# Lab 4: *k*NN and Evaluation

## Practical exercises

Consider the following data:

|        | input |       | output |       |
|--------|-------|-------|--------|-------|
|        | $y_1$ | $y_2$ | $y_3$  | $y_4$ |
| $x_1$  | 1     | 1     | A      | 1.4   |
| $x_2$  | 2     | 1     | B      | 0.5   |
| $x_3$  | 2     | 3     | B      | 2     |
| $x_4$  | 3     | 3     | B      | 2.2   |
| $x_5$  | 1     | 0     | A      | 0.7   |
| $x_6$  | 1     | 4     | A      | 1.2   |

1. Assuming a *k*-nearest neighbor with *k*=3 applied within a leave-one-out schema:

    a) Let $y_3$ be the output variable (*categoric*). Classify $x_1$ when considering uniform weights and:

        i. Euclidean (*l*2) distance (real input variables)

        ii. Hamming distance (categorical input variables)

    b) Let $y_4$ be the output variable (*numeric*). Considering cosine similarity, provide the mean regression estimate for $x_1$.

    c) Consider a weighted-distance *k*-nearest neighbor with Euclidean (*l*2) distance, identify the:

        i. weighted modeestimate of $x_1$ for the $y_3$outcome
        ii. weighted mean estimate of $x_1$ for the $y_4$outcome

2. Let $x_j$ be the measurement on variable $y_j$ for a given observation $x$.

Given the learnt regression model $\hat{x}_4 = 1 - 0.8x_1 + 0.2x_2{}^2 + 0.2x_1x_2$:

a) Compute the $y_4$ regression estimates for the observations of the aforementioned dataset

b) Compute the training Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE)

c) Perform a residue analysis to assess the presence of systemic biases against $y_1$ and $y_2$

**3.** Consider the probabilistic outcome of a classifier for the given six observations to be

$$p(y_3 = A \,|\, x) = \big[p(y_3 = A \,|\, x_1), \ldots, p(y_3 = A \,|\, x_6)\big] = [0.45, 0.4, 0.3, 0.6, 0.8, 0.4]$$

**a)** Draw the training ROC curve

**b)** Compute the training AUC

**c)** Would you change the default 0.5 probability threshold for this classifier in order to maximize training F1?

# Programming quest

1. Consider the accuracy estimates collected under a 5-fold CV for two predictive models M1 and M2, $acc_{M1}$=(0.7,0.5,0.55,0.55,0.6) and $acc_{M2}$=(0.75,0.6,0.6,0.65,0.55).

   Using **scipy**, assess whether the differences in predictive accuracy are statistically significant.

   *Resource*: https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.ttest_rel.html

4. Consider the *housing* dataset available at https://web.ist.utl.pt/~rmch/dscience/data/housing.arff and the *Regression* notebook available at the course's webpage. Using a 10-fold cross-validation:

   a) Assess the MAE of a kNN regressor for $k \in \{1, 5, 9\}$ (remaining parameters as default)

   b) Compare the RMSE of the default kNN and decision tree regressors