# Module 2 Lab:
# Data Annotation in Active Learning

## OBJECTIVE

**Data annotation** is to label collected data, which is necessary to build machine learning models, especially for supervised learning models in different applications such as Image Processing, Natural Language Processing, and Smart Grid.

**Active learning** is a category of machine learning models in which a learning algorithm can interactively query a user to label samples with the predefined label. The goal of this lab is to build a pipeline that implements data annotation in the active learning.

## PREREQUISITES

You have to complete all submodule 2.4 lecture slides and install packages below.

*sklearn, modAL, numpy, IPython, matplotlib*

## INSTRUCTIONS

- Import ***MNIST*** dataset (http://yann.lecun.com/exdb/mnist/) with *sklearn.datasets*
- Split the dataset into training and testing datasets $D_{training}$ and $D_{testing}$
- Select 100 training samples to build a subset $D_s$ from $D_{training}$, which is used to build a classifier for active learning
- Build a data pool of unlabeled data $D_{pool} = D_{training} - D_s$
- Initializing a classifier *C* based on **Logistic Regression** with $D_s$ via for active learning
- Implement data annotation on $D_{pool}$ by querying users to labeled samples that are from $D_{pool}$, extend $D_s$ with the annotation results to build to $D'_s$, and retrain *C* with $D'_s$ and check the performance, where the performance is evaluated with accuracy