

## CS 536 : Final Project - Data Completion and Interpolation

16:198:536

For your final project, you will take a data set of the form  $\{\underline{x}^1, \underline{x}^2, \dots, \underline{x}^m\}$ , and construct a system for predicting or interpolating missing features (frequently given as empty or NA) from the present features in new records. Notice that in this case, instead of singling out a single feature or factor for prediction/regression/classification, *any* feature of  $\underline{X}$  might be missing and need to be predicted or reconstructed from the features that are present - and the features that are present may vary from instance to instance.

Before getting into the specifics of the data, here are some problems that you will have to address when coming up with a solution:

- How to represent or process the data. Data features may contain a number of diverse data types (real values, integer values, categorial or binary values, ordered categorial values, open/natural language responses). How can you represent these for processing and prediction?
- How to model the problem of interpolation. What are the inputs, what are the outputs? An important if subtle question to consider here - what does it mean to predictor or interpolate a missing feature?
- Model selection. What kind of model or models do you want to consider?
- Quantifying loss or error. How can you quantify how good a model is, how to measure its loss/error? This is important not only in terms of evaluating your model, but in terms of training as well - how can you refine and improve your model without a way of comparing them?
- Training. What kind of training algorithm can you apply to your model(s)? What design choices do you have to make here?
- Feature selection. It is frequently useful in learning problems to focus on specific features and exclude others, to try to eliminate spurious features and focus on what matters. How can that be applied here?
- Validation. How can you prevent or avoid over-fitting? Can you apply the usual training/testing/validation paradigm to this problem? How do you choose the training or testing data? Note that a record won't need to be complete to still be useful, potentially, in interpolation. Can cross-validation be applied here? This can be especially important when the data set is not overwhelmingly large and data must be used carefully.
- Evaluation. How good is your final model? How can you evaluate this? What are the limits and strengths of your model - how many features does a new record actually need to be able to interpolate well?

These are important questions to answer when building and evaluating machine learning systems, and I expect thorough treatments of each in your final writeup.

## 1 The Data

For this project, you may choose one of two datasets to try to model. Both come from psychology studies in the Many Labs series, an attempt to test the replicability or generalizability of psychological effects. In both cases, multiple labs attempted the same experiments, to determine the extent the results of those experiments generalized and could be relied on in diverse circumstances. Being psychology studies, the data sets are records of personal answers and reports from subjects on a variety of questions relating to topics like perception and mental biases. Additionally, the data sets include demographic information about the subjects and information about the researchers and experimental environments. All are potentially useful for the problem of prediction and interpolation.

You may choose either data set for this project. In either case, as described, given a new or partial record, you want to be able to predict or infer the values of missing features (any missing feature!) based on the features that are present. The following are summaries of the two projects, and descriptions of the provided project files (also available through the websites).

- **The Original Many Labs Project:** The original Many Labs project ([project website](#)) attempted to replicate 28 psychological studies, across 60 different labs, trying to determine to what extent the originally studied effect was reproducible. Questions given to subjects touched on a diverse array of topics from nationalism, to the perceptions of numbers, to feelings about art and mathematics.
  - **The Main Data Set:** Tab.delimited.Cleaned.dataset.WITH.variable.labels.csv *Note - as the title suggests, this is a tab-delimited file.*
  - **An Explanation of the Variables and Values:** Codebook.xlsx
  - **The Survey Questions:** DetailedCodebook.HTML.pages.zip
  - **Extraneous Data:** Datasets.zip
- **Many Labs 3:** Many psychology studies depend on drawing their subjects from a pool of undergraduate students willing (or compelled) to take part in the studies. Many Labs 3 ([project website](#)) attempts to determine the quality of these subjects for experimentation, by looking at how attitudes and mental states vary across the semester. If there is a significant time-of-semester effect, this would imply that the results of your experiment could be influenced by when in the semester you chose to do it. Questions touched on here include various self reports on attitudes about work, challenges, and self-perception.
  - **The Main Data Set:** ML3AllSites.csv
  - **An Explanation of the Variables and Values:** ML3 Variable Codebook.xlsx
  - **The Survey Questions:** ML3\_Computer\_Scripts\_by\_collection\_site.zip
  - **Extraneous Data:** ML3 Final Data.zip

*A cautionary note for this data set: this data set seems to have a lot of variability in how people were able to answer questions.*

## 2 Requirements

You must submit any code written to complete the assignment (and any supplemental model files, etc). The code should be well-commented and clear. Additionally, the code should be your own - any learning algorithms you implement should be your own work, and not rely on other people's algorithms or external libraries. You may use existing libraries or frameworks to help *represent* your model and handle it in code, but you are not allowed to use any built in training or optimization algorithms, or any pre-coded machine learning algorithms. Your system must learn based on your own algorithms.

Additionally, you must submit a writeup that includes and addresses the following points, quantifying relationships and plotting data where appropriate:

- **Describe your Model:** What approach did you take? What design choices did you make, and why? How did you represent the data? How can you evaluate your model for goodness of fit? Did you make an effort to identify and exclude irrelevant variables? How did you handle missing data?

- **Describe your Training Algorithm:** Given your model, how did you fit it to the data? What design choices did you make, and why? Were you able to train over all features? What kinds of computational tradeoffs did you face, and how did you settle them?
- **Describe your Model Validation:** How did you try to avoid overfitting the data? How did you handle the modest (in ML terms) size of the data set?
- **Evaluate your Model:** Where is your model particularly successful, where does it lack? Does it need a certain amount of features in order to interpolate well? Are there some features it is really good at predicting and some it is really poor at predicting? Why do you think that is?
- **Analyze the Data:** What features were particularly valuable in predicting/interpolating? What features weren't particularly useful at all? Were there any features about the researcher/experimental environment that were particularly relevant to predicting answers - and if so, what conclusions can you draw about the replicability of those effects?
- **Generate Data:** Use your system to try to generate realistic data, and compare your generated data to the real data. How good does it look? What does it mean for it to 'look good'?

*A Note on Data Types:* In both data sets (though moreso in Many Labs 3) there are a number of features that are given in the form of natural language answers. Utilizing natural language data can be difficult, so you may ignore these data points for the purpose of this project.

### 3 Bonus

This project has three potential bonuses available.

- **Bonus 1:** Build your system to include processing and analysis of natural language answers, to try to use them to inform prediction of other features. How can you represent natural language answers in a useful way? What language processing is necessary to be able to use them? Do they actually provide useful information for the prediction problem?
- **Bonus 2:** Build your system to include prediction about the natural language answers. Based on available features, what can you predict about a person's answers to natural language questions (if not the actual answers themselves)? How can you assess the quality of these predictions?
- **Bonus 3:** Build your system to include generation of natural language answers to questions, based on available features. How can you model this generation problem, and how can you evaluate the quality of the answers?