2. 3101282	C85 NaN	No
4	3.0	0	26.0	1				

```
In [72]: 1 # Check frame structure  
2 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 893 entries, 0 to 892  
Data columns (total 8 columns):  
Pclass      892 non-null object  
Sex         892 non-null object  
Age         715 non-null object  
Survived    893 non-null object  
Name        893 non-null object  
Ticket      893 non-null object  
Cabin       206 non-null object  
Selected    893 non-null object  
dtypes: object(8)  
memory usage: 55.9+ KB
```

```
In [74]: 1 # Check number of rows and columns  
2 data.shape
```

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3



```
In [74]: 1 # Check number of rows and columns  
2 data.shape
```

```
Out[74]: (893, 8)
```

```
In [55]: 1 # Slice the frame into selected features and target variable  
2 data=data.loc[1:,'Pclass':'Survived']  
3 data.head(5)
```

```
Out[55]:
```

	Pclass	Sex	Age	Survived
1	NaN	NaN	NaN	class
2	3.0	1	22.0	0
3	1.0	0	38.0	1
4	3.0	0	26.0	1
5	1.0	0	35.0	1

```
In [56]: 1 # Check null values  
2 data.isnull().sum()
```

```
Out[56]: Pclass      1  
Sex         1  
Age       178  
Survived     0  
dtype: int64
```

```
In [57]: 1 # Clean up by dropping the null entries  
2 dataset = data.dropna()
```

```
In [58]: 1 # Ensure it is cleaned.
```



```
In [58]: 1 # Ensure it is cleaned.  
2 dataset.isnull().sum()
```

```
Out[58]: Pclass      0  
Sex         0  
Age         0  
Survived    0  
dtype: int64
```

```
In [59]: 1 # Check first 5 entries  
2 dataset.head(5)
```

```
Out[59]:
```

Pclass	Sex	Age	Survived	
2	3.0	1	22.0	0
3	1.0	0	38.0	1
4	3.0	0	26.0	1
5	1.0	0	35.0	1
6	3.0	1	35.0	0

```
In [60]: 1 # Define x (features) and y (label/classifier)  
2 x = dataset.iloc[:, :-1].values  
3 y = dataset['Survived']
```

```
In [61]: 1 # Split the dataset into training and testing sets (20% for testing)  
2 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.20, random_state=27)
```

```
In [75]: 1 # Training/fitting the model with the Decision Tree Classifier  
2 tree_model = tree.DecisionTreeClassifier()  
3 tree_model.fit(x_train,y_train)
```

File Edit View Insert Cell Kernel Widgets Help

Trusted Python 3

In [75]:

```
1 # Training/fitting the model with the Decision Tree Classifier
2 tree_model = tree.DecisionTreeClassifier()
3 tree_model.fit(x_train,y_train)
```

Out[75]:

```
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False,
random_state=None, splitter='best')
```

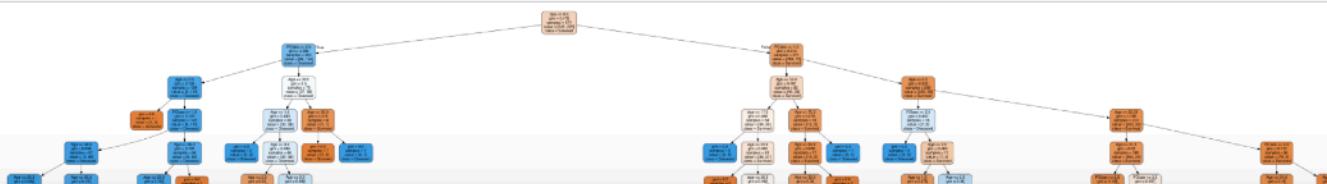
In [76]:

```
1 # Define the feature and target variable names for displaying in the tree
2 titanic_feature_names = ['PClass',
3 'Sex',
4 'Age']
5 titanic_target_names = ['Survived', 'Diseased']
```

In [77]:

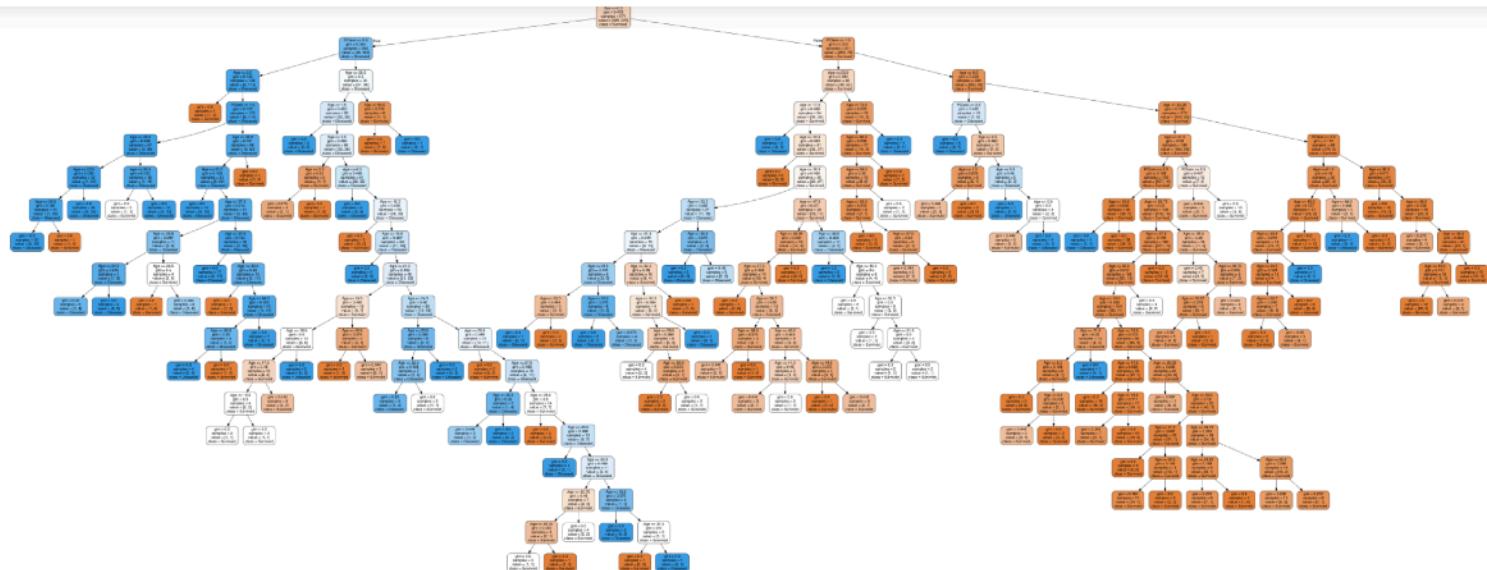
```
1 # Generate the tree
2 dot_data = tree.export_graphviz(tree_model, out_file=None,
3                                 feature_names=titanic_feature_names,
4                                 class_names=titanic_target_names,
5                                 rounded=True, filled=True)
6 # Draw graph
7 graph = pydotplus.graph_from_dot_data(dot_data)
8 # Show graph
9 Image(graph.create_png())
```

