

Modulo IV.1 - Tecniche di osservazione e valutazione delle risorse intelligenti : introduzione al Data Mining

1. Data Mining
2. Tipi di apprendimento nel Data Mining
3. Algoritmi di classificazione
4. Algoritmi di clustering
5. Algoritmi di regressione
6. KNIME
7. Materiali supplementari: utilizzo di KNIME

Bibliografia



1. Data Mining

Il Data Mining è il processo di ricerca e analisi di grandi database per trovare informazioni utili al processo decisionale.

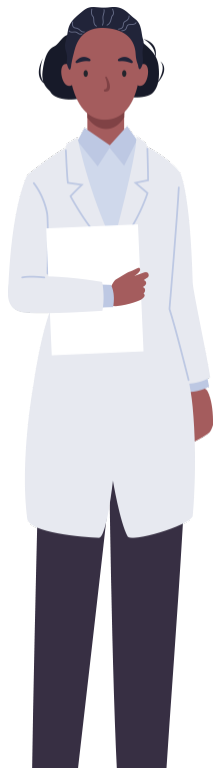
Esistono numerose tecniche di DM che impiegano l'analisi matematica per dedurre i modelli e le tendenze esistenti nei dati. In genere, questi schemi non possono essere rilevati con l'esplorazione tradizionale dei dati perché le relazioni sono troppo complesse o perché il volume di dati da analizzare è troppo grande.

Attualmente il Data Mining viene utilizzato continuamente per l'analisi di grandi quantità di dati in vari campi della conoscenza, come l'istruzione, l'economia, gli affari, l'ambiente, ...

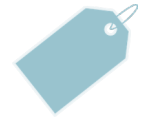




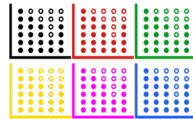
1.1. Data Mining. Nozioni di base



È il processo di estrazione di informazioni utili da grandi volumi di dati.
 I dati possono provenire da fonti diverse, come sensori, social media, transazioni, ecc.
 L'obiettivo è identificare modelli, tendenze e relazioni nascoste nei dati.



È un'attività che si svolge in un ambiente di lavoro dove si utilizzano strumenti e tecniche
 per analizzare i dati e prendere decisioni basate sui risultati.



È un'attività che si svolge in un ambiente di lavoro dove si utilizzano strumenti e tecniche
 per analizzare i dati e prendere decisioni basate sui risultati.



È un'attività che si svolge in un ambiente di lavoro dove si utilizzano strumenti e tecniche
 per analizzare i dati e prendere decisioni basate sui risultati.



1.2. Data Mining. Processo di applicazione di tecniche di data mining

Definizione del problema

Questa è la prima fase in cui un problema specifico viene tradotto in un problema di data mining, in cui vengono sollevati gli obiettivi dell'analisi e le domande di ricerca.

Preparazione e raccolta dei dati

È la fase più estesa del processo, poiché la qualità dei dati è una delle sfide più importanti del data mining. I dati grezzi devono essere identificati, puliti e archiviati in un formato predefinito.

Modellazione e valutazione

In questa fase vengono selezionate e applicate diverse tecniche di modellazione dei dati (algoritmi), per poi stabilire i parametri e i valori ottimali di queste tecniche.

Distribuzione

È l'ultima fase in cui i risultati del data mining vengono organizzati e presentati tramite grafici e report.



1.2. Data Mining. Processo di applicazione di tecniche di data mining



È importante notare che ogni processo di data mining è un processo iterativo, il che significa che il processo non si ferma quando una particolare soluzione viene impiegata. Può essere solo una nuova voce per un altro processo di data mining (Rodríguez-Arribas, 2021). In altre parole, in molte occasioni l'applicazione delle tecniche di DM richiede diverse iterazioni e l'uso di diversi algoritmi per poter estrarre i risultati finali della ricerca che stiamo conducendo.





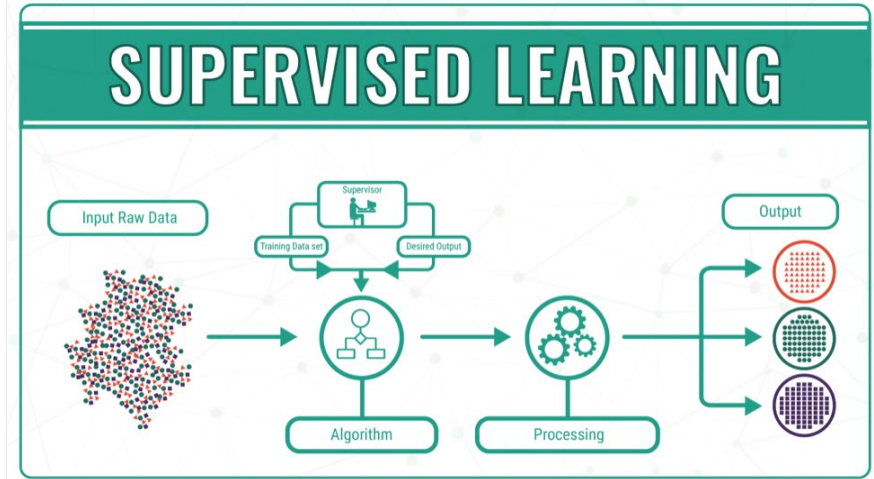
2.1. Tipi di apprendimento nel Data Mining

APPRENDIMENTO SUPERVISIONATO

L'obiettivo fondamentale dell'apprendimento supervisionato è la creazione di un modello in grado di prevedere i valori corrispondenti agli oggetti in ingresso dopo aver familiarizzato con una serie di esempi, i dati di addestramento .

Questa tecnica si compone di due fasi fondamentali :

1. una fase di addestramento, in cui si utilizza un insieme di dati etichettati, che contengono i dati di input e i risultati desiderati per quei dati di addestramento con un algoritmo che permette di dedurre una funzione dai dati grezzi (raw data) che stiamo fornendo all'algoritmo ;
2. la fase di test, in cui la funzione ottenuta nella fase precedente viene utilizzata per generare nuove previsioni con nuovi set di dati .



Il processo è noto come apprendimento supervisionato, poiché conoscendo le risposte di ciascun esempio dell'insieme di formazione, è possibile correggere la funzione generata dall'algoritmo . La preparazione dell'algoritmo viene supervisionata correggendo i suoi parametri, in base ai risultati ottenuti iterativamente .

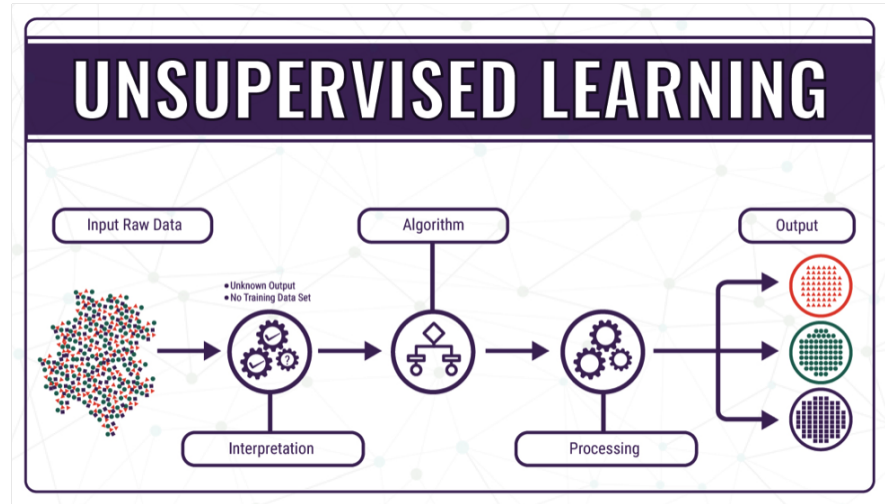


APPRENDIMENTO NON -SUPERVISIONATO

Questo tipo di apprendimento è l'altro approccio di base al Machine Learning (ML). L'apprendimento non supervisionato prevede dati non etichettati che l'algoritmo deve cercare di comprendere da solo.

L'obiettivo di questo tipo di apprendimento è far sì che la macchina impari senza l'aiuto o le indicazioni dei data scientist, cioè senza supervisione e senza un set di dati di addestramento. Inoltre, la macchina stessa aggiusterà i risultati e i raggruppamenti quando ci saranno risultati più adatti, permettendo alla macchina di comprendere i dati e di elaborarli nel modo migliore.

L'apprendimento non-supervisionato viene utilizzato per esplorare dati sconosciuti e non etichettati. Può rivelare modelli che potrebbero essere stati trascurati o esaminare grandi insiemi di dati che sarebbero troppo impegnativi per una singola persona.





