

Course Project

CS 348 - Winter 2020

(Total weight 25%)

Project Description

Important Dates

- Project mixer (in class): Thursday, Jan 23
- Milestone 0 (preparation): Thursday, Jan 30 [Not Graded]
- Milestone 1 (proposal): Thursday, Feb 06
- Milestone 2 (mid-term report): Thursday, Mar 05
- Early in-class demo: Thursday, April 02
- Other demo slots: April 02-03 (TBD)
- Final report and code submission: due by the scheduled demo slot; updatable until Thur April 09

Overview

The project is to build a database-driven application from the ground up. There will be some examples and instructions to help you get started. No prior experience in developing such applications is assumed. This document also describes a number of possible ideas.

Submission and Grading

There will be three milestones and a final project demo; see Important Dates above for due dates. You will find the details of what to submit for each checkpoint later in this document under respective sections. Only one member needs to submit through MarkUs (<https://www.student.cs.uwaterloo.ca/~markus/>) on behalf of the entire project team; make sure all members are added to the submission. Because of the open-ended nature of course projects, certain instructions may not apply to your particular project; when in doubt, consult the instructor.

Each project will be graded on a scale of 0-100 points. A breakdown is as follows:

- 1) 30 points for submitting the required work at three checkpoints (milestone 1, milestone 2, and demo). (milestone 0 will not be graded);
- 2) 40 points for completing the proposed work;
- 3) 30 points for the quality of the work.

Out of the 70 points (2 & 3) for completeness/quality, 5 to 10 points are reserved for impressive and/or innovative work beyond what is expected. In other words, meeting the expectation will ensure a project grade of at least A-, but A and A+ will require exceptional work.

What the “required work” means may evolve over the course of the project. Milestone 0 ensures that each student finds a team and get familiar with the environment setup. Milestone 1 proposal helps you get a feel for the amount of work involved, and work with you to ensure that it meets the minimum requirements of depth and scope for a course project.

Teamwork

The project should be completed in (4 or 5)-person teams. Any other team size requires explicit approval from the instructor (please make your request through Piazza private questions); team sizes below 4 or above 5 are strongly discouraged. Regardless of the team size, an equal amount of work is expected and the same grading scale will be applied. All members in a team will receive identical grades for the project. You are required to report each team member's effort and progress in the milestone and final reports. If there is any problem working with your team members that you cannot resolve by yourself, bring it to the instructor's attention as soon as possible. Last-minute complaints of the form "my partner did nothing" will not be entertained. There will be a "project mixer" class (see Important Dates above) where you can pitch your project ideas by giving a short presentation and the number of additional team members you are looking for. Students looking for projects can talk those giving presentations to see if there is a match. Even if you do not have a concrete idea by the time of the mixer, you can still give a short presentation saying what domains and types you are interested in, and look for teams. Besides the mixer, you are also welcome to use Piazza to share ideas and recruit team members.

Platform Issues

To develop the project, there are two main options for platforms: (i) DB2 on the servers (e.g. `linux.student.cs.uwaterloo.ca`) from Assignment 1; (ii) cloud-based platforms (e.g. Google Cloud Platform). We provide some sample codes and instructions for you to start with one of two options. If you wish, you may use other languages, tools, or application development frameworks, or run servers on your own machines. Setting up the whole application/database stack is non-trivial and can be a rewarding experience. However, the course staff can only support the technologies used by the provided examples.

Tasks Overview

The project is to build a database-driven application. Specifically, you will need to complete the following tasks over the course of this semester. Note that different members of a team can work on some of these tasks concurrently.

1. Pick your favorite data management application. It should be relatively substantial, but not too enormous. Several project ideas are described at the end of this document, but you are encouraged to come up with your own. When picking an application, keep the following questions in mind:
 - a. How do you plan to acquire the data to populate your database? Use of real datasets is highly recommended. You may use program-generated "fake" datasets if real ones are too difficult to obtain.
 - b. How are you going to use the data? What kind of queries do you want to ask? How is the data updated? Your application should support both queries and updates.
2. Design the database schema. Start with an E/R diagram and convert it to a relational schema. Identify any constraints that hold in your application domain, and code them as database constraints. If you plan to work with real datasets, it is important to go over some samples of real data to validate your design (in fact, you should start Task 7 below as early as possible, in parallel to Tasks 3-6). Do not forget to apply database design theory and check for redundancies.
3. Create a sample database using a small dataset. You may generate this small dataset by hand. You will find this sample database very useful in testing, because large datasets make debugging difficult. It is a good idea to write some scripts to create/load/destroy the sample database automatically; they will save you lots of typing when debugging.