Out[77]:





Out[77]:



```
In [66]: 1 # Evaluating the model by using the trained model against the testing dataset  
2 clf_prediction = tree_model.predict(x_test)
```

```
In [67]: 1 # Display the accuracy score  
2 print(accuracy_score(clf_prediction, y_test))
```

0.7972027972027972

Adding the other models

```
In [79]: 1 # Import Pandas, Scikit-learn, Graphviz modules
2 from sklearn import tree
3 from sklearn.model_selection import train_test_split
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.tree import export_graphviz
6 from sklearn.metrics import classification_report
7 from sklearn.metrics import accuracy_score
8 import pandas as pd
9 import pydotplus
10 from IPython.display import Image
11 from sklearn.neighbors import KNeighborsClassifier # K-Nearest Neighbour Algorithm
12 from sklearn.naive_bayes import GaussianNB          # Gaussian Naive Bayes Algorithm
13 from sklearn.svm import SVC                         # Support Vector Machine Algorithm
14 from sklearn.ensemble import RandomForestClassifier # Random Algorithm
```

```
In [80]: 1 # Import Titanic dataset
2 data = pd.read_csv('http://localhost/datasets/titanic.csv')
3 # Check the first 5 entries
4 data.head(5)
```

Out[80]:

	Pclass	Sex	Age	Survived	Name	Ticket	Cabin	Selected
0	continuous	discrete	continuous	discrete	string	string	string	discrete
1	NaN	NaN	NaN	class	meta	meta	meta	meta
2	3.0	male	22.0	0	Braund, Mr. Owen Harris	A/5 21171	Nan	No
3	1.0	female	38.0	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	PC 17599	C85	No

```
In [80]: 1 # Import Titanic dataset
2 data = pd.read_csv('http://localhost/datasets/titanic.csv')
3 # Check the first 5 entries
4 data.head(5)
```

Out[80]:

	Pclass	Sex	Age	Survived	Name	Ticket	Cabin	Selected
0	continuous	discrete	continuous	discrete		string	string	string
1	NaN	NaN	NaN	class		meta	meta	meta
2	3.0	male	22.0	0	Braund, Mr. Owen Harris	A/5 21171	Nan	No
3	1.0	female	38.0	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	PC 17599 STON/O2. 3101282	C85 NaN	No
4	3.0	female	26.0	1				

```
In [81]: 1 # Convert Male and Female into binary values of 0 and 1.
2 data.Sex[data.Sex == 'male'] = 1
3 data.Sex[data.Sex == 'female'] = 0
4 print(data)
```

```
Pclass      Sex      Age  Survived \
0  continuous  discrete  continuous  discrete
1        NaN      NaN       NaN    class
2        3.0      1.0     22.0      0
3        1.0      0.0     38.0      1
4        3.0      0.0     26.0      1
..        ...
888       2.0      1.0     27.0      0
889       1.0      0.0     19.0      1
890       3.0      0.0      NaN      0
891       1.0      1.0     26.0      1
892       3.0      1.0     32.0      0
```



```
In [82]: 1 # Check the first 5 items.  
2 data.head(5)
```

Out[82]:

	Pclass	Sex	Age	Survived	Name	Ticket	Cabin	Selected
0	continuous	discrete	continuous	discrete	string	string	string	discrete
1	NaN	NaN	NaN	class	meta	meta	meta	meta
2	3.0	1	22.0	0	Braund, Mr. Owen Harris	A/5 21171	Nan	No
3	1.0	0	38.0	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	PC 17599 STON/O2. 3101282	C85 NaN	No
4	3.0	0	26.0	1				

```
In [83]: 1 # Slice the frame into selected features and target variable  
2 data=data.loc[1:,'Pclass':'Survived']  
3 data.head(5)
```

Out[83]:

	Pclass	Sex	Age	Survived
1	NaN	NaN	NaN	class
2	3.0	1	22.0	0
3	1.0	0	38.0	1
4	3.0	0	26.0	1
5	1.0	0	35.0	1

```
In [84]: 1 # Check null values  
2 data.isnull().sum()
```

Out[84]: Pclass 1
Sex 1



```
In [84]: 1 # Check null values  
2 data.isnull().sum()
```

```
Out[84]: Pclass      1  
Sex         1  
Age       178  
Survived     0  
dtype: int64
```

```
In [85]: 1 # Ensure it is cleaned.  
2 dataset.isnull().sum()
```

```
Out[85]: Pclass      0  
Sex         0  
Age         0  
Survived    0  
dtype: int64
```

```
In [86]: 1 # Check first 5 entries  
2 dataset.head(5)
```

```
Out[86]:  
Pclass  Sex  Age  Survived  
2      3.0   1   22.0      0  
3      1.0   0   38.0      1  
4      3.0   0   26.0      1  
5      1.0   0   35.0      1  
6      3.0   1   35.0      0
```

```
In [87]: 1 # Define x (features) and y (label/classifier)
```



```
In [87]: 1 # Define x (features) and y (label/classifier)
2 x = dataset.iloc[:, :-1].values
3 y = dataset['Survived']
```

```
In [88]: 1 # Split the dataset into training and testing sets (20% for testing)
2 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.20, random_state=27)
```

```
In [90]: 1 # Training/Fitting with Different Classifiers Setup
2 # Decision Tree Classifier
3 tree_model = tree.DecisionTreeClassifier()
4 # Support Vector Machine Classifier
5 SVC_model = SVC(gamma='scale')
6 # KNN model requires you to specify n_neighbors,
7 # the number of points the classifier will look at to determine what class a new point belongs to
8 KNN_model = KNeighborsClassifier(n_neighbors=5)
9 # KNN_model = KNeighborsClassifier(algorithm='scale', leaf_size=30, metric='minkowski',
10 #                                     metric_params=None, n_jobs=None, n_neighbors=5, p=2)
11 # Gaussian Naive Baye model
12 GNB_model = GaussianNB()
13 # Random Forest
14 # Create the model with 100 trees
15 RF_model = RandomForestClassifier(n_estimators=100,
16                                   bootstrap=True,
17                                   max_features = 'sqrt')
```

```
In [91]: 1 # Start training/fitting
2 tree_model.fit(x_train, y_train)
3 SVC_model.fit(x_train, y_train)
4 KNN_model.fit(x_train, y_train)
5 GNB_model.fit(x_train, y_train)
6 RF_model.fit(x_train, y_train)
```

```
In [93]: 1 # Training completed. Use models to predict with the test data and store the outcome  
2 SVC_prediction = SVC_model.predict(x_test)  
3 KNN_prediction = KNN_model.predict(x_test)  
4 GNB_prediction = GNB_model.predict(x_test)  
5 TREE_prediction = tree_model.predict(x_test)  
6 RF_prediction = RF_model.predict(x_test)
```

