# Lecture 4
# How to Conduct Data Understanding Phase?

Kim, Yang Sok

Dept. of MIS, Keimyung University

- **Introduction**
- **Tasks & Outputs**
  - Gathering data
  - Describing Data
  - Exploring Data
  - Verifying Data Quality
- **Exercises**
  - Exercise 1: Write a data description report
  - Exercise 2: Loading Dataset with Rapidminer
  - Exercise 3: Exploring Data with Rapidminer
- **Conclusion**

# Introduction

- **In the Data Understanding phase of CRISP-DM, you obtain data and verify that it is appropriate for your needs.**
  - You might identify issues that cause you to return to business understanding and revise your plan.
  - You may even discover flaws in your business understanding, another reason to rethink goals and plans.

- **The Data Understanding phase includes four tasks. These are**
  - Gathering data
  - Describing data
  - Exploring data
  - Verifying data quality

# Tasks & Outputs

- **The Data Understanding phase requires you to acquire the data listed in the project resources.**

- **This initial collection includes <span style="color:red">data loading</span>, if this is necessary for data understanding.**

  - For example, if you use a specific tool for data understanding, it makes perfect sense to load your data into this tool. <u>In this lecture, we will use Rapidminer, so it is necessary load data with Rapidminer in this phase.</u>

  - If you acquire **<span style="color:red">multiple data sources</span>** then you need to consider how and when you're going to **<span style="color:red">integrate these</span>**.

- **Just one deliverable exists for this task: the <u>initial data collection report</u>.**

- **You may do a lot of work to assemble the data you need before you can write this report. First, you will make your plan, as follows:**

  - **Outline data requirements**: Create a list of the types of data necessary to address the data mining goals. Expand the list with details such as the required time range and data formats.

  - **Verify data availability**: Confirm that the required data exists, and that you can use it. If some of the data you want is unavailable, decide how you will address that issue. Consider alternatives such as

    - Substituting with an alternative data source

    - Narrowing the scope of the project

    - Gathering new data

  - **Define selection criteria**: Identify the specific data sources (databases, files, documents, and so on.) you will use. Within those sources, specify the tables, fields, and case ranges that are relevant to this project.

**Tasks & Outputs**

Gathering Data - Outputs

- **Once you've gone through these steps, you must actually obtain the data.**

- **At this stage, import the data into the data-mining platform you'll be using for the project to confirm that it is possible to do so and that you understand the process.**

- **In the course of this trial you may discover software (or hardware) limitations you had not anticipated, such as**
  - Limits on the number of cases or fields, or on the amount of memory you may use
  - Inability to read the data formats of your sources
  - Difficulty dealing with imperfections in the data (for example, you might encounter products that won't import or analyze incomplete datasets)

**Tasks & Outputs**

**Gathering Data - Outputs**

- **Finally, summarize the gathering process in a report.**

- **The report should describe your requirements, and explain in some detail exactly what data you have gathered and from what sources.**

- **Here you confirm that you have actually obtained the data and that it is compatible with your data-mining platform.**
  - If you have run into difficulties, you'll explain what they were and how you have addressed them (using alternative sources, revising plans, changing formats).

**Tasks & Outputs**

**Describing Data - Tasks**

- **Now that you have data, prepare a general description of what you have.**

- **Examine the "gross" or "surface" properties of the acquired data and report on the results.**

**Task & Outputs**

**Describing Data - Outputs**

- **Data description report**
  - Describe the data that has been acquired including its format, its quantity (for example, the number of records and fields in each table), the identities of the fields and any other surface features which have been discovered.

  - Evaluate whether the data acquired satisfies your requirements.

**Tasks & Outputs**

Exploring Data - Tasks

- **During this stage you'll address data mining questions using querying, data visualization and reporting techniques. These may include:**
  - Distribution of key attributes (for example, the target attribute of a prediction task)
  - Relationships between pairs or small numbers of attributes
  - Results of simple aggregations
  - Properties of significant sub-populations
  - Simple statistical analyses

- **These analyses may directly address your data mining goals. They may also contribute to or refine the data description and quality reports, and feed into the transformation and other data preparation steps needed for further analysis.**

**Tasks & Outputs**

Exploring Data - Outputs

- The deliverable for this task is the <span style="color:red">data exploration report.</span>

- This report describe results of your data exploration, including first findings or initial hypothesis and their impact on the remainder of the project.

- If appropriate you could include graphs and plots here to indicate data characteristics that suggest further examination of interesting data subsets.

- This report should include a more detailed description of the data than the data description report, including distributions, summaries, and any signs of data quality problems.

- **You have the data and you've examined it, and now you have to determine whether it's good enough to support your goals.**

- **You will often have some quality problem to address yet still be able to move forward, but at times the data quality is so poor that it cannot support your plan and you'll have to look for alternatives.**

- **Some of the worst data problems would include**
  - The data you need doesn't exist. (Did it never exist, or was it discarded? Can this data be collected and saved for future use?)
  - It exists, but you can't have it. (Can this restriction be overcome?)
  - You find severe **data quality issues**.
    - Is the data complete (does it cover all the cases required)?
    - Is it correct, or does it contain errors and, if there are errors, how common are they?
    - Are there missing values in the data? If so, how are they represented, where do they occur, and how common are they?

**Tasks & Outputs**

Verifying Data Quality - Outputs

- **The deliverable for this task is the <span style="color:red">data quality report</span>.**

- **This summarizes the data that you have, minor and major quality issues that you have found, and possible remedies for quality problems or alternatives (such as using an alternative data resource).**

- **If you are facing any really serious data quality issues and can't identify an adequate solution, you may have to recommend reconsidering goals or plans.**

# Exercise 1: Write a data description report

**Exercise - Write a Data Description Report**

**Template**

| | |
|---|---|
| Name | Name of dataset |
| Contributor | Dataset creator |
| Business Objective | Explain business objectives |
| **Data Mining Objectives** | **Explain data mining objectives** |
| **Data Mining Tasks** | **Prediction, classification, clustering, etc.** |
| **Performance Measures** | **Accuracy, Precision, Recall, F-Measure, RMSE, MAE, R2, …** |
| Dataset Characteristics | Multivariate, univariate, sequential, time-series, text, other |
| Number of Examples | Number of examples collected |
| Area | Life Sciences, Physical Sciences, CS / Engineering, Social Sciences, Business, Game, Other |
| Attribute Characteristics | Numerical, categorical, text, mixed |
| Number of Attributes | Number of attributes |
| Missing Values? | Yes/No |
| Version | 1.0, … |
| Date of Creation | 2019-08-05 |

| Attribute Definition | | | | |
|---|---|---|---|---|
| Name | Value Type | Role | Description | Allowable Values |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

These three items need knowledge on data mining. So now you may skip these three items

17

- **Univariate data**
  - This type of data consists of only one variable.
  - The analysis of univariate data is thus the simplest form of analysis since the information deals with only one quantity that changes.
  - It does not deal with causes or relationships and the main purpose of the analysis is to describe the data and find patterns that exist within it.
  - The example of a univariate data can be height.

- **Bivariate data**
  - This type of data involves two different variables.
  - The analysis of this type of data deals with causes and relationships and the analysis is done to find out the relationship among the two variables.
  - Example of bivariate data can be temperature and ice cream sales in summer season.

- **Multivariate data**
  - When the data involves three or more variables, it is categorized under multivariate.
  - Example of this type of data is suppose an advertiser wants to compare the popularity of four advertisements on a website, then their click rates could be measured for both men and women and relationships between variables can then be examined.

- ## Sequential Data:
  - The order of the data matters, but the time stamp is irrelevant or it doesn't matter. (Example: DNA sequence. As you see the concept of time is irrelevant, so the order is not temporal.)

- ## Temporal Sequence:
  - In addition to the order of data, the time stamp also matters. (Example: Data collected from customers' shopping behavior, considering their transaction time stamp as the temporal dimension.)

- ## Time Series:
  - The data is in order, with a fixed time-difference between occurrence of successive data points. (Example: Time series of the temperature of a surface being recorded every 120 seconds.)

