# Knime Hands-on
# World Indicators Dataset

Prof. Dr. Daniel A. Keim
University of Konstanz, Germany

2nd ACM Europe Summer School in Data Science 12 – 18 July 2018, Athens, Greece
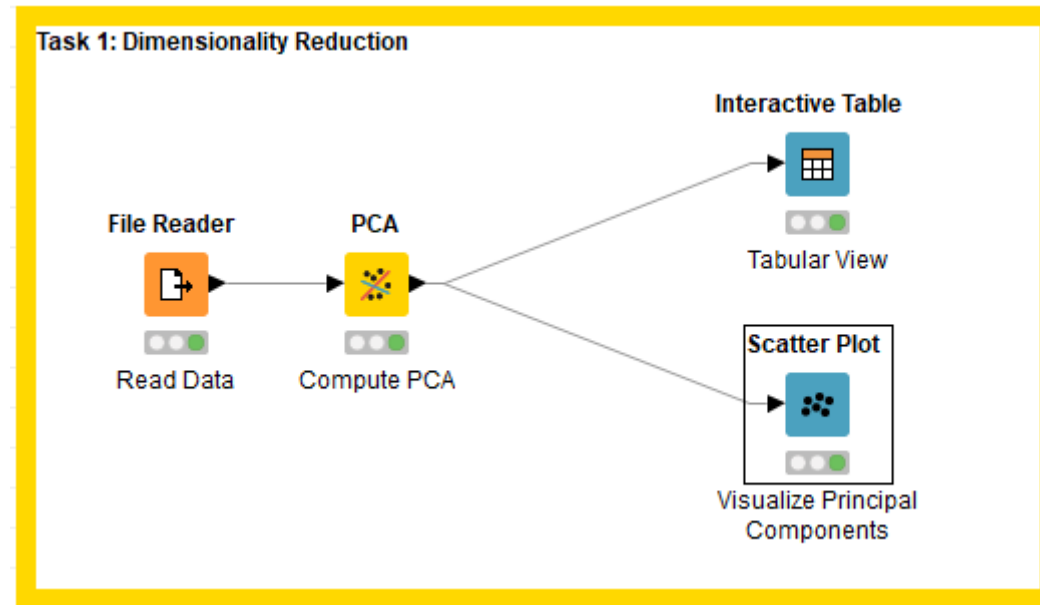
# Stock Dataset

- `Close/Open`: stock price at open and close

- `High/Low`: max and min price during the day

- `Volume`: trading volume

- `Number of records`: the number of contained records.

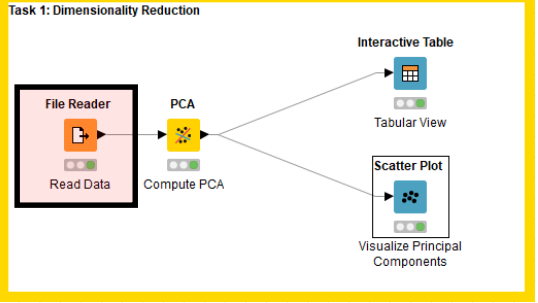| Close | Company | Date | High | Low | Number of Records | Open | Volume |
|---|---|---|---|---|---|---|---|
| 133.9 | Amazon | 1/4/2010 | 136.61 | 133.14 | 1 | 136.25 | 7600543 |
| 30.57 | Apple | 1/4/2010 | 30.64 | 30.34 | 1 | 30.49 | 123432050 |
| 53.64 | Biogen Idec | 1/4/2010 | 53.97 | 53.6 | 1 | 53.97 | 2469662 |
| 134.69 | Amazon | 1/5/2010 | 135.48 | 131.81 | 1 | 133.43 | 8856456 |
| 30.63 | Apple | 1/5/2010 | 30.8 | 30.46 | 1 | 30.66 | 150476004 |
| 53.38 | Biogen Idec | 1/5/2010 | 55 | 53 | 1 | 54.72 | 4899370 |
| 132.25 | Amazon | 1/6/2010 | 134.73 | 131.65 | 1 | 134.6 | 7180977 |
| 30.14 | Apple | 1/6/2010 | 30.75 | 30.11 | 1 | 30.63 | 138039594 |
| 53.43 | Biogen Idec | 1/6/2010 | 53.7 | 52.8 | 1 | 53.1 | 5555723 |
| 130 | Amazon | 1/7/2010 | 132.32 | 128.8 | 1 | 132.01 | 11030124 |
| 30.08 | Apple | 1/7/2010 | 30.29 | 29.86 | 1 | 30.25 | 119282324 |
| 52.99 | Biogen Idec | 1/7/2010 | 53.5 | 52.46 | 1 | 53.23 | 3659834 |

# Stock Dataset – Task 1

- Find similar stocks in terms of the Dimensions *open*, *close*, *high*, and *low.*

- Apply a dimension reduction technique.

