**Prediction for the countries/areas which are most likely to be affected by the virus in the next decade**

The Problem approach and the dataset description

Assumptions:

1. There has been lots of missing data in the 27 lines .Not sure what is the purpose of that data. I am not considering that data for the time being. Assume some of the missing data can be filled by the regions/countries likely effected
2. Not considering the Data with year format XXXX-XXXX which represents duration of year.
3. There is no description on the "STATUS" column
4. The output to be in an **csv** file having list of countries and effected region coordinated where occurrence time is less than 5 years

The Dataset

1. (42042, 15), 15 Columns of data
2. `Categorical values`
   `['VECTOR','SOURCE_TYPE','COUNTRY_ID','infection_source',`
   `'infection_time','X','Y']`

`The Problem:`

`1. Since there is no clear distinction about a value to be predicted this is clearly an` **`Unsupervised problem.`** `Since we need to classify and segment the data as we do not know what value to predict and based on that we can provide the probability of region to be in which region.`

`2.We can hence create cluster of the region (segment the regions) based on the High and Low probability of infection time.`

`3. If the "infection time" is less than 5 years high likely to be in the High risk cluster and more than 5 will be in Low risk cluster.`

`4.Using K-means (algorithm approach) we can find the no of best suited clusters and we can classify as "Risk region"`

`Solution approach:`

1. `The Vector, infection-time ,infection-source has a direct relation with the disease occurrence as from the visualization graphs and how they have effected over the time.`
   *Aedes aegypti* - Dengue
   *Ae. Albopictus –Chickungunya*

Feature heatmap

Countries effected by Source of Vectors

Countries effected by Vectors

The Elbow Method

Clusters of Countries