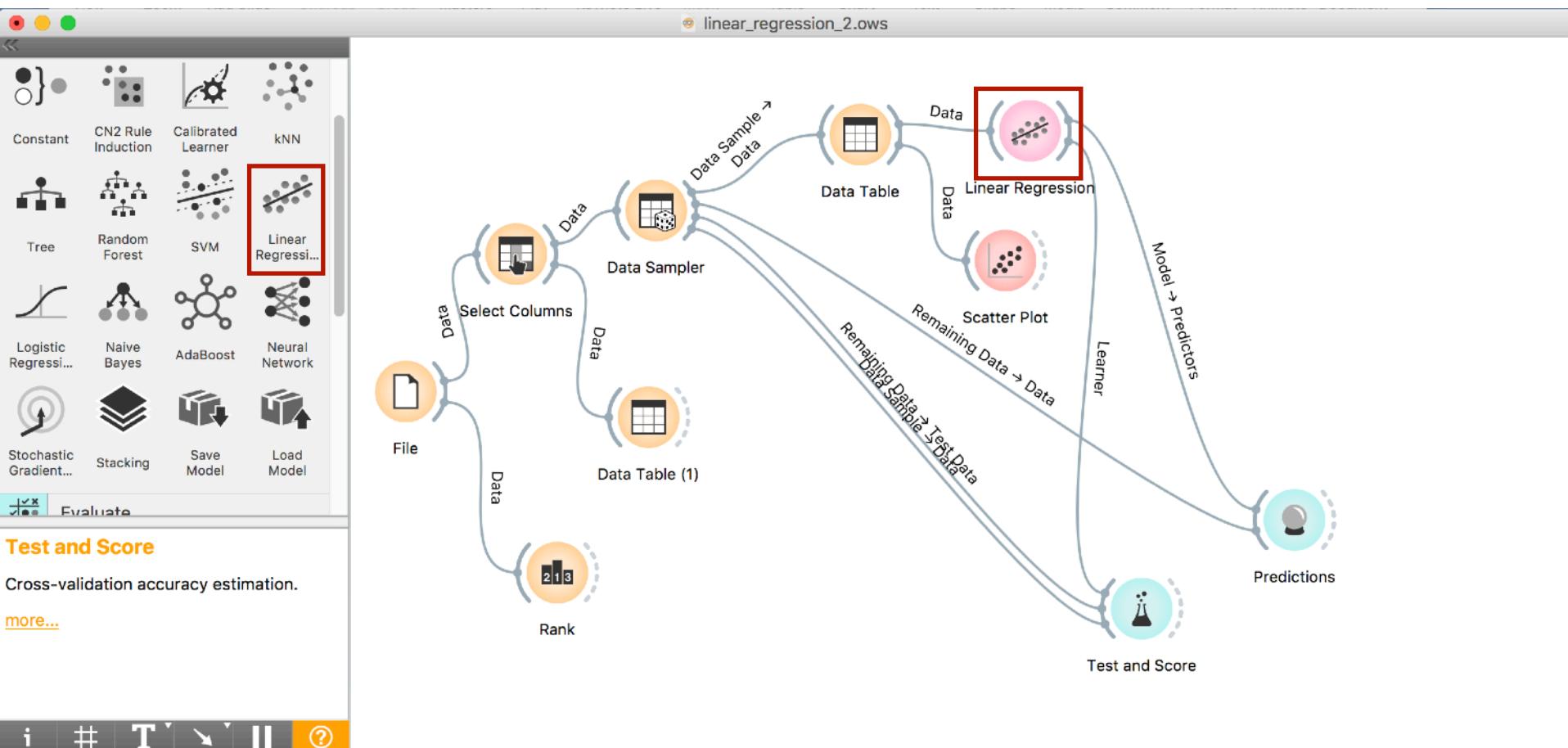
```
In [95]: 1 # Accuracy score is the simplest way to evaluate  
2 print("SVC")  
3 print(accuracy_score(SVC_prediction, y_test))  
4 print("KNN")  
5 print(accuracy_score(KNN_prediction, y_test))  
6 print("GNB")  
7 print(accuracy_score(GNB_prediction, y_test))  
8 print("TREE")  
9 print(accuracy_score(TREE_prediction, y_test))  
10 print("RF")  
11 print(accuracy_score(RF_prediction, y_test))
```

```
SVC  
0.6083916083916084  
KNN  
0.8041958041958042  
GNB  
0.7272727272727273  
TREE  
0.7972027972027972  
RF  
0.8181818181818182
```