- Inspect the data in a table and visualize the projected dataset in a scatter plot. Identify similar stocks visually.

# Task 1: Knime Workflow

# Task 1: Read CSV

# Task 1: PCA Configuration

# Task 1: PCA Tabular View



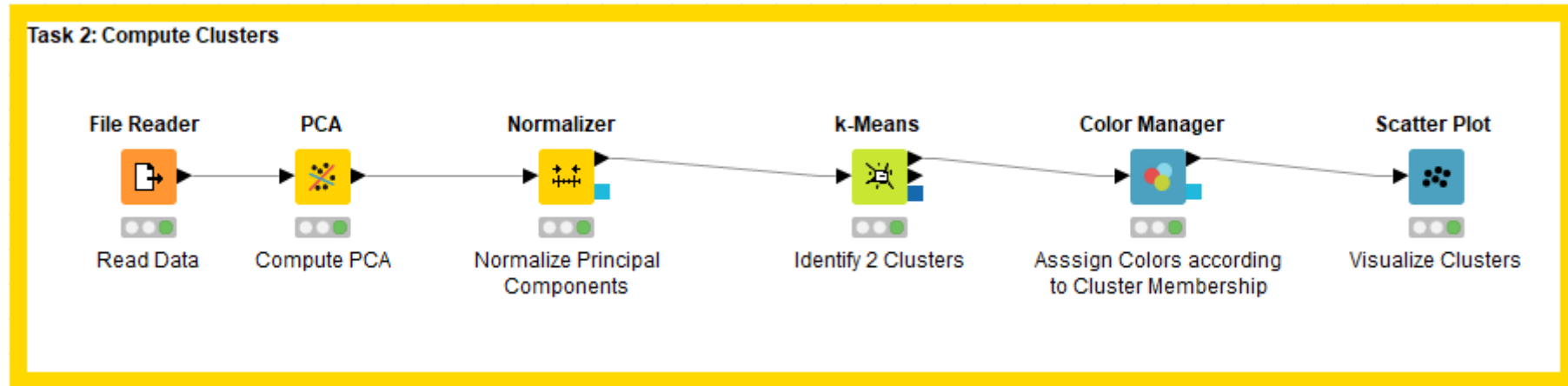| Row ID | Close | Company | Date | High | Low | Number... | Open | Volume | PCA dime... | PCA di... | PCA di... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Row0 | 133.9 | Amazon | 1/4/2010 | 136.61 | 133.14 | 1 | 136.25 | 7600543 | -32,600,482.... | 99.457 | -1.765 |
| Row1 | 30.57 | Apple | 1/4/2010 | 30.64 | 30.34 | 1 | 30.49 | 123432050 | 83,231,024.866 | 102.78 | 0.076 |
| Row2 | 53.64 | Biogen Idec | 1/4/2010 | 53.97 | 53.6 | 1 | 53.97 | 2469662 | -37,731,363.... | 270.939 | -0.331 |
| Row3 | 134.69 | Amazon | 1/5/2010 | 135.48 | 131.81 | 1 | 133.43 | 8856456 | -31,344,569.... | 99.472 | 0.729 |
| Row4 | 30.63 | Apple | 1/5/2010 | 30.8 | 30.46 | 1 | 30.66 | 150476004 | 110,274,978.... | 54.517 | 0.04 |
| Row5 | 53.38 | Biogen Idec | 1/5/2010 | 55 | 53 | 1 | 54.72 | 4899370 | -35,301,655.... | 266.155 | -1.2 |
| Row6 | 132.25 | Amazon | 1/6/2010 | 134.73 | 131.65 | 1 | 134.6 | 7180977 | -33,020,048.... | 103.539 | -1.733 |
| Row7 | 30.14 | Apple | 1/6/2010 | 30.75 | 30.11 | 1 | 30.63 | 138039594 | 97,838,568.866 | 77.052 | -0.331 |
| Row8 | 53.43 | Biogen Idec | 1/6/2010 | 53.7 | 52.8 | 1 | 53.1 | 5555723 | -34,645,302.... | 266.532 | 0.08 |
| Row9 | 130 | Amazon | 1/7/2010 | 132.32 | 128.8 | 1 | 132.01 | 11030124 | -29,170,901.... | 101.753 | -1.551 |
| Row10 | 30.08 | Apple | 1/7/2010 | 30.29 | 29.86 | 1 | 30.25 | 119282324 | 79,081,298.866 | 110.925 | -0.122 |
| Row11 | 52.99 | Biogen Idec | 1/7/2010 | 53.5 | 52.46 | 1 | 53.23 | 3659834 | -36,541,191.... | 270.322 | -0.341 |
| Row12 | 133.52 | Amazon | 1/8/2010 | 133.68 | 129.03 | 1 | 130.56 | 9833829 | -30,367,196.... | 102.04 | 1.802 |
| Row13 | 30.28 | Apple | 1/8/2010 | 30.29 | 29.87 | 1 | 30.04 | 111969081 | 71,768,055.866 | 123.908 | 0.154 |
| Row14 | 54.08 | Biogen Idec | 1/8/2010 | 54.33 | 52.81 | 1 | 53 | 2996438 | -37,204,587.... | 270.477 | 0.54 |
| Row15 | 130.31 | Amazon | 1/11/2010 | 132.8 | 129.21 | 1 | 132.62 | 8786668 | -31,414,357.... | 104.83 | -1.77 |
| Row16 | 30.02 | Apple | 1/11/2010 | 30.43 | 29.78 | 1 | 30.4 | 115557365 | 75,356,339.866 | 117.461 | -0.298 |
| Row17 | 53.9 | Biogen Idec | 1/11/2010 | 54.35 | 53.53 | 1 | 54.1 | 1813289 | -38,387,736.... | 271.751 | -0.288 |
| Row18 | 127.35 | Amazon | 1/12/2010 | 129.82 | 126.55 | 1 | 128.99 | 9098190 | -31,102,835.... | 110.394 | -1.263 |
| Row19 | 29.67 | Apple | 1/12/2010 | 29.97 | 29.49 | 1 | 29.88 | 148614774 | 108,413,748.... | 59.59 | -0.107 |