4. Design a user interface for your application (e.g. command line-based on school serve, or web-based). Think about how a typical user would interact with the database. Optionally, it might be useful to build a “canned” demo version of the interface first (i.e., with hard-coded rather than dynamically generated responses), while you brush up your user interface design skills at the same time. Do not spend too much time on refining the look of your interface; you just need to understand the basic “flow” in order to figure out what database operations are needed in each step of the user interaction.
5. Write SQL queries that will supply dynamic contents for the interface you designed for Task 4. Also write SQL code that modifies the database on behalf of the user. You may hard-code the query and update parameters. Test these SQL statements in the sample database.
6. Choose an appropriate platform for your application. JDBC? Python or PHP? JavaScript or plain HTML? Start by implementing a “hello world” type of simple database-driven application, deploy it in your development environment, and make sure that all parts are working together correctly. The course website will provide pointers to working examples.
7. Acquire the large “production” dataset, either by downloading it from a real data source or by generating it using a program. Make sure the dataset fits your schema. For real datasets, you might need to write programs/scripts to transform them into a form that is appropriate for loading into a database. For program-generated datasets, make sure they contain enough interesting “links” across rows of different tables, or else all your join queries may return empty results. Keep in mind that the school servers has limited capacity: for larger databases, you may need to consider cloud-based platforms.
8. Test the SQL statements you developed for Task 5 in the large database. Do you run into any performance problems? Try creating some additional indexes to improve performance.
9. Implement and debug the application and the user interface. Test your application with the smaller sample database first. You may need to iterate the design and implementation several times in order to correct any unforeseen problems.
10. Test your application with the production dataset. Resolve any performance problems.
11. Polish the user interface. You may add as many features as you like, though they are optional because they are not the main focus of this course.

Milestone 0: Project Preparation

You should have formed a team and started thinking about task 1-6. In particular, your team should get familiar with one of the platforms and be able to run “hello world” type of simple database-drive applications (task 6). Some examples (ProjectNote-GCP.pdf, ProjectNote-JDBC-DB2.pdf) are provided on Learn under Course Project to help you start. The goal of milestone 0 is to ensure that you have found a team, and your team has already tried out the chosen programming platform and environment (as switching platforms later will be painful and time consuming).

Submit the following electronically under “Project Milestone 0” via MarkUs by **Thur, Jan 30, 11pm**

1. A text file **members.txt**, listing the members of your team, and for each member, a description of effort and progress made by this member to date.
2. A progress report **report.pdf** should at least contain:
 - A brief description of your application: which dataset do you plan to use, who are the users of the application, etc.

- A brief description of your choice of platform and user interface. How users interact with your application (e.g. command line or web, etc.)?
3. A .zip or .tar.gz archive of your source code **code.zip**. The source code directory should at least contain:
- README.txt to describe how to create and load your sample database to your chosen platform. You don't have to use the datasets described in the report. Toy datasets (e.g. a single table) can be used.
 - Source code for your "hello world" type of simple database-driven applications with your sample database (e.g. the application allows user to connect to the database, to select some rows from a table, etc.)

Milestone 1: Project Proposal

You should have completed Tasks 1-5 and have started thinking about 6 and 7. If you plan to work with real data, you should also have made significant progress on Task 7 (you should at least ensure that it is feasible to obtain the real dataset, transform it, and load it into your database). **Submit** the following electronically under "Project Milestone 1" via MarkUs by **Thur, Feb 06, 11pm**:

1. A text file **members.txt**, listing the members of your team, and for each member, a description of effort and progress made by this member to date.
2. A progress report **report.pdf** should at least contain:
 - A brief description of your application.
 - A plan for getting the data to populate your database, as well as some sample data.
 - A list of assumptions that you are making about the data being modeled.
 - An E/R diagram for your database design.
 - A list of database tables with keys declared.
 - A description of the user interface. You can write a brief English description of how users interact with the interface (e.g., "the user selects a car model from a pull-down menu, clicks on the 'go' button, and a new page will display all cars of this model that are available for sale"). Or, instead, you can submit a canned demo version of the website.
3. A .zip or .tar.gz archive of your source code **code.zip**. The source code directory should at least contain:
 - A README file describing how to create and load your sample database.
 - Files containing the SQL code used for creating tables, constraints, stored procedures and triggers (if any).
 - A file test-sample.sql containing the SQL statements you wrote for Task 5.
 - A file test-sample.out showing the results of running test-sample.sql over your sample database.
 - Code implementing a simple but working database-driven application on your chosen platform (1-2 simple tasks/ can be run, e.g. adding a row), which can serve as a starting point for completing your project.
 - If applicable, any code for downloading/scraping/transforming real data that you have written for Task 7 so far.