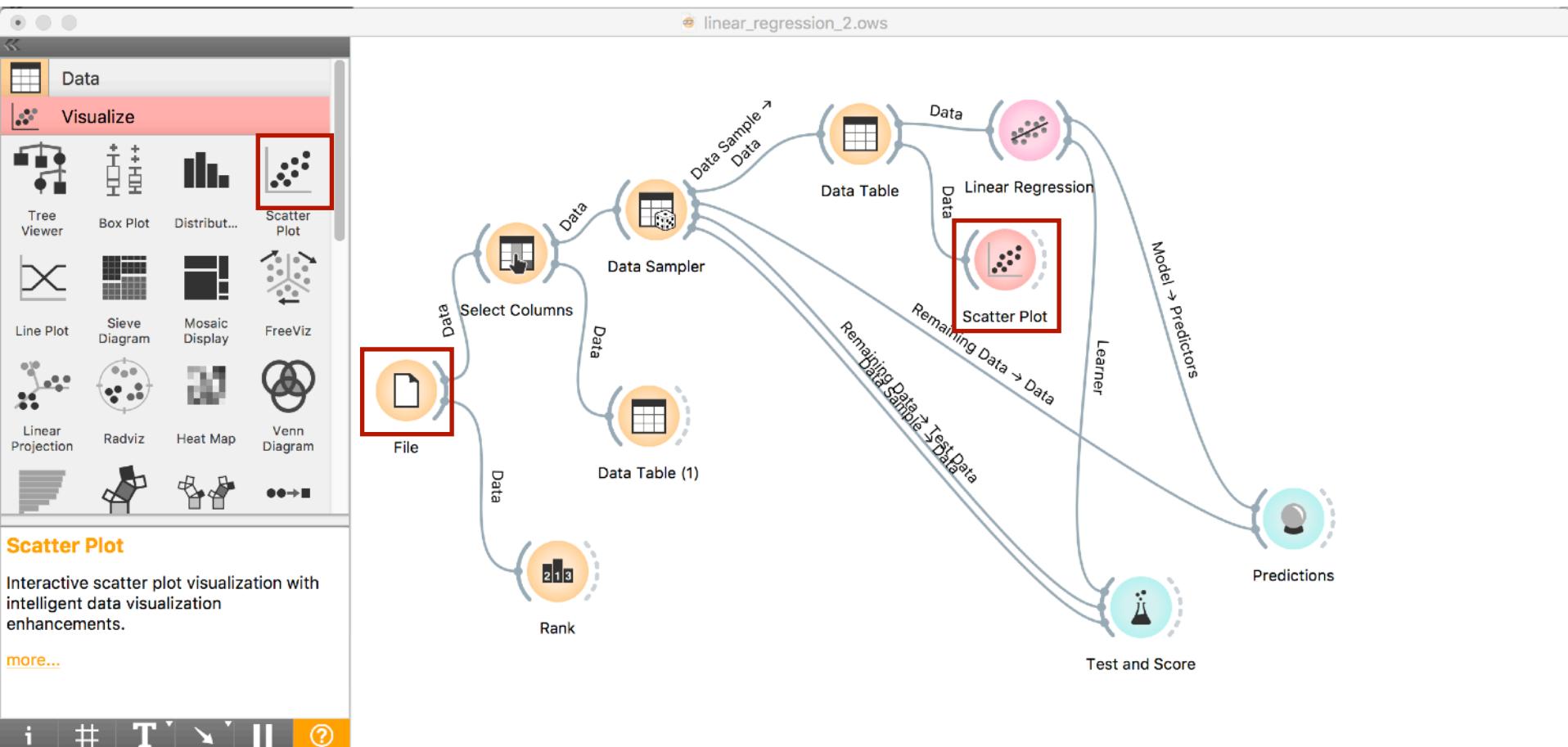
```
In [ ]:
```

```
1
```

The same logic applies to Regression models.



Same pipeline concept.



File

File: **housing.tab** Reload

URL:

Info

Housing dataset
Data collected by the U.S Census Service concerning housing in Boston.

506 instance(s)
13 feature(s) (no missing values)
Regression; numerical class (no missing values)
0 meta attribute(s)

Columns (Double click to edit)

	Name	Type	Role	Values
1	CRIM	N numeric	feature	
2	ZN	N numeric	feature	
3	INDUS	N numeric	feature	
4	CHAS	N numeric	feature	
5	NOX	N numeric	feature	
6	RM	N numeric	feature	

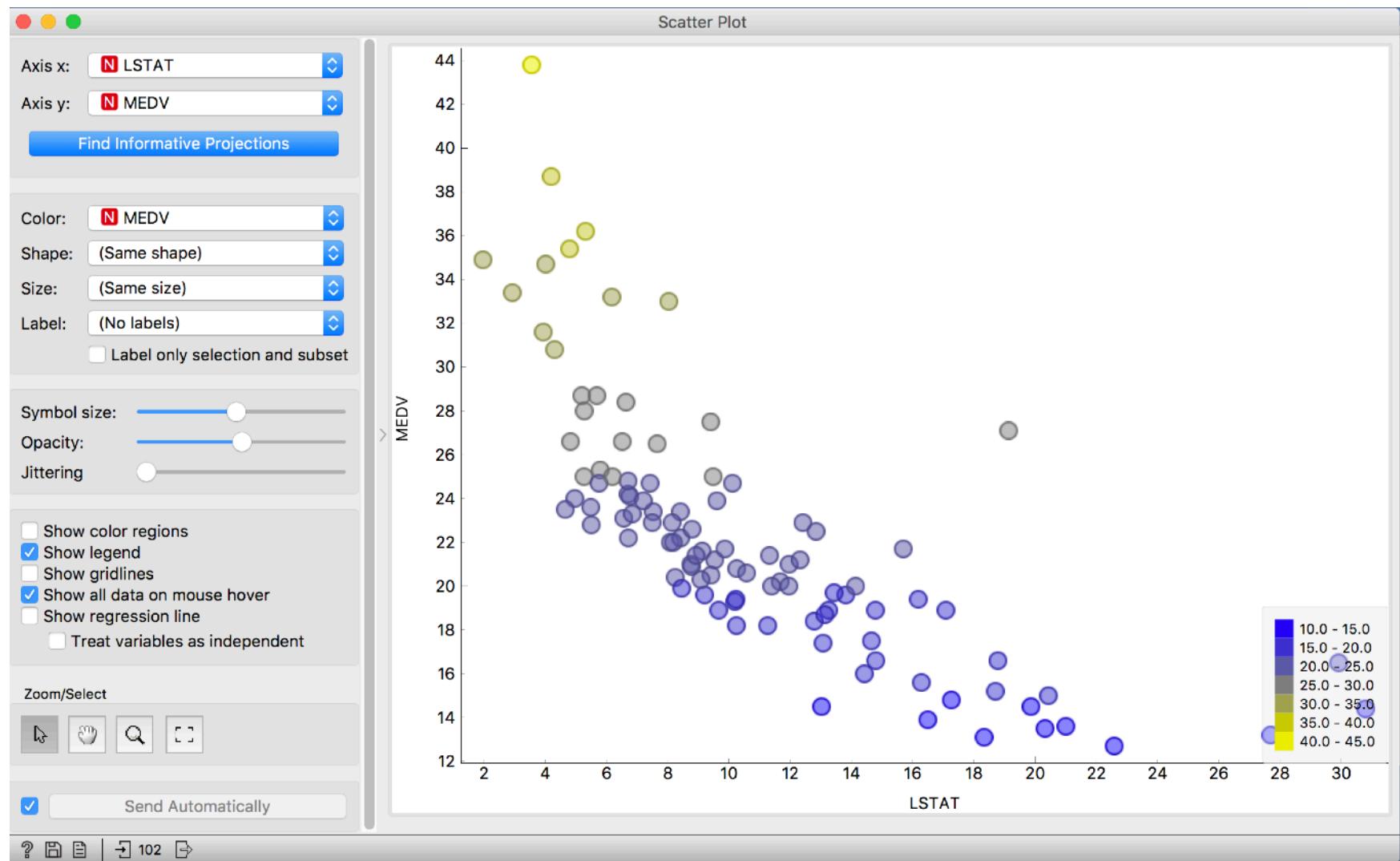
Browse documentation datasets Reset Apply