APPRENDIMENTO SEMI-SUPERVISIONATO

Attualmente sono in corso numerose ricerche sui metodi di apprendimento semi-supervisionato . Queste tecniche di apprendimento automatico utilizzano dati di addestramento sia etichettati che non etichettati : in genere, una piccola quantità di dati etichettati insieme a una grande quantità di dati non etichettati (Zhu e Goldberg ,2009).

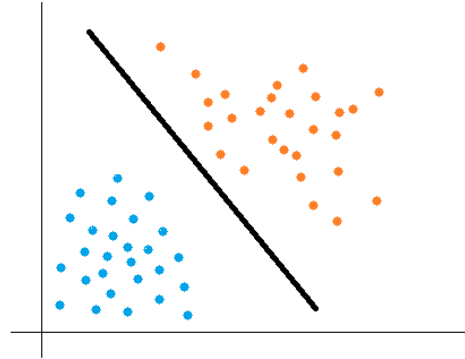
In altre parole, cercano di migliorare i modelli di previsione ottenuti utilizzando esclusivamente dati etichettati, esplorando le informazioni strutturali contenute nei dati non etichettati . Possiamo dire che l'apprendimento semi -supervisionato cerca di combinare i due approcci tradizionali del data mining (apprendimento supervisionato e apprendimento non supervisionato) per mantenere il meglio di ciascuno di essi.





3. Algoritmi di classificazione

d H ¼ P c t = | Ó | r k g g e r ¼ w w s g u C L s H P = s L i H w w ¼ =
 C L ¼ O c e l l f ¼ ¼ i s g g f l s i r = s i i ¼ Ó p r c s i ¼ d u r f l i H
 C L ¼ O t c g C g i ¼ ¼ ¼ O = ¼ O ¼ Ó | c | r s d r ¼ g | i c w ¼ ¼ H H i s c u O |
 L = P g t = s e = t O c e l l f ¼ | C g g t H
 à ¼ r k g g e r ¼ w w s f = t g t f ¼ ¼ C c s g g Ó | ¼ C c s = O t s = s
 Ó s f C s c g s C r = T C g g s Ó | = f ¼ r ¼ C ¼ r t a Ó | r k g g e r ¼ c s
 r t H c | ¼ t = ¼ H C ¼ s i | r f i l i i r b r = s r | r f = O ¼

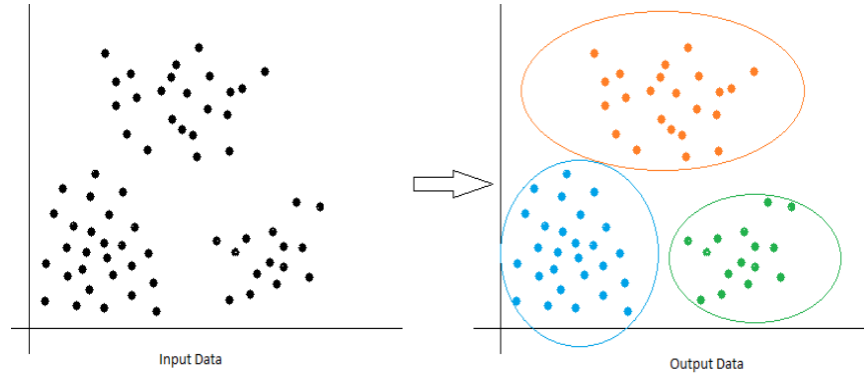


Questi algoritmi lavorano generalmente sulle informazioni fornite da un insieme di campioni, modelli, esempi o prototipi di addestramento che vengono presi come rappresentanti delle classi e conservano un'etichetta di classe corretta. Questo insieme di prototipi correttamente etichettati è chiamato insieme di addestramento e costituisce la conoscenza disponibile per la classificazione di nuovi campioni.

L'obiettivo della classificazione supervisionata è quello di determinare, in base a ciò che è noto, a quale classe dovrebbe appartenere un nuovo campione, considerando le informazioni che possono essere estratte.

4. Algoritmi di clustering

Gli algoritmi di clustering sono responsabili del raggruppamento degli oggetti in un set di dati in base alle loro somiglianze. In questo modo, gli oggetti che si trovano all'interno di un cluster o di un gruppo hanno più somiglianze tra loro che differenze. Questi algoritmi lavorano con dati non etichettati, quindi è l'algoritmo stesso che analizza i dati per trovare il numero ottimale di raggruppamenti per l'insieme di dati in ingresso, poiché non abbiamo una conoscenza preliminare delle caratteristiche dei dati e delle loro classi.



I raggruppamenti effettuati dagli algoritmi possono essere di due tipi:

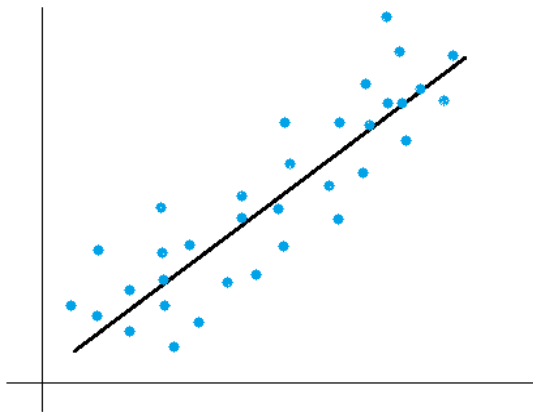
- hard cluster** : ogni dato appartiene esclusivamente a un gruppo;
- soft (diffuse) cluster** : i dati possono appartenere a più gruppi in misura diversa, cioè gli stessi dati possono avere un grado di appartenenza del 60% al gruppo 1 e del 40% al gruppo 2.



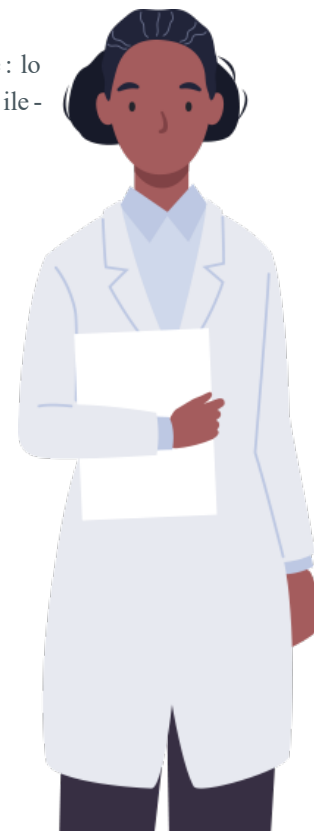


Gli algoritmi di regressione sono un sottocampo dell'apprendimento supervisionato il cui obiettivo è stabilire un metodo per la relazione tra un certo numero di caratteristiche e una variabile obiettivo continua.

Si tratta di algoritmi che stabiliscono una traiettoria per fornire la tendenza di un insieme di dati, cioè: lo scopo di questi algoritmi è di mettere in relazione un certo numero di caratteristiche e una variabile-obiettivo continua.



Questa tecnica è utile per prevedere esiti che sono valori continui, il che significa che la risposta alla domanda di ricerca è presentata da una quantità che può essere determinata in modo flessibile in base agli input del modello, piuttosto che essere limitata a un insieme finito di etichette come nel caso della classificazione.





6. KNIME

KNIME è un'applicazione open-source che ci permette di applicare ai nostri set di dati o a set di dati campione :

1. metodi statistici,
2. algoritmi di data mining o Machine Learning (apprendimento automatico),
3. tecniche di visualizzazione .

È costruito sulla piattaforma Eclipse ed è programmato in Java. Essere un software open-source ha molti vantaggi : il suo codice appartiene alla comunità di utenti e sviluppatori, il che garantisce che sarà sempre uno strumento libero che può essere scaricato e utilizzato gratuitamente secondo i termini della licenza GPLv3. Permette inoltre di incorporare codice sviluppato in R o Python .

È uno strumento di "programmazione visiva". L'analisi dei dati può essere effettuata in modo intuitivo, impostando il processo con un semplice clic del mouse . I "nodi" di cui abbiamo bisogno vengono posizionati, senza aver bisogno di conoscerne il nome o la configurazione, poiché è sempre possibile ricevere aiuto .

È uno strumento progettato per essere semplice da usare . Il concetto più importante nell'uso dello strumento è quello di flusso di lavoro (workflow). Un flusso di lavoro è una sequenza di passi configurati dall'utente . Formalmente è un insieme di nodi uniti da frecce che rappresentano il flusso di dati da un nodo all'altro . Un nodo racchiude diversi lavori che possono essere eseguiti con i dati ; ci sono nodi per molti compiti .

