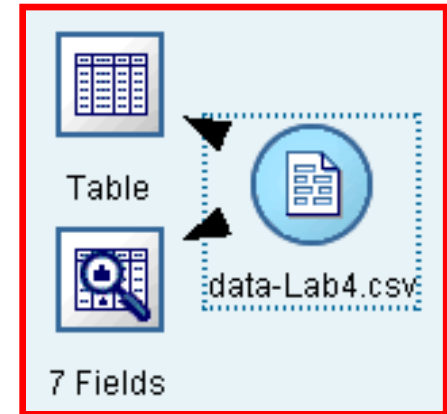# Data Mining – LAB 4

## Classification – Neural Network

# Data File

- Download from:
  - *www.comp.polyu.edu.hk/~csamak/data/data-Lab4.csv*
- For Virtualbox image, it is placed at:
  - `C:\Data\data-Lab4.csv`

# Data Understanding

- Load the data file (data-Lab4.csv) into PASW



**Think about these**:

1. How many attributes are being used in the dataset?
2. How many records are stored in the dataset?
3. Are there problems with the dataset?

# Data Preparation - Transformation

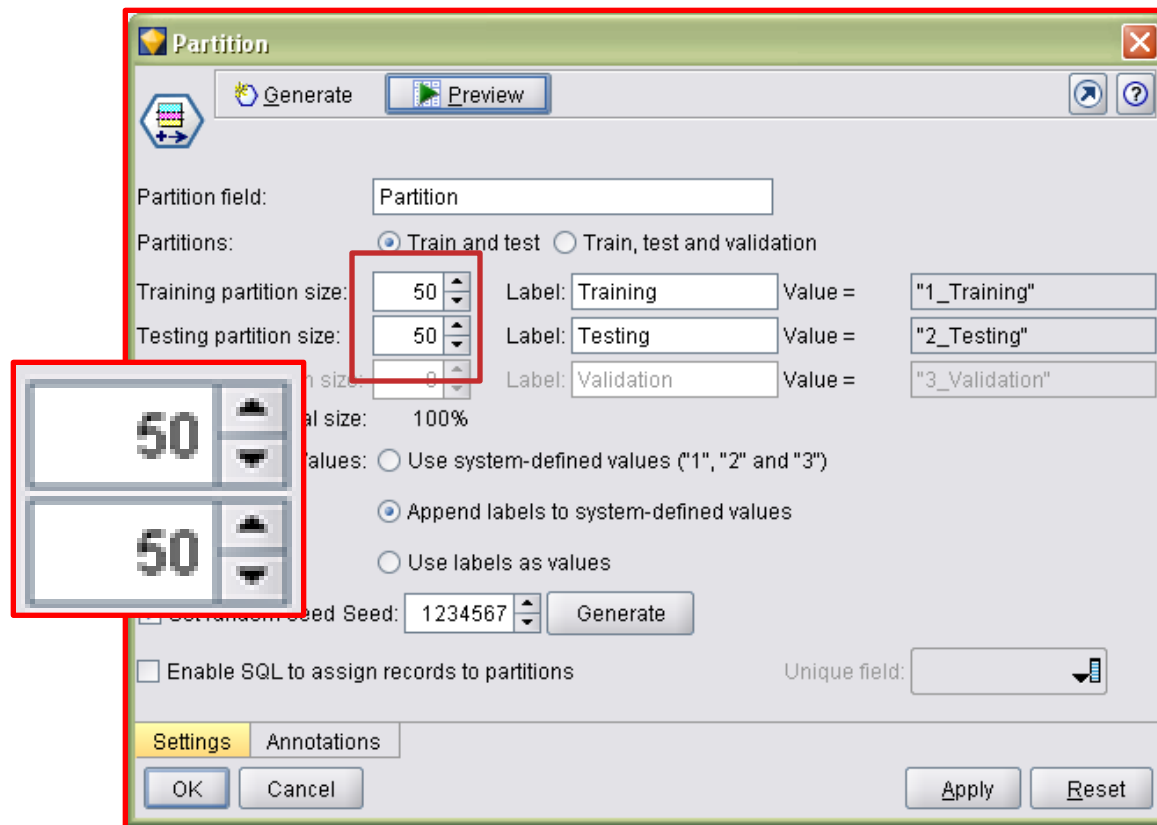- "Derive" a new field "Na_to_K" (ratio of Na to K, ie Na/K)