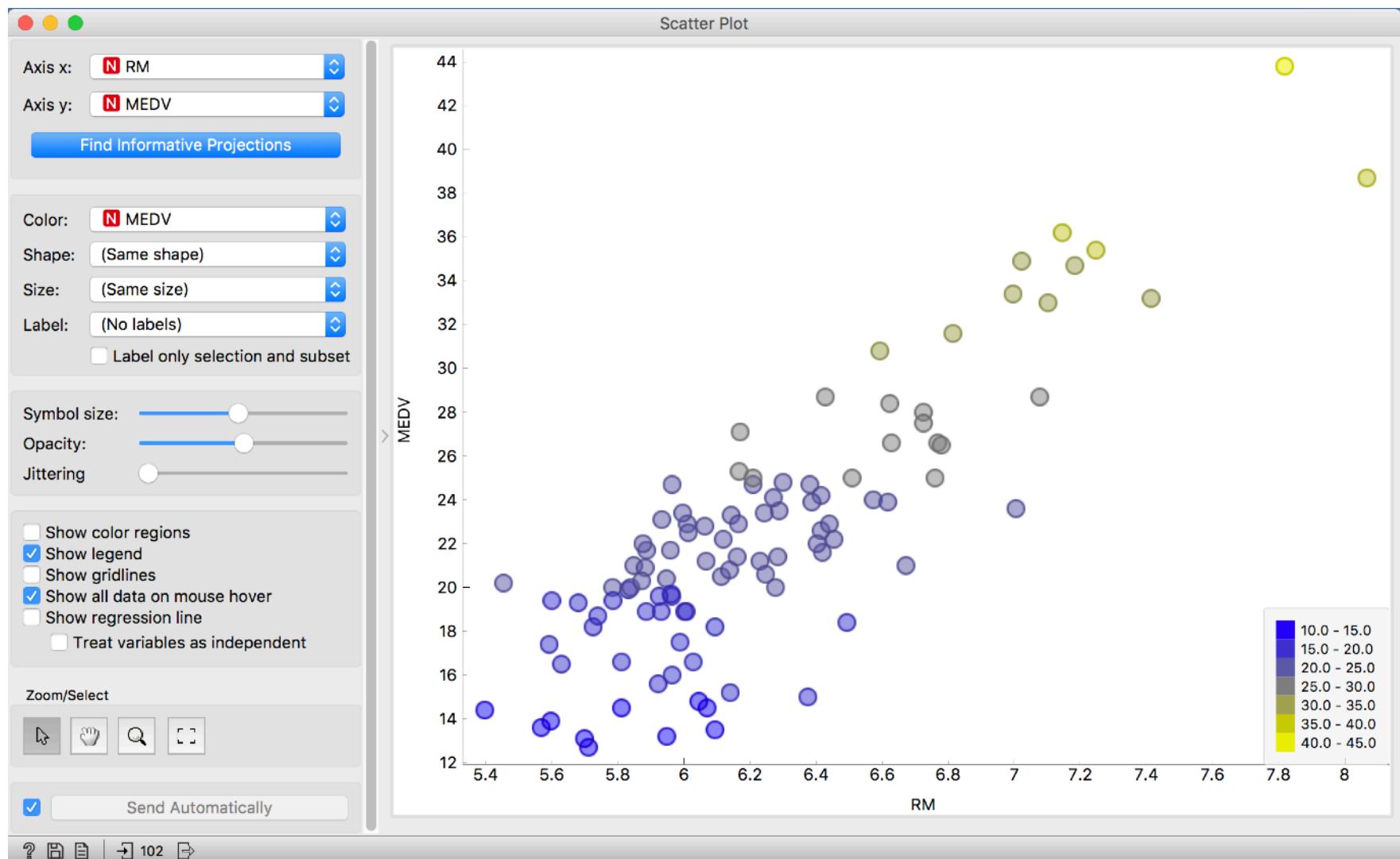
File Read send more

Tree View Line Project

```

graph LR
    Data((Data)) --> LR((Linear Regression))
    LR --> SP((Scatter Plot))
    LR --> Model((Model → Predictor's))
    Model --> Predictions((Predictions))
    Predictions --> TSS((Test and Score))
    TSS --> TestData((Test Data))
    TestData --> RemainingData((Remaining Data → Data))
    RemainingData --> Learner((Learner))
    Learner --> Model
    
```





linear_regression_2.ows

Show probabilities for

Linear Regression

	MEDV	RM	PTRATIO	LSTAT	AGE	DIS	
1	21.4	18.6	6.405	20.9	10.63	85.4	2.7147
2	18.6	19.3	6.137	20.9	13.44	87.4	2.7147
3	18.8	20.1	6.167	20.9	12.33	90.0	2.4210
4	15.2	19.5	5.851	20.9	16.47	96.7	2.1069
5	15.0	19.5	5.836	20.9	18.66	91.9	2.2110
6	18.5	20.4	6.127	20.9	14.09	85.2	2.1224
7	21.1	19.8	6.474	20.9	12.27	97.1	2.4329
8	18.9	19.4	6.229	20.9	15.55	91.2	2.5451
9	21.0	21.7	6.195	20.9	13.00	54.4	2.7778
10	25.3	22.8	6.715	17.8	10.16	81.6	2.6775
11	17.0	18.8	5.913	17.8	16.21	92.9	2.3534
12	18.3	18.7	6.092	17.8	17.09	95.4	2.5480
13	21.2	18.5	6.254	17.8	10.45	84.2	2.2565
14	17.5	18.3	5.928	17.8	15.76	88.2	2.4631
15	21.0	21.2	6.176	17.8	12.04	72.5	2.7301
16	19.3	19.2	6.021	17.8	10.30	82.6	2.7474
17	17.9	20.4	5.872	17.8	15.37	73.1	2.4775

Model MSE RMSE MAE R2

Linear Regression 43.035 6.560 4.365 0.556

Restore Original Order

i # T ? 404

linear_regression_2.ows

Test and Score

Sampling

Cross validation
Number of folds: 5
 Stratified

Cross validation by feature

Random sampling
Repeat train/test: 10
Training set size: 70 %
 Stratified

Leave one out

Test on train data

Test on test data

Evaluation Results

Model	MSE	RMSE	MAE	R2
Linear Regression	43.035	6.560	4.365	0.556

Model → Predictors

Predictions

Score

?

