Peter McKay, Daniel Walinsky
2015-02-19

# Project Proposal

For our final project, we intend to tackle the following Kaggle problem: the Forest Cover Type Prediction. In brief, we wish to, given a collection of geographic and cartographic variables that apply to a given 30 x 30 meter chunk of forest, determine what kind of tree coverage will be most prevalent in that region.

Examples of data fields we have access to include the soil type (Cryaquolis, Granile, Rock outcrop), the elevation, and average slope.

We are faced, then, with a multi-class classification problem. From such data fields in our testing file, we wish to discover whether a particular cell can be categorized as Spruce/Fir, Lodgepole Pine, Ponderosa Pine, Cottonwood/Willow, Aspen, Douglas-fir, or Krummholz.

Our first step must be to acquire the data and examine it, in order to discover what sort of model would best suit our data. To that end, we must first settle upon a set of tools to assist us in this matter. We plan to carry out most of our research using Python. More specifically, Python 3.4.2, with the help of a number of libraries designed to assist in machine learning and data analysis.

**pandas: Python Data Analysis Library:** We will make use of data structures and input mechanisms introduced in pandas. These data structures, while containing a great deal of complex functionality, extend iterables with sufficient cleanliness to allow for compatibility with most other python libraries.

**matplotlib:** The first of these such libraries that sounds helpful is matplotlib. matplotlib extends python with some matlab-like functionality, allowing us to generate graphs and examine the data in question. Should we need to scale the data, or remove a number of corrupted entries, matplotlib should give us an idea of where to start

**numpy:** numpy provides a collection of very useful tools for analysis, which will come in handy as we examine pairs of

**scikit-learn:** In an effort to avoid duplicating the wheel (and running out of RAM while we try to get that wheel to load), we will be making use of scikit-learn to actually carry out our machine algorithms.

Random forests and decision trees would be poetically appropriate, but possibly less than optimal.

intro

related work

methodology

results

reflection

conclusion