**Exercise - Write a Data Description Report**

**Template – Attribute Definition**

- **Name is a unique identifier, so names <span style="color:red">should not duplicate.</span>**
- **Value type can be classified into four types**

Why I need to concern data value types?

```
                          Data
      ┌────────────────────┼──────────────┬─────────────┐
  Categorical          Numerical        Text       Date-Time
   ┌─────┴─────┐      ┌─────┴─────┐
 Nominal    Ordinal  Discrete  Continuous
 ┌───┴────┐
Binominal  Polynomial
```

- **Attribute is one of following roles: (general) attribute, label, id, weight, batch, cluster, prediction, outlier, cost, base value**
  - (general) attribute is an attribute to describe the example.
    - When data set load by read operators, Rapidminer regards all attributes as this type.
  - id is an attribute used to identify individual example.
  - label is an attribute that should be predicted by the model.
    - For predicting problems (classification and regression), an example must have a label attribute.

- **Allowable Values are defined by min and max values for the numerical attribute value type and by a list of values for the categorical attribute value type.**

- **Go to https://archive.ics.uci.edu/ml/datasets/wine+quality**



- **Use information in the web page, the paper, and the dataset in order to write the data description report for the Wine Quality Prediction project.**
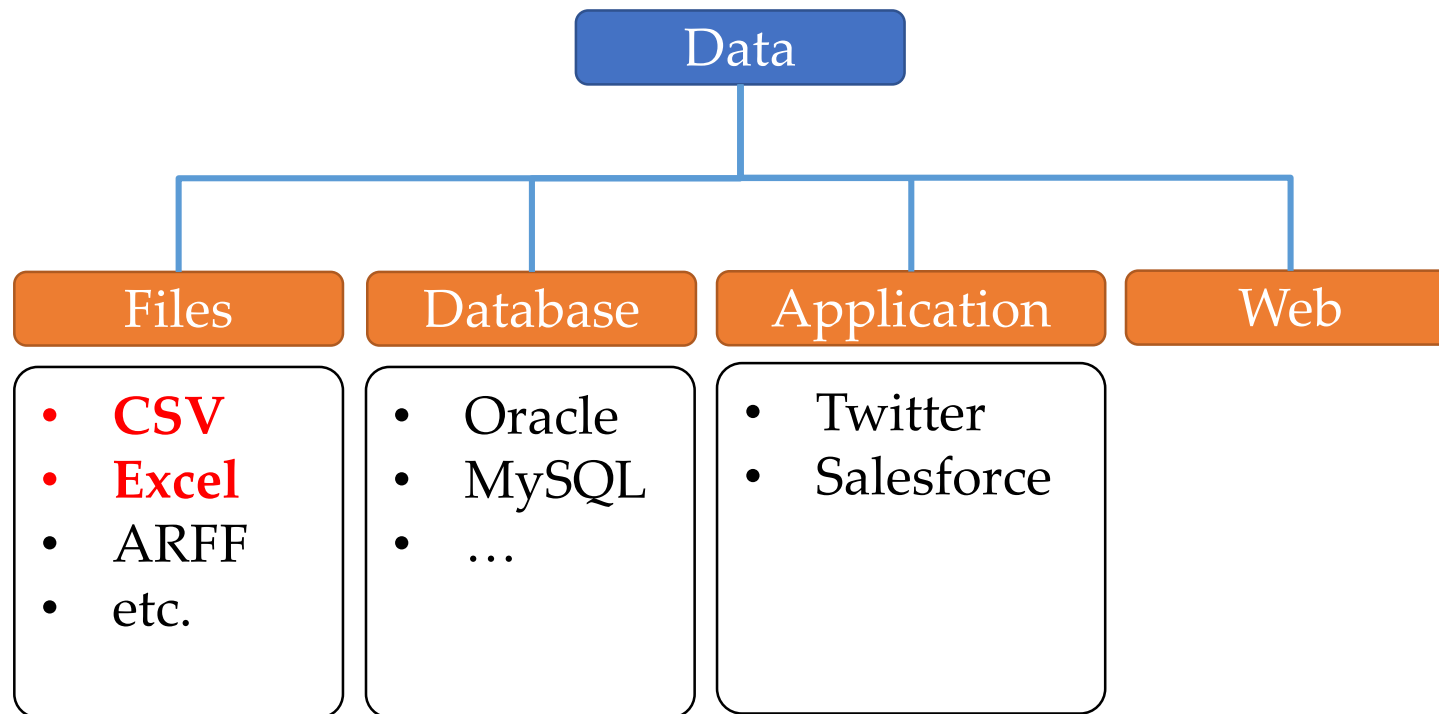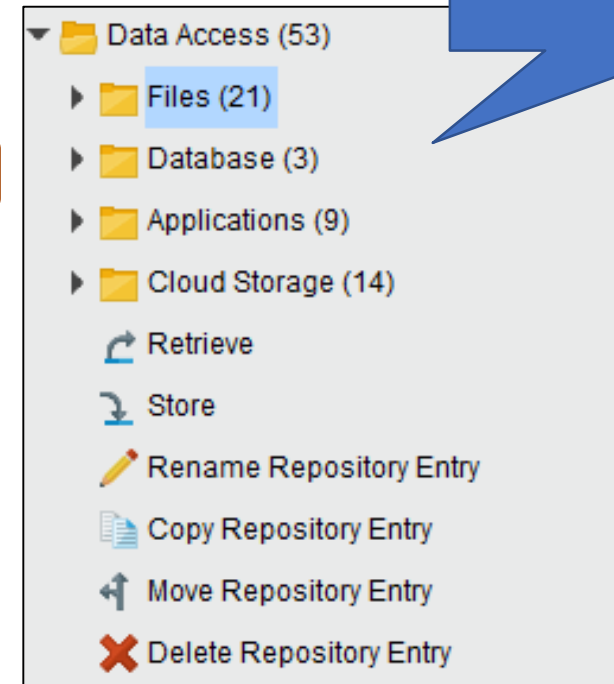
# Exercise 2: Loading Dataset with Rapidminer

· **Data can be stored into different sources with different formats:**

Rapidminer provides various operators to access different data sources



Data
- Files
  - **CSV**
  - **Excel**
  - ARFF
  - etc.
- Database
  - Oracle
  - MySQL
  - …
- Application
  - Twitter
  - Salesforce
- Web

Data Access (53)
- Files (21)
- Database (3)
- Applications (9)
- Cloud Storage (14)
- Retrieve
- Store
- Rename Repository Entry
- Copy Repository Entry
- Move Repository Entry
- Delete Repository Entry

**Exercise: Working with Rapidminer to Load Dataset**

**Download Wine Datasets**

- **Go to https://archive.ics.uci.edu/ml/datasets/wine+quality**



- **Click "Data Folder" link**

## Index of /ml/machine-learning-databases/wine-quality

- Parent Directory
- winequality-red.csv
- winequality-white.csv
- winequality.names

Create "Dataset" folder in "Desktop" and save these two dataset into "Dataset" folder

