**CIND 820 XJH Final Project by Sylvia Pereira**
Supervisor: Dr. Tamer Abdou tamer.abdou@torontomu.ca

**Data Preparation**

The data preparation was performed in Alteryx Designer Desktop software.

<u>Glossary</u>
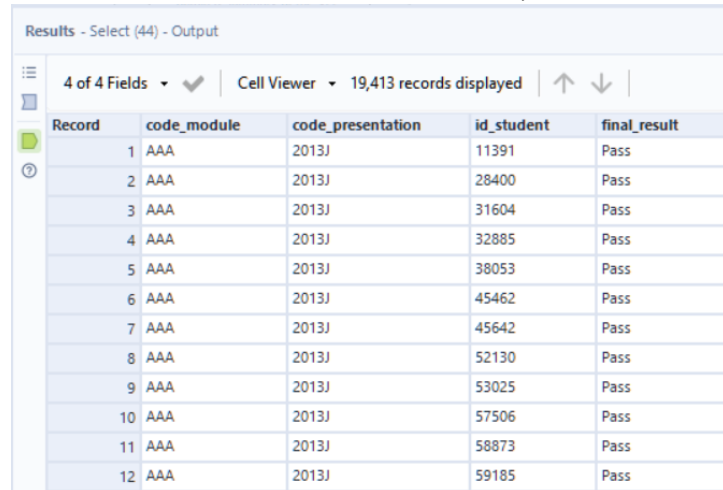VLE: Virtual Learning Environment

Below you can follow the steps I used to prepare the dataset for the Exploratory Data Analysis (EDA) report.

Original datasets can be found at: https://analyse.kmi.open.ac.uk/open_dataset

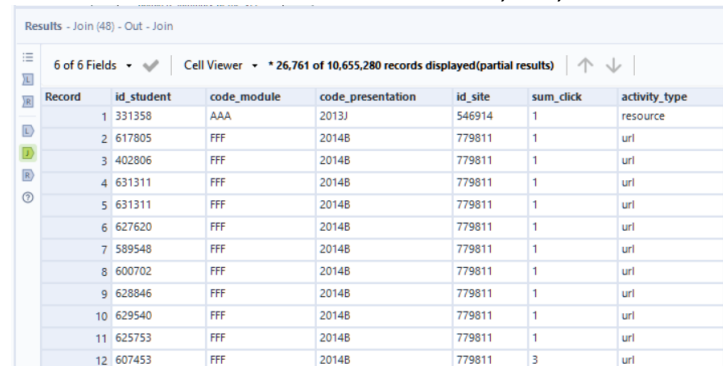| | |
|---|---|
| The first step in the process was to filter students that Passed or Failed the courses and remove undesired columns. Other results, such as withdrawal and distinction, were dropped from the dataset. Table: **StudentInfo.csv** | **Result:** Table with four columns and 19,413 records.  |
| **StudentVle.csv** and **vle.csv** tables were combined to summarize the VLE interactions per student. Duplicates and undesired fields were removed. | **Result:** Table with six columns and 10,655,280 records.  |