# Data Preparation - Transformation

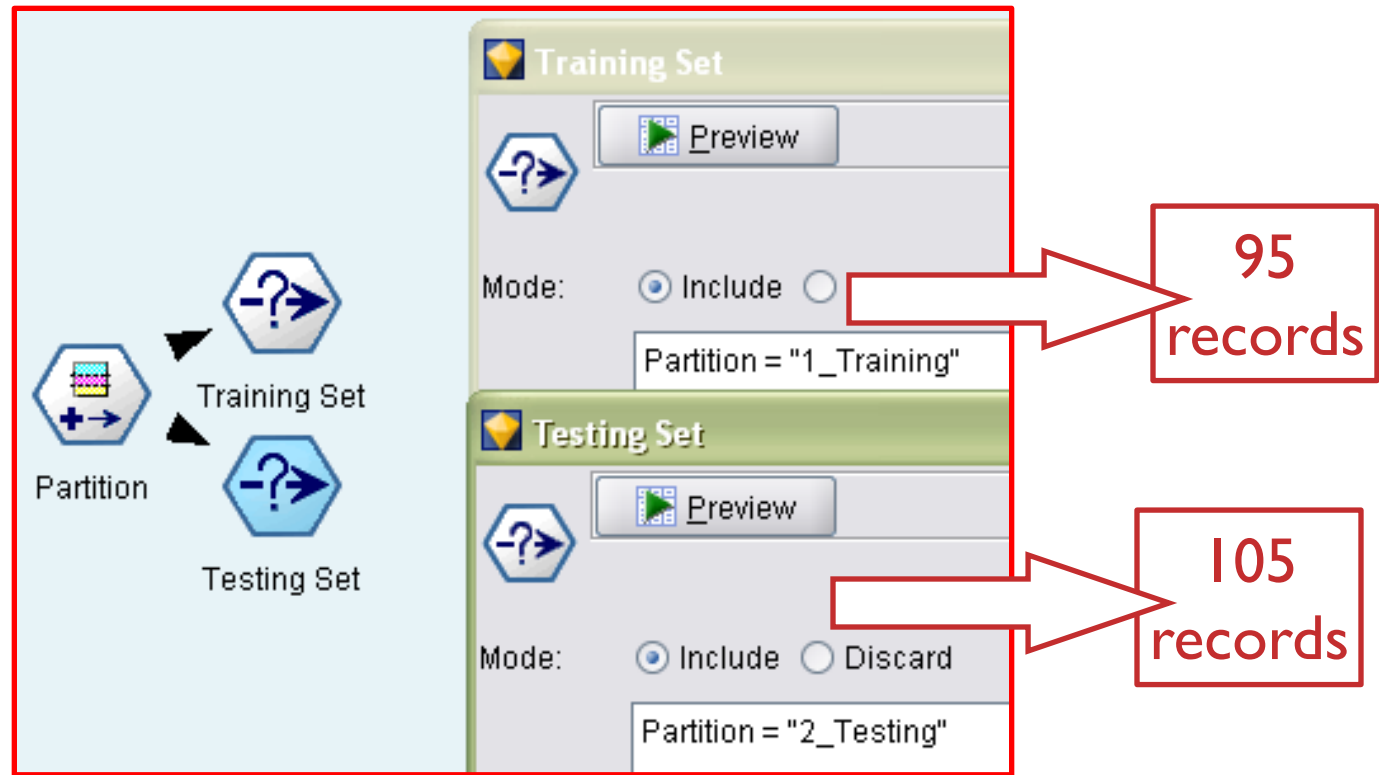- Use "Filter" node to discard the fields "Na" and "K"

# Data Preparation

- Add "Partition" node to divide the dataset into two, Training and Testing, in 50/50.

# Data Preparation

- Use two "Select" nodes to get the records

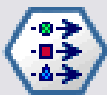# Prepare to Build the Model

- Use a "Type" node to refresh fields' value in application's memory
- Set the **INPUT** ➔ Age, Sex, BP, Cholesterol, Na_to_K
- Set the **OUTPUT** ➔ Drug
- Select "NONE" for the field, Partition

# Prepare to Build the Model

# Neural Net Model

Classification for attribute, "Drug" – Add the node, "Neural Net".

# Neural Net Model

- Link the model to "Testing Set" and add a node, "Analysis" to check the model



Do you know why you have a different figure?

# Things to Note



Six training methods for building model.

Randomly splits the data into separate training and testing sets for purposes of model building

"**Stop on**" – Stopping criteria
**Default** – the network stops training when the network appears to have reached its optimally trained state.

# The Six Training Methods

- **Quick**. This method uses rules of thumb and characteristics of the data to choose an appropriate shape (topology) for the network.

- **Dynamic**. This method creates an initial topology but modifies the topology by adding and/or removing hidden units as training progresses.

- **Multiple**. This method creates several networks of different topologies (the exact number depends on the training data). These networks are then trained in a pseudo-parallel fashion. At the end of training, the model with the lowest RMS error is presented as the final model.

- **Prune**. This method starts with a large network and removes (prunes) the weakest units in the hidden and input layers as training proceeds. This method is usually slow, but it often yields better results than other methods.

- **RBFN**. The radial basis function network (RBFN) uses a technique similar to k-means clustering to partition the data based on values of the target field.

- **Exhaustive prune**. This method is related to the Prune method. It starts with a large network and prunes the weakest units in the hidden and input layers as training proceeds. With Exhaustive Prune, network training parameters are chosen to ensure a very thorough search of the space of possible models to find the best one. This method is usually the slowest, but it often yields the best results. Note that this method can take a long time to train, especially with large datasets.

# Other Stop Options

- **Accuracy** (%). With this option, training will continue until the specified accuracy is attained. This may never happen, but you can interrupt training at any point and save the net with the best accuracy achieved so far.
- **Cycles**. With this option, training will continue for the specified number of **cycles** (passes through the data).
- **Time** (mins). With this option, training will continue for the specified amount of time (in minutes). Note that training may go a bit beyond the specified time limit in order to complete the current cycle.

# Advanced Settings



**Hidden layers**. Select the number of hidden layers for the neural network. More hidden layers can help neural networks learn more complex relationships, but they also increase training time.

**Layer 1, 2, 3**. For each layer, specify the number of hidden units to include. More hidden units per layer can also help in learning complex tasks, but as with additional hidden layers, they also increase training time.

**Persistence**. Specify the number of cycles for which the network will continue to train if no improvement is seen. Higher values can help networks escape local minima, but they also increase training time.

**Alpha and Eta**. These parameters control the training of the network.

# Try Yourself

- Separate the dataset into 70/30 for the Training and Testing.

- Check the results.

- What are the differences?