

## CS 327E Milestone 5 due Sunday, 03/01.

This is the first of two milestones that makes use of Apache Beam for cleansing your main dataset.

1. Review your latest tables in your `<source>_modeled` dataset. Make a list of the remaining standardization and normalization problems present in your data. For example, duplicate records in table `XYZ` caused by non-conforming dates in column `abc`. Add those issues to the `TRANSFORMS.txt` file.
2. Choose one of the tables you identified in `TRANSFORMS.txt`. Write a short Beam pipeline that cleanses the data in this table. The pipeline should satisfy the following requirements:
  - use the Direct Runner to execute the pipeline
  - run a BigQuery query that contains a `limit` clause over the table(s) in your `<source_modeled>` dataset
  - make an input `PCollection` from the BigQuery results
  - write the input `PCollection` to a local file named `input.txt`
  - apply one or more custom `DoFns` through a `ParDo`
  - write the output `PCollection` to a local file named `output.txt`
  - write the output `PCollection` to a new BigQuery table in your `<source>_modeled` dataset
  - execute the pipeline from your `<source>_modeled.ipynb` notebook.
3. Verify that the BigQuery result table from the previous step contains a primary key. If it's a child table, it must also have a foreign key. Run the SQL statements to verify these constraints from your `<source>_modeled.ipynb` notebook.

### Coding Conventions:

- The Beam pipeline should be in a file named `<table>_beam.py` where `<table>` is the name of the table that is being transformed.
- The BigQuery result table should be named `<table>_Beam` and reside in your `<source_modeled>` dataset.
- The `DoFn` code should be commented sufficiently to understand the main logic of the `transform(s)`.

CS 327E Milestone 5 Rubric

**Due Date: 03/01/20**

<p>Create a file <code>&lt;table&gt;_beam.py</code> that takes in data from your <code>&lt;source&gt;_modeled</code> dataset, performs a DoFn transform on the data, and writes it back out into another table. Sufficiently comment the code to show understanding of the Apache Beam pipeline.</p> <p>In addition, a <code>TRANSFORMS.txt</code> file should now be present for all groups. If a transformation could not be found, please refer to the TAs for assistance.</p> <ul style="list-style-type: none"> <li>-100 missing <code>&lt;table&gt;_beam.py</code> from repository</li> <li>-50 code does not implement the DoFn transform</li> <li>-50 code does not pull from or write back to your dataset</li> <li>-40 code does not write to two output files <code>input.txt</code> and <code>output.txt</code> (these text files need not be pushed to your repo)</li> <li>-50 code does not write to output table <code>&lt;table&gt;_Beam</code></li> <li>-30 code missing comments</li> <li>-40 missing pipeline run call from <code>&lt;source&gt;_modeled.ipynb</code></li> <li>-40 missing <code>TRANSFORMS.txt</code></li> <li>-20 missing or incorrect primary key verification on output table <code>&lt;table&gt;_Beam</code> in <code>&lt;source&gt;_modeled.ipynb</code></li> <li>-20 missing or incorrect foreign key verification on output table <code>&lt;table&gt;_Beam</code> in <code>&lt;source&gt;_modeled.ipynb</code> if output table is a child table</li> </ul>	<p>100</p>
<p><code>submission.json</code> submitted into Canvas. Your project <b>will not</b> be graded without this submission. The file should have the following schema:</p> <pre>{   "commit-id": "your most recent commit ID from Github",   "project-id": "your project ID from GCP" }</pre> <p>Example:</p> <pre>{   "commit-id": "dab96492ac7d906368ac9c7a17cb0dbd670923d9",   "project-id": "some-project-id" }</pre>	<p>Required</p>
<p><b>Total Credit:</b></p>	<p><b>100</b></p>