

## Preparación de Datos.

Tras la investigación de los datos, hemos encontrado que hay ciertos DataSets (**df\_distance\_km**, **df\_distance\_min** y **df\_historic\_order\_demand**) que requerían de ciertas imputaciones de datos.

Para '**df\_distance\_km**', vemos que hay distancias de 0 entre diferentes clientes cuando la distancia es mayor. Para solventar este problema, hemos decidido hacer uso de la **fórmula de la distancia euclidiana**, que obtiene una distancia entre 2 coordenadas (lat, long) aproximada en kilómetros.

Para '**df\_distancia\_min**', tenemos tiempos de 0, dándose el mismo "fallo" que en el '**df\_distance\_km**', por lo que hemos imputado los nuevos datos gracias a las nuevas distancias corregidas. Para calcular los tiempos necesitamos la velocidad, como no tenemos una velocidad exacta y puede variar dependiendo de los tramos, usaremos la **fórmula de la velocidad implícita** para sacar una "velocidad media". Con las distancias en km divididas entre la velocidad implícita, podremos obtener una distancia en minutos aproximada, ya que la velocidad es media y puede variar y los datos de las distancias de km también son aproximadas.

Para '**df\_historic\_order\_demand**', hay ciertos valores nulos en diferentes datos de la columna '**order\_demand**', por lo que hemos imputado los valores medios de ese cliente en el mes indicado.