| | |
|---|---|
| Using the Join tool, the results of **StudentInfo.csv** and the combined dataset from **StudentVle.csv** + **vle.csv** were joined again to form the working dataset. | **Result:** Table with seven columns and 8,754,616 records.<br><br>Results - Join (41) - Out - Join<br>7 of 7 Fields ✓ Cell Viewer ▾ * 23,996 of 8,754,616 records displayed(partial results)<br><table><tr><th>Record</th><th>id_student</th><th>course_name</th><th>code_presentation</th><th>id_site</th><th>activity_type</th><th>sum_click</th><th>final_result</th></tr><tr><td>1</td><td>349182</td><td>DDD</td><td>2013J</td><td>674010</td><td>subpage</td><td>1</td><td>Fail</td></tr><tr><td>2</td><td>349182</td><td>DDD</td><td>2013J</td><td>674449</td><td>url</td><td>1</td><td>Fail</td></tr><tr><td>3</td><td>349182</td><td>DDD</td><td>2013J</td><td>674100</td><td>subpage</td><td>2</td><td>Fail</td></tr><tr><td>4</td><td>349182</td><td>DDD</td><td>2013J</td><td>674100</td><td>subpage</td><td>2</td><td>Fail</td></tr><tr><td>5</td><td>349182</td><td>DDD</td><td>2013J</td><td>674100</td><td>subpage</td><td>1</td><td>Fail</td></tr><tr><td>6</td><td>349182</td><td>DDD</td><td>2013J</td><td>673740</td><td>oucontent</td><td>11</td><td>Fail</td></tr><tr><td>7</td><td>349182</td><td>DDD</td><td>2013J</td><td>673740</td><td>oucontent</td><td>4</td><td>Fail</td></tr><tr><td>8</td><td>349182</td><td>DDD</td><td>2013J</td><td>673740</td><td>oucontent</td><td>2</td><td>Fail</td></tr><tr><td>9</td><td>349182</td><td>DDD</td><td>2013J</td><td>673740</td><td>oucontent</td><td>5</td><td>Fail</td></tr><tr><td>10</td><td>349182</td><td>DDD</td><td>2013J</td><td>673740</td><td>oucontent</td><td>1</td><td>Fail</td></tr><tr><td>11</td><td>349182</td><td>DDD</td><td>2013J</td><td>673727</td><td>oucontent</td><td>2</td><td>Fail</td></tr><tr><td>12</td><td>349182</td><td>DDD</td><td>2013J</td><td>673727</td><td>oucontent</td><td>1</td><td>Fail</td></tr></table> |
| Using the Find Replace tool, the fields course_name and code_presentation were replaced by one field called course_year_month to clarify the information.<br>The old fields course_name and code_presentation were dropped from the dataset as well as the id_site field since they are not relevant to this analysis. | **Result:** Table with five columns and 8,754,616 records.<br><br>Results - Select (5) - Output<br>5 of 5 Fields ✓ Cell Viewer ▾ * 23,689 of 8,754,616 records displayed(partial results)<br><table><tr><th>Record</th><th>id_student</th><th>course_year_month</th><th>activity_type</th><th>sum_click</th><th>final_result</th></tr><tr><td>1</td><td>349182</td><td>DDD 2013_October</td><td>subpage</td><td>1</td><td>Fail</td></tr><tr><td>2</td><td>349182</td><td>DDD 2013_October</td><td>url</td><td>1</td><td>Fail</td></tr><tr><td>3</td><td>349182</td><td>DDD 2013_October</td><td>subpage</td><td>2</td><td>Fail</td></tr><tr><td>4</td><td>349182</td><td>DDD 2013_October</td><td>subpage</td><td>2</td><td>Fail</td></tr><tr><td>5</td><td>349182</td><td>DDD 2013_October</td><td>subpage</td><td>1</td><td>Fail</td></tr><tr><td>6</td><td>349182</td><td>DDD 2013_October</td><td>oucontent</td><td>11</td><td>Fail</td></tr><tr><td>7</td><td>349182</td><td>DDD 2013_October</td><td>oucontent</td><td>4</td><td>Fail</td></tr><tr><td>8</td><td>349182</td><td>DDD 2013_October</td><td>oucontent</td><td>2</td><td>Fail</td></tr><tr><td>9</td><td>349182</td><td>DDD 2013_October</td><td>oucontent</td><td>5</td><td>Fail</td></tr><tr><td>10</td><td>349182</td><td>DDD 2013_October</td><td>oucontent</td><td>1</td><td>Fail</td></tr><tr><td>11</td><td>349182</td><td>DDD 2013_October</td><td>oucontent</td><td>2</td><td>Fail</td></tr><tr><td>12</td><td>349182</td><td>DDD 2013_October</td><td>oucontent</td><td>1</td><td>Fail</td></tr></table> |
| To reduce the number of records, the fields: id_student, course_year_month, and activity_type were grouped while the function "sum" was used for sum_clicks.<br>Then, the table was sorted by ID_Student field. | **Result:** Table with five fields and 167,933 records<br><br>Results - Sort (19) - Output<br>5 of 5 Fields ✓ Cell Viewer ▾ * 23,947 of 167,933 records displayed(partial results)<br><table><tr><th>Record</th><th>ID_Student</th><th>course_year_month</th><th>Activity_Type</th><th>Sum_click</th><th>final_result</th></tr><tr><td>1</td><td>6516</td><td>AAA 2014_October</td><td>homepage</td><td>497</td><td>Pass</td></tr><tr><td>2</td><td>6516</td><td>AAA 2014_October</td><td>forumng</td><td>451</td><td>Pass</td></tr><tr><td>3</td><td>6516</td><td>AAA 2014_October</td><td>subpage</td><td>143</td><td>Pass</td></tr><tr><td>4</td><td>6516</td><td>AAA 2014_October</td><td>dataplus</td><td>21</td><td>Pass</td></tr><tr><td>5</td><td>6516</td><td>AAA 2014_October</td><td>resource</td><td>31</td><td>Pass</td></tr><tr><td>6</td><td>6516</td><td>AAA 2014_October</td><td>oucontent</td><td>1505</td><td>Pass</td></tr><tr><td>7</td><td>6516</td><td>AAA 2014_October</td><td>url</td><td>143</td><td>Pass</td></tr><tr><td>8</td><td>11391</td><td>AAA 2013_October</td><td>oucontent</td><td>553</td><td>Pass</td></tr><tr><td>9</td><td>11391</td><td>AAA 2013_October</td><td>url</td><td>5</td><td>Pass</td></tr><tr><td>10</td><td>11391</td><td>AAA 2013_October</td><td>forumng</td><td>193</td><td>Pass</td></tr><tr><td>11</td><td>11391</td><td>AAA 2013_October</td><td>homepage</td><td>138</td><td>Pass</td></tr><tr><td>12</td><td>11391</td><td>AAA 2013_October</td><td>resource</td><td>13</td><td>Pass</td></tr></table> |
| In the final step, using the CrossTab tool, I pivoted the table's orientation by moving vertical data onto the horizontal axis.<br>The data corresponding to the learner's trajectories through the VLE are shown on which each row | **Result:** Table with 23 fields and 19,077 records |