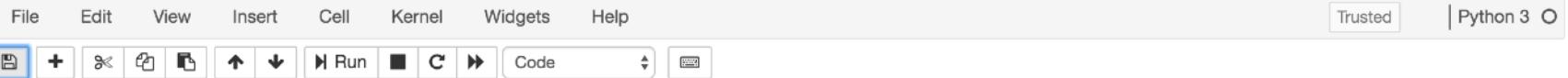
Test and Score

Cross-validation a

more...

i # T ↵ || ?

**The same coding logic applies to building
Regression models in Sckit-learn**

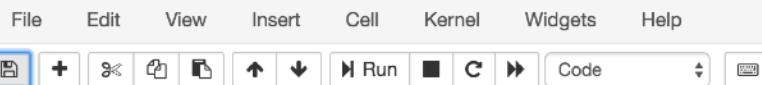


Building and Visualizing a Linear Regression Model

```
In [96]: 1 # Import Pandas, Matplotlib, Numpy, and Scikit-learn modules
2 from matplotlib import pyplot as plt
3 import pandas as pd
4 import numpy as np
5 from sklearn import linear_model
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import classification_report
8 from sklearn.metrics import confusion_matrix
9 from sklearn.metrics import accuracy_score
```

```
In [97]: 1 # Import dataset
2 column_names = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT', 'MEDV']
3 data = pd.read_csv('http://localhost/datasets/housing.csv', header=None, delimiter=r"\s+", names=column_names)
4 # Check the first 5 entries
5 print(data.head(5))
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222.0				
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222.0				



Trusted

Python 3

```
In [98]: 1 # Display number of rows and columns
2 print(np.shape(data))
```

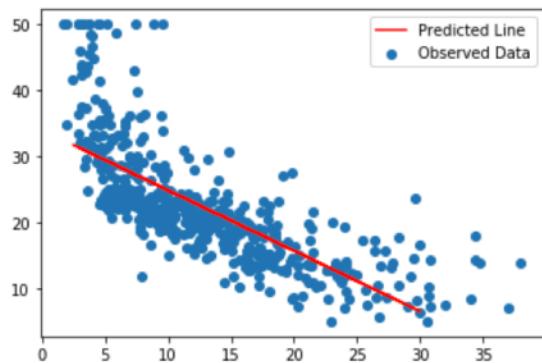
(506, 14)

```
In [99]: 1 # Descriptive statistics of features
2 print(data.describe())
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	\
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	
mean	3.613524	11.363636	11.136779	0.069170	0.554695	6.284634	
std	8.601545	23.322453	6.860353	0.253994	0.115878	0.702617	
min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000	
25%	0.082045	0.000000	5.190000	0.000000	0.449000	5.885500	
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.208500	
75%	3.677082	12.500000	18.100000	0.000000	0.624000	6.623500	
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000	
	AGE	DIS	RAD	TAX	PTRATIO	B	\
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	
mean	68.574901	3.795043	9.549407	408.237154	18.455534	356.674032	
std	28.148861	2.105710	8.707259	168.537116	2.164946	91.294864	
min	2.900000	1.129600	1.000000	187.000000	12.600000	0.320000	
25%	45.025000	2.100175	4.000000	279.000000	17.400000	375.377500	
50%	77.500000	3.207450	5.000000	330.000000	19.050000	391.440000	
75%	94.075000	5.188425	24.000000	666.000000	20.200000	396.225000	
max	100.000000	12.126500	24.000000	711.000000	22.000000	396.900000	
	LSTAT	MEDV					\
count	506.000000	506.000000					
mean	12.653063	22.532806					
std	7.141062	9.197104					
min	1.730000	5.000000					
25%	6.250000	17.025000					

