# Analysing anomaly in king-country house process data

#Last amended: 31st Dec, 2023
#My folder: C:\Users\Ashok\OneDrive\Documents\king_country
# Kaggle: https://www.kaggle.com/datasets/harlfoxem/housesalesprediction

# Ref: h2o:

    a. https://github.com/h2oai/h2o-tutorials/tree/master/best-practices/anomaly-detection
    b. https://github.com/h2oai/h2o-tutorials/tree/master/best-practices

## Steps:

    a. Import data in h2o.ai flow
    b. Do not split frame
    c. Ignore *id* and *date* features
    d. Build autoencoder model as given at the end of this document.
    e. Predict Reconstruction error of original data.
    f. Download predictions of MSE after combining with original data frame
    g. Take the csv file to Windows 10 and open it in Excel
    h. Perform plotting
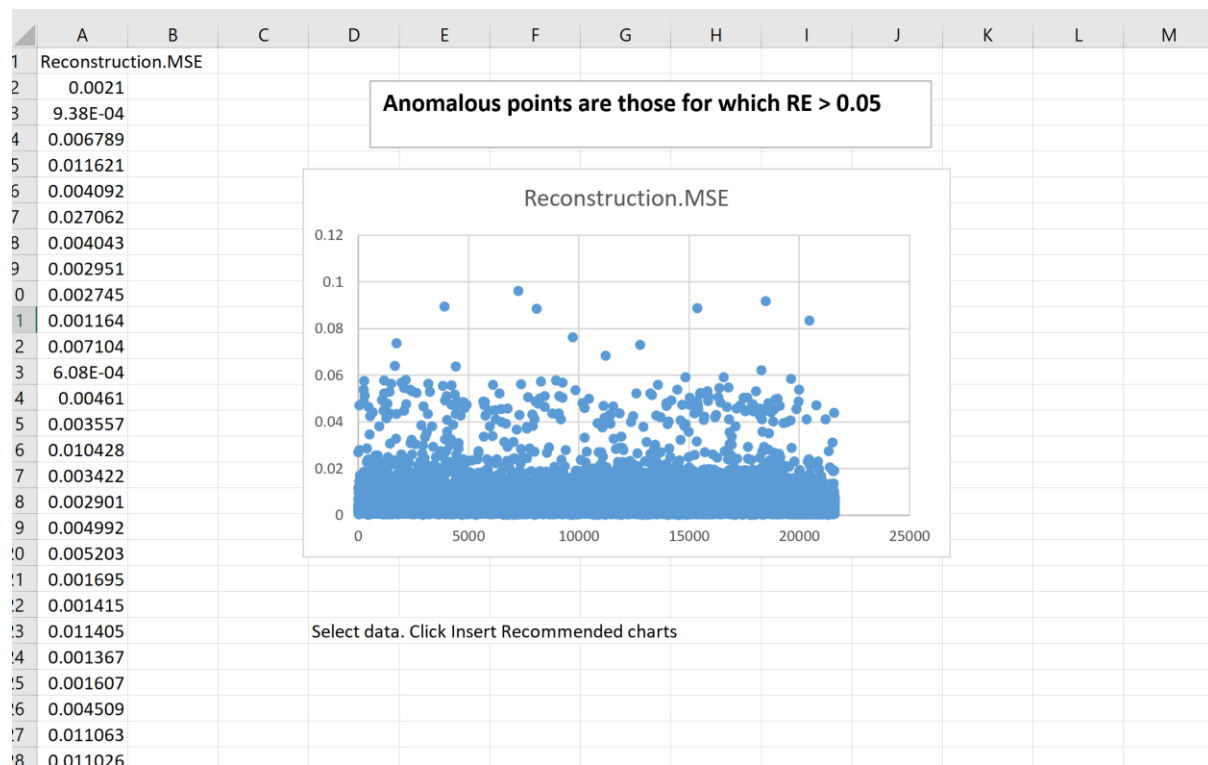


*Figure 1: Reconstruction error plot*

Excel sheet with reconstruction error (RE) and other features. Column '*tags*' has been derived from RE. RE > 0.06 is 1 else 0.