**Exercise: Working with Rapidminer to Load Dataset**

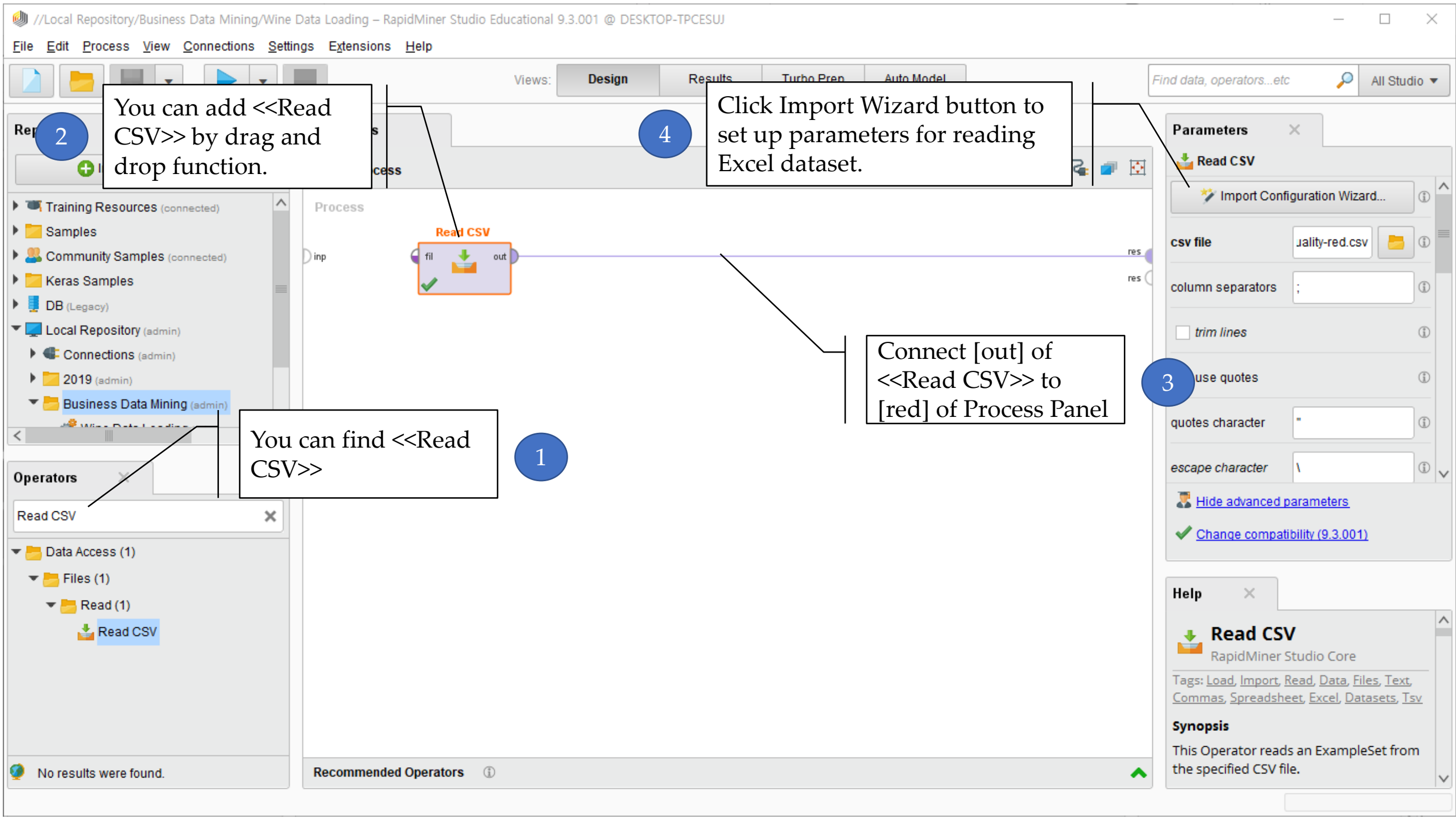**Load Dataset Using <<Read Excel>>**

- **Tasks**
  - Read "winequality-red.csv" stored in CSV format

- **Steps**
  1. Find <<Read CSV>> operator from Operator Panel
  2. Add <<Read CSV>> to Process Panel
  3. Set up parameters of <<Read CSV>>
  4. Connect [exa] port of <<Read CSV>> to [res] port of Process Panel
  5. Execute the analysis process
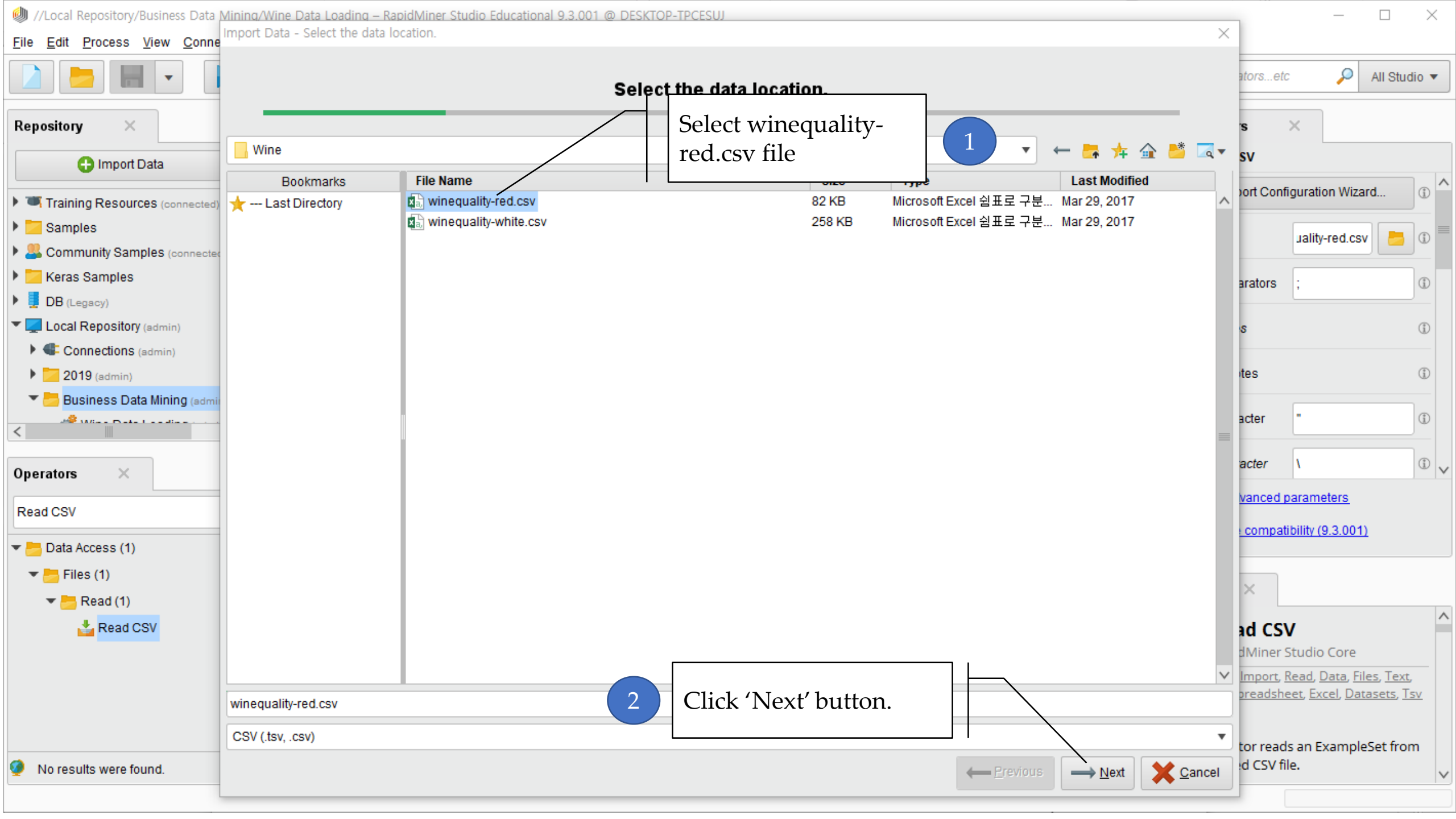  6. Save the process in to 'Exercise' repository

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   **Design**   Results   Turbo Prep   Auto Model

Find data, operators...etc   🔍   All Studio ▼

**You can add <<Read CSV>> by drag and drop function.**

(2)

**Click Import Wizard button to set up parameters for reading Excel dataset.**

(4)

**Parameters**   ✕

⬇ Read CSV

⚡ Import Configuration Wizard...   ⓘ

csv file   ᴜality-red.csv   📁   ⓘ

column separators   ;   ⓘ

☐ *trim lines*   ⓘ

**Process**

Rep...

▶ 🦉 Training Resources (connected)
▶ 📁 Samples
▶ 👥 Community Samples (connected)
▶ 📁 Keras Samples
▶ 🗄 DB (Legacy)
▼ 🖥 Local Repository (admin)
   ▶ 🔌 Connections (admin)
   ▶ 📁 2019 (admin)
   ▼ 📁 Business Data Mining (admin)
       ⚙ Wine Data Loading

Read CSV

inp   ◗ fil   ⬇ out ◖   res
                    ✓           res

(1)

**You can find <<Read CSV>>**

**Connect [out] of <<Read CSV>> to [red] of Process Panel**

(3)

use quotes   ⓘ

quotes character   "   ⓘ

escape character   \   ⓘ

👥 Hide advanced parameters

✓ Change compatibility (9.3.001)

**Operators**   ✕

Read CSV   ✕

▼ 📁 Data Access (1)
   ▼ 📁 Files (1)
       ▼ 📁 Read (1)
           ⬇ Read CSV

🌐 No results were found.

**Recommended Operators**   ⓘ   ⌃

**Help**   ✕

⬇ **Read CSV**
RapidMiner Studio Core

Tags: Load, Import, Read, Data, Files, Text, Commas, Spreadsheet, Excel, Datasets, Tsv

**Synopsis**

This Operator reads an ExampleSet from the specified CSV file.

Import Data - Select the data location.