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
In [104]: 1 # Display scatter plot to fit the line
2 a=reg_model.predict(x_test[:None])
3 plt.scatter(x,y)
4 plt.plot(x_test,a,'r')
5 plt.legend(['Predicted Line','Observed Data'])
6 plt.show()
```

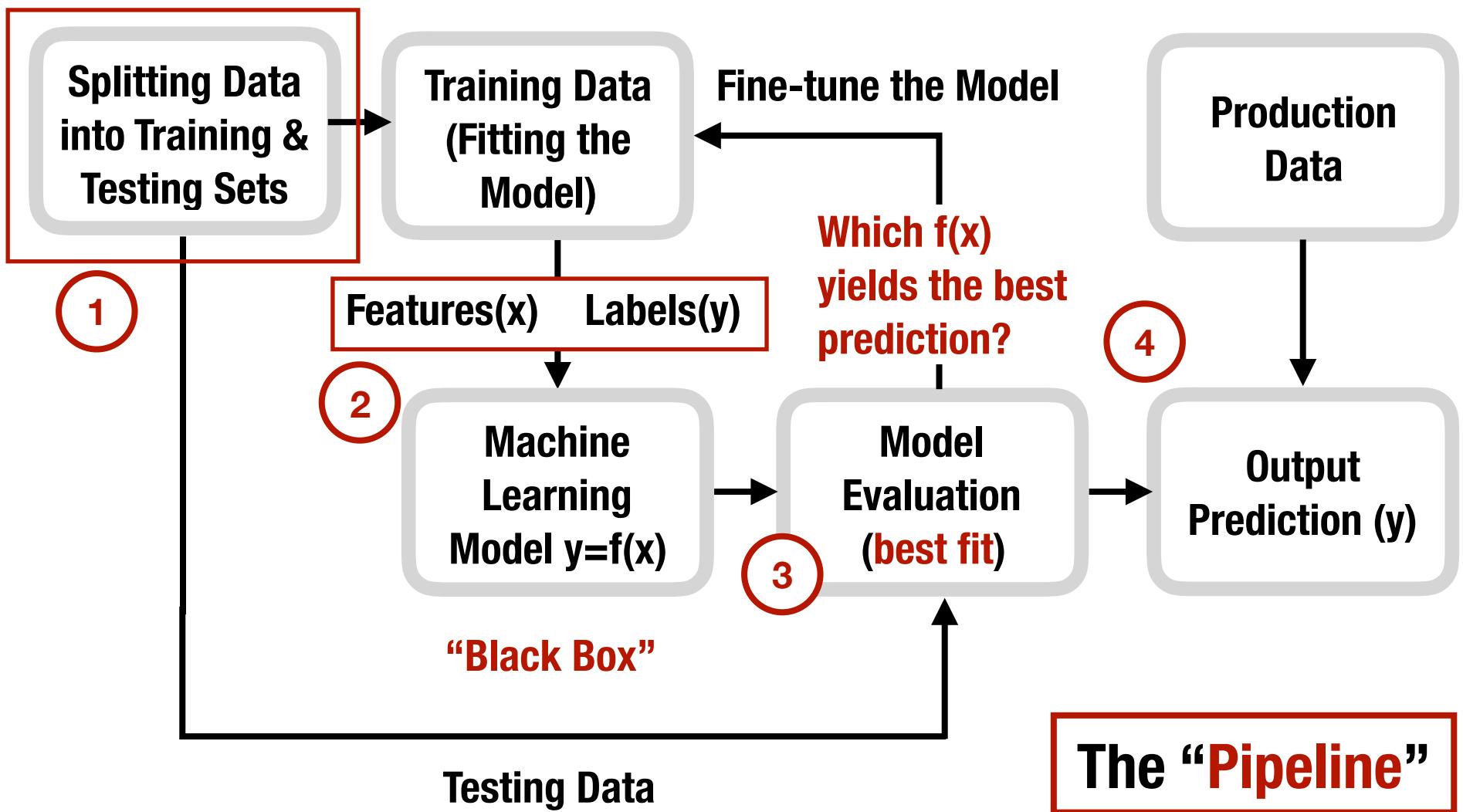


```
In [105]: 1 # Mean Square Error (MSE)
2 np.mean((a-y_test)**2)
```

Out[105]: 40.13019759337232

In []:

Mastering the **pipeline workflow is
the key to get started.**



An end-to-end open source machine learning platform

[TensorFlow](#)[For JavaScript](#)[For Mobile & IoT](#)[For Production](#)

The core open source library to help you develop and train ML models. Get started quickly by running Colab notebooks directly in your browser.

[Get started with TensorFlow](#)

<https://www.tensorflow.org/>

TensorFlow Core

[Overview](#) [Tutorials](#) [Guide](#) [TF 1](#)

TensorFlow is an end-to-end open source platform for machine learning

TensorFlow makes it easy for beginners and experts to create machine learning models. See the sections below to get started.

[See tutorials](#)[See the guide](#)

Tutorials show you how to use TensorFlow with complete, end-to-end examples.

Guides explain the concepts and components of TensorFlow.



<https://www.tensorflow.org/overview/>

beginner.ipynb

File Edit View Insert Runtime Tools Help

Share   

Table of contents



+ Code + Text 

RAM 
Disk 

 Editing 

Copyright 2019 The TensorFlow Authors.

TensorFlow 2 quickstart for beginners

Section

TensorFlow 2 quickstart for beginners



[View on TensorFlow.org](#)



[Run in Google Colab](#)



[View source on GitHub](#)



[Download notebook](#)

This short introduction uses [Keras](#) to:

1. Build a neural network that classifies images.
2. Train this neural network.
3. And, finally, evaluate the accuracy of the model.

This is a [Google Colaboratory](#) notebook file. Python programs are run directly in the browser—a great way to learn and use TensorFlow. To follow this tutorial, run the notebook in Google Colab by clicking the button at the top of this page.

1. In Colab, connect to a Python runtime: At the top-right of the menu bar, select *CONNECT*.
2. Run all the notebook code cells: Select *Runtime > Run all*.

[Install](#)[Learn](#)[API](#)[Resources](#)[More](#) Search

English

[GitHub](#)[Sign in](#)[TensorFlow tutorials](#)[Quickstart for beginners](#)[Quickstart for experts](#)**BEGINNER**[ML basics with Keras](#)[Load and preprocess data](#)[Estimator](#)**ADVANCED**[Customization](#)[Distributed training](#)[Images](#)[Text](#)[Structured data](#)[Generative](#)[TensorFlow](#) > [Learn](#) > [TensorFlow Core](#) > [Tutorials](#)

TensorFlow 2 quickstart for beginners

[Run in Google Colab](#)[View source on GitHub](#)[Download notebook](#)

This short introduction uses [Keras](#) to:

1. Build a neural network that classifies images.
2. Train this neural network.
3. And, finally, evaluate the accuracy of the model.

This is a [Google Colaboratory](#) notebook file. Python programs are run directly in the browser—a great way to learn and use TensorFlow. To follow this tutorial, run the notebook in Google Colab by clicking the button at the top of this page.

1. In Colab, connect to a Python runtime: At the top-right of the menu bar, select *CONNECT*.
2. Run all the notebook code cells: Select *Runtime > Run all*.

Download and install TensorFlow 2. Import TensorFlow into your program:

⚠ Note: Upgrade pip to install the TensorFlow 2 package. See the [install guide](#) for details.

https://www.youtube.com/watch?v=OXNC_sefxi4

Deploying ML+ Applications in Flask

Predict Fake News

Enter the news url

Predict

Training the Model

Building a Fake News Classifier & Deploying it Using Flask

[Article](#) [Github link](#) [Demo Site](#)

1. [1st Fake News](#)
2. [1st Fake News](#)
3. [1st Real News](#)

```
In [108]: 1 # Importing the libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.metrics import classification_report, confusion_matrix
5 from sklearn.pipeline import Pipeline
6 from sklearn.model_selection import train_test_split
7 from sklearn.naive_bayes import MultinomialNB
8 from sklearn.feature_extraction.text import TfidfVectorizer
9 import pickle
10
11 #Importing the cleaned file containing the text and label
12 news = pd.read_csv('news.csv')
13 X = news['text']
14 y = news['label']
```

```
In [109]: 1 #Splitting the data into train
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
3
4 #Creating a pipeline that first creates bag of words(after applying stopwords) & then applies Multinomial Naive Bay
5 pipeline = Pipeline([('tfidf', TfidfVectorizer(stop_words='english')),
6                      ('nbmodel', MultinomialNB())])
```



