# HW #2 Grading Sheet

Student ID: _____123_____  Student Name: _____Nick Noteworthy_____

**Problem 1a (20 pts)** - Structure of the MovieLens data.

| Questions | Pts | Comments |
|---|---|---|
| Is the schema homogeneous?<br>- Movies - *yes*<br>- Ratings - *yes*<br>- Tags - *yes*<br>- Links - *yes* | 4 | Most everyone got these. |
| How are fields accessed (by name or position)?<br>- All files - *by position* | 2 | CSV's are kind of both - but most said name-based (due to header, which was a perfectly good answer. |
| How are records delimited?<br>- All files - *newline* | 2 | Most everyone got this. |
| How are fields delimited?<br>- All files - *comma* | 2 | Ditto. |
| How are field values encoded?<br>- All files - *plain text* | 2 | Ditto. |
| How complex (primitive values vs. arrays)? What about relationships (e.g., are files 1NF or not)?<br>- Movies - all *primitive* except genres (which is an *array*)<br>- Ratings - all *primitive*<br>- Tags - all *primitive*<br>- Links - all *primitive* | 2<br><br>1<br>1<br>1 | Most everyone got these; a few failed to mention genres being a pipe-delimited list-valued field. |
| What are the semantics of the data, and are validity checks possible?<br>- Could check that the *links* to other movie databases are *valid* (similar to FK checks)<br>- Could check *movie ids* when used as *FKs*<br>- Could check *user ids* when used as *FKs*<br>- Could check *sensibility of timestamp values* (e.g., in a year range after movies first appeared in the world) | 3 | There were a number of thoughts possible and given here; in grading we were mostly just looking for some attention being paid to this question, including things that could be wrong and be checked (or not). |
| **Total** | 20 | |

**Problem 1b (10 pts)** - Granularity of the MovieLens data.

| Questions | Pts | Score | Comments |
|---|---|---|---|
| What kind of thing (person, object, etc.) do the records represent? <br> - Movies - *movies* <br> - Ratings - *user movie ratings* <br> - Tags - *user annotations* <br> - Links - *cross-references to two other movie databases* | <br><br>2<br>2<br>2<br><br>2 | | Most everyone got these. |
| Are the records homogeneous? <br> - All files - *yes* | <br>1 | | Most everyone got this. |
| Any alternative interpretations of the records? <br> - *No fixed tag interpretation(s)* | <br><br>1 | | Here we just wanted something to be said about this question (showing that it was thought about). |
| **Total** | 10 | | |

**Problem 1c (20 pts)** - Scope of the MovieLens data.

| Questions | Pts | Score | Comments |
|---|---|---|---|
| What characteristics of the things are captured by the record fields? <br> - Movies - *title, genre(s)* <br> - Ratings - *user, rating, tstamp* <br> - Tags - *tag value, tstamp* <br> - Links - *remote movie ids* | <br><br>4<br>6<br>4<br>2 | | Most everyone got this, though a few forgot to mention the timestamps (which were the rating and tag timestamps, not movie timestamps). |
| Are related fields consistent? | (0) | | Here we looked at comments but didn't grade the commentary. |
| Can additional characteristics be deduced from the characteristics? | (0) | | Ditto. |
| Do records all have the same fields? <br> - All files - *yes* | <br>1 | | Most everyone got this. |
| Do the records represent the entire populations of their things? <br> - Subset of *all movies* | <br><br>1 | | Here the sought-after point was that not every movie ever made, or perhaps every user, was in the data set. |
| Are there multiple records for the same thing (requiring de-dup'ing)? <br> - All files - *no* | <br><br>1 | | Here we just looked for some thoughts. |
| Any heterogeneous sets of records? <br> - All files - *no* | <br>1 | | Most everyone got this. |
| **Total** | 20 | | |

**Problem 2 (50 pts)** - Data wrangling.

| Questions | Pts | Score | Comments |
|---|---|---|---|
| a. Misplaced value to extract?<br>  -  Movie year<br>  -  Could use string regex | 5<br>5 | | Here we wanted both identifying the value (year) and mentioning how it might be extracted. |
| b. Lifting of one->several fields?<br>  -  IMDB: primary name<br>  -  Could separate: first + last | 5<br>5 | | Here we were thinking of IMDB actor names, but we also took other answers that were "lifting-like" as valid answers. |
| c. Aggregation - ratings example?<br>  -  Could aggregate ratings:<br>    *movieId,avgRating*<br>    31,2.75<br>    66,4.23<br>    ... | 5<br><br>5 | | Most folks got the idea here. We wanted both a written example and a data example of what it would look like. |
| d. Movies file problems/issues?<br>  -  Non-1NF field (genres) should be normalized<br>  -  Misplaced value (year) should be extracted<br>    *movieId,title,year:*<br>    1,Toy Story,1995<br>    ...<br>    *movieId,genre:*<br>    1,Adventure<br>    1,Animation<br>    ... | 5<br><br>5 | | Here we wanted all of the issues to be mentioned - both the misplaced year and the fact that genres being non-1NF would be a pain for queries/analysis. |
| e. How could reviews be combined across IMDB and MovieLens?<br>  -  Need to ***union*** the ratings to get them into one dataset.<br>  -  Must first ***standardize*** their values (e.g., to a 1-10 scale, multiplying MovieLens values by 2 after averaging them by movie):<br>    *movieId,avgRating*<br>    31,3.50<br>    66,6.37<br>    ... | 4<br><br>2<br><br>2<br><br><br>2 | | ATTENTION: Most folks did NOT get this, at least not in the way we wanted. To deal with scores across databases, scores on different scales, some work is needed. One DB had averages, so first the other DB had to be averaged. Also, the score scales were different (by 2x), so one would need to equalize the scales. Then - and only then - one could "union" the results, as in, combine the scores TOGETHER to have a score that reflects both. Aggregation and normalization were the big asks here. |
| **Total** | 50 | | |

**Total Score: _____ ( = _____ + _____ + _____ + _____ )**