Select the data location.

**Select winequality-red.csv file** ①

| Wine | | | | |
|---|---|---|---|---|

File · Edit · Process · View · Conne

Repository

Import Data

Training Resources (connected)
Samples
Community Samples (connected)
Keras Samples
DB (Legacy)
Local Repository (admin)
Connections (admin)
2019 (admin)
Business Data Mining (admin)

| Bookmarks | File Name | Size | Type | Last Modified |
|---|---|---|---|---|
| --- Last Directory | winequality-red.csv | 82 KB | Microsoft Excel 쉼표로 구분... | Mar 29, 2017 |
| | winequality-white.csv | 258 KB | Microsoft Excel 쉼표로 구분... | Mar 29, 2017 |

Operators

Read CSV

Data Access (1)
Files (1)
Read (1)
Read CSV

**Click 'Next' button.** ②

winequality-red.csv

CSV (.tsv, .csv)

No results were found.

Previous · Next · Cancel

port Configuration Wizard...

uality-red.csv

arators

ites

acter

acter

vanced parameters

compatibility (9.3.001)

ad CSV

dMiner Studio Core

Import, Read, Data, Files, Text,
preadsheet, Excel, Datasets, Tsv

tor reads an ExampleSet from
d CSV file.

File   Edit   Process   View   Conne

**Import Data - Specify your data format**

# Specify your data format

| | Header Row | 1 | File Encoding | x-windows-949 ▾ | | Use Quotes | " |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Start Row | 1 | Escape Character | \ | | Trim Lines | |
| | Column Separator | Semicolon ";" ▾ | Decimal Character | . | | Skip Comments | # |

| 1 | fixed aci... | volatile a... | citric acid | residual ... | chlorides | free sulf... | total sulf... | density | pH | sulphates | alcohol | quality |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 2 | 7.4 | 0.7 | 0 | 1.9 | 0.076 | 11 | 34 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 |
| 3 | 7.8 | 0.88 | 0 | 2.6 | 0.098 | 25 | | | | | 9.8 | 5 |
| 4 | 7.8 | 0.76 | 0.04 | 2.3 | 0.092 | 15 | | | | | 9.8 | 5 |
| 5 | 11.2 | 0.28 | 0.56 | 1.9 | 0.075 | 17 | | | | | 9.8 | 6 |
| 6 | 7.4 | 0.7 | 0 | 1.9 | 0.076 | 11 | | | | | 9.4 | 5 |
| 7 | 7.4 | 0.66 | 0 | 1.8 | 0.075 | 13 | | | | | 9.4 | 5 |
| 8 | 7.9 | 0.6 | 0.06 | 1.6 | 0.069 | 15 | | 0.9964 | 3.3 | 0.46 | 9.4 | 5 |
| 9 | 7.3 | 0.65 | 0 | 1.2 | 0.065 | 15 | 21 | 0.9946 | 3.39 | 0.47 | 10 | 7 |
| 10 | 7.8 | 0.58 | 0.02 | 2 | 0.073 | 9 | 18 | 0.9968 | 3.36 | 0.57 | 9.5 | 7 |
| 11 | 7.5 | 0.5 | 0.36 | 6.1 | 0.071 | 17 | 102 | 0.9978 | 3.35 | 0.8 | 10.5 | 5 |
| 12 | 6.7 | 0.58 | 0.08 | 1.8 | 0.097 | 15 | 65 | 0.9959 | 3.28 | 0.54 | 9.2 | 5 |
| 13 | 7.5 | 0.5 | 0.36 | 6.1 | 0.071 | 17 | 102 | 0.9978 | 3.35 | 0.8 | 10.5 | 5 |
| 14 | 5.6 | 0.615 | 0 | 1.6 | 0.089 | 16 | 59 | 0.9943 | 3.58 | 0.52 | 9.9 | 5 |
| 15 | 7.8 | 0.61 | 0.29 | 1.6 | 0.114 | 9 | 29 | 0.9974 | 3.26 | 1.56 | 9.1 | 5 |
| 16 | 8.9 | 0.62 | 0.18 | 3.8 | 0.176 | | | | | 0.88 | 9.2 | 5 |
| 17 | 8.9 | 0.62 | 0.19 | 3.9 | 0.17 | | | | | 0.93 | 9.2 | 5 |

System automatically setup parameters for data format. You can change them if necessary.

**1**

**2**   Click 'Next' button.

✓ no problems.

Previous   →Next   ✕ Cancel

SV

port Configuration Wizard...

uality-red.csv

arators   ;

vanced parameters

compatibility (9.3.001)

**ad CSV**
dMiner Studio Core
Import, Read, Data, Files, Text,
preadsheet, Excel, Datasets, Tsv

tor reads an ExampleSet from
ed CSV file.

**Repository**

Import Data

▸ Training Resources (connected)
▸ Samples
▸ Community Samples (connected)
▸ Keras Samples
▸ DB (Legacy)
▾ Local Repository (admin)
  ▸ Connections (admin)
  ▸ 2019 (admin)
  ▾ Business Data Mining (admin)

**Operators**

Read CSV

▾ Data Access (1)
  ▾ Files (1)
    ▾ Read (1)
      Read CSV

No results were found.

File   Edit   Process   View   Conne

Repository

Import Data

- Training Resources (connected)
- Samples
- Community Samples (connected)
- Keras Samples
- DB (Legacy)
- Local Repository (admin)
  - Connections (admin)
  - 2019 (admin)
  - Business Data Mining (admin)
    - Wine Data Loading

Operators

Read CSV

- Data Access (1)
  - Files (1)
    - Read (1)
      - Read CSV

No results were found.

**Import Data - Format your columns.**

# Format your columns.

Date format  [ Enter value... ▼ ]        ☐ Replace errors with missing values ⓘ

| fixed acidity ⚙ ▼ | volatile aci... ⚙ ▼ | citric acid ⚙ ▼ | residual su... ⚙ ▼ | chlorides ⚙ ▼ | free sulfur ... ⚙ ▼ | total sulfur ... ⚙ ▼ | density |
|---|---|---|---|---|---|---|---|
| real | real | real | real | real | integer | integer | real |
| 1 | 7.400 | 0.700 | 0.000 | 1.900 | 0.076 | 11 | 34 | 0.998 |
| 2 | 7.800 | 0.880 | 0.000 | 2.600 | 0.098 | 25 | 67 | 0.997 |
| 3 | 7.800 | 0.760 | 0.040 | 2.300 | 0.092 | 15 | 54 | 0.997 |
| 4 | 11.200 | 0.280 | 0.560 | 1.900 | | | | 0.998 |
| 5 | 7.400 | 0.700 | 0.000 | 1.900 | | | | 0.998 |
| 6 | 7.400 | 0.660 | 0.000 | 1.800 | | | | 0.998 |
| 7 | 7.900 | 0.600 | 0.060 | 1.600 | | | | 0.996 |
| 8 | 7.300 | 0.650 | 0.000 | 1.200 | 0.065 | 15 | 21 | 0.995 |
| 9 | 7.800 | 0.580 | 0.020 | 2.000 | 0.073 | 9 | 18 | 0.997 |
| 10 | 7.500 | 0.500 | 0.360 | 6.100 | 0.071 | 17 | 102 | 0.998 |
| 11 | 6.700 | 0.580 | 0.080 | 1.800 | 0.097 | 15 | 65 | 0.996 |
| 12 | 7.500 | 0.500 | 0.360 | 6.100 | 0.071 | 17 | 102 | 0.998 |
| 13 | 5.600 | 0.615 | 0.000 | 1.600 | 0.089 | 16 | 59 | 0.994 |
| 14 | 7.800 | 0.610 | 0.290 | 1.600 | 0.114 | 9 | 29 | 0.997 |
| 15 | 8.900 | 0.620 | 0.180 | 3.800 | 0.176 | 52 | 145 | 0.999 |
| 16 | 8.900 | 0.620 | 0.190 | 3.900 | 0.170 | 51 | 148 | 0.999 |
| 17 | 8.500 | 0.280 | 0.560 | 1.800 | | | 103 | 0.997 |

**① System automatically setup attributes' value types and role. But you can change them manually.**

**② Click 'Finish' button.**

✓ no problems.