```
In [109]: 1 #Splitting the data into train
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
3
4 #Creating a pipeline that first creates bag of words(after applying stopwords) & then applies Multinomial Naive Bay
5 pipeline = Pipeline([('tfidf', TfidfVectorizer(stop_words='english')),
6                      ('nbmodel', MultinomialNB())])
7
8 #Training our data
9 pipeline.fit(X_train, y_train)
10
11 #Predicting the label for the test data
12 pred = pipeline.predict(X_test)
```

```
In [110]: 1 splitting the data into train
2 train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
3
4 reating a pipeline that first creates bag of words(after applying stopwords) & then applies Multinomial Naive Bayes
5 peline = Pipeline([('tfidf', TfidfVectorizer(stop_words='english')),
6                      ('nbmodel', MultinomialNB())])
7
8 raining our data
9 peline.fit(X_train, y_train)
10
11 redicting the label for the test data
12 ed = pipeline.predict(X_test)
```

```
In [111]: 1 #Checking the performance of our model
2 print(classification_report(y_test, pred))
3 print(confusion_matrix(y_test, pred))
4
5 #Saving the file
```

```
In [110]: 1 #Splitting the data into train
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
3
4 #Creating a pipeline that first creates bag of words(after applying stopwords) & then applies Multinomial Naive Bay
5 pipeline = Pipeline([('tfidf', TfidfVectorizer(stop_words='english')),
6 ('nbmodel', MultinomialNB())])
7
8 #Training our data
9 pipeline.fit(X_train, y_train)
10
11 #Predicting the label for the test data
12 pred = pipeline.predict(X_test)
```

```
In [111]: 1 #Checking the performance of our model
2 print(classification_report(y_test, pred))
3 print(confusion_matrix(y_test, pred))
4
5 #Serialising the file
6 with open('model.pickle', 'wb') as handle:
7     pickle.dump(pipeline, handle, protocol=pickle.HIGHEST_PROTOCOL)
```

	precision	recall	f1-score	support
FAKE	0.98	0.70	0.82	655
REAL	0.75	0.98	0.85	612
accuracy			0.84	1267
macro avg	0.87	0.84	0.83	1267
weighted avg	0.87	0.84	0.83	1267

```
[[458 197]
 [ 10 602]]
```

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
In [1]: 1 #!pip install newspaper3k  
2 # import nltk  
3 # nltk.download('punkt')
```

```
In [*]: 1 #Importing the Libraries  
2 import numpy as np  
3 from flask import Flask, request, render_template  
4 from flask_cors import CORS  
5 import os  
6 from sklearn.externals import joblib  
7 import pickle  
8 import newspaper  
9 from newspaper import Article  
10 import urllib  
11  
12 #Loading Flask and assigning the model variable  
13 app = Flask(__name__)  
14 CORS(app)  
15  
16 with open('model.pickle', 'rb') as handle:  
17     model = pickle.load(handle)  
18  
19 @app.route('/')  
20 def main():  
21     return render_template('main.html')  
22  
23 #Receiving the input url from the user and using Web Scrapping to extract the news content  
24 @app.route('/predict', methods=['GET', 'POST'])  
25 def predict():  
26     url = request.get_data(as_text=True)[5:]  
27     url = urllib.parse.unquote(url)  
28     article = Article(str(url))  
29     article.download()
```

Running the Flask App



```
In [*]: 1 #Importing the Libraries
2 import numpy as np
3 from flask import Flask, request, render_template
4 from flask_cors import CORS
5 import os
6 from sklearn.externals import joblib
7 import pickle
8 import newspaper
9 from newspaper import Article
10 import urllib
11
12 #Loading Flask and assigning the model variable
13 app = Flask(__name__)
14 CORS(app)
15
16 with open('model.pickle', 'rb') as handle:
17     model = pickle.load(handle)
18
19 @app.route('/')
20 def main():
21     return render_template('main.html')
22
23 #Receiving the input url from the user and using Web Scrapping to extract the news content
24 @app.route('/predict', methods=['GET', 'POST'])
25 def predict():
26     url = request.get_data(as_text=True)[5:]
27     url = urllib.parse.unquote(url)
28     article = Article(str(url))
29     article.download()
30     article.parse()
31     article.nlp()
32     news = article.summary
33     #Passing the news article to the model and returing whether it is Fake or Real
```

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3



```
17     model = pickle.load(handle)
18
19 @app.route('/')
20 def main():
21     return render_template('main.html')
22
23 #Receiving the input url from the user and using Web Scrapping to extract the news content
24 @app.route('/predict',methods=['GET','POST'])
25 def predict():
26     url = request.get_data(as_text=True)[5:]
27     url = urllib.parse.unquote(url)
28     article = Article(str(url))
29     article.download()
30     article.parse()
31     article.nlp()
32     news = article.summary
33     #Passing the news article to the model and returing whether it is Fake or Real
34     pred = model.predict([news])
35     return render_template('main.html', prediction_text='The news is "{}".format(pred[0]))
36
37 if __name__ == '__main__':
38     from werkzeug.serving import run_simple
39     run_simple('localhost', 9019, app)
40 # if __name__ == '__main__':
41 #     app.run(debug = True)
```

```
/Users/yssuen/opt/anaconda3/lib/python3.7/site-packages/scikit-learn/externals/joblib/__init__.py:15: DeprecationWarning:
sklearn.externals.joblib is deprecated in 0.21 and will be removed in 0.23. Please import this functionality directly
from joblib, which can be installed with: pip install joblib. If this warning is raised when loading pickled models,
you may need to re-serialize those models with scikit-learn 0.21+.
    warnings.warn(msg, category=DeprecationWarning)
* Running on http://localhost:9019/ (Press CTRL+C to quit)
127.0.0.1 - - [20/Apr/2020 18:08:33] "GET / HTTP/1.1" 200 -
```

THANK YOU FOR YOUR TIME!

Popular Algorithms

Unsupervised Learning	Supervised Learning	Reinforcement Learning
Principal Component Analysis	K-Nearest Neighbour	Q-Learning
Hierarchical Clustering	Logical Regression	Markov Decision Processes
K-Means Clustering	Linear Regression	
Association Rules	Decision Tree	
Collaborative Filtering	Random Forest	
	Support Vector Machine	
	Naive Baynes	
	Artificial Neural Network	