# Deploying predictive models with the Actor framework

Brian Gawalt

October 14, 2015

**Abstract**

The majority of data science and machine learning tutorials focus on generating models: assembling a dataset; splitting the data into training, validation, and testing subsets; building the model; and demonstrating its generalizability. But when it's time to repeat the analogous steps when using the model in production, issues of high latency can arise. To an end user, the cost of too much time spent featurizing raw data and evaluating a model over features can wind up erasing any gains a smarter prediction can offer. Exposing concurrency in these model-usage steps, and then capitalizing on that concurrency, can cut down on that latency cost. This work describes how the Actor framework can be used to bring a predictive model to a real-time setting. Two case-study examples are described: a simple text classifier (with accompanying code), and a live deployment built for the freelancing platform Upwork.

## 1 Introduction

tk tk