Figure 2: Note the IF condition in the IInd column. The above table is sorted by tag values.

Sort the Excel sheet by *tags* and select top 250 points for plotting in [batchgeo.com](http://batchgeo.com). Free version of batchgeo.com can take at most 250 points as a csv file—just drag and drop csv file to plot.

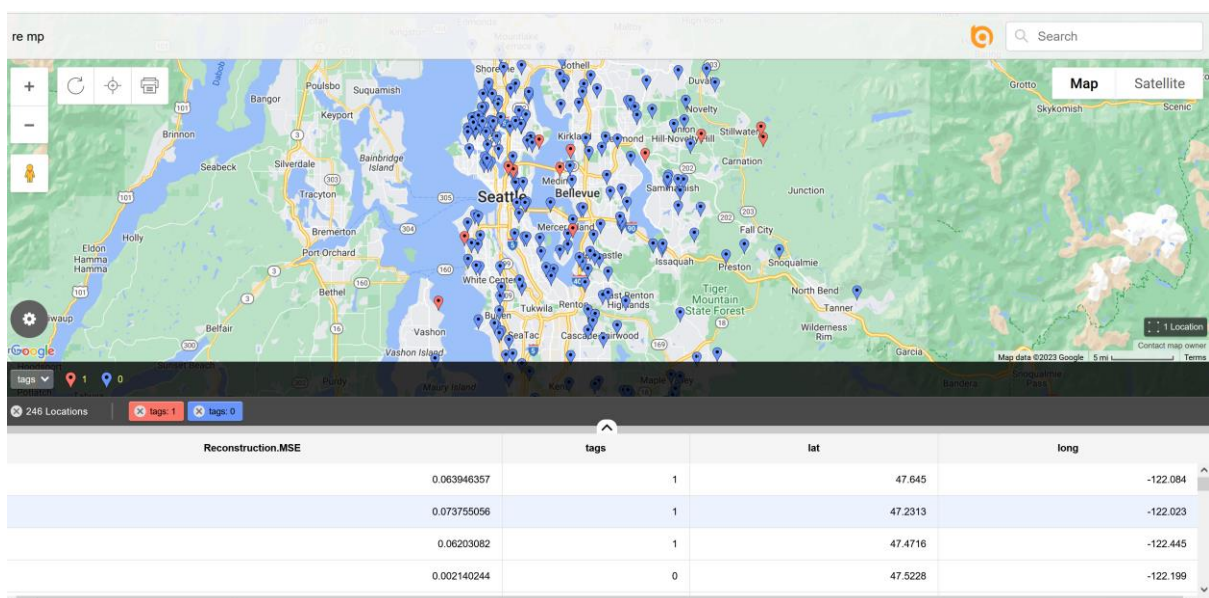Map displaying high Reconstruction Error (points are in red):