and each column corresponds to a particular resource within the VLE.

Results - Data Cleansing (36) - Output26

23 of 23 Fields ▾ ✓  | Cell Viewer ▾  * 16,347 of 19,077 records displayed(partial results) | ↑ ↓ |

| Record | ID_Student | course_year_month | final_result | dataplus | dualpane | externalquiz | folder | forumng |
|---|---|---|---|---|---|---|---|---|
| 1 | 606606 | EEE 2013_October | Fail | 0 | 0 | 0 | 0 | 19 |
| 2 | 622480 | GGG 2014_February | Fail | 0 | 0 | 0 | 0 | 108 |
| 3 | 695877 | GGG 2014_October | Fail | 0 | 0 | 0 | 0 | 67 |
| 4 | 446761 | DDD 2014_October | Pass | 0 | 0 | 28 | 0 | 210 |
| 5 | 2481765 | BBB 2013_February | Pass | 0 | 0 | 0 | 0 | 450 |
| 6 | 1618172 | EEE 2014_February | Pass | 0 | 1 | 0 | 0 | 246 |
| 7 | 693629 | FFF 2014_October | Fail | 0 | 0 | 0 | 0 | 7 |
| 8 | 696415 | DDD 2014_October | Pass | 0 | 0 | 16 | 0 | 805 |
| 9 | 372345 | FFF 2014_October | Pass | 0 | 1 | 0 | 0 | 158 |
| 10 | 551836 | BBB 2013_October | Pass | 0 | 0 | 0 | 0 | 267 |
| 11 | 696804 | BBB 2014_October | Pass | 0 | 0 | 0 | 0 | 7 |
| 12 | 630156 | GGG 2014_February | Fail | 0 | 0 | 0 | 0 | 0 |