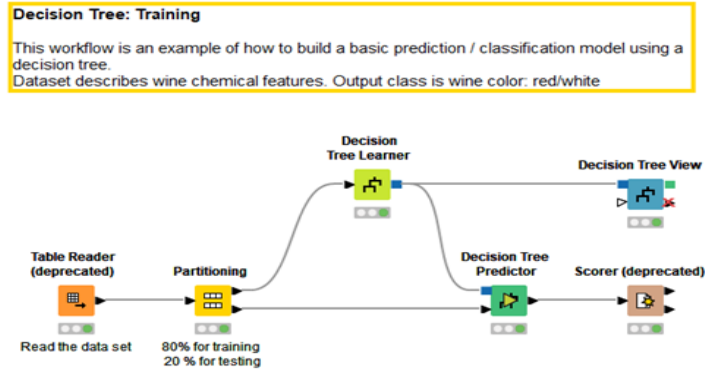


Sono presenti nodi per:

- caricare i dati da file o database.
- Creare, modificare o eliminare righe o colonne dal set di dati con cui stiamo lavorando.
- Calcolare statistiche medie, percentili, correlazioni, ecc.
- Combinare dati provenienti da fonti diverse.
- Costruire e valutare modelli di Machine Learning come: classificazione, regressione o clustering.
- Visualizzare i dati utilizzando grafici a barre, a torta, a dispersione e altri tipi di grafici più avanzati.
- Generare report.

Un flusso di lavoro potrebbe avere un nodo per caricare un set di dati da un file Excel, poi un nodo per selezionare gli attributi (colonne) da quel set di dati e infine un altro nodo per visualizzare le statistiche degli attributi selezionati.

Qui di seguito, un esempio di «albero decisionale»



Modulo IV.1. Materiali supplementari : utilizzo di KNIME

1. KNIME, installazione

2. KNIME, flussi di lavoro

2.1. Nodi

2.2. L'area di lavoro

3. Esempio generico: classificazione delle specie floreali

4. Esempio con i dati di un intervento terapeutico intelligente (EarlyCare)

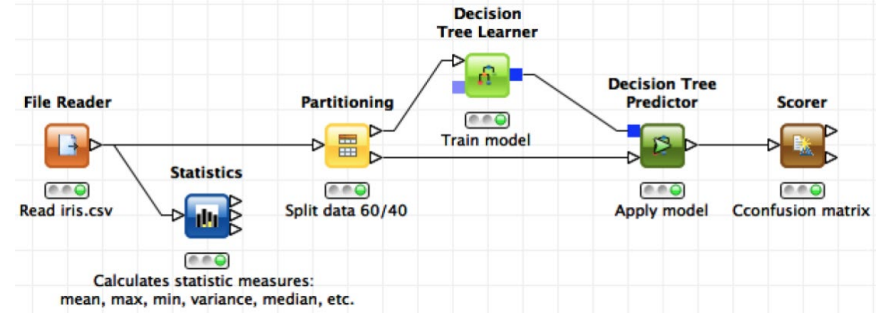
Riferimenti dal web





2. KNIME, flussi di lavoro

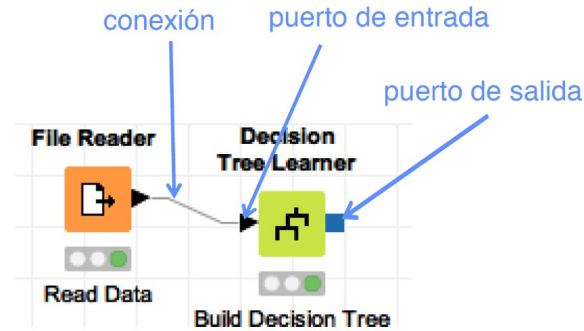
Questa sezione è stata generata automaticamente da un software di traduzione e potrebbe contenere errori di grammatica o di sintassi. Il contenuto non è stato verificato e potrebbe non essere accurato. Si consiglia di leggere attentamente il testo originale e di verificare la correttezza delle informazioni presentate.



2.1. Nodi

I nodi (nell'immagine accanto, i nodi di entrata e di uscita, oltre alla connessione) incapsulano gli algoritmi che implementano le azioni che possono essere eseguite sui dati:

- manipolazione di righe, colonne, ecc.
- Creazione di modelli di data mining.
- Valutazione dei modelli.
- Applicazione dei modelli a nuovi dati.
- Processi ETL (Extract, Load, Transform).
- Creazione di report personalizzati.





2.2. L'area di lavoro

L'area di lavoro è la cartella o directory del nostro computer in cui sono memorizzati tutti i progetti realizzati con KNIME. Sarà necessario scegliere un'area di lavoro prima di avviare il programma (si può anche lasciare la cartella che appare di default al momento dell'installazione).



The screenshot displays the KNIME Analytics Platform interface with a workflow titled "Building a Simple Classifier". The workflow consists of the following nodes:

- Partitioning**: A "Random drawing" node with parameters "80% upper part" and "20% lower part". It splits the data into a "training set" and a "test set".
- Decision Tree Learner**: A node that receives the "training set" and is configured to "Train to predict class 'income'".
- Decision Tree Predictor**: A node that receives the "test set" and is configured to "Apply decision tree model to test set".
- Interactive Table**: A node that receives the output from the predictor and is configured to "Display table of the entire data".
- Interactive Table (local)**: A node that receives the output from the predictor and is configured to "Show entire data as table".

The interface also includes several side panels:

- KNIME Explorer**: Shows the project structure, including "My-KNIME-Hub", "EXAMPLES", and "LOCAL (Local Workspace)".
- Workflow Coach**: Lists "Recommended Nodes" such as "Decision Tree Predictor" (85%), "Decision Tree To Image" (5%), and "PMML Writer" (3%).
- Node Repository**: A tree view of available nodes categorized by function like "IO", "Manipulation", "Views", "Analytics", "DB", etc.
- Outline**: Provides a high-level overview of the workflow structure.
- Console**: Displays the KNIME startup log, including the version "v4.0.1.v201908131317" and a warning: "WARN Color Manager 0:2 Column 'income' has no nominal values".
- Description**: Provides a detailed explanation of the "Decision Tree Learner" node, stating it induces a classification decision tree in main memory and handles both nominal and numerical target attributes.



2.2.1. L'area di lavoro: KNIME Explorer

Y | % | \$ | ¼ | ¶ | P | L | | W | ¶ | ↑ | ¶ | ¶ | ↑ | s | g | i | ¶ | | | C | ¶ | ↑ | s | i | ¶ | s | | ¼ | L | g | g | | Ö | | ¶ | ¶ | ¶ | ¶ | g | ¼ | ¶ | ¶ | | Ö | ¶ | W | g | | ¶ | C | ¶ | ¶ | ¶ | ¼ | ¶ | ¶ | ¶ | ¶ | s | g | C | ¶ | ¶ | ¶ | ¼ | ¶ | ¶ | ¶ | ¶ | ¼ | L | g | g | | Ö | | ¶ | ¶ | ¶ | ¶



KNIME Explorer

- EXAMPLES (knime-guest@http://publ
- LOCAL (Local Workspace)
 - Example Workflow

Welcome to KNIME Analytics Platform

0: Example Wo

This Example Workflow uses a **File Reader** node to import the Iris dataset (included). It then assigns some basic statistics with a **Statistics** node. The data is split into training and testing fractions with a predictive model in PMML from the training fraction which is then applied to the test fraction using the **Scorer** node, which is applied after the **Decision Tree Predictor**. Finally, errors can be highlighted certain classes of errors which can then be visualized using a **Scatter Plot** node.

Workflow Coach

Decision Tree Lear



2.2.2. L'area di lavoro: flusso di lavoro, editor

Questa è l'area di lavoro principale, dove i nodi vengono trascinati, collegati e il flusso di lavoro viene configurato.

The screenshot displays the KNIME Analytics Platform interface. The main workspace is highlighted with a red border and contains a workflow diagram. A yellow box highlights a text description of the workflow. The workflow consists of the following nodes:

- File Reader**: Read iris.csv
- Color Manager**: Assign colors
- Statistics**: Calculate statistic measures: mean, max, min, variance, median, etc.
- Partitioning**: Split data 60/40
- Decision Tree Learner**: Train model
- Decision Tree Predictor**: Apply model
- Scatter Plot**: View test data
- Scorer**: Compute confusion matrix
- Interactive Table**: Explore test data

The interface also shows a Node Repository on the left with categories like IO, Manipulation, Views, Analytics, Database, and Other Data Types. The top bar includes a zoom level of 75% and a toolbar with various icons. The bottom bar shows an Outline and Console panel.