← Previous     ⚑ Finish     ✖ Cancel

File  Edit  Process  View  Connections  Settings  Extensions  Help

Views:  **Design**  Results  Turbo Prep  Auto Model

Find data, operators...etc    All Studio ▼

**Repository** ✕

⊕ **Import Data**    ≡ ▼

▶ 🖳 Training Resources (connected)
▶ 📁 Samples
▶ 👥 Community Samples (connected)
▶ 📁 Keras Samples
▶ 🗄 DB (Legacy)
▼ 🖥 Local Repository (admin)
  ▶ 🔌 Connections (admin)
  ▶ 📁 2019 (admin)
  ▼ 📁 Business Data Mining (admin)
    👤 Wine Data Loading

**Operators** ✕

Read CSV ✕

▼ 📁 Data Access (1)
  ▼ 📁 Files (1)
    ▼ 📁 Read (1)
      ⬇ Read CSV

🌐 No results were found.

**Pro**    **1**

Click "Run" icon to execute the analysis process.

100% 🔍 🔍 🔍

**Process**

) inp

**Read CSV**

fil  ⬇  out

res
res

**Recommended Operators** ⓘ

**Parameters** ✕

⬇ Read CSV

✨ Import Configuration Wizard... ⓘ

csv file        uality-red.csv  📁  ⓘ

column separators    ;    ⓘ

☐ trim lines    ⓘ

☑ use quotes    ⓘ

quotes character    "    ⓘ

escape character    \    ⓘ

👤 Hide advanced parameters

✔ Change compatibility (9.3.001)

**Help** ✕

⬇ **Read CSV**
   RapidMiner Studio Core

Tags: Load, Import, Read, Data, Files, Text, Commas, Spreadsheet, Excel, Datasets, Tsv

**Synopsis**

This Operator reads an ExampleSet from the specified CSV file.

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   Design   **Results**   Turbo Prep   Auto Model

Find data, operators...etc       All Studio ▼

Result History        ▦ **ExampleSet (Read CSV)**       ✕

Open in    [⚡] Turbo Prep    [🤖] Auto Model                        ter (1,599 / 1,599 examples):  all ▼

**"Result" view is activated.**  ①

**"ExampleSet" tab is activated. This tab shows examples of the loaded dataset**  ②

| Row No. | fixed acidity | volatile acidity | citric acid | residual sug... | chlorides | free sulfur d... | total sulfur d... | density | pH | sulphates | alcohol | quality |
|---------|---------------|------------------|-------------|-----------------|-----------|------------------|-------------------|---------|-------|-----------|---------|---------|
| 1 | 7.400 | 0.700 | 0 | 1.900 | 0.076 | 11 | 34 | 0.998 | 3.510 | 0.560 | 9.400 | 5 |
| 2 | 7.800 | 0.880 | | | | 25 | 67 | 0.997 | 3.200 | 0.680 | 9.800 | 5 |
| 3 | 7.800 | 0.760 | | | | 15 | 54 | 0.997 | 3.260 | 0.650 | 9.800 | 5 |
| 4 | 11.200 | 0.280 | | | | 60 | | 0.998 | 3.160 | 0.580 | 9.800 | 6 |
| 5 | 7.400 | 0.700 | | | | 11 | 34 | 0.998 | 3.510 | 0.560 | 9.400 | 5 |
| 6 | 7.400 | 0.660 | 0 | 1.800 | 0.075 | 13 | 40 | 0.998 | 3.510 | 0.560 | 9.400 | 5 |
| 7 | 7.900 | 0.600 | 0.060 | 1.600 | 0.069 | 15 | 59 | 0.996 | 3.300 | 0.460 | 9.400 | 5 |
| 8 | 7.300 | 0.650 | 0 | 1.200 | 0.065 | 15 | 21 | 0.995 | 3.390 | 0.470 | 10 | 7 |
| 9 | 7.800 | 0.580 | 0.020 | 2 | 0.073 | 9 | 18 | 0.997 | 3.360 | 0.570 | 9.500 | 7 |
| 10 | 7.500 | 0.500 | 0.360 | 6.100 | 0.071 | 17 | 102 | 0.998 | 3.350 | 0.800 | 10.500 | 5 |
| 11 | 6.700 | 0.580 | 0.080 | 1.800 | 0.097 | 15 | 65 | 0.996 | 3.280 | 0.540 | 9.200 | 5 |
| 12 | 7.500 | 0.500 | 0.360 | 6.100 | 0.071 | 17 | 102 | 0.998 | 3.350 | 0.800 | 10.500 | 5 |
| 13 | 5.600 | 0.615 | 0 | 1.600 | 0.089 | 16 | 59 | 0.994 | 3.580 | 0.520 | 9.900 | 5 |
| 14 | 7.800 | 0.610 | 0.290 | 1.600 | 0.114 | 9 | 29 | 0.997 | 3.260 | 1.560 | 9.100 | 5 |
| 15 | 8.900 | 0.620 | 0.180 | 3.800 | 0.176 | 52 | 145 | 0.999 | 3.160 | 0.880 | 9.200 | 5 |
| 16 | 8.900 | 0.620 | 0.190 | 3.900 | 0.170 | 51 | 148 | 0.999 | 3.170 | 0.930 | 9.200 | 5 |

ExampleSet (1,599 examples, 0 special attributes, 12 regular attributes)

Select "Save Process as …" menu form File menu.

1

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   **Design**   Results   Turbo Prep   Auto Model

Find data, operators...etc   All Studio ▼

**Repository**   ✕

➕ Import Data   ≡ ▼

▶ 🎓 Training Resources (connected)
▶ 📁 Samples
▶ 👥 Community Samples (connected)
▶ 📁 Keras Samples
▶ 📦 DB (Legacy)
▼ 🖥 Local Repository (admin)
  ▶ 🔌 Connections (admin)
  ▶ 📁 2019 (admin)
  ▼ 📁 Business Data Mining (admin)
    ⚙ Wine Data Loading (admin – v1, 8/3/
  ▶ 📁 data (admin)

**Operators**   ✕

Read Excel   ✕

▼ 📁 Data Access (4)
  ▼ 📁 Files (4)
    ▼ 📁 Read (3)
      📄 Read CSV
      📄 Read Excel
      📄 Read Excel with Format
    ▼ 📁 Write (1)
      📄 Write Excel
▼ 📁 Extensions (1)

**Process**

🔵 Process

100% 🔍 🔍 🔍 ⊞ 🔗 🗔 ⊡

Process

**Read CSV**

)inp   fil   📥   out

✔

You can find that new process added into this repository folder   **1**

Recommended Operators   ⓘ

**Parameters**   ✕

🔲 Process

logverbosity            init ▼   ⓘ

logfile                  [        ] 📁   ⓘ

resultfile               [        ] 📁   ⓘ

random seed             2001        ⓘ

send mail               never ▼     ⓘ

encoding                SYSTEM ▼    ⓘ

👤 Hide advanced parameters

✔ Change compatibility (9.3.001)

**Help**   ✕

🔲 **Process**
   RapidMiner Studio Core

**Synopsis**
The root operator which is the outer most operator of every process.

# Exercise 2: Changing Meta Data Information

**Exercise: Changing Meta Data Information**

**Task & Process**

- **Tasks**
  - Modify meta data information after loading dataset.

- **Steps**
  1. Select <<Read CSV>> operator
  2. Click 'Edit List' button beside 'data set meta data information' parameter
  3. Change meta data information

# Exercise 3: Exploring Data with Rapidminer

- **You can understand data by reviewing  descriptive statistics, which include the followings:**

- **Central value**
  - Numerical value
    - Average(Mean)
    - Median
    - Mode
  - Categorical
    - Least / Most  Frequent

- **Dispersion**
  - Range
  - Min & Max
  - Variance
  - Standard deviation

File  Edit  Process  View  Connections  Settings  Extensions  Help

Views:    Design    **Results**    Turbo Prep    Auto Model

Find data, operators...etc    All Studio ▼

Result History    ExampleSet (Read CSV)    ✕

| Name | Type | Missing | Statistics |
|------|------|---------|-----------|

Filter (12 / 12 attributes):    Search for Attributes

**Data**

| | | | | | |
|---|---|---|---|---|---|
| **Label** **quality** | Polynominal | 0 | Least 3 (10) | | Values ...[4 more] |

📋 Copy statistics to clipboard

⚠ Toggle   Copy the statistics of all attributes to the clipboard.

