

OBJECTIVE:

- Familiarize with the Jupyter Notebook environment, Python commands and basic libraries.
- Explore the *MNIST* dataset
- Create a kNN Classifier for *MNIST*

PREPARATION & REFERENCES:

- Complete IMLP (Muller/Guido book) Ch.1, working through ch01-introduction code & notebook.
- Sklearn datasets [Sklearn datasets](#)

TASK:

WRITE Python code, using basic libraries (*sklearn*, *numpy*, *matplotlib*) to perform the tasks below.

Place the solution to each number below in a unique cell; use comment notation to include the number for reference, and a typed response to any explicit questions (in bold).

Load the Data

- Using Sklearn.datasets [fetch\\_openml](#) loader, load the 'mnist\_784', version 1 dataset.  
**What type of object is returned?** \_\_\_\_\_

Analyze the Data

Provide Python statement(s) to support your answers to the questions below:

- Display the object's attributes (*keys*). **How many are there?** \_\_\_\_\_
- Display a partial description of the dataset, just the first 200 characters.
- How many samples are in the dataset?** \_\_\_\_\_
- How many features does each data sample have?** \_\_\_\_\_
- What is the target shape?** \_\_\_\_\_
- What is the difference between the storage structure for the data and that for the target?** \_\_\_\_\_

Create a custom Function – just for practice...

- Write a function that calculates an index value based on what alpha characters are passed to it. The function takes in two alpha characters and returns the summation of their numeric placement in the alphabet {'a'=1, 'b'=2,...'z'=26}. Example: *fnct('c','y')* would return 28.
- Invoke the function with your initials, and display the **resulting index value.** ? \_\_\_\_\_

Plot your DIGIT from the Data – just for fun!

Each sample represents a single image as 28x28 pixels and each pixel holds a value 0-255 (the pixel's intensity from white to black). That's why there are 784 total features per sample – 1 per pixel.

- You will first need to grab your digit's feature vector from the data, using your custom index value (#9 above).
- Next, you will need to reshape it for plotting. Simply invoked the *reshape* method with your feature vector.
- Now you are ready to plot it using [matplotlib pyplot's imshow](#) method. You need only the first two arguments: your reshaped array, and colormap 'binary'. **What is your digit?** \_\_\_\_\_

OK, back to the main task - Let's do some TRAINING

FIRST – separate your dataset into two sets, one for training and one for testing. We won't use the sklearn train\_test\_split this time, because the MNIST dataset has already been carefully separated into separate training and test sets. More importantly, processing it could take a long time, so we are just going to use a fraction of it.

- Assign the first 6,000 samples as your **training set** (X\_train, y\_train), and the last 1,000 samples as your **test set** (X\_test, y\_test). Remember, X comes from your *data*, and y comes your *targets*.

NEXT – let's build a [kNN classifier](#), with initial **k-value =3** [tip: enable GPUs (*Runtime->Change runtime type*)]

- What is your test *score* when k is 3? \_\_\_\_\_
- Does it improve when you set algorithm to 'distance'? (yes/no) \_\_\_\_\_ **set weights to 'distance'**
- Can you find a k that performs better? \_\_\_\_\_ If yes, what k? \_\_\_\_\_

DELIVERABLES

- *xyLab2.py* - Python source file, where *xy* are your initials; fully documented, including a general header & comments
- *xyLab2.pdf* – which SHOWS ALL ... code, outputs, embedded comments - generated from a full run of matching .py file  
**[CLEAN:** Before you print .pdf, be sure all cells and outputs are displayed, close help windows, remove unnecessary cells]