# Berkeley Data Analytics Final Project: Predictive Model for Building Project Bids

**Class Presentation by:**
Loc Nguyen
Ashok Ramaswami

Speaker notes - Ashok:
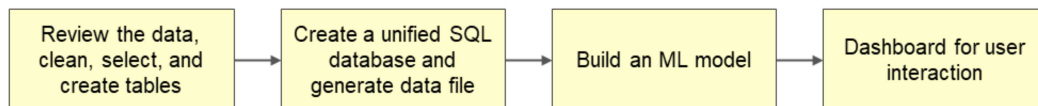- We're working with real world data from a company that provides a variety of services for buildings
- The company has a very successful track record and reputation with over 1300 completed projects
- An aspect of the work that the company would like to improve is their bid evaluation and pricing
- This is because about 50% of their projects have been over budget
- The question we'd like to address is whether an ML model can predict if a given project is likely to be over or under budget given the type, scope, and other details of a project
- Gaining some insight on this could help the company bid more accurately and manage better the budget for projects that fit certain patterns
- Our goal: Build a model that predicts the likelihood of a project being under or over budget

Speaker notes - Ashok:
Our methodology consists of four steps, broadly. We'll expand on this in subsequent slides.
- 1. Clean up data: using python with pandas and numpy libraries. Then use sqlalchemy to load the data into SQL Server.
- 2. Create a database: the ETL returns 4 tables to load on our server, then we use SQL to merge them before extracting the cleaned data.
- 3. Processing the cleaned data, then build and train a ML model.
- 4. Build a deliverable dashboard for the users to apply that model

## ETL

```
Run the Master function to import data into SQL Database:

▷ ▶≡ M↓

    master_ETL(building_path, project_path, hour_path, emp_path)

Start cleaning
 Clean building complete!
 Clean project complete!
 Clean hour complete!
 Clean emp complete!
Finish importing all df into SQL database.
Done. Full ETL ran successully in 9.968253374099731 seconds.
```
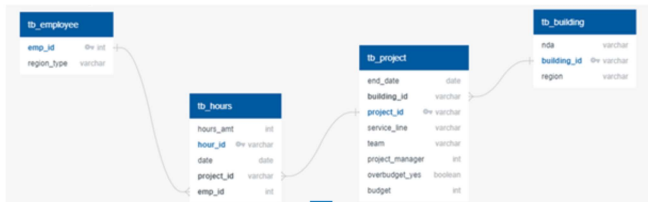
ETL - Loc:
Step 1 was to load all the raw data, cleaned them up (remove or fill na), kept just the columns that we needed, convert the datatypes, and rename them into something meaningful. Then we built a few functions to have all of them auto-loaded the next time we have a new raw dataset.
The snippet showed it took about 10 seconds to finish the ETL with 4 tables.

Speaker notes - Loc:
- After the ETL process, we move on to the 2nd step to work on our server. Our original data is in 4 tables: Building, project, hours, and employees.
- We use SQL to:
  - Merged the project table with the building table in order to extract the building information for each project,
  - Merged the hours table with employees table to get the employee region for each hour entry.
  - Completed data exploration to remove features that may overfit the model. For e.g: after merging the hours table, we found out that almost everyone was in the same region. So the employee region will not be a good feature anymore. Thus, we removed it in our cleaned df.

Speaker notes - Ashok:
- Our data was fairly clean but still needed some preprocessing. This included:
    Summing the number of hours from the hour table to the project
    Calculating the project month duration using the first date/last date of the
    hours table and the project end date
    Counting the number of employees working on each project.

    Following this cleanup process, we prepared the dataset for one-hot-coding,
    scaling, and generating a model. We have 10 features for the ML model.

# Training and Test of ML models

| MODEL | ACCURACY |
|---|---|
| SVM | 0.729 |
| Random Forest | 0.681 |
| Logistic Regression | 0.492 |

Speaker notes - Ashok:
We created a dataset based on the feature engineering that we described in earlier slides
- Specifically we created ML models with three approaches:
  - ○ SVM
  - ○ Random Forest
  - ○ Logistic Regression
- We applied the standard form for separating training and test data
- We also created a version of the ML model using binning for one of the features
- We then reviewed the model accuracy for each approach

As an additional step, we ran a cross validation to evaluate model accuracy under different test samples

Speaker notes - Loc:
- Our goal is to create a simple dashboard that would take some user input parameters (such as project budget, duration, number of employees, estimated project manager, estimated hours, etc.) and predict the target
- To implement the dashboard we:
    Download the ML model using pickle library
    Developed a Flask app with python to take and process user input
    Developed UI code to display user input fields and display a prediction

## Demo

- We'll take 2-3 minutes for a demo of our code

Loc:
- Over-budget: 'Commissioning', 'SB', Burke, 110700, 'SF', True, 1200.0, 50, 25
- Under-budget: 'WELL', 'SB', Ashok, 56221, 'SF', True, 200, 15, 6

## Suggested Future Work

- Predictive value
  - Build a linear regression model to return a predicted budget instead of logistic regression.
  - Analyze how much impact each feature contribute to the output result, thus allow end-user to have better adjustment on their input.

- Improved data
  - We believe that developing a richer dataset, in terms of additional features and completeness, will enhance the model accuracy.

- Enhanced dashboard
  - We recommend additional work to build a highly interactive dashboard that can be integrated with workflows for pricing and project management to manage outcomes more proactively.
  - The model is vulnerable to outlier, so we need to find that outlier range to better guide the user input.
  - The dashboard is limited to local computer only and not shareable yet.

Loc:
1. Predict how much should we bid for a particular project, similar to a pricing tool.
   Consider adding/reducing other features to further increase the accuracy, but careful not to overfit the model.
2. The more data, the better our model.
   Make sure to collect only clean data.
3. - The model is vulnerable to outlier (such as budget as low as $1, or 1000+ employees, but still return "Under-budget").
   - The end-user is usually non-tech people, so this dashboard has many rooms to improve. Such as: include some pie chart/report, save/export multiple predictions/scenarios, provide some suggestion on what should the user do to have the expected outcome (including limiting the range of input to avoid outliers).
   - The flask app is not shareable yet.