**Statistics**

| | | | Min | Max | Average |
|---|---|---|-----|-----|---------|
| **fixed acidity** | Real | 0 | 4.600 | 15.900 | 8.320 |
| **volatile acidity** | Real | 0 | 0.120 | 1.580 | |

On the table header, click right mouse button and execute "Copy statistics to clipboard" to copy statistics    ① 1

**Visualizations**

| | | | Min | Max | Average |
|---|---|---|-----|-----|---------|
| **citric acid** | Real | 0 | 0 | 1 | |
| **residual sugar** | Real | 0 | 0.900 | 15.500 | 2.539 |
| **chlorides** | Real | 0 | 0.012 | 0.611 | 0.087 |

**Annotations**

| | | | Min | Max | Average |
|---|---|---|-----|-----|---------|
| **free sulfur dioxide** | Integer | 0 | 1 | 72 | 15.875 |
| **total sulfur dioxide** | Integer | 0 | 6 | 289 | 46.468 |
| **density** | Real | 0 | 0.990 | 1.004 | 0.997 |

Showing attributes 1 - 12    Examples: 1,599   Special Attributes: 1   Regular Attributes: 11

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   Design   **Results**   Turbo Prep   Auto Model

Find data, operators...etc        All Studio ▾

Result History        **ExampleSet (Read CSV)**   ✕

Name          |—|   Type        Missing        Statistics        Filter (12 / 12 attributes):   Search for Attributes

**Data**

① Click 'Details' link.

Values

Label
**quality**          Polynominal     0

Least        Most
3 (10)       5 (681)

5 (681), 6 (638),
7 (199), 4 (53),
...[2 more]

**Statistics**

Open visualizations        Details...

② Sort table based on Nominal value

Min          Max          Average
**fixed acidity**    Real        0        4.600        15.900       8.320

**Visualizations**

③ Select data and copy (Ctrl+C) them.

🪟 Nominal values                                         ✕

**citric acid**    Re

**Annotations**

| Index | Nominal value ↑ | Absolute count | Fraction |
|-------|-----------------|----------------|----------|
| 6 | 3 | 10 | 0.006 |
| 4 | 4 | 53 | 0.033 |
| 1 | 5 | 681 | 0.426 |
| 2 | 6 | 638 | 0.399 |
| 3 | 7 | 199 | 0.124 |
| 5 | 8 | 18 | 0.011 |

**residual sugar**    Re

**chlorides**      Re

**free sulfur dioxide**    Int

                                                    ✖ Close

Showing attributes 1 - 12          Examples: 1,599   Special Attributes: 1   Regular Attributes: 11

[1] Process 04:10:20

| attribute | value type | number of missing values | value | frequency | ratio | average | stdev | min | max |
|---|---|---|---|---|---|---|---|---|---|
| quality | Polynominal | 0 | 3 | 10 | 0.006 | | | | |
| | | | 4 | 53 | 0.033 | | | | |
| | | | 5 | 681 | 0.426 | | | | |
| | | | 6 | 638 | 0.399 | | | | |
| | | | 7 | 199 | 0.124 | | | | |
| | | | 8 | 18 | 0.011 | | | | |
| fixed acidity | Real | 0 | | | | 8.320 | 1.741 | 4.600 | 15.900 |
| volatile acidity | Real | 0 | | | | 0.528 | 0.179 | 0.120 | 1.580 |
| citric acid | Real | 0 | | | | 0.271 | 0.195 | 0.000 | 1.000 |
| residual sugar | Real | 0 | | | | 2.539 | 1.410 | 0.900 | 15.500 |
| chlorides | Real | 0 | | | | 0.087 | 0.047 | 0.012 | 0.611 |
| free sulfur dioxide | Integer | 0 | | | | 15.875 | 10.460 | 1.000 | 72.000 |
| total sulfur dioxide | Integer | 0 | | | | 46.468 | 32.895 | 6.000 | 289.000 |
| density | Real | 0 | | | | 0.997 | 0.002 | 0.990 | 1.004 |
| pH | Real | 0 | | | | 3.311 | 0.154 | 2.740 | 4.010 |
| sulphates | Real | 0 | | | | 0.658 | 0.170 | 0.330 | 2.000 |
| alcohol | Real | 0 | | | | 10.423 | 1.066 | 8.400 | 14.900 |

Update table using copied data

1

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   Design   **Results**   Turbo Prep   Auto Model

Find data, operators...etc   All Studio ▼

Result History   **ExampleSet (Read CSV)**   ✕

| Name | | Type | Missing | Statistics | | | |
|---|---|---|---|---|---|---|---|

**Data**

Click "Details.." link to see detailed example distribution for each value

Filter (12 / 12 attributes):   Search for Attributes   ▼

**Statistics**

Values

| Label ∧ **quality** | Polynominal | 0 | [histogram: 5, 6, 7, 4, 8] Open visualizations | Least 3 (10) | Most 5 (681) | 5 (681), 6 (638), 7 (199), 4 (53), ...[2 more] Details... |

**Visualizations**

| ∨ **fixed acidity** | Real | 0 | Min 4.600 | Max 15.900 | Average 8.320 |

| ∨ **volatile acidity** | | 0 | Min 0.120 | Max 1.580 | Average 0.528 |

Click 'quality' row, then you can see "Open visualization" and "Details" links   ①

**Annotations**

| ∨ **citric acid** | | 0 | M 0 | | |

Click "Open visualization" to see chart   ②

| ∨ **residual sugar** | Real | 0 | Min 0.900 | Max 15.500 | Average 2.539 |

| ∨ **chlorides** | Real | 0 | Min 0.012 | Max 0.611 | Average 0.087 |

| ∨ **free sulfur dioxide** | Integer | 0 | Min 1 | Max 72 | Average 15.875 |

| ∨ **total sulfur dioxide** | Integer | 0 | Min 6 | Max 289 | Average 46.468 |

③

Showing attributes 1 - 12

Examples: 1,599   Special Attributes: 1   Regular Attributes: 11

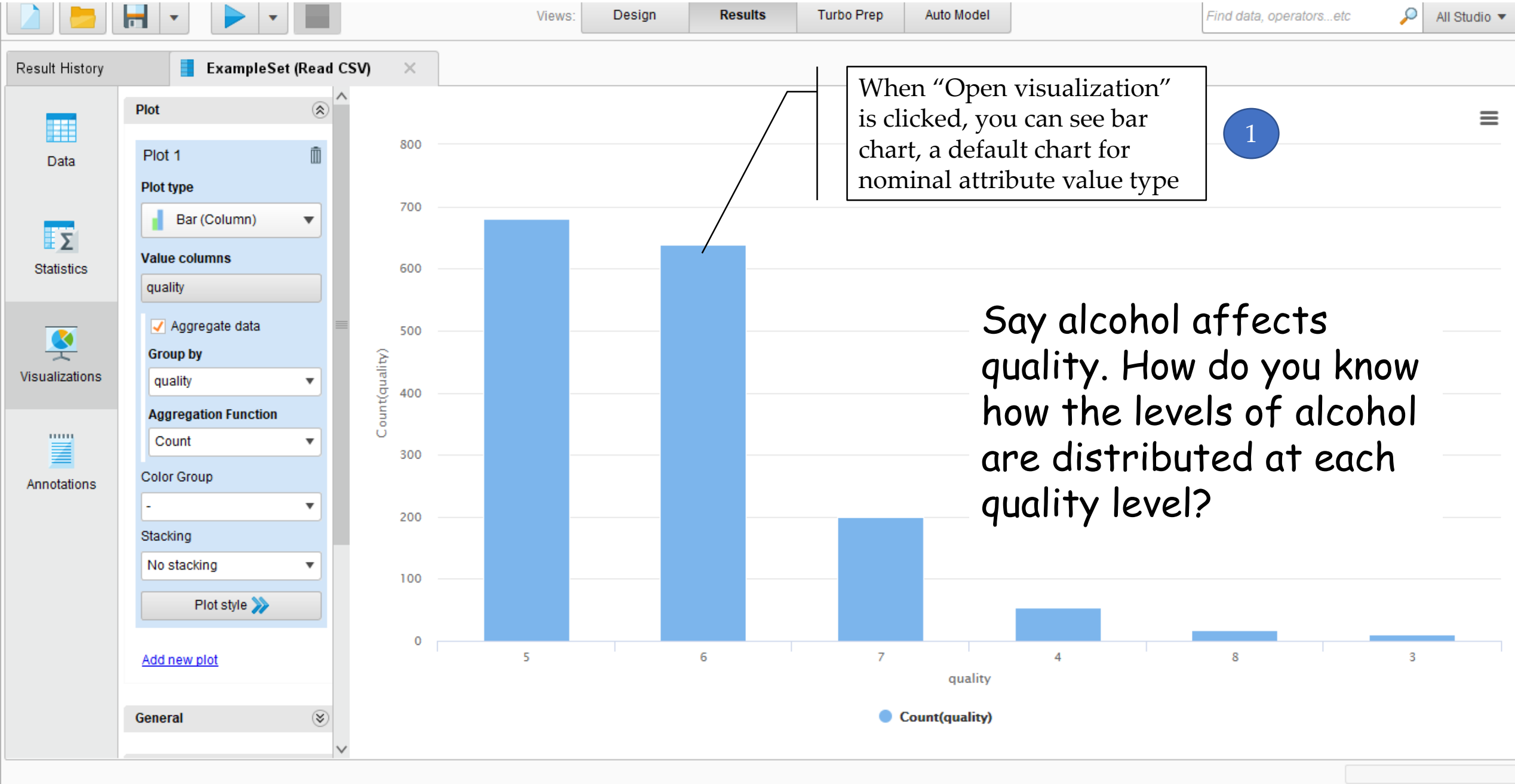**Exercise: Exploring Data with Rapidminer**