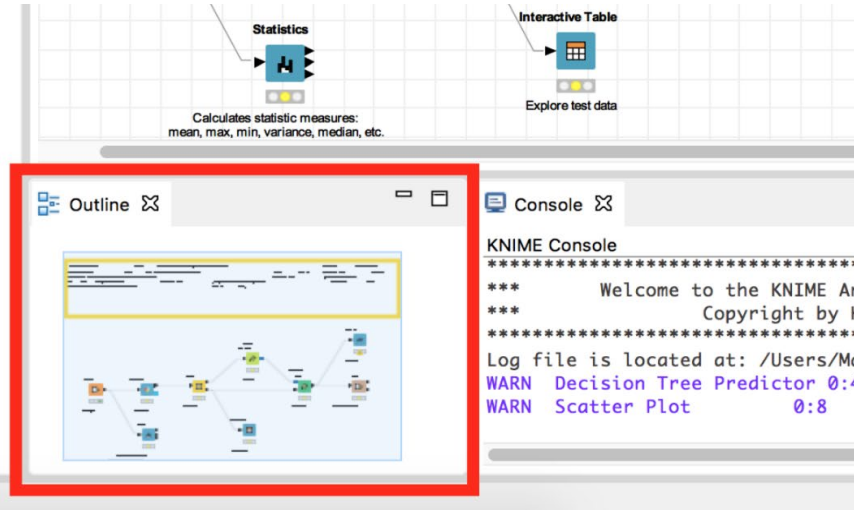




2.2.3. L'area di lavoro: schema

Il diagramma mostra un workflow di lavoro in KNIME con i seguenti componenti:

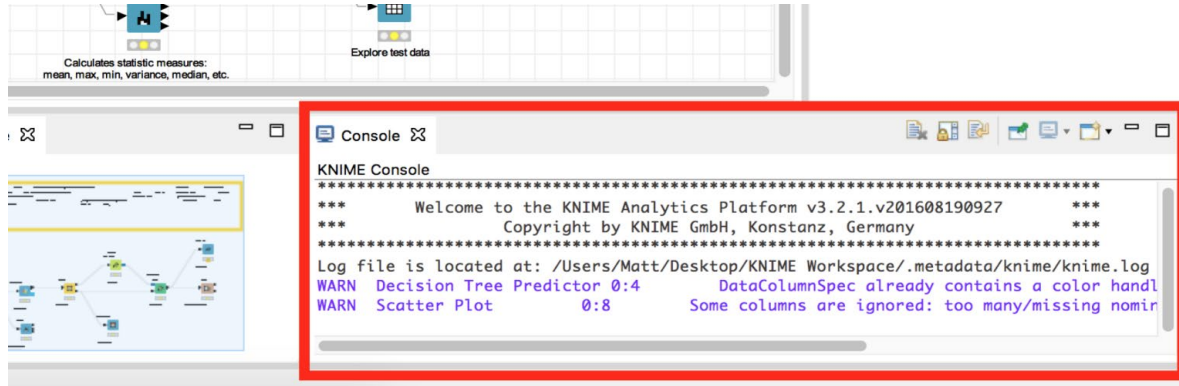
- ▶ IO
- ▶ Manipulation
- ▶ Views
- ▶ Analytics
- ▶ Database
- ▶ Other Data Types
- ▶ Structured Data
- ▶ Scripting
- ▶ Tool Integration
- ▶ Community Nodes
- ▶ KNIME Labs
- ▶ Workflow Control
- ▶ Social Media
- ▶ Reporting
- ▶ Chemistry
- ▶ ChemAxon / Infocom





2.2.4. L'area di lavoro: la console

È un campo di output di testo che visualizza gli avvisi e gli errori che si verificano durante l'esecuzione del flusso di lavoro. Visualizza anche informazioni rilevanti sul processo di esecuzione.





2.2.5. L'area di lavoro: archivio dei nodi

[Illegible text, possibly a corrupted or stylized header]



The screenshot displays the KNIME Workflow Coach interface. On the left, the **Node Repository** panel is highlighted with a red border, showing a tree view of node categories: IO, Manipulation, Views, Analytics, Database, Other Data Types, Structured Data, Scripting, Tool Integration, Community Nodes, KNIME Labs, Workflow Control, Social Media, Reporting, Chemistry, and ChemAxon / Infocom. The main workspace shows a workflow with nodes: File Reader (Read iris.csv), Color Manager (Assign colors), Statistics (Calculates statistic measures: mean, max, min, variance, median, etc.), Partitioning (Split data 60/40), Decision Tree Learner (Train model), Interactive Table (Explore test data), Decision Tree Predictor (Apply model), Scorer (Compute confusion matrix), and Scatter Plot (View test data). The bottom right shows the **Console** with KNIME version information and a warning message: `WARN Decision Tree Predictor 0:4 DataColumnSpec already contains a color hc`.

2.2.6. L'area di lavoro: coach del flusso di lavoro

Se abbiamo dato il permesso di raccogliere i nostri dati, questa sezione ci suggerisce quali nodi è più probabile che dobbiamo utilizzare in un determinato momento.



Workflow Coach

Recommended Nodes	Community
Decision Tree Predictor	85%
Decision Tree To Image	5%
Decision Tree to Ruleset	3%
PMML Writer	3%
Decision Tree View	1%
PMML To Cell	<1%
Boosting Learner Loop End	<1%
Model Writer	<1%
Model Loop End	<1%

Try this:

KNIME's Interactive Visualizations:

- 1) Execute the workflow
- 2) Open the Scorer node view
- 3) Hit a cell in the confusion matrix
- 4) Open the Interactive Table view
- 5) Select "Hitte"->"Filter"->"Show Hitlled Only"

This shows only the misclassified data rows.



2.2.7. L'area di lavoro: descrizione dei nodi

Il workflow di esempio utilizza un nodo File Reader per importare il dataset Iris (incluso). Assegna le proprietà visive con un nodo Color Manager e calcola alcune statistiche di base con un nodo Statistics. I dati vengono divisi in frazioni di addestramento e test con un nodo Partitioning. Il nodo Decision Tree Learner genera un modello predittivo in PMML dalla frazione di addestramento, che viene poi applicato alla frazione di test utilizzando il nodo Decision Tree Predictor. Le prestazioni del modello vengono valutate con il nodo Scorer, che viene applicato dopo il nodo Decision Tree Predictor. Infine, gli errori possono essere esplorati interattivamente utilizzando un nodo Interactive Table per evidenziare alcune classi di errori, che possono poi essere visualizzati utilizzando un nodo Scatter Plot.

The screenshot displays the KNIME Analytics Platform interface. On the left, a workflow is visible with nodes: File Reader (Read iris.csv), Column Filter (Petal ONLY), Color Manager (Assign colors), Statistics (Calculates statistic measures: mean, max, min, variance, median, etc.), Partitioning (Split data 60/40), Decision Tree Learner (Train model), Decision Tree Predictor (Apply model), Interactive Table (Explore test data), Scorer (Compute confusion matrix), and Scatter Plot (View test data). A yellow text box highlights the workflow description. On the right, the 'Node Description' panel for the 'File Reader' node is shown, containing instructions on how to interact with the table header and a list of output ports.

File Reader

Click on the table header

If the column header in the preview table is clicked, a new dialog opens where column properties can be set: name and type can be changed (and will be fixed then). A pattern can be entered that will cause a "missing cell" to be created when it's read for this column. Additionally, possible values of the column domain can be updated by selecting "Domain". And, you can choose to skip this column entirely, i.e. it will not be included in the output table then.