| Row20 | 53.25 | Biogen Idec | 1/12/2010 | 54.03 | 52.94 | 1 | 53.66 | 3242899 | -36,958,126.... | 270.212 | -0.458 |
| Row21 | 129.11 | Amazon | 1/13/2010 | 129.71 | 125.75 | 1 | 127.9 | 10727856 | -29,473,169.... | 107.617 | 0.634 |
| Row22 | 30.09 | Apple | 1/13/2010 | 30.13 | 29.16 | 1 | 29.7 | 151472335 | 111,271,309.... | 54.479 | 0.26 |
| Row23 | 54.08 | Biogen Idec | 1/13/2010 | 54.33 | 52.96 | 1 | 53.33 | 2131312 | -38,069,713.... | 271.773 | 0.32 |
| Row24 | 127.35 | Amazon | 1/14/2010 | 130.38 | 126.4 | 1 | 129.14 | 9788435 | -30,412,590.... | 108.884 | -1.434 |
| Row25 | 29.92 | Apple | 1/14/2010 | 30.07 | 29.86 | 1 | 30.02 | 108288411 | 68,087,385.866 | 130.748 | -0.066 |
| Row26 | 54.36 | Biogen Idec | 1/14/2010 | 54.36 | 53.59 | 1 | 53.81 | 1770910 | -38,430,115.... | 271.707 | 0.24 |
| Row27 | 127.14 | Amazon | 1/15/2010 | 129.65 | 127.06 | 1 | 129.18 | 15382763 | -24,818,262.... | 99.083 | -1.456 |
| Row28 | 29.42 | Apple | 1/15/2010 | 30.23 | 29.41 | 1 | 30.13 | 148584065 | 108,383,039.... | 59.552 | -0.488 |
| Row29 | 54.18 | Biogen Idec | 1/15/2010 | 54.85 | 53.78 | 1 | 54.06 | 2997202 | -37,203,823.... | 269.153 | -0.077 |
| Row30 | 127.61 | Amazon | 1/19/2010 | 128 | 124.33 | 1 | 126.2 | 8900116 | -31,300,909.... | 114.029 | 0.796 |
| Row31 | 30.72 | Apple | 1/19/2010 | 30.74 | 29.61 | 1 | 29.76 | 182501620 | 142,300,594.... | -1.479 | 0.701 |
| Row32 | 55.23 | Biogen Idec | 1/19/2010 | 55.31 | 54.29 | 1 | 54.29 | 2941572 | -37,259,453.... | 268.127 | 0.497 |
| Row33 | 125.78 | Amazon | 1/20/2010 | 129.2 | 125.08 | 1 | 127.13 | 9081533 | -31,119,492.... | 113.178 | -1.128 |
| Row34 | 30.25 | Apple | 1/20/2010 | 30.79 | 29.93 | 1 | 30.7 | 153037892 | 112,836,866.... | 50.406 | -0.307 |
| Row35 | 54.72 | Biogen Idec | 1/20/2010 | 55.37 | 54.01 | 1 | 55.04 | 2634095 | -37,566,930.... | 268.66 | -0.427 |
| Row36 | 126.62 | Amazon | 1/21/2010 | 128.15 | 125 | 1 | 127.26 | 9976146 | -30,224,879.... | 111.676 | -0.568 |
| Row37 | 29.72 | Apple | 1/21/2010 | 30.47 | 29.6 | 1 | 30.3 | 152038565 | 111,837,539.... | 52.969 | -0.397 |
| Row38 | 53.21 | Biogen Idec | 1/21/2010 | 54.53 | 52.67 | 1 | 54.53 | 4731164 | -35,469,861.... | 267.034 | -1.181 |
| Row39 | 121.43 | Amazon | 1/22/2010 | 127.67 | 120.76 | 1 | 125.6 | 11577818 | -28,623,207.... | 114.592 | -3.39 |
| Row40 | 28.25 | Apple | 1/22/2010 | 29.64 | 28.17 | 1 | 29.54 | 220441872 | 180,240,846.... | -66.217 | -0.847 |
| Row41 | 53.06 | Biogen Idec | 1/22/2010 | 53.71 | 52.6 | 1 | 53.19 | 3167791 | -37,033,234.... | 271.005 | -0.266 |
| Row42 | 120.31 | Amazon | 1/25/2010 | 122.28 | 118.12 | 1 | 122.1 | 12031053 | -28,169,972.... | 120.131 | -1.517 |