**Visualize Data with Charts**

- **You can use charts to easily identify patterns within data.**

- **Many charts are available, but we only focus on the following four charts**

|  | **Single Attribute** | **Multiple Attributes** |
|---|---|---|
| **Categorical Values** | Bar Charts | Stacked Bar chart |
| **Numerical Values** | Histogram | Scatter Plot |

# Visualize single nominal attribute with a bar chart



When "Open visualization" is clicked, you can see bar chart, a default chart for nominal attribute value type

Say alcohol affects quality. How do you know how the levels of alcohol are distributed at each quality level?

# Visualize single nominal attribute with a bar chart

# Visualize single numerical attribute with a histogram



As we know distribution of 'alcohol' using histogram, we want to convert this into categorical one.

# Use <<Discretize by Binning>> to generate categorized attribute.



Add <<Discretize by Binning>>.

Find <<Discretize by Binning>> to split data into bins.

Set 'attribute filter type' as 'single' and attribute is 'alcohol'

Set 'number of bins' as '5'.

Set 'range name type' as 'interval'.

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   Design   **Results**   Turbo Prep   Auto Model

Find data, operators...etc        All Studio ▼

Result History   |   ▦ **ExampleSet (Discretize)**   ✕

Open in   ▦ Turbo Prep   🤖 Auto Model

Filter (1,599 / 1,599 examples):   all ▼

You can find alcohol has been changed into a categorical format

| Row No. | quality | alcohol | fixed acidity | volatile acidity | citric acid | residual sug... | chlorides | free sulfur d... | total sulfur d... | density | pH | sulphates |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 5 | [-∞ - 9.7] | | | | | 0.076 | 11 | 34 | 0.998 | 3.510 | 0.560 |
| 2 | 5 | [9.7 - 11.0] | | | | | 0.098 | 25 | 67 | 0.997 | 3.200 | 0.680 |
| 3 | 5 | [9.7 - 11.0] | | | | | 0.092 | 15 | 54 | 0.997 | 3.260 | 0.650 |
| 4 | 6 | [9.7 - 11.0] | 11.200 | 0.280 | 0.560 | 1.900 | 0.075 | 17 | 60 | 0.998 | 3.160 | 0.580 |
| 5 | 5 | [-∞ - 9.7] | 7.400 | 0.700 | 0 | 1.900 | 0.076 | 11 | 34 | 0.998 | 3.510 | 0.560 |
| 6 | 5 | | | | 0 | 1.800 | 0.075 | 13 | 40 | 0.998 | 3.510 | 0.560 |
| 7 | 5 | | | | | 1.600 | 0.069 | 15 | 59 | 0.996 | 3.300 | 0.460 |
| 8 | 7 | | | | 0 | 1.200 | 0.065 | 15 | 21 | 0.995 | 3.390 | 0.470 |
| 9 | 7 | [-∞ - 9.7] | 7.800 | 0.580 | 0.020 | 2 | 0.073 | 9 | 18 | 0.997 | 3.360 | 0.570 |
| 10 | 5 | [9.7 - 11.0] | 7.500 | 0.500 | 0.360 | 6.100 | 0.071 | 17 | 102 | 0.998 | 3.350 | 0.800 |
| 11 | 5 | [-∞ - 9.7] | 6.700 | 0.580 | 0.080 | 1.800 | 0.097 | 15 | 65 | 0.996 | 3.280 | 0.540 |
| 12 | 5 | [9.7 - 11.0] | 7.500 | 0.500 | 0.360 | 6.100 | 0.071 | 17 | 102 | 0.998 | 3.350 | 0.800 |
| 13 | 5 | [9.7 - 11.0] | 5.600 | 0.615 | 0 | 1.600 | 0.089 | 16 | 59 | 0.994 | 3.580 | 0.520 |
| 14 | 5 | [-∞ - 9.7] | 7.800 | 0.610 | 0.290 | 1.600 | 0.114 | 9 | 29 | 0.997 | 3.260 | 1.560 |
| 15 | 5 | [-∞ - 9.7] | 8.900 | 0.620 | 0.180 | 3.800 | 0.176 | 52 | 145 | 0.999 | 3.160 | 0.880 |
| 16 | 5 | [-∞ - 9.7] | 8.900 | 0.620 | 0.190 | 3.900 | 0.170 | 51 | 148 | 0.999 | 3.170 | 0.930 |

Data

Statistics

Visualizations

Annotations

Click 'Run' icon to execute the analysis process.   ①

Click 'Statistics' menu.   ②

ExampleSet (1,599 examples, 1 special attribute, 11 regular attributes)

File    Edit    Process    View    Connections    Settings    Extensions    Help

Views:    Design    **Results**    Turbo Prep    Auto Model

Find data, operators...etc    All Studio

Result History    **ExampleSet (Discretize)**

| Name | Type | Missing | Statistics |
| --- | --- | --- | --- |

Filter (12 / 12 attributes):    Search for Attributes

**Data**

**Statistics**

Label
**quality**    Polynominal    0

Least
3 (10)

5 (681)

5 (681), 6 (638),
7 (199), 4 (53),
...[2 more]

Open visualizations

Details...

Click 'Open visualizations' menu.    2

**Visualizations**

| | Least | Most | Values |
| --- | --- | --- | --- |
| **alcohol**    Nominal    0 | [13.6 - ∞] (8) | [9.7 - 11.0] (639) | [9.7 - 11.0] (639), [-∞ - 9.7] (552), ...[3 more] |

**Annotations**

Click 'quality' row.    1

| | Min | Max | Average |
| --- | --- | --- | --- |
| **fixed acidity** | | 15.900 | 8.320 |
| **volatile acidity** | | 1.580 | 0.528 |
| **citric acid**    Real    0 | 0 | 1 | 0.271 |
| **residual sugar**    Real    0 | 0.900 | 15.500 | 2.539 |
| **chlorides**    Real    0 | 0.012 | 0.611 | 0.087 |
| | Min | Max | Average |

# Visualize multiple nominal attributes with a stacked bar chart

Views: Design | **Results** | Turbo Prep | Auto Model

Find data, operators...etc | All Studio ▼

Result History | 🔲 **ExampleSet (Discretize)** ✕

**Plot** ⊗

**Plot 1** 🗑

**Plot type**

📊 Bar (Column) ▼

**Value column**

quality ▼

☑ Aggregate data

**Group by**

quality ▼

**Aggregation Function**

Count ▼

Color Group

alcohol ▼

Stacking

No stacking ▼

Plot style »

Add new plot

**General** ⊗

ExampleSet

Set 'Color Group' as 'alcohol'. ①

Set 'Stacking' as 'No stacking'. ②

Count(quality) axis: 0, 50, 100, 150, 200, 250, 300, 350, 400

quality axis: 5, 6, 7, 4, 8, 3

● alcohol: [−∞ − 9.7]   ● alcohol: [9.7 − 11.0]   ● alcohol: [12.3 − 13.6]   ● alcohol: [13.6 − ∞]   ● alcohol: [11.0 − 12.3]

