# Assignment: Interpretability and Explainability in Machine Learning
## Concrete data

Pedro Delicado, Universitat Politècnica de Catalunya - Barcelona TECH

2023-12-20

## Concrete Dataset

UC Irvine Machine Learning Repository, Concrete Dataset.

**Abstract:** Concrete is the most important material in civil engineering. The concrete compressive strength is a highly nonlinear function of age and ingredients. These ingredients include cement, blast furnace slag, fly ash, water, superplasticizer, coarse aggregate, and fine aggregate.

**Data Characteristics:** The actual concrete compressive strength (MPa) for a given mixture under a specific age (concretes) was determined from laboratory. Data is in raw form (not scaled).

**Summary Statistics:**

- Number of instances (observations): 1030
- Number of Attributes: 9
- Attribute breakdown: 8 quantitative input variables, and 1 quantitative output variable
- Missing Attribute Values: None

---

**Variable Information:** Given is the variable name, variable type, the measurement unit and a brief description.

*Input Variables:*

- Cement (component 1) – quantitative – kg in a m3 mixture

- Blast Furnace Slag (component 2) – quantitative – kg in a m3 mixture
- Fly Ash (component 3) – quantitative – kg in a m3 mixture

- Water (component 4) – quantitative – kg in a m3 mixture

- Superplasticizer (component 5) – quantitative – kg in a m3 mixture
- Coarse Aggregate (component 6) – quantitative – kg in a m3 mixture
- Fine Aggregate (component 7) – quantitative – kg in a m3 mixture
- Age – quantitative – concrete (1~365)

*Response variable:* - Concrete compressive strength – quantitative – MPa – Output Variable

```r
library(readxl)
concrete <- as.data.frame(read_excel("Concrete_Data.xls"))
DescVars <- names(concrete)
names(concrete) <- c("Cement","Slag","FlyAsh","Water","Superplast",
"CoarseAggr","FineAggr","Age","Strength")
```

**Data processing: Creating training and test sets**

Create a training sample choosing 700 data at random. The non-chosen data will be the test set.

## 1. Fit a Random Forest

    a. Compute the *Variable Importance* by the reduction of the **impurity** at the splits defined by each variable.

    b. Compute the Variable Importance by out-of-bag random permutations.

    c. Do a graphical representation of both Variable Importance measures.

    d. Compute the Variable Importance of each variable by Shapley Values.

## 2. Fit a linear model and a gam model.

    a. Summarize, numerically and graphically, the fitted models.

    b. Compute the Variable Importance by Shappley values in the linear and gam fitted models. Compare your results with what you have learned before.

## 3. Relevance by Ghost Variables

Compute the relevance by ghots variables in the three fitted modls.

## 4. Global Importance Measures and Plots using the library DALEX

    a. Compute Variable Importance by Random Permutations

    b. Do the Partial Dependence Plot for each explanatory variable.

    c. Do the Local (or Conditional) Dependence Plot for each explanatory variable.

## 5. Local explainers with library DALEX

Choose two instances in the the test set, the prediction for which we want to explain:

- The data with the lowest value in Strength.
- The data with the largest value in Strength.

For these two instances, do thefollowing tasks for the fitted random forest.

    a. Explain the predictions using SHAP.

    b. Explain the predictions using Break-down plots.

    c. Explain the predictions using LIME.

    d. Do the Individual conditional expectation (ICE) plot, or ceteris paribus plot

    e. Plot in one graphic the Individual conditional expectation (ICE) plot for variable **Age** for eachcase in the test sample. Add the global Partial Depedence Plot.