

# ETL Process

## Cities.csv

To fix **encoding errors** in the variable *city\_code*:

- Created a list of correct city codes.
- Used fuzzy matching to find the best match for each erroneous city code using the fuzzywuzzy library.
- Created a new DataFrame 'cities\_cleaned' with the updated 'city\_code' variable.

## Products.csv

Products data has been cleaned:

- Created a new DataFrame `product\_cleaned` based on the original data.
- Converted hierarchy columns ('hierarchy1\_id', 'hierarchy2\_id', 'hierarchy3\_id', 'hierarchy4\_id', 'hierarchy5\_id') to categorical dtype.

## Sales.csv

Sales data has been cleaned:

- Created a new DataFrame `sales\_cleaned` based on the original data.
- Converted 'date' column to datetime format.
- Converted promotional columns ('promo\_type\_1', 'promo\_bin\_1', 'promo\_type\_2', 'promo\_bin\_2', 'promo\_discount\_2', 'promo\_discount\_type\_2') to categorical dtype.

## Forecast Revenue.csv

- Renamed column 'WEEK' to 'Week' in fore\_rev DataFrame
- Converted 'Date' column to datetime dtype and formatted dates as '%d/%m/%Y' in fore\_rev\_cleaned DataFrame