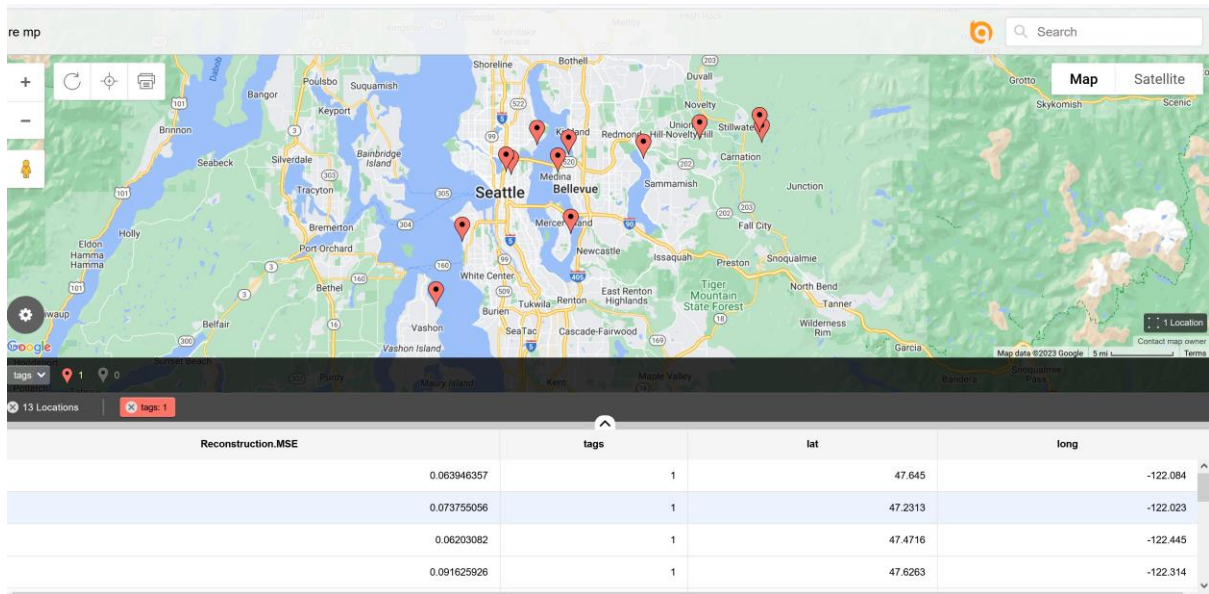Figure 3: tags 1 and tags 0 both getting displayed.

*Figure 4 Only tag 1 getting displayed*

Most of the tag 1 points are near waterfront.

```
# H2o: king-country autoencoder model

buildModel 'deeplearning', {
                    "model_id":"deeplearning-autoencoder",
                    "training_frame":"kc_house_data.hex", <== Full frame, no split
                    "nfolds":0,
                    "ignored_columns":["id","date"],      <==
                    "ignore_const_cols":true,
                    "activation":"Tanh",                  <==
                    "hidden":[5],                         <==
                    "epochs":100,                         <==
                    "variable_importances":false,         <==
                    "score_each_iteration":false,
                    "standardize":true,
                    "train_samples_per_iteration":-2,
                    "adaptive_rate":true,
                    "input_dropout_ratio":0,
                    "l1":0,
                    "l2":0,
                    "loss":"Automatic",
                    "distribution":"AUTO",
                    "quantile_alpha":0.5,
                    "huber_alpha":0.9,
                    "score_interval":5,
                    "score_training_samples":10000,
                    "score_validation_samples":0,
                    "score_duty_cycle":0.1,
                    "stopping_rounds":5,
                    "stopping_metric":"AUTO",
                    "stopping_tolerance":0,
                    "max_runtime_secs":0,
                    "autoencoder":true,
                    "categorical_encoding":"AUTO",
                    "auc_type":"AUTO",
                    "gainslift_bins":-1,
                    "overwrite_with_best_model":true,
                    "target_ratio_comm_to_comp":0.05,
                    "seed":-1,
                    "rho":0.99,
                    "epsilon":1e-8,
                    "nesterov_accelerated_gradient":true,
                    "max_w2":3.4028235e+38,
                    "initial_weight_distribution":"UniformAdaptive",
                    "classification_stop":0,
                    "regression_stop":0.000001,
                    "score_validation_sampling":"Uniform",
```

```
        "diagnostics":true,
        "fast_mode":true,
        "force_load_balance":true,
        "single_node_mode":false,
        "shuffle_training_data":true,                    <==
        "missing_values_handling":"MeanImputation",
        "quiet_mode":false,
        "sparse":false,
        "col_major":false,
        "average_activation":0,
        "sparsity_beta":0,
        "max_categorical_features":2147483647,
        "reproducible":false,
        "export_weights_and_biases":false,
        "mini_batch_size":1,
        "elastic_averaging":false
}
```

#############