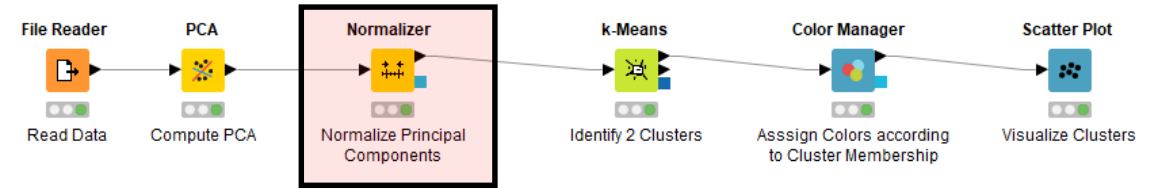# Task 1: PCA Scatter Plot (pc0, pc1)

# Task 2 – Compute Clusters

- Cluster the projected dataset using k-means.

- Assign colors corresponding to cluster membership.

- Visualize the cluster memberships.

# Task 2 – Knime Workflow

# Normalize PCs



Task 2: Compute Clusters

| File Reader | PCA | Normalizer | k-Means | Color Manager | Scatter Plot |
|---|---|---|---|---|---|
| Read Data | Compute PCA | Normalize Principal Components | Identify 2 Clusters | Assign Colors according to Cluster Membership | Visualize Clusters |



Dialog - 3:22 - Normalizer (Normalize Principal)

File

Methods | Flow Variables | Memory Policy

◉ Manual Selection  ○ Wildcard/Regex Selection

**Exclude**

Column(s): [          ]  [Search]
☐ Select all search hits

D Close
D High
D Low
I Number of Records
D Open
I Volume

◉ Enforce exclusion

**Select**

[ add >> ]

[ add all >> ]

[ << remove ]

[ << remove all ]

**Include**

Column(s): [          ]  [Search]
☐ Select all search hits

D PCA dimension 0
D PCA dimension 1
D PCA dimension 2

○ Enforce inclusion

**Settings**

◉ Min-Max Normalization          Min: 0.0
                                 Max: 1.0

○ Z-Score Normalization (Gaussian)

○ Normalization by Decimal Scaling

[ OK ]  [ Apply ]  [ Cancel ]  (?)

# k-means Clustering

# Assign Colors to Clusters

# Visualize Cluster Memberships

# Task 3 – Interactive Cluster Separation

- Use a *better suited* clustering technique

- Assign colors and visualize the cluster memberships

- Open two tables to show the data from at least two different clusters.

# Task 3: Knime Workflow

# Task 3: Distance Computations

# Task 3: DBScan Clustering

# Task 3: DBScan Meta Data



Task 3: Interactive Cluster Separation

| Row ID | Count |
|---|---|
| Noise | 96 |
| Cluster_0 | 5 |
| Cluster_1 | 7 |
| Cluster_2 | 5 |
| Cluster_3 | 7 |
| Cluster_4 | 6 |
| Cluster_5 | 7 |
| Cluster_6 | 7 |
| Cluster_7 | 5 |
| Cluster_8 | 5 |
| Cluster_9 | 1125 |
| Cluster_10 | 2500 |

Table View - 3:28 - Interactive Table (Cluster Count)
File  Hilite  Navigation  View  Output

# Task 3: Select Instances of Cluster 9

# Task 3: Inspect Instances of Cluster 9