

# **COURSE MEDIAN PREDICTION VIA SYLLABI ANALYSIS**

CORALIE PHANORD, GRAESON MCMAHON, KELSEY JUSTIS

## **1. PROBLEM STATEMENT**

College students often find median grades helpful in the course selection process. Knowledge of a course's median grade provides insight on the difficulty of their courses enabling an individual to set a well-balanced schedule. In this project, supervised machine-learning algorithms are used to predict median grades given a course syllabi, or more specifically the identifiable features present in the syllabi.

## **2. METHOD**

We are considering two regression-based approaches for our problem. Currently, we are leaning towards decision tree algorithms such as Ross Quinlan's C4.5, as they excel at handling the mixture of numerical and categorical data which characterizes our feature set. We are also considering artificial neural networks.

## **3. DATA**

The data used provides training examples in the form of syllabus-median grade pairs; the median grades serve as points to be fitted by algorithms learning features that are present in the corresponding syllabi. Important characteristics of the data are detailed below:

### **3.1 Sources**

Data publically available via department/college websites for Dartmouth and other colleges and from the Registrar's Office/department databases (requests currently pending).

### **3.2 Training Set**

Syllabi-median grade pairs from random (partitioning process to be determined) quarters.

### **3.3. Testing Set**

Syllabi-median grade pairs from random (partitioning process to be determined) quarters.

Although not present before the project due date, we hope to also test on present term's data.

### **3.4 Number of Examples**

An exact figure is currently being determined; modest estimates (without approved department requests) suggest 500+ examples are obtainable via Dartmouth websites alone.

### **3.5 Number of Features**

This is expected to vary throughout the feature selection process as we test and discover features existing in the full text of the syllabi.

### **3.6 Types of Features Present**

The text each syllabi contains offers uncountably many data points, because of this, we will test and discover how document formatting/aesthetics, language rhetoric, word connotation, and word frequency factor into a course's median grade. Tentative features include:

- Syllabus length
- Proportion of bolded, italicized, and underlined words
- Department origin
- Course Number
- Occurrence of words: "quiz," "exam," "midterm," "test," and "project", "labs"
- Occurrence of terms suggesting course rigor such as: "late", "penalized", "mandatory", and "prerequisite"
- Frequency of negative words such as: "no" and "not"

- Number of X-hours
- Number of Teaching Assistants
- Number of “%” symbols
- Class timeslot

#### 4. MILESTONE GOAL & PROJECT TIMELINE

##### 4.1 Collect and format data for processing

*Accomplished by January 29*

Data source and format choice are finalized; data is ready to process.

##### 4.2 Flexible parser to extract features from syllabi

*Accomplished by February 12*

Code is developed and capable of scanning syllabi for features in raw text and document formatting.

##### 4.3 Initial development of chosen algorithm

*Started by February 17 (Milestone due date)*

Algorithm used to learn features from syllabi via parser is chosen. Early work is started with emphasis on code structure; outside implementations of chosen algorithm are examined for best practices.

#### 5. REFERENCES

**Data for Dartmouth College Course Medians**

<http://www.dartmouth.edu/~reg/transcript/medians/>

**Sample Data for department course syllabi**

**Biology:** <https://biology.dartmouth.edu/undergraduate/courses-and-syllabi>

**Computer Science:** <http://web.cs.dartmouth.edu/undergraduate/undergraduate-courses>