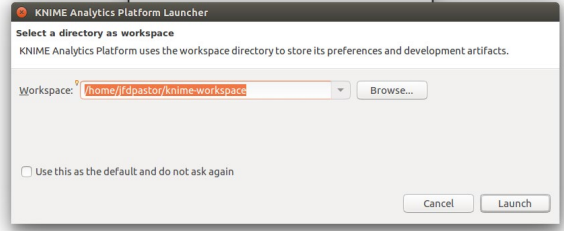
Ports

Output Ports

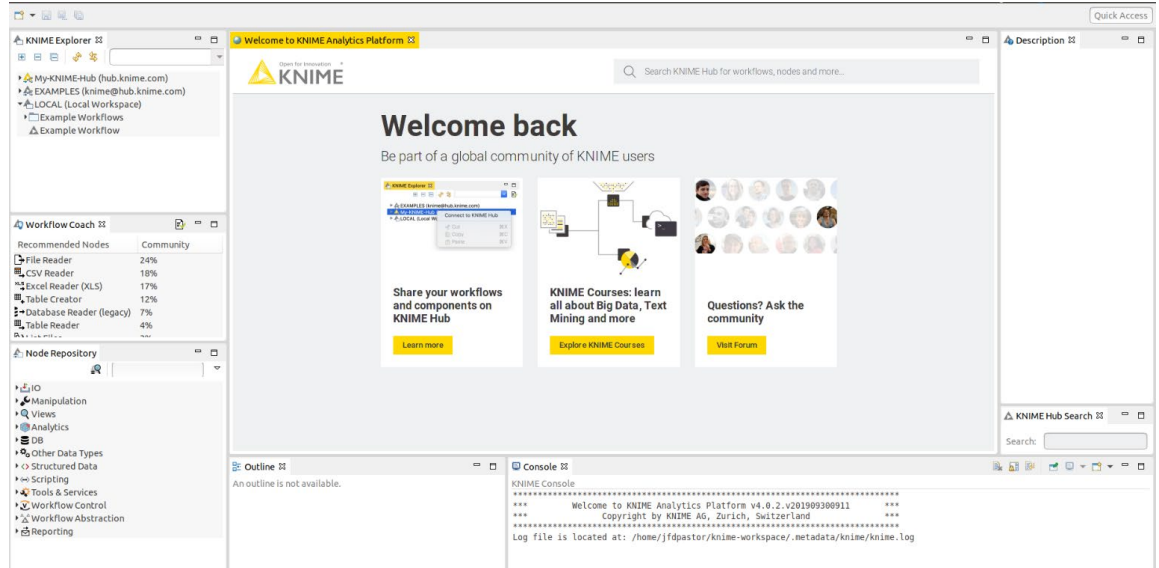
- 0 Datable just read from the file



3. Esempio generico: classificazione delle specie floreali



Supponendo che KNIME sia già stato installato, si deve andare nella cartella in cui si trova ed eseguirlo facendo doppio clic sulla sua icona. All'apertura, ci chiederà la cartella "Workspace". Questa è la cartella in cui si troveranno tutti i nostri progetti.

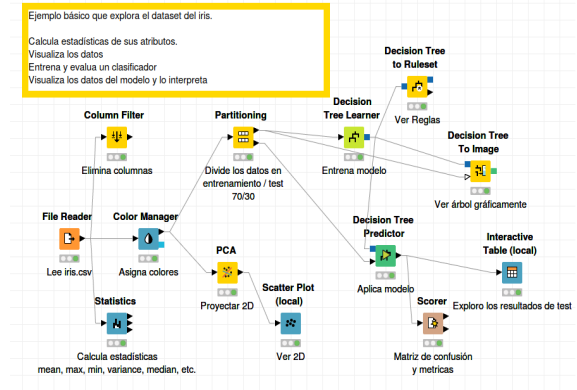


3.1. Esempio generico: classificazione delle specie floreali

- Calcolare le statistiche degli attributi.
- visualizzare i dati,
- addestrare e valutare un classificatore.

Si tratta di un flusso di lavoro di base, che lavora con il set di dati di Iris. L'Iris è costituito da 150 esempi appartenenti a 3 diverse specie di fiori. Ogni esempio ha 4 attributi che descrivono il fiore: lunghezza del sepal, lunghezza del petalo, larghezza del sepal e larghezza del petalo. Con questo set di dati, ci occuperemo di:

- calcolare le statistiche degli attributi.
- visualizzare i dati,
- addestrare e valutare un classificatore.





3.2. Esempio generico: classificazione delle specie floreali

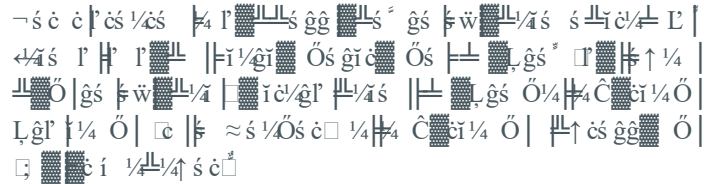
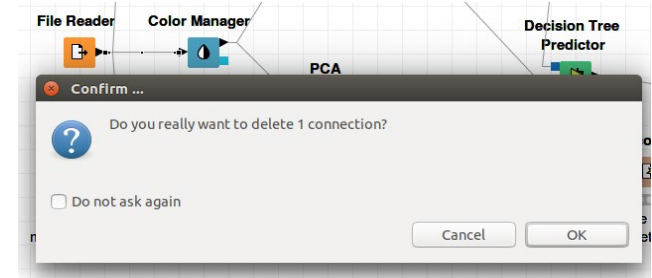
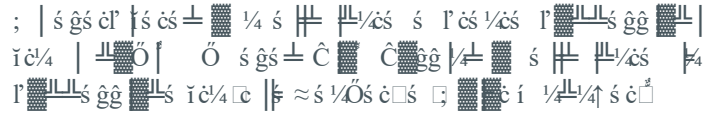
Nell'editor si può vedere una serie di nodi interconnessi. In questo editor, i nodi vengono trascinati, uniti, configurati ed eseguiti per eseguire operazioni e analisi sui dati.

Dispone di strumenti di navigazione come lo zoom in/out (ingrandimento o rimpicciolimento) e permette di aggiungere commenti.



È possibile eseguire ogni nodo o l'intero flusso di lavoro con pulsanti simili a "play".

È necessario eseguire i nodi successivi ogni volta che viene apportata una modifica a un nodo. In altre parole, se si modifica un parametro di un nodo che si trova all'inizio del flusso di lavoro, è necessario premere il pulsante play con le due frecce bianche per eseguire nuovamente tutti i nodi successivi.



3.3. Esempio generico: classificazione delle specie floreali

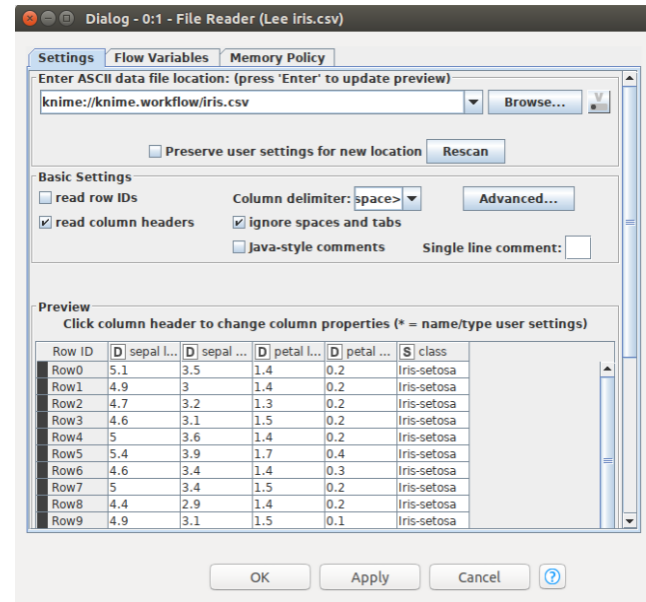
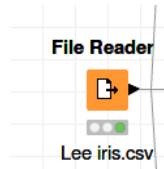
Loading data

Il nodo "File Reader" è il nodo utilizzato per caricare i set di dati (leggere i dati da qualsiasi punto siano memorizzati). È possibile caricare i dati da un url (Internet) o dal disco rigido del computer.

Y C̄ḡḡ |̄ |̄ P̄ē↑ L̄c̄/4s̄ |̄ P̄Ȫ ↔4̄ s̄ P̄Ȫ ȪC̄C̄ |̄ P̄ |̄
ḡL̄ Ȫ |̄ s̄ ḡḡ

B L̄c̄/4s̄ |̄ |̄ P̄ē↑ L̄c̄/4w̄ |̄ s̄ Φ C̄ḡḡ |̄ |̄ |̄ C̄ḡī/4s̄
|̄ |̄ s̄ ḡī/4w̄ |̄ s̄ |̄ |̄ Ȫs̄ |̄ |̄ |̄ Ȫs̄ |̄ ē |̄ P̄=̄s̄
Ȫs̄ ḡ |̄ s̄ c̄/4=̄ |̄ P̄c̄/4s̄

f L̄s̄ ḡī Φ |̄ s̄ P̄s̄ ḡḡ/4c̄ C̄s̄ c̄P̄=̄T̄ 1/4 V̄|̄ s̄ |̄ ē |̄ Ȫ |̄ Ȫ/4 |̄
ḡ |̄ ḡs̄ C̄/4/4 |̄ Ȫ/4 W̄|̄ |̄ 1/4 V̄|̄ |̄ Ȫ/4 C̄L̄ |̄ |̄ s̄ W̄|̄ |̄
1/4 V̄|̄ |̄ Ȫ/4ī/4L̄ |̄ W̄ |̄ |̄ s̄ P̄





, * [ð ± Ć [s L s c [[gg e l ¼ w [L s Ó s H g Ć s l [¼ c s ¼ H

Colorazione dei dati

Il nodo "Colour Manager" consente di colorare il dataset in base ai valori di uno dei suoi attributi.

