# Week 2 Practical Lab: Data Pre-processing and Statistical Analysis using KNIME

**Useful resources:**
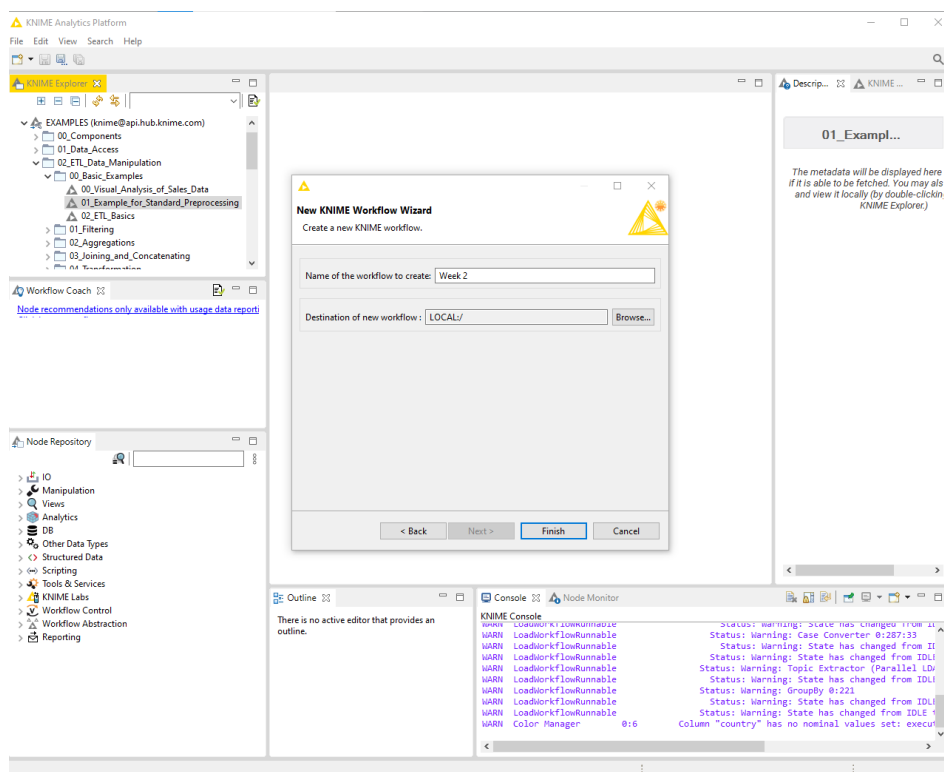
- **KNIME Analytics Platform** download link:
  https://www.knime.com/downloads/download-knime
  (It is recommender to download the installer version 561 MB for Windows)
- **Getting Started Guide** https://www.knime.com/getting-started-guide
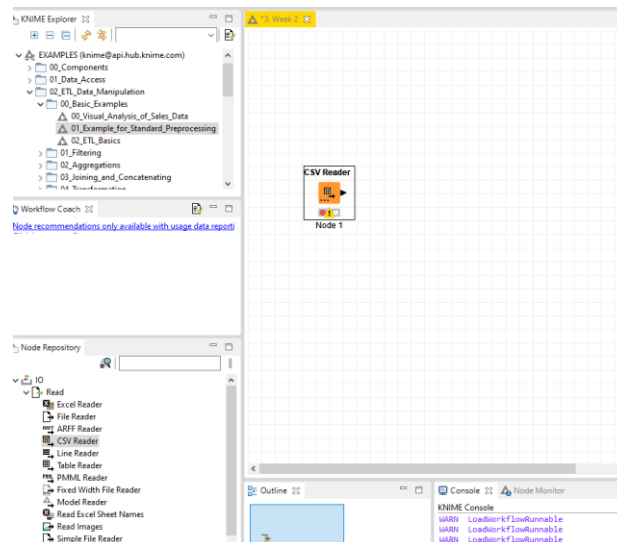- **KNIME Learning** https://www.knime.com/learning

**Objective:**

- Perform data preprocessing using KNIME
- Perform Statistical analysis using KNIME
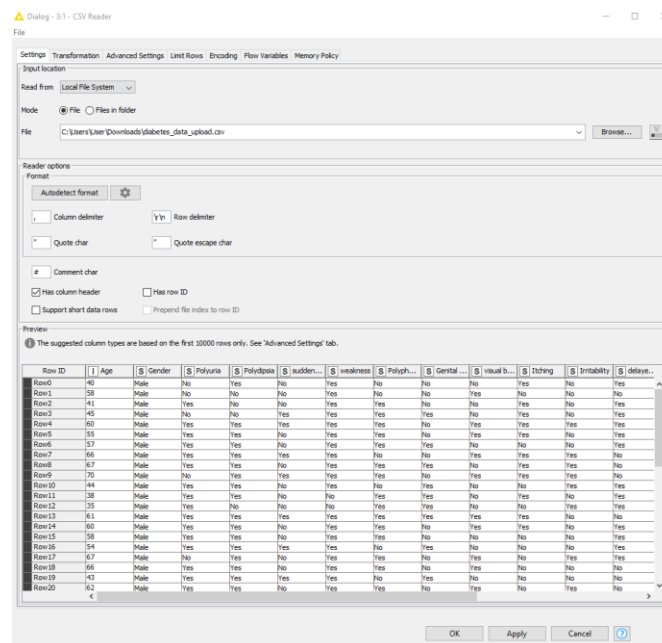
**Instructions:**

1. Open the KNIME software installed on your computer.
2. Create a new workflow from File-> New -> **New KNIME Workflow** -> Next -> *Enter a project name* -> Finish.

