

Drivendata.org is running a competition, to support the Predict the Next Pandemic Initiative, for the creation of a model that can predict the number of dengue fever cases reported each week in San Juan, Puerto Rico and Iquitos, Perú using data obtained from the US government. The data for this project can be found at [drivendata.org/competitions/44/dengai-predicting-disease-spread/](https://drivendata.org/competitions/44/dengai-predicting-disease-spread/). According to the CDC about 2.5 billion people live in places where they are at risk of contracting dengue fever and an estimate of 50 to 100 million cases occur yearly. Cases of dengue have been increasing in recent years and will most likely continue to as the global climate becomes hotter. Accurately predicting dengue fever outbreaks is crucial for public health and can help advise the planning of preventative measures. Models created for this competition will be used to by the government to help aid in predictive analysis. Dengue is spread by mosquitoes and thus correlates to environmental factors. The data set will include features like climate and precipitation. In total the data set contains twenty two features. There are multiple models that can be used to create a predictive model for this problem. The CDC uses C-algorithms for modeling disease outbreak and that will be one of the approaches I plan on using. I also plan on using random forest to create a second model because of the flexibility allowed and the tendency for creating a model that fits well. Since there is a separate set of test data provided it should help reduce any overfitting. After creating those two predictive models and I will evaluate any shortcomings and see if any other machine learning algorithms will be able to remedy them. I plan on delivering the annotated code with a slide deck for this project.