Il risultato è una tabella in cui ogni riga è colorata in base al valore dell'attributo scelto.

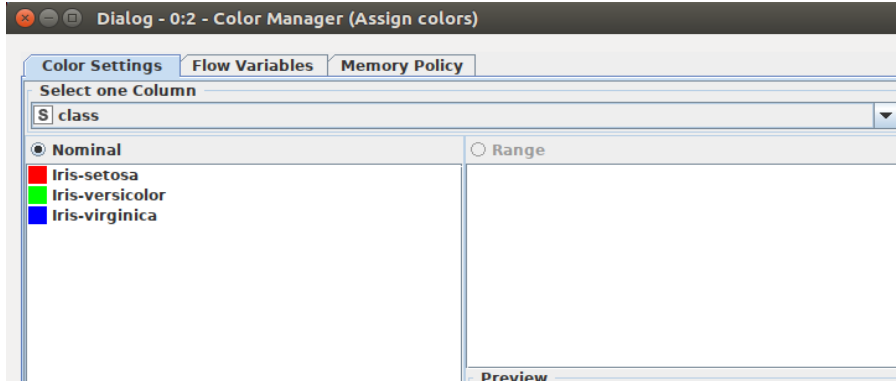
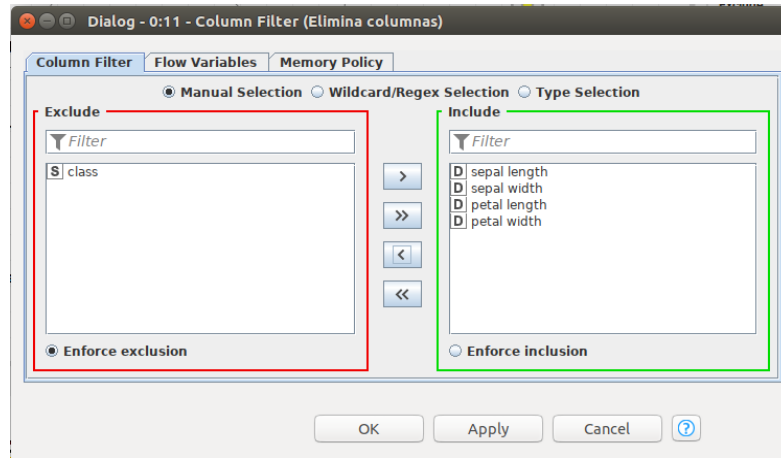
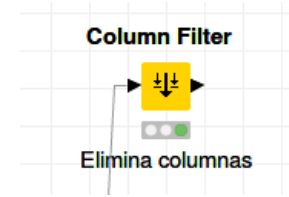


Table with Colors - 0:2 - Color Manager (Asigna colores)					
File Hilite Navigation View					
Table "default" - Rows: 150		Spec - Columns: 5		Properties Flow	
Row ID	D sepal l...	D sepal ...	D petal l...	D petal ...	S class
Row24	4.8	3.4	1.9	0.2	Iris-setosa
Row25	5	3	1.6	0.2	Iris-setosa
Row26	5	3.4	1.6	0.4	Iris-setosa
Row27	5.2	3.5	1.5	0.2	Iris-setosa
Row28	5.2	3.4	1.4	0.2	Iris-setosa
Row29	4.7	3.2	1.6	0.2	Iris-setosa
Row30	4.8	3.1	1.6	0.2	Iris-setosa
Row31	5.4	3.4	1.5	0.4	Iris-setosa
Row32	5.2	4.1	1.5	0.1	Iris-setosa
Row33	5.5	4.2	1.4	0.2	Iris-setosa
Row34	4.9	3.1	1.5	0.2	Iris-setosa
Row35	5	3.2	1.2	0.2	Iris-setosa
Row36	5.5	3.5	1.3	0.2	Iris-setosa
Row37	4.9	3.6	1.4	0.1	Iris-setosa
Row38	4.4	3	1.3	0.2	Iris-setosa
Row39	5.1	3.4	1.5	0.2	Iris-setosa
Row40	5	3.5	1.3	0.3	Iris-setosa
Row41	4.5	2.3	1.3	0.3	Iris-setosa
Row42	4.4	3.2	1.3	0.2	Iris-setosa
Row43	5	3.5	1.6	0.6	Iris-setosa
Row44	5.1	3.8	1.9	0.4	Iris-setosa
Row45	4.8	3	1.4	0.3	Iris-setosa
Row46	5.1	3.8	1.6	0.2	Iris-setosa
Row47	4.6	3.2	1.4	0.2	Iris-setosa
Row48	5.3	3.7	1.5	0.2	Iris-setosa
Row49	5	3.3	1.4	0.2	Iris-setosa
Row50	7	3.2	4.7	1.4	Iris-versic...
Row51	6.4	3.2	4.5	1.5	Iris-versic...
Row52	6.9	3.1	4.9	1.5	Iris-versic...
Row53	5.5	2.3	4	1.3	Iris-versic...
Row54	6.5	2.8	4.6	1.5	Iris-versic...
Row55	5.7	2.8	4.5	1.3	Iris-versic...
Row56	6.3	3.3	4.7	1.6	Iris-versic...
Row57	4.8	3.4	3.2	1	Iris-versic...

3.5. Esempio generico: classificazione delle specie floreali

Rimuovere colonne

c | f | r | s | l | w | p | class |
g | l | p | s | g | w | l | w | p | class |
r | l | w | p | class |
l | w | p | class |
c | f | r | s | l | w | p | class |



Nell'esempio, la colonna "Classe", che contiene il nome della specie a cui appartiene il fiore descritto in ogni esempio, verrà eliminata.

L'operazione viene eseguita solo per causare un errore nel flusso di lavoro.

Saper identificare i tipi di errore è fondamentale per utilizzare uno strumento come KNIME.





Comando di filtraggio delle colonne per il nodo Column Filter.

Errori e notifiche

- Rimuovere il collegamento tra "File Reader" e "Color Manager".
- Configurare "Column Filter" per rimuovere la classe.
- Collegiamo "Column Filter" con "Color Manager".

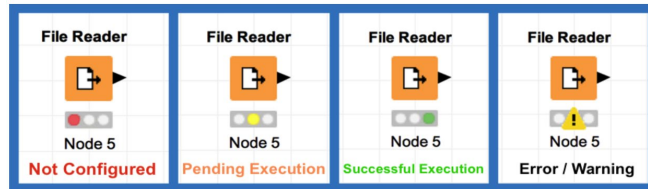
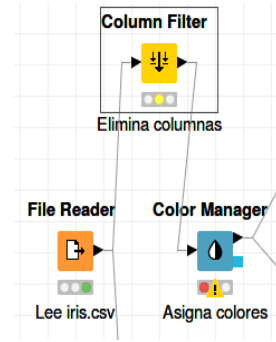
Ora si ottiene un errore in "Color Manager", perché questo nodo ha usato la classe per dare colore agli esempi. Per continuare, ristabiliamo la connessione tra "File Reader" e "Color Manager".

Un nodo può trovarsi in 4 stati diversi.

- Non configurato. È necessario fare doppio clic su di esso e scegliere qualche parametro importante che lo strumento non può scegliere per noi.
- In attesa. Il pulsante di esecuzione non è ancora stato premuto.
- Eseguito.
- Errore/Avvertimento. Non può essere eseguito.

(Come sopra, quando si elimina una colonna utilizzata da un nodo successivo).

(Colonna utilizzata da un nodo successivo).



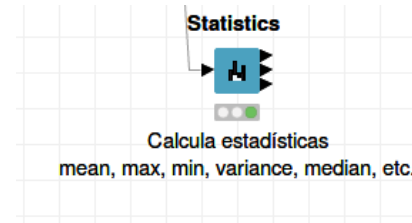
3.7. Esempio generico: classificazione delle specie

floreali



Calcolare statistiche

Il nodo "Statistiche" consente di ottenere statistiche da una tabella di dati. Selezionando il nodo e facendo clic su "Visualizzazione statistiche", si ottiene una tabella con le statistiche per ciascun attributo.