Milestone 2: Project Mid-term Report

You should have completed Tasks 1-8 and have made good progress on 9. **Submit** the following electronically under “Project Milestone 2” via MarkUs by **Thur, Mar 05, 11pm**:

1. A text file **members.txt**, listing the members of your team, and for each member, a description of effort and progress made by this member since the last milestone.
2. A progress report **report.pdf** should at least contain:
 - New assumptions, E/R diagram, and list of tables (if they have changed since Milestone 1).
 - A brief description of the platform you chose in Task 6.
 - Changes you made to the database during performance tuning in Task 8, e.g., additional indexes created.
3. A .zip or .tar.gz archive of your source code **code.zip**. The source code directory should at least contain:
 - A README file describing how to generate the “production” dataset and load it into your database. Do not submit the production dataset itself through if it is too big; instead, submit the URL where you download/scrape the raw data (if applicable), and the code that extracts and transforms (or generates) the production dataset.
 - A file test-production.sql containing the SQL statements you wrote for Task 5. You may wish to modify some queries to return only the top 10 result rows instead of all result rows (there might be lots for large datasets).
 - A file test-production.out showing the results of running test-production.sql over the production dataset.
 - Code implementing a simple but working database-driven application on your chosen platform. This code includes 1-2 more tasks than the version in milestone 1.

Project Demo & Final Report and Code

At the end of the semester, you will need to present a working demo of your system. Instructions on how to sign up for the demo will be given during the second to last week of the class. Prior to your demo, **submit** the following electronically under “Project Final”:

- A text file **members.txt**, listing the members of your team, and for each member, a description of effort made by this member throughout the semester, highlighting the effort since the last milestone.
- A final project report **report.pdf**, including a brief description of your application, the E/R diagram for your database design, assumptions that you are making about the data being modeled, and the list of database tables with descriptions.
- A .zip or .tar.gz archive of all your source code **code.zip**. The source code directory should also contain a README file describing how to set up your servers and database, and how to compile and deploy your application.

The code and report are updatable until **Thur, April 09, 11pm**.

Project Ideas

Below is a list of possible project ideas for which high-quality datasets exist. Of course, you are welcome to come up with your own.

Entertainment, sports, or financial websites

Examples include those that allow visitors to explore information about movies, music, sports, games, stocks, etc. There are already many commercial offerings for such purposes. While there is less room for innovation, there are plenty of examples of what a good website would look like, as well as high-quality, well-formatted datasets. For example, IMDb makes their movie database available <http://www.imdb.com/interfaces>; historical stock quote can be downloaded and scraped from many sites such as Yahoo! and Google Finance. This project is well-suited for those who just want to learn how to build database-backed websites as beginners. You can always spice things up by adding features that you wish those websites had (e.g., different ways for summarizing, exploring, and visualizing the data).

Websites providing access to datasets of public interest

If you are interested in doing some good to society while learning databases, this project is for you. There are many interesting datasets “available” to the public, but better ways for accessing and analyzing them are still sorely needed. Here are some examples:

- Data.gov (<http://www.data.gov/>) has a huge compilation of data sets produced by the US government.
- The Supreme Court Database (<http://scdb.wustl.edu/data.php>) tracks all cases decided by the US Supreme Court.
- US government spending data (<https://www.usaspending.gov/>) has information about and database downloads of government contracts and awards.
- Federal Election Commission (<https://www.fec.gov/data/>) has campaign finance data to download as well as nice interfaces for exploring the data.
- GovTrack.us (<http://www.govtrack.us/developers>) tracks all bills through the Congress and all votes casted by its members. The Washington Post has a nice (albeit outdated) website (<http://projects.washingtonpost.com/congress/113/>) for exploring this type of data (in predefined ways), but you can be creative with additional and/or more flexible exploration and analysis options. Vote Smart (<https://votesmart.org/>) has a wealth of additional, useful information on votes, such as issue tags, synopses and highlights.
- Each state legislature maintains its own voting records. For example, you can find North Carolina’s here: <http://www.ncleg.net/Legislation/voteHistory/voteHistory.html>. Some states provide records in already structured formats, but for others, you may need to scrape their websites.
- The Washington Post maintains a list of datasets (<http://www.washingtonpost.com/wp-srv/metro/data/datapost.html>) that have been used to generate investigative news pieces. Most of these datasets hide behind some interface and may need to be scraped. Use this list for examples of what datasets are “interesting” and how to present data to the public effectively.
- Stanford Journalism Program maintains a list of curated transportation-related datasets (<http://www.datadrivenstanford.org/>).

- National Institute for Computer-Assisted Reporting maintains a list of datasets of public interest (<http://www.ire.org/nicar/database-library/>). Use this list for examples of what datasets are “interesting”—they are generally not available to the public, but there may be alternative ways to obtain them.

Your task would be to take one of such datasets, design a good relational schema, clean up/restructure the data, and build an application for the public to explore the dataset. You are welcome to come up with your own.