TEQSA: PRV14311
CRICOS: 03836J

Australia Advance Education Group Pty Ltd. trading as
Sydney International School of Technology and Commerce
ABN 74 613 055 440 |ACN 613 055 440
Level 14/233 Castlereagh Street, Sydney NSW 2000

P a g e  | 1

**3.** Select a data set of your choice from the given link (https://archive.ics.uci.edu/ml/index.php ) and open the file using appropriate node reader. (In this example, the following data set is chosen: https://archive.ics.uci.edu/ml/datasets/Early+stage+diabetes+risk+prediction+dataset. )

    a. Drag and drop the appropriate node from the Node reposity in the workflow (for this example, CSV reader.)



    b. To configure the node, right click on the node and select **Configure….**

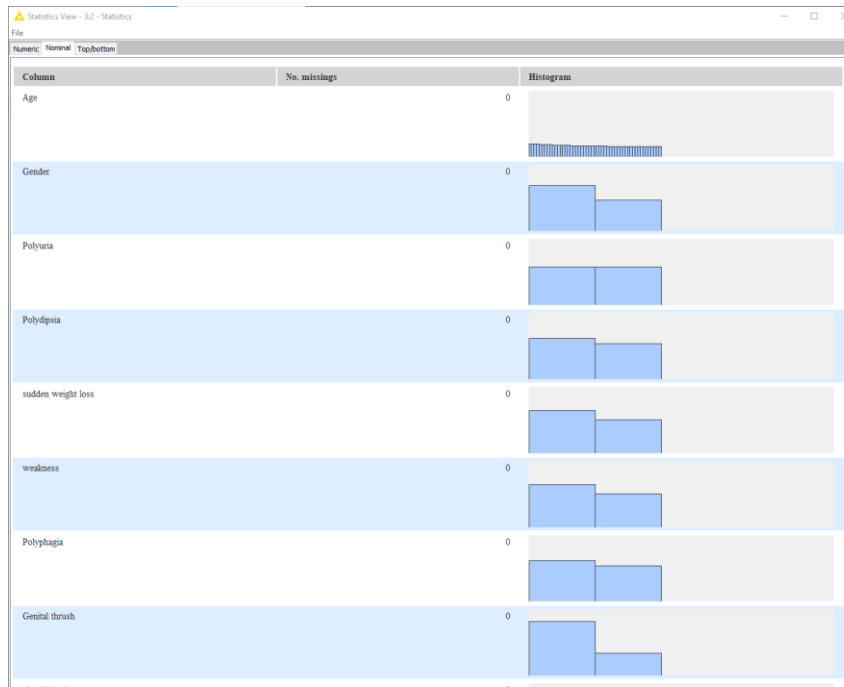    c. Browse and select your chosen file and click Okay if the data preview is accurate.



    d. Execute the node by clicking right on the node and then select **Execute**.

TEQSA: PRV14311          Australia Advance Education Group Pty Ltd. trading as          P a g e  | 2
CRICOS: 03836J          Sydney International School of Technology and Commerce
ABN 74 613 055 440 |ACN 613 055 440
Level 14/233 Castlereagh Street, Sydney NSW 2000

4. Select the **Statistics** node from the **Node Reposity** and connect with the **File Reader node** to understand on which columns or rows you need to perform pre-precessing. Then execute the node.



5. After executing the node, right click and select "View: Statistics View" to explore the statistical measures.



6. If there is any missing value in your data set, then select the "Missing Value" node as part of the pre-processing and connect it with the File Reader node.
    a. Go to Configuration mode for the missing value node and select appropriate rules for different data types.
    b. For column specific rules, Select the Column Settings tab.

TEQSA: PRV14311
CRICOS: 03836J

Australia Advance Education Group Pty Ltd. trading as
Sydney International School of Technology and Commerce
ABN 74 613 055 440 |ACN 613 055 440
Level 14/233 Castlereagh Street, Sydney NSW 2000

P a g e  | 3

7. Explore other pre-processing nodes (Convert & Replace, Filter, Split & Combine, Normaliser etc.) and discuss the suitability with your Lecturer.

TEQSA: PRV14311
CRICOS: 03836J

Australia Advance Education Group Pty Ltd. trading as
Sydney International School of Technology and Commerce
ABN 74 613 055 440 |ACN 613 055 440
Level 14/233 Castlereagh Street, Sydney NSW 2000

P a g e  | 4