Statistics View - 0:9 - Statistics (Calculates statistic measures)

File

Numeric Nominal Top/bottom

Column	Min	Mean	Median	Max	Std. Dev.	Skewness	Kurtosis	No. Missing	No. +∞	No. -∞	Histogram
sepal length	4.3	5.8433	?	7.9	0.8281	0.3149	-0.5521	0	0	0	
sepal width	2	3.0573	?	4.4	0.4359	0.319	0.2282	0	0	0	
petal length	1	3.758	?	6.9	1.7653	-0.2749	-1.4021	0	0	0	
petal width	0.1	1.1993	?	2.5	0.7622	-0.103	-1.3406	0	0	0	

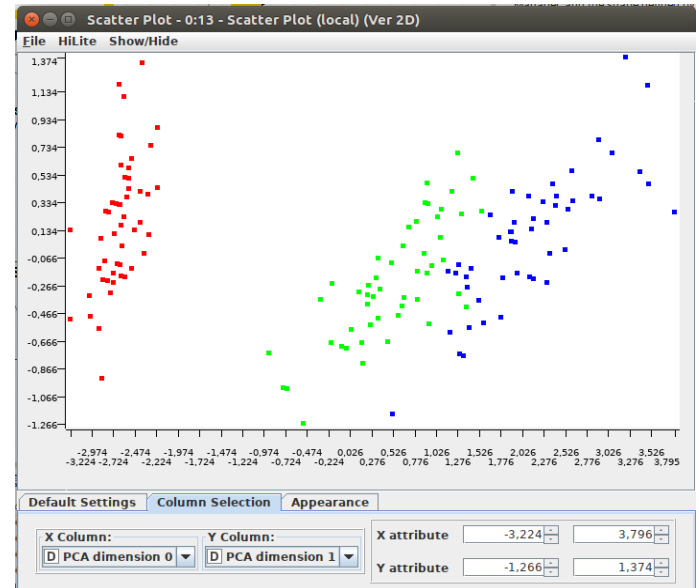
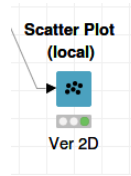


3.9. Esempio generico: classificazione delle specie floreali

Grafico di dispersione

Il grafico di dispersione (scatter plot) è uno strumento fondamentale per visualizzare i dati multivariati. In questo caso, le prime due componenti principali (PCA dimension 0 e 1) vengono utilizzate per ridurre la dimensionalità dei dati, consentendo di visualizzare la relazione tra diverse variabili in un piano bidimensionale. I punti sono colorati in base alle diverse specie floreali, permettendo di osservare eventuali cluster o tendenze nel dataset.

Nell'esempio vengono visualizzate le prime due componenti principali; questa tecnica consente di riassumere e visualizzare (con due attributi) un insieme di dati con diverse colonne.



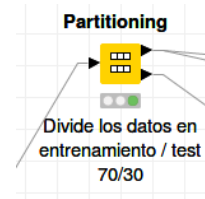
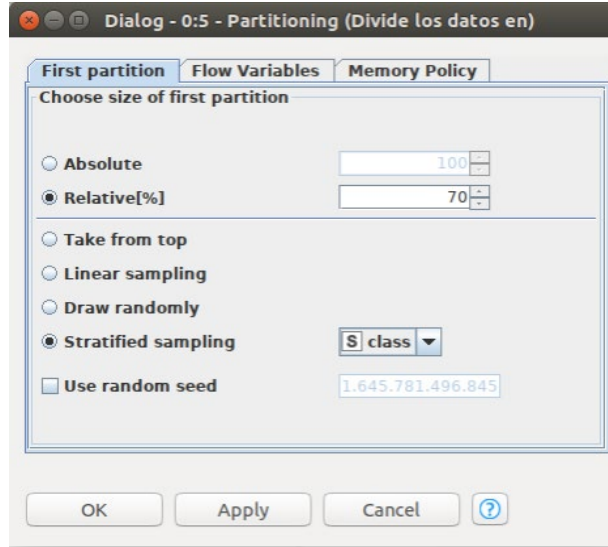
3.10. Esempio generico: classificazione delle specie

floreali



Partizionamento dei dati

Il dataset viene diviso in due parti: una per l'addestramento (70%) e una per la validazione (30%).



Nel nodo si può scegliere la percentuale di istanze da utilizzare per imparare e la percentuale da utilizzare per il test.

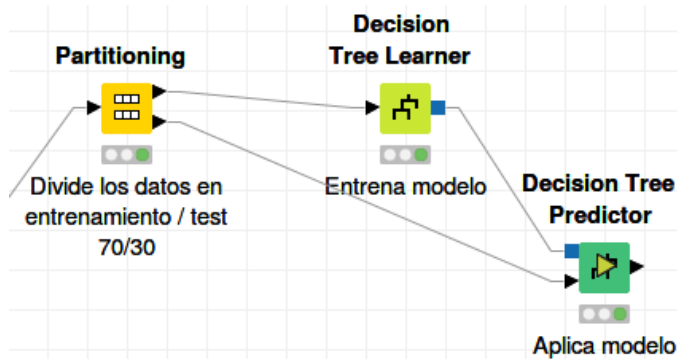
Esistono opzioni per rendere questo partizionamento completamente casuale o "stratificato", in modo da mantenere la proporzione degli spazi di lavoro.

3.11. Ejemplo genérico: clasificación de especies florales

Creación de modelos de datos en línea

El flujo de trabajo comienza con la división de los datos en conjuntos de entrenamiento y prueba.

- Se divide el conjunto de datos en un 70% de entrenamiento y un 30% de prueba.
- El modelo se entrena con los datos de entrenamiento.
- El modelo se aplica a los datos de prueba para evaluar su rendimiento.





; cs ¼wlls Ó ± Ös ||Ó |Ó¼¼± ±±±↑

L'esempio mostra l'interno del nodo KNIME "Learner" di un albero di classificazione.

Permette di configurare, tra l'altro, la misura di qualità degli attributi, se ha un pruning o meno, ecc.



Dialog - 0:10 - Decision Tree Learner (Entrena modelo)

Options PMMLSettings Flow Variables Memory Policy

General

- Class column: class
- Quality measure: Gini index
- Pruning method: No pruning
- Reduced Error Pruning
- Min number records per node:
- Number records to store for view:
- Average split point
- Number threads:
- Skip nominal columns without domain information

Root split

- Force root split column
- Root split column: petal width

Binary nominal splits

- Binary nominal splits
- Max #nominal:
- Filter invalid attribute values in child nodes

OK Apply Cancel ?



3.12. Esempio generico: classificazione delle specie floreali

Visualizzazione dei risultati

È possibile utilizzare un nodo interattivo di tipo tabella per visualizzare i valori della classe effettiva e della classe prevista, per tutti gli esempi di test. In questo modo è possibile vedere gli esempi mal classificati.

