

## INSTRUCTIONS

1. The assignment contains four questions. A few bonus questions are mentioned.
2. This assignment is due on **6th Feb, 23:59 \*\***(No Further extensions\*\*).
3. Assignment must be implemented in Python 3 only.
4. You are allowed to use libraries for data preprocessing (numpy, pandas etc) and for evaluation metrics, data visualization (matplotlib etc.).
5. You will be evaluated not just on the overall performance of the model and also on the experimentation with hyper parameters, data preprocessing techniques etc.
6. The report file must be a well documented jupyter notebook, explaining the experiments you have performed, evaluation metrics and corresponding code. The code must run and be able to reproduce the accuracies, figures/graphs etc.
7. For all the questions, you must create a train-validation data split and test the hyperparameter tuning on the validation set. Your jupyter notebook must reflect the same.
8. Any attempts at **plagiarism will be penalized heavily**.
9. Make sure you run and save your notebooks before submission.
10. For question 3 of the Decision Trees section, output your model's depth first traversal into outputimp.txt and submit it along with the ipynb file.
11. Naming convention for the ipynb file is <roll\_number>\_assign1.ipynb
12. Compress your submission files into a zip file with the naming convention: <roll\_number>\_assign1.zip and submit in the portal.

## ▼ 1) REGRESSION

Please find the Diamond Price Prediction Data set

<https://drive.google.com/drive/folders/1qE1tm3Ke3uotTyv6SUqru09t-AkcwRK?usp=sharing>.

"description.txt" contains the feature description of data, "diamonds.csv" has the data.

```
# To read data from diamonds.csv
import pandas as pd
headers = ["carat", "cut", "color", "clarity", "depth", "table", "price", "x", "y", "z"]
data = pd.read_csv('diamonds.csv', na_values='?',
                  header=None, names = headers)
data = data.reset_index(drop=True)
data = data.iloc[1:]
data.describe()
#print(data)
```

# This is formatted as code

## KNN Regression [Diamond Price Prediction Dataset]

1. a) Build a knn regression algorithm [using only python from scratch] to predict the price of diamonds.

# code for knn regression

1. b) Do we need to normalise data? [If so Does it make any difference?].

# give proper explanation

2. Experiment with different distance measures[Euclidean distance, Manhattan distance, Hamming Distance] to handle categorical attributes.

# show all the experiments

3. Report Mean Squared Error(MSE), Mean-Absolute-Error(MAE), R-squared (R2) score in a tabular form.

# report a table

4. a) Choose different K values (k=2,3,5,7,11,16) and experiment. Plot a graph showing R2 score vs k.

# plot

4. b) Are the R-squared scores the same? Why / Why not? How do we identify the best K? Suggest a computational procedure, with a logical explanation.

# Explanation

5. a) Also, report the performance of scikit-learn's kNN regression algorithm.

# scikit-learn KNN Regressor

5. b) Compare it with the algorithm you built. [ you can use complexities, R2 score etc..]

```
# Comparison
```

6. From the above experiments, what do you think are advantages and disadvantages of the knn regression algorithm?

```
# report this along with the experiments
```

## ▼ 2) Linear Regression

Dataset - same as above (Diamond Price Detection)

2a) Implement a Linear Regression model (from the scratch) taking suitable independent variables from the dataset.

Report and Calculate the error obtained.

2b) What are the best suitable features you used to predict the price of the dataset and Why?

Idea: Use Correlation to get the suitable features and Report the values accordingly.

```
#code for Correlation between features and the Diamond Price.
```

Explanation for 2b) -

2c) Use the module Linear Regression from sklearn to predict the price of diamonds(considering the same attributes as before) and compare the result obtained with the above.

```
# import sklearn model
```

2d) Now, using the whole dataset, predict the price of the Diamonds using the module of Linear Regression from sklearn. Report the changes you have observed compared to before? Adding extra features did it make the prediction better or worse. Comment?

2e) Now, compare the algorithms KNN regression and Linear Regression. What are the differences you have observed? Which is better and why. Your statements should be backed up with statistics.