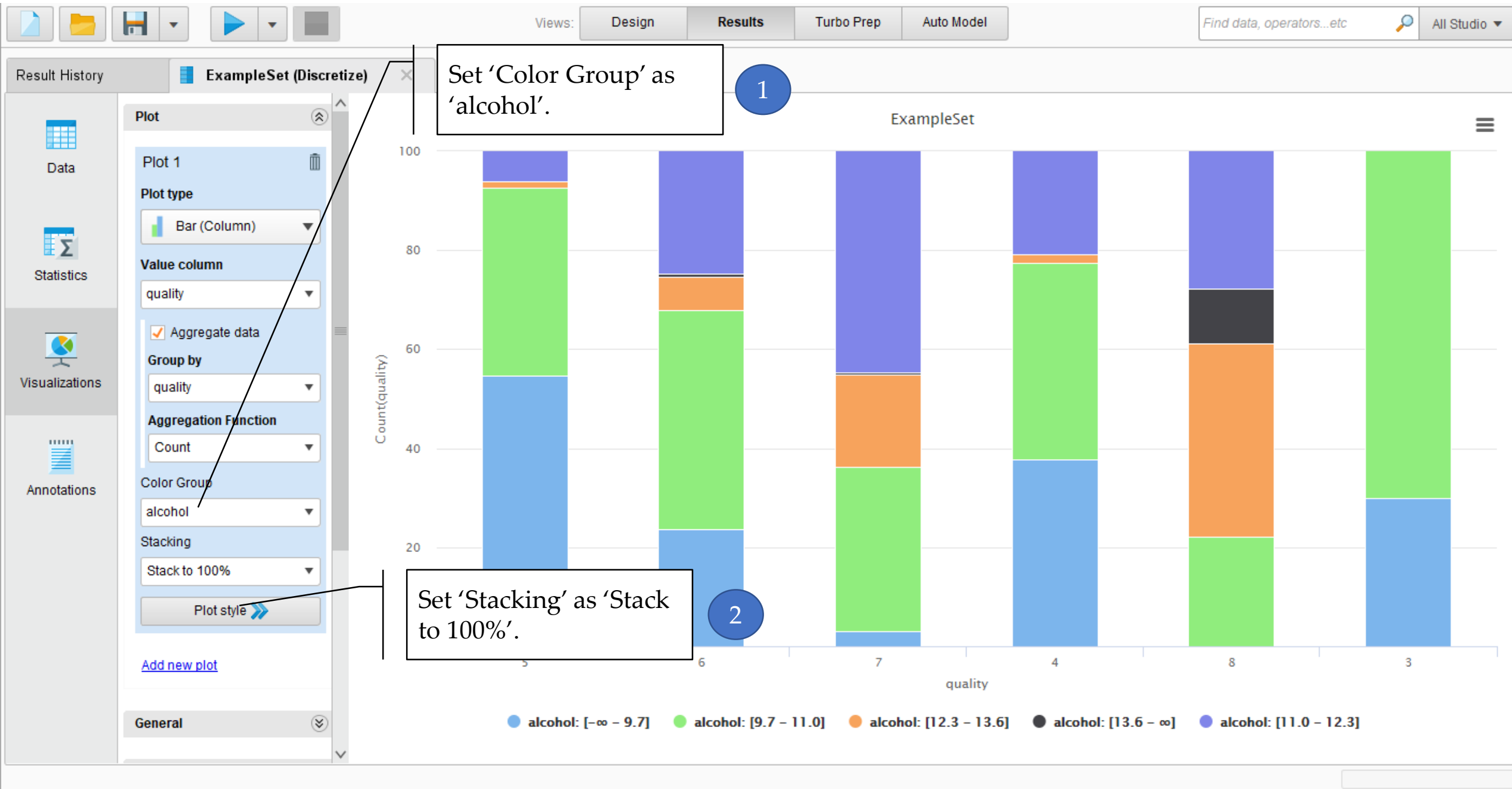# Visualize multiple nominal attributes with a stacked bar chart
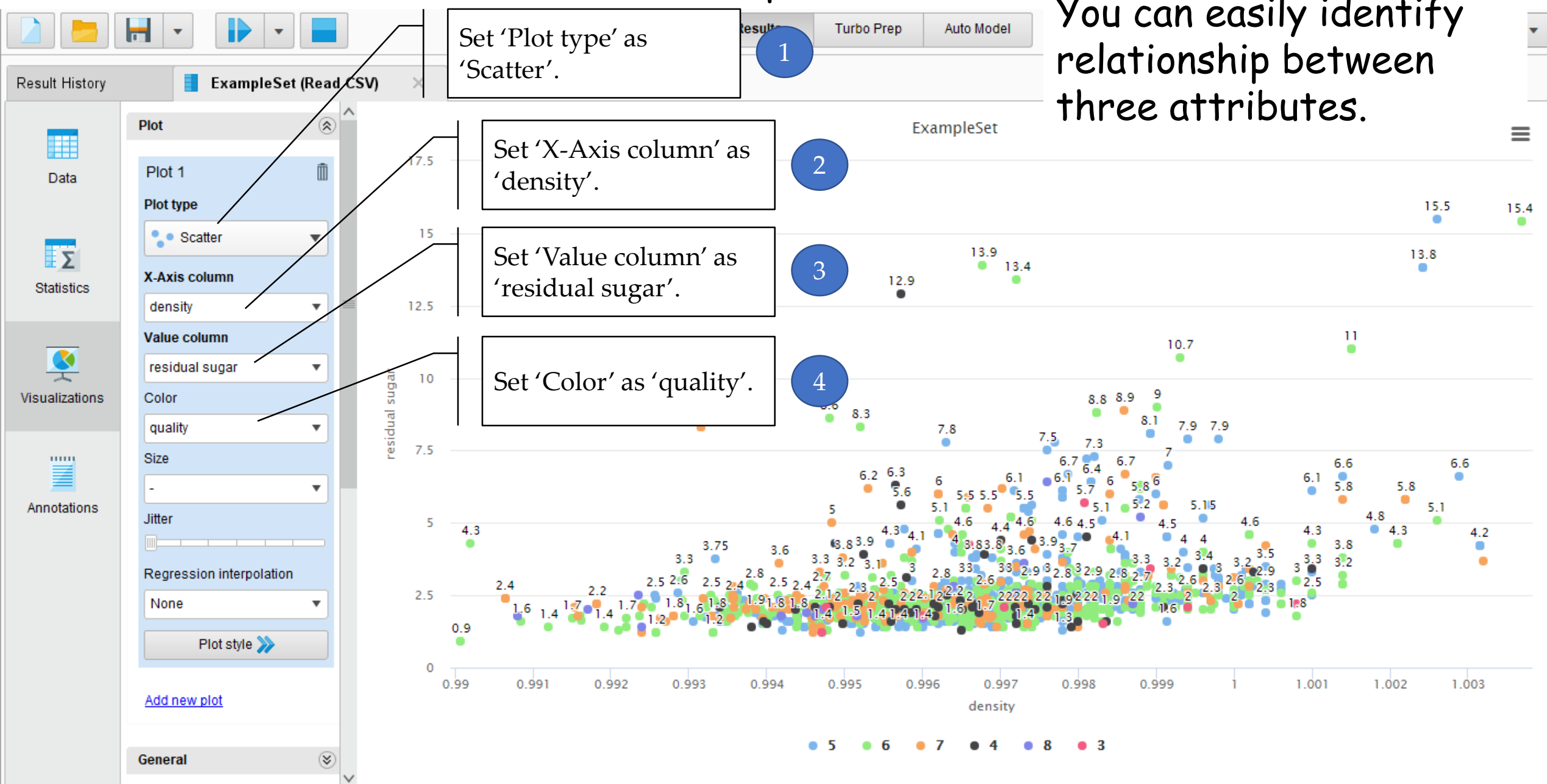
# Visualize multiple nominal attributes with a stacked bar chart

# Visualize three attributes with a scatter plot

You can easily identify relationship between three attributes.

Set 'Plot type' as 'Scatter'. **1**

Set 'X-Axis column' as 'density'. **2**

Set 'Value column' as 'residual sugar'. **3**

Set 'Color' as 'quality'. **4**



[1] Process  03:29:26

# Conclusion

## Conclusion

- **In Data Understanding phase, it is necessary to gather, describe and explore data and verify data quality.**

- **Through exercises, we learn how to load dataset and how to explore dataset with Rapidminer**

- **Now you need to have knowledge on data to be analyzed and go to next phase of CRISP-DM.**

# QUESTIONS?