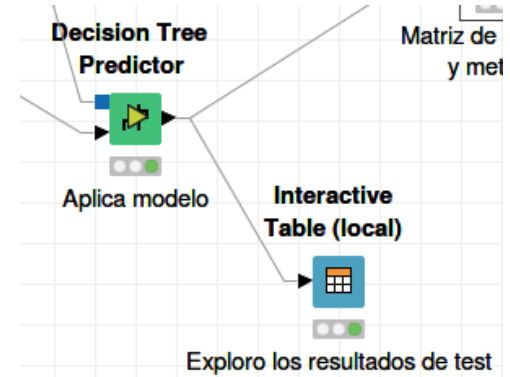


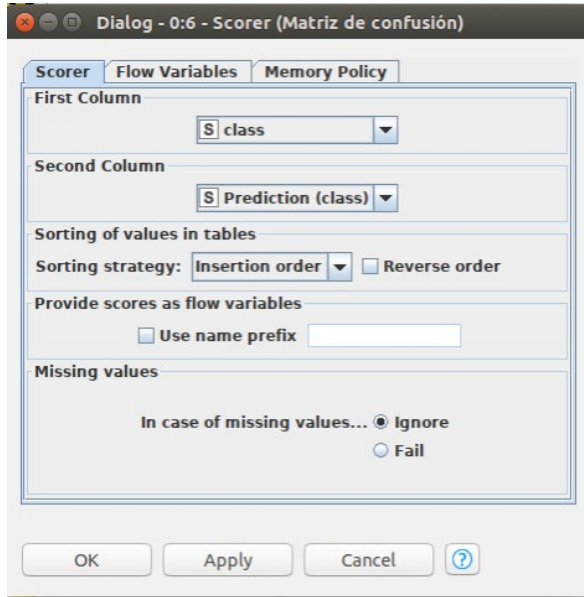
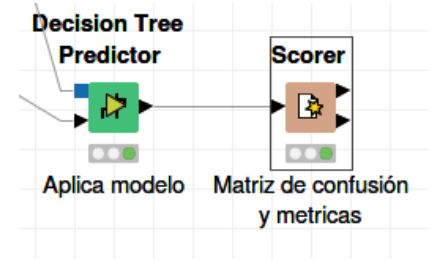
Table View - 0:7 - Interactive Table (local) (Exploro los resultados de test)						
File Hilite Navigation View Output						
Row ID	D sepal l...	D sepal ...	D petal l...	D petal ...	S class	S Prediction (class)
Row0	5.1	3.5	1.4	0.2	Iris-setosa	Iris-setosa
Row6	4.6	3.4	1.4	0.3	Iris-setosa	Iris-setosa
Row16	5.4	3.9	1.3	0.4	Iris-setosa	Iris-setosa
Row17	5.1	3.5	1.4	0.3	Iris-setosa	Iris-setosa
Row18	5.7	3.8	1.7	0.3	Iris-setosa	Iris-setosa
Row24	4.8	3.4	1.9	0.2	Iris-setosa	Iris-setosa
Row27	5.2	3.5	1.5	0.2	Iris-setosa	Iris-setosa
Row28	5.2	3.4	1.4	0.2	Iris-setosa	Iris-setosa
Row30	4.8	3.1	1.6	0.2	Iris-setosa	Iris-setosa
Row32	5.2	4.1	1.5	0.1	Iris-setosa	Iris-setosa
Row37	4.9	3.6	1.4	0.1	Iris-setosa	Iris-setosa
Row39	5.1	3.4	1.5	0.2	Iris-setosa	Iris-setosa
Row40	5	3.5	1.3	0.3	Iris-setosa	Iris-setosa
Row44	5.1	3.8	1.9	0.4	Iris-setosa	Iris-setosa
Row46	5.1	3.8	1.6	0.2	Iris-setosa	Iris-setosa
Row52	6.9	3.1	4.9	1.5	Iris-versicolor	Iris-virginica
Row59	5.2	2.7	3.9	1.4	Iris-versicolor	Iris-versicolor
Row60	5	2	3.5	1	Iris-versicolor	Iris-versicolor



Creazione di modelli di data mining



Il processo di creazione di modelli di data mining inizia con la raccolta e la pulizia dei dati. Successivamente, si procede alla selezione delle variabili e alla scelta del modello di machine learning appropriato. Infine, il modello viene valutato e ottimizzato per garantire la massima accuratezza e performance.



Per configurare il nodo "Scorer" dobbiamo definire quale sia la colonna che rappresenta la classe reale e quale quella della classe prevista dal modello.



Visualizzazione dei risultati

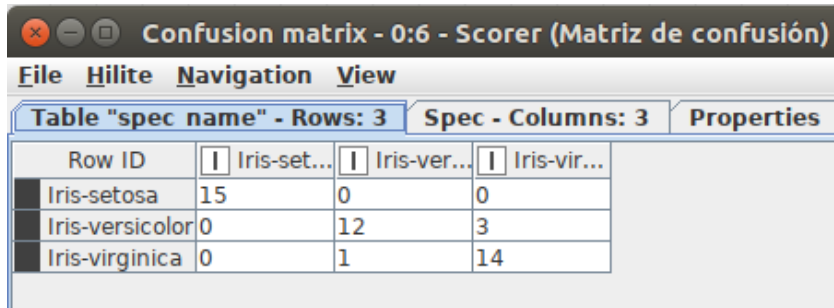
Nello stesso nodo si può ottenere anche la matrice di confusione, che mette in relazione la classe reale con la classe prevista, per osservare i tipi di errori nel modello.

In questa tabella:

- righe: sono le classi reali.
- colonne: sono le classi previste.

Nell'esempio:

- 3 esempi di iris-versicolor sono stati classificati erroneamente come iris-virginia.
- 1 esempio di iris-virginia è stato erroneamente classificato come iris-versicolor.



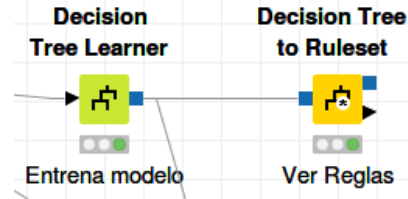
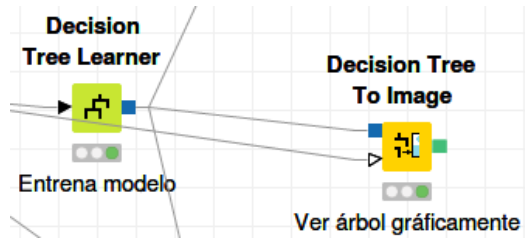
Row ID	Iris-set...	Iris-ver...	Iris-vir...
Iris-setosa	15	0	0
Iris-versicolor	0	12	3
Iris-virginica	0	1	14



3.14 . Ejemplo generico: clasificacion de las especies floreali

Visualizaci3n del modelo

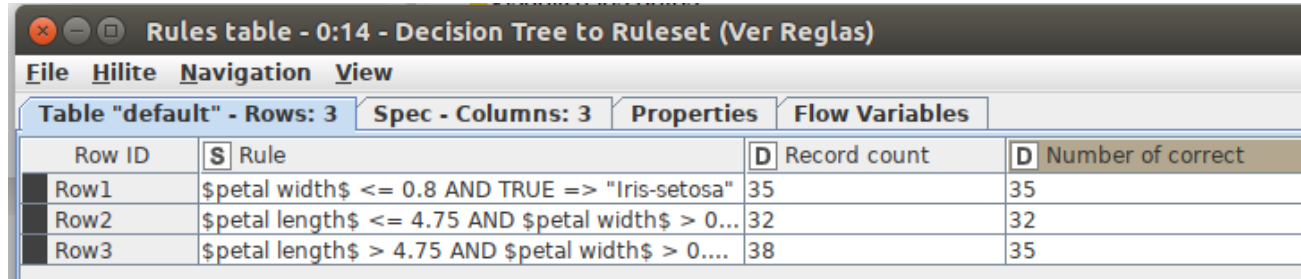
El modelo de aprendizaje de un 3rbol de decisi3n puede ser visualizado de diferentes maneras. En la imagen superior se muestra un ejemplo de un 3rbol de decisi3n con nodos y ramas, representando la estructura del modelo. El 3rbol muestra una jerarqu3a de decisiones basadas en variables de entrada, que finalmente conducen a una clasificaci3n de las especies floreali.



Visualizzazione del modello



Un albero più grande sarebbe poco pratico da visualizzare graficamente, quindi può essere tradotto in un insieme di regole, in una rappresentazione più compatta.



Row ID	Rule	Record count	Number of correct
Row1	\$petal width\$ <= 0.8 AND TRUE => "Iris-setosa"	35	35
Row2	\$petal length\$ <= 4.75 AND \$petal width\$ > 0...	32	32
Row3	\$petal length\$ > 4.75 AND \$petal width\$ > 0...	38	35

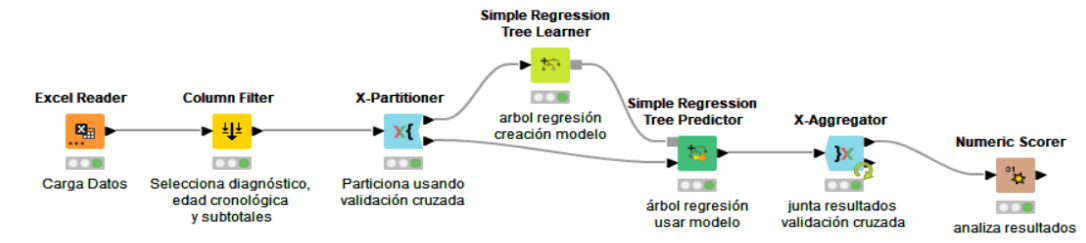


* [Illegible text]

- È necessario scaricare il file "eEarlyCare.knwf".
- In KNIME Explorer fare clic con il tasto destro del mouse e poi "Importa flusso di lavoro KNIME ...".
- Quindi l'opzione "Seleziona file" → "Sfoglia".
- Selezionarlo e cliccare su "Ok".



Si tratta di un flusso di lavoro che utilizza gli item della scala eEarlyCare, l'età cronologica e il sesso come variabili indipendenti e la diagnosi principale come variabile dipendente. Esplora il flusso di lavoro.



Ejemplo de regresión con los datos de eEarlyCare





Riferimenti dal web

<https://www.knime.com>