Explanation -

2f) Plot the predicted values from KNN regression, Linear Regression and Actual Diamond Price.

```
#plot
```

## ▼ KNN Classifier

In this problem you are required to train, test and validate a K-Nearest Neighbor Classifier on the famous CIFAR-10 dataset. The next few cells will guide you through the process. Follow along...

```
# Run some setup code for this notebook.
```

```
import random
import numpy as np
import matplotlib.pyplot as plt
import os
```

```
from __future__ import print_function
```

```
from builtins import range
from six.moves import cPickle as pickle
from imageio import imread
import platform
```

```
# This is a bit of magic to make matplotlib figures appear inline in the notebook
# rather than in a new window.
```

```
%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'
```

```
# Some more magic so that the notebook will reload external python modules;
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```
# Download the dataset
!wget http://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz -O cifar-10-python.ta
!tar -xvzf cifar-10-python.tar.gz
!rm cifar-10-python.tar.gz
```

--2021-01-25 17:02:02-- <http://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz>

```
Resolving www.cs.toronto.edu (www.cs.toronto.edu)... 128.100.3.30
Connecting to www.cs.toronto.edu (www.cs.toronto.edu)|128.100.3.30|:80... con
HTTP request sent, awaiting response... 200 OK
Length: 170498071 (163M) [application/x-gzip]
Saving to: 'cifar-10-python.tar.gz'
```

```
cifar-10-python.tar 100%[=====>] 162.60M 71.2MB/s in 2.3s
```

```
2021-01-25 17:02:04 (71.2 MB/s) - 'cifar-10-python.tar.gz' saved [170498071/1
```

```
cifar-10-batches-py/
cifar-10-batches-py/data_batch_4
cifar-10-batches-py/readme.html
cifar-10-batches-py/test_batch
cifar-10-batches-py/data_batch_3
cifar-10-batches-py/batches.meta
cifar-10-batches-py/data_batch_2
cifar-10-batches-py/data_batch_5
cifar-10-batches-py/data_batch_1
```

```
# This cell loads the training and testing dataset. Please note the variables at
# the end of the cell as you would require them to access the train/test data
# and labels throughout the assignment
```

```
def load_pickle(f):
    version = platform.python_version_tuple()
    if version[0] == '2':
        return pickle.load(f)
    elif version[0] == '3':
        return pickle.load(f, encoding='latin1')
    raise ValueError("invalid python version: {}".format(version))

def load_CIFAR_batch(filename):
    """ load single batch of cifar """
    with open(filename, 'rb') as f:
        datadict = load_pickle(f)
        X = datadict['data']
        Y = datadict['labels']
        X = X.reshape(10000, 3, 32, 32).transpose(0,2,3,1).astype("float")
        Y = np.array(Y)
        return X, Y

def load_CIFAR10(ROOT):
    """ load all of cifar """
    xs = []
    ys = []
    for b in range(1,6):
        f = os.path.join(ROOT, 'data_batch_{}'.format(b))
        X, Y = load_CIFAR_batch(f)
        xs.append(X)
        ys.append(Y)
    Xtr = np.concatenate(xs)
    Ytr = np.concatenate(ys)
    del X, Y
    Xte, Yte = load_CIFAR_batch(os.path.join(ROOT, 'test_batch'))
    return Xtr, Ytr, Xte, Yte
```

```

cifar10_dir = 'cifar-10-batches-py'

# Cleaning up variables to prevent loading data multiple times
try:
    del X_train, y_train
    del X_test, y_test
    print('Clear previously loaded data.')
except:
    pass

X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

# As a sanity check, we print out the size of the training and test data.
print('Training data shape: ', X_train.shape)
print('Training labels shape: ', y_train.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)

Training data shape: (50000, 32, 32, 3)
Training labels shape: (50000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)

```

Next we visualize the CIFAR-10 dataset. Although these functions are being written for you, we highly recommend you go through the code and make yourself familiar as these are things you will be required to do very often when working on AI/ML projects

```

# Visualize some examples from the dataset.
# We show a few examples of training images from each class.
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship',
num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y)
    idxs = np.random.choice(idxs, samples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()

```



In the next cell we flatten each image into a single dimensional vector so that it is easy to process. You should be able to reason about the dimensions comfortable.

```
# Subsample the data for more efficient code execution in this exercise
num_training = 5000
mask = list(range(num_training))
X_train = X_train[mask]
y_train = y_train[mask]

num_test = 500
mask = list(range(num_test))
X_test = X_test[mask]
y_test = y_test[mask]

# Reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
print(X_train.shape, X_test.shape)

(5000, 3072) (500, 3072)
```

In the next cell you are going to implement the main KNearestNeighbor class and keep adding functions to it as and when required in the subsequent steps.

```
# You will keep coming back to this cell to add more functions as and when
# required. Right now it is very simple!
class KNearestNeighbor(object):
    """ a kNN classifier with L2 distance """

    def __init__(self):
        pass
```

Go back to KNearestNeighbor class and add a method to train the classifier. Your function will be called as below. In KNN classifier, this step is a simple memorization of the training data.

```
# Make sure to add train() in the classifier class before executing this cell
classifier = KNearestNeighbor()
classifier.train(X_train, y_train)
```

Go back to KNearestNeighbor class and add a method to compute distances between each pair of test image and train image. You can use two loops to do this. Remember we are using standard L-2 distance metric. Precisely your method should return a distance matrix( $D$ ) where  $D(i, j) = \text{L-2 distance between the } i^{th} \text{ test image and the } j^{th} \text{ train image}$ . Your function will be called as below

```
# Make sure to add compute_distances_two_loops() in the classifier class
# before executing this cell
dists = classifier.compute_distances_two_loops(X_test)
print(dists.shape)
```

Next implement the function to predict labels. Again go back to the KNearestNeighbor class cell. Your function will be called as below

```
# Make sure to add predict_labels() in the classifier class
# before executing this cell
y_test_pred = classifier.predict_labels(dists, k=1)

# Compute and print the fraction of correctly predicted examples
num_correct = np.sum(y_test_pred == y_test)
accuracy = float(num_correct) / num_test
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))

# Use this cell to compute accuracies for k = 3, 5, 7
```

Now we have a basic classifier ready but it is extremely inefficient. In Machine Learning writing vectorised code is one of the most important skills.

Now you have to again go back to the KNearestNeighbor class cell and add functions to compute the distance matrix using:

1. single loop
2. no loops

From here onwards we won't do as much hand holding as done before. You are expected to write all the code from scratch

```
# Use this cell to call the functions
```



Next, you are required to tabulate the time taken to compute the distance matrix using each of the three above methods implemented for  $k = 1, 3, 5, 7$ .

# Use this cell to present your timing results for computing distance matrix

As you can see that in KNN algorithm it is important to tune the hyperparameter  $K$ . We will do this using Cross Validation Leave One Out approach. The idea is to split the train set into few folds(here we recommend you to set the fold number to 5). Then, for each value of  $K$ , we leave one of these folds out to evaluate performance but use the others to train. Repeat this by leaving each fold out once. You will get 5 accuracies in this case, one for leaving out each fold.

Depending on the average accuracy decide on the optimal value of  $K$ . Report the same. We also expect you to present plots showing the average accuracy and standard deviation for each value of  $K$ .

# Use this cell to implement Cross Validation

In this last part of the notebook you are expected to do the following and present appropriate reports/results:

1. Try subtracting the mean from each pixel in each image. This mean is computed across all pixels in all images. Report accuracies on the test set for  $k = 1, 3, 5, 7$ . Is there any change in the accuracy? Why?
2. Try subtracting the mean from each pixel in each image. This mean is computed across all pixels at that position from all images. Report accuracies on the test set for  $k = 1, 3, 5, 7$ . Is there any change in the accuracy? Why?
3. What is the time complexity of training using KNN classifier? What is the time complexity while testing? Is KNN a linear classifier or can it learn any boundary?
4. Bonus: Explore what image features you can use to better the performance of KNN classifier.

## ► Decision Trees

The Wisconsin Breast Cancer Dataset(WBCD) can be found here(<https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.data>)

This dataset describes the characteristics of the cell nuclei of various patients with and without breast cancer. The task is to classify a decision tree to predict if a patient has a benign or a malignant tumour based on these features.

Attribute Information:

#	Attribute	Domain
-----		
1.	Sample code number	id number
2.	Clump Thickness	1 - 10
3.	Uniformity of Cell Size	1 - 10
4.	Uniformity of Cell Shape	1 - 10
5.	Marginal Adhesion	1 - 10
6.	Single Epithelial Cell Size	1 - 10
7.	Bare Nuclei	1 - 10
8.	Bland Chromatin	1 - 10
9.	Normal Nucleoli	1 - 10
10.	Mitoses	1 - 10
11.	Class:	(2 for benign, 4 for malignant)

[